# A PROBABILISTIC ANALYSIS OF LOW-RANK APPROXIMATIONS IN OPTIMIZATION PROBLEMS WITH ELLIPSOIDAL CONSTRAINTS

A Dissertation

Presented to the Faculty of the Graduate School

of Cornell University

in Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

by

Spyridon Schismenos

August 2009

A PROBABILISTIC ANALYSIS OF LOW-RANK APPROXIMATIONS IN

OPTIMIZATION PROBLEMS WITH ELLIPSOIDAL CONSTRAINTS

Spyridon Schismenos, Ph.D.

Cornell University 2009

Low-rank matrix approximations have been used in many applications, because they provide compact representations of the data and reveal the underlying structure. This dissertation is concerned with applications of low-rank approximations in optimization problems. Motivation comes from a recent effort in designing radiotherapy treatment plans for patients with cancer. The problem was formulated as a second-order cone program. Due to the size of the problem, low-rank matrices were used in order to create a computationally tractable approximation. This work is an attempt to theoretically explain the success of low-rank approximations in such problems. The main vehicle for this analysis is a stylized optimization problem with randomly sampled ellipsoidal constraints. We consider two different matrix approximations, one based on the Singular Value Decomposition and one based on column sampling, and apply them to the matrices in the stylized problem. We provide results about the probability distributions of the optimal values of these problems as well as their relative difference. Since the focus is on problems with large number of constraints, we provide asymptotic results, when the number of constraints tends to infinity. We finally compare the performance of the two approximations and discuss the implications of our results.

## BIOGRAPHICAL SKETCH

Spyros was born and raised in Panaitólion, a village in western Greece. While in high school, he decided that he did not want to become a dentist, like his parents and his brother, and he chose to study Applied Mathematics at the National Technical University in Athens. In 2004, he left Greece and came to Ithaca, NY to pursue a Ph.D. in Operations Research.

After graduation, he will work for JPMorgan Chase in London, UK.

To my family.

# ACKNOWLEDGEMENTS

**TABLE OF CONTENTS**

# LIST OF FIGURES

# CHAPTER 1

## INTRODUCTION

## 1.1 Low-rank matrix approximations

In many applications where the data can be formulated as a matrix with a large number of columns and rows, it is of interest to find compact representations of the data. The concept of the rank of a matrix is fundamentally connected to finding such compact representations. For an $m \times n$ matrix $A$, the rank is defined as the dimension of the range of $A$, $\text{ran}(A) = \{y \in \mathbb{R}^m | y = Ax \text{ for some } x \in \mathbb{R}^n\}$. If the rank of $A$ is equal to $k$, then all columns of $A$ can be written as linear combinations of a subset of the columns of size $k$.

In many cases though, the rank of $A$ is not significantly smaller than $\min\{m, n\}$. One can then try to find an $m \times n$ matrix of low rank that is an approximation to $A$. A natural way to define such an approximation is as the solution to the problem

$$\min \ \|A - X\|_F$$

$$\text{subject to } \ X \in \mathbb{R}^{m \times n}, \ \text{rank}(X) \leq k, \tag{1.1}$$

where $\|\cdot\|_F$ denotes the Frobenius norm. The solution to this problem is related to the Singular Value Decomposition (SVD) of $A$ and is described in the following theorem, see [23].

**Theorem 1.1.1** *Let $U^T A V = \Sigma = diag(\sigma_1, ..., \sigma_p) \in \mathbb{R}^{m \times n}, p = \min\{m, n\}$ be the SVD of the $m \times n$ matrix $A$, where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices. Then a*

*solution to Problem (1.1) is given by the matrix*

$$A^{(k)} = U\Sigma^{(k)}V^T,$$

*where $\Sigma^{(k)} = diag(\sigma_1, ..., \sigma_k, 0, ..., 0)$. Furthermore,*

$$\|A - A^{(k)}\|_F = \sqrt{\sum_{j=k+1}^{p} \sigma_j^2}.$$

This fundamental result shows that the singular values of the matrix *A* determine how close it is to a low-rank matrix and gives an algorithm for computing such an approximation based on the SVD.

Computing a low-rank approximation using the SVD is appealing from a theoretical point of view, since it provides the closest matrix with a given rank. For many applications where the data matrix is large, calculating the SVD can be impractical since it requires a large number of operations and it has large memory requirements. Recent research has thus focused on algorithms that are suboptimal, in the sense that the low-rank matrices that they calculate are not the closest possible to the original matrix, [18, 36, 22, 19, 20, 2]. The advantage of using such algorithms is that they are faster than the SVD-based algorithm and need less memory, making them much more suitable for large scale applications.

Low-rank approximations have found numerous applications in various fields. Examples include Latent Semantic Indexing, [7, 31], Support Vector Machine training, [22, 36], Computer Vision, [30], and Web Search models, [26]. In these applications, the data consist of a matrix that, although not of low rank, can be approximated well by a low-rank matrix. Calculating such a low-rank approximation can reveal the underlying structure of the data and allow for fast computations.

## 1.2 An application in optimization

The topic of this dissertation is the application of low-rank approximations in optimization problems. More specifically, we are concerned with the problem of approximating quadratic constraints using low-rank matrices. Motivation comes from a recent optimization problem, [13, 12], arising in designing Intensity Modulated Radiation Therapy (IMRT) treatment plans. The problem is formulated as a second-order cone program, but due to the size of the problem, low-rank approximations were used in order to make the problem computationally tractable. Computational results suggest that the approximation was successful.

Briefly, in [13, 12], an optimization formulation for the IMRT problem was proposed, which we present in more detail in Section 2.1. The analysis takes into account various uncertainties that affect the radiation dose delivered to a patient. These uncertainties are due to position uncertainties, or motion and deformation of the patient or of the inner organs between daily treatments. Constraints are imposed on the dose delivered to each part of the body, that assure high dose with high probability delivered to tumor parts and low dose with high probability delivered to healthy parts of the body. Each such constraint is a quadratic constraint in $\mathbb{R}^d$, of the form

$$\sqrt{x^T W_i x} \leq a_i + b_i^T x,$$

where $W_i$ is a covariance matrix.

The number of quadratic constraints of this type is around 10000, while the dimension $d$ of the problem is around 1000, making the problem computationally intractable due to large memory requirements. In order to make this prob-

lem more tractable, a simple approximation scheme was used in [13, 12] that essentially replaced the covariance matrices $W_i$ with sample covariance matrices $C_i$ based on a small number of scenarios. The problem was then formulated as an optimization problem with a linear objective function.

The optimization problem was solved numerically in a patient case. The resulting treatment plan was satisfactory, achieving sparing of the healthy tissue and delivering a high dose to the target volume, while accounting for uncertainties. In order to further test the quality of the solution obtained using this scenario-based approximation, the problem was solved using various numbers of scenarios, [13]. The solutions were found to be of similar quality.

This empirical phenomenon is striking. The $d \times d$ covariance matrices $W_i$ were substituted by matrices $C_i$ of rank in the order of 10, reducing the memory requirements significantly and making the problem computationally tractable. On the other hand, the solution was satisfactory, suggesting that it was close to the solution that would have been obtained without the approximation.

## 1.3   Dissertation layout

Motivated by the optimization problem that we present in Section 2.1, we perform a theoretical analysis of low-rank approximations in optimization problems with quadratic constraints. Our first goal is to explore mathematically the success of the low-rank approximation in [13]. Furthermore, through our theoretical study we provide insight to the properties of low-rank approximations in optimization and compare the performance of different kinds of approximations.

In order to perform our analysis, we construct a stylized optimization problem that is based on the robust IMRT problem and has a feasible region that is the intersection of ellipsoids of the same shape, randomly rotated around the origin. We argue in Section 2.2 that this is a reasonable approximation to the constraints of the IMRT problem. We then approximate the constraints by using two different approximations: one based on the SVD and one based on column sampling. We focus on the optimal values of these problems and especially on their relative difference. Since low-rank approximations are particularly attractive when the number of constraints is large, our study involves an asymptotic analysis of the optimal values, when the number of constraints tends to infinity.

In Chapter 2 we present the robust IMRT problem formulation in more detail. We then construct the stylized optimization problem. The stylized problem is a sampled optimization problem, with a linear objective function and constraints drawn independently from a specific distribution. Sampled problems have been used recently in order to approximate the chance-constrained problem, which is more difficult to solve in practice. We give a brief introduction to the relevant theory and then present the results that connect chance-constrained and sampled problems. We then capitalize on them in order to describe a method than can be used to give properties of the optimal value of any sampled problem with independent constraints.

In Chapter 3 we present the approximate optimization problem that results when the matrices in the constraints of the original problem are approximated with low-rank matrices calculated using the SVD. Using such an approximation is a natural choice since it approximates the matrices in an optimal way and it has a nice geometric interpretation in our setting. We then apply the technique

that we developed in Chapter 2 to the optimal values of the two problems. More specifically, we explicitly solve the corresponding chance-constrained optimization problem by exploiting symmetry in the distribution that we have chosen for generating constraints in our model. We conclude this chapter with an asymptotic result about the the optimal value of the chance-constrained problem that will be used later.

We present an alternative approach on approximating a matrix using column sampling in Chapter 4. This algorithm requires time that is linear in the dimension of the matrix. Thus, it is faster than the SVD-based algorithm and is more attractive in a practical setting. After describing in detail the algorithm and related work, we give the form of the optimization problems in our model under this approximation. In order to give some insight into the behavior of the approximating problem, we present an analysis when the constraints of the original problem are nearly spherical. The analysis is related to the well-known coupon-collector's problem from probability theory. We then apply the technique from Chapter 2 to the approximating problem. The analysis is more involved than in the SVD-based approximation. We are not able to explicitly solve the corresponding chance-constrained optimization problem, but we provide a bound on its optimal value.

In Chapter 5, we use our fundamental result about sampled problems from Chapter 2 in order to characterize the asymptotic behavior of the optimal value of sampled problems when the number of constraints tends to infinity. The asymptotic behavior depends on a certain asymptotic behavior of the optimal value of the corresponding chance-constrained problem.Applying our methods, we get asymptotic results related to the optimal values and the errors of the

approximations when the number of constraints tends to infinity. Finally, we combine the results that we have for the approximations in order to discuss the implications of our analysis for using low-rank approximations in practice and to compare the two methods that we have studied.

## 2.1   The robust IMRT problem

Intensity Modulated Radiation Therapy (IMRT) is a method for planning and delivering radiation therapy to patients with cancer. The objective of IMRT is to shape the distribution of the dose delivered to the targeted tissues while sparing healthy tissues, [10, 24, 29]. Briefly, dose distributions are formed by the superposition of a series of beamlets intersecting the target from many directions.

From an optimization point of view, one wants to find the optimal intensity assigned to each beamlet in order to hit the targeted tissues with a high dose while delivering a low dose to healthy tissues. In addition to that, the optimization formulation has to address uncertainties about the exact locations of the targeted areas. These uncertainties are unavoidable and are due to either position uncertainties, or motion and deformation of the patient or of the inner organs during, or between daily treatments.

There exist many ways to formulate this problem as an optimization problem that addresses uncertainties. The approach in [13, 12] can be interpreted either from a probabilistic, or from a robust optimization point of view. We present here, almost verbatim from [12], the probabilistic derivation.

Using the conventional modelling approach of using voxels (volume elements) as sample points on a grid where we measure amounts of dose absorbed, let $x$ be the beamlet intensity vector and $D_i(x)$ be the total dose delivered to voxel $i$ over the course of all $N$ treatments. Then $D_i(x)$ is viewed as a random

variable, since the exact dose in each treatment depends on the exact position of the voxel, which is random and as argued in [13] can be considered to be normally distributed. We denote by $\mu_i$ and $\sigma_i^2$ the mean and the variance of $D_i(x)$ respectively.

If the voxel $i$ is part of a healthy region, then we do not want the dose in that voxel to exceed some threshold $M_1 \geq \mu_i$. We require that the dose exceeds $M_1$ with probability at most $\delta$, where $\delta$ is a small constant, e.g. 0.05. So, we require that $P(D_i(x) > M_1) \leq \delta$ which is equivalent to

$$P\left(\frac{D_i(x) - \mu_i}{\sigma_i} > \frac{M_1 - \mu_i}{\sigma_i}\right) \leq \delta. \tag{2.1}$$

Since the random variable $(D_i(x) - \mu_i)/\sigma_i$ is approximately distributed as a standard normal, (2.1) becomes

$$\sigma_i \leq \frac{M_1 - \mu_i}{\Phi^{-1}(1 - \delta)}, \tag{2.2}$$

where $\Phi(\cdot)$ is the cumulative distribution of a standard normal random variable.

In complete analogy, if voxel $i$ comes from the target volume, then we do not want the dose to fall below a certain threshold $M_2 < \mu_i$. A similar argument leads to the constraint

$$\sigma_i \leq \frac{\mu_i - M_2}{\Phi^{-1}(1 - \delta)}. \tag{2.3}$$

From the i.i.d. assumption we get that $\mu_i = NE[D_{i1}(x)]$ and $\sigma_i^2 = N\text{Var}[D_{i1}(x)]$, so we need to calculate the mean and the variance of the dose for a single treatment. Let $Y_i$ denote a random column vector, indexed by beamlets, representing the dose delivered to voxel $i$ from each beamlet, if the beamlets have unit intensity. Then the dose to voxel $i$ in a single treatment is given by

$$D_{i1}(x) = Y_i^T x.$$

This immediately implies that we have

$$\mu_i = NE[Y_i^T]x$$

and

$$\sigma_i^2 = Nx^T\text{Cov}(Y_i)x.$$

Constraint (2.2) becomes then

$$\sqrt{Nx^T\text{Cov}(Y_i)x} \leq \frac{M_1 - NE[Y_i^T]x}{\Phi^{-1}(1-\delta)} \qquad (2.4)$$

and (2.3) becomes

$$\sqrt{Nx^T\text{Cov}(Y_i)x} \leq \frac{NE[Y_i^T]x - M_2}{\Phi^{-1}(1-\delta)}. \qquad (2.5)$$

The next question that arises is how one can estimate the covariance matrices $\text{Cov}(Y_i)$ for each voxel $i$. Estimating the covariance matrices assuming some probability distribution for the movement of the voxels and some model for the dose deposition is possible, but the biggest challenge is storing these matrices. In practice, $\text{Cov}(Y_i)$ is a $d \times d$ matrix, where $d$ is the number of beamlets,which is on the order of one thousand. There can easily be tens of thousands of voxels, so storing a covariance matrix for each voxel requires too much memory for this formulation to be tractable. In order to make the problem tractable, the following model of the random dose was adopted in [13, 12].

Suppose that on any single treatment, one of $m$ possible scenarios $s_1, ..., s_m$ can occur with probabilities $p_1, .., p_m$ respectively. Let $a_{ij}$ be a column vector, indexed by beamlets, giving the dose delivered to voxel $i$ in scenario $j$ from each beamlet, when the beamlets have unit intensity.

Let $p$ be the column vector of scenario probabilities, and

$$A_i = \begin{bmatrix} a_{i,1}^T \\ a_{i,2}^T \\ \vdots \\ a_{i,m}^T \end{bmatrix}$$

be a matrix where the $j$-th row contains the vector giving the dose to voxel $i$ in scenario $j$ from each beamlet. Then the mean of the dose to voxel $i$ is

$$E[D_{i1}(x)] = p^T A_i x.$$

Let $e$ denote an $m \times 1$ vector where each element is equal to 1, let $I$ denote the $m \times m$ identity matrix and $P$ denote the diagonal matrix where $P_{jj} = p_j$. Then

$$\text{Var}[D_{i1}(x)] = \left[A_i x - e(p^T A_i x)\right]^T P \left[A_i x - e(p^T A_i x)\right]$$

$$= \|RA_i x\|^2,$$

where $R = P^{1/2}(I - ep^T)$.

Putting it all together, we see that constraint (2.2) can be written as

$$\|RA_i x\|_2 \leq \frac{M_1 - N p^T A_i x}{\Phi^{-1}(1 - \delta) \sqrt{N}}. \tag{2.6}$$

Constraint (2.3) can be expressed in a similar way. Finally, further constraints are imposed on the maximum total dose to target voxels with high probability, the minimum dose per scenario to target voxels and dose-volume constraints of the form "no more than a specific percentage of a healthy structure may receive more than a given dose". The objective function is then written as the sum of weighted penalties, penalizing the violation of the constraints.

This formulation relies on approximating the dose distribution, which is in general continuous, by a discrete distribution. This discrete distribution depends on the scenarios that we choose and the probability assigned to each of

them. It is clear that the proposed approximation to the original formulation is much more tractable. Instead of storing a $d \times d$ covariance matrix for each constraint, the approximation requires only an $m \times d$ matrix.

In each constraint, the original covariance matrix $W_i = \text{Cov}(Y_i)$ is replaced by the covariance matrix $C_i = A_i^T R^2 A_i$ with respect to the discrete distribution. Given the fact that we have $m$ scenarios, $C_i$ is of rank at most $m$. In the patient case that was solved in [13, 12], $m$ was on the order of ten. From a geometric point of view this means that for each $i$, the convex quadratic constraint defined by the matrix $W_i$ is replaced by a convex quadratic constraint that is defined by the matrix $C_i$. Intuitively, the quality of the optimal solution depends on how close the two constraints are, or in other words by how well $W_i$ can be approximated by a low-rank approximation.

In [13], the problem was solved using various numbers of scenarios, and the quality of the solution was evaluated with a Dose-Expected Volume Histogram (DEVH). Briefly, a DEVH plots the expected volume of a structure that receives at least some dose. All solutions were found to be of similar quality, achieving sparing of the healthy tissue, while delivering a high dose to the target volume.

## 2.2 Stylized model

The success of the low-rank formulation in the robust IMRT problem motivates us to study the effect of low-rank approximations to optimization problems of the same form. Our goal is to study the relation between the optimal values of the original and the approximating problem. More specifically, we want to analyze theoretically how this relation depends on the shape of the constraints,

the objective function, the type of the low-rank approximation that is used and the number of constraints in the problem.

The reasons for using a stylized model are twofold. First, we want to work with an optimization problem that is simple enough to be analyzed, while sharing the fundamental properties of the IMRT problem. Second, using a stylized model allows us to illuminate how the various parameters of the problem, such as the dimension, the singular values of the matrices, the number of constraints and the objective function, affect the behavior of the low-rank approximation.

Before presenting the model, we take a closer look at the constraints of the original problem in the patient case studied in [13]. As we saw in Section 2.1, the constraints are of the form

$$\sqrt{x^T W_i x} \le a_i + b_i^T x.$$

First, we examine the covariance matrices $W_i$ of the doses in 100 voxels of size 1cm from the area of interest. We estimate the covariance matrices by sampling 3000 shifts from a uniform distribution in a cube centered at the origin, with edges parallel to the axes and with edge length equal to 2cm. Shifts in this context correspond to rigid body deformations, i.e., we assume that the entire patient is moved this much. The resulting matrices have a small number of significant singular values. Furthermore, as one can see from Figure 2.1, the matrices have singular values that decay in a similar way.

Using the same sampled shifts, we also examine the means of the doses in each voxel. For each voxel, most entries of the mean dose vector are near zero, except a few positive entries. A plot of the mean dose vector with respect to beamlets for a single voxel is given in Figure 2.2.

Figure 2.1: Plot of the 20 largest singular values of the constraint matrices in the IMRT problem.

The first assumption that we introduce concerns the singular values of the matrices $W_i, i = 1, ..., n$ in the constraints. We assume in our model that the matrices $W_i, i = 1, ..., n$ have the same singular values. This seems reasonable, since the estimated matrices in our experimental result have singular values that decay in a similar way. In addition, this assumption allows us to emphasize the dependence of the behavior of the low-rank approximations in our model on the singular values of the constraint matrices.

Our second assumption concerns the linear term $b_i^T x$ in the constraints. We assume that $b_i = 0$, $i = 1, ..., n$. Although this is not a close approximation of the constraints in the IMRT problem, we introduce this assumption in order to keep the model simple and to make the results and conclusions transparent. This as-

Figure 2.2: Plot of the average dose to a specific voxel from each beamlet

sumption implies that all constraints in our model are ellipsoidal cylinders centered at 0. The methodology that we develop can be used for analyzing models of the same form where the centers of the ellipsoidal cylinders are sampled randomly uniformly in a ball centered at the origin. In that case, the formulas are similar to the ones that we provide here, but more involved.

Under these assumptions, the stylized optimization problem takes the form

$$\max \ c^T x \tag{2.7}$$

$$\text{subject to} \ \ x^T W_i x \leq a_i^2, i = 1, ..., n.$$

Let $W_i = Q_i^T A Q_i$ be the SVD of the covariance matrices in the constraints for $i = 1, ..., n.$, where $A$ does not depend on $i$ because we have assumed that the matrices $W_i$, $i = 1, ..., n$ have the same singular values. Without loss of generality we assume that the largest diagonal entry of $A$ is equal to 1 and we consider

15

matrices $A$ such that the diagonal elements are in decreasing order, i.e., $1 = A_{11} \geq A_{22} \geq \cdots \geq A_{dd} \geq 0$. Furthermore, we assume that $a_i = 1, i = 1, ..., n$, so that all constraints are ellipsoidal cylinders of the same shape, rotated around the origin. Without loss of generality we choose $c \in \mathbb{R}^d$ such that $\|c\|_2 = 1$.

The final assumption concerns the orthogonal matrices $Q_i, i = 1, ..., n$. For our stylized model, we assume that $Q_i, i = 1, ..., n$ are independent random matrices and they follow the orthogonally invariant distribution in the set $O^d$ of $d \times d$ orthogonal matrices. This implies that the constraints are ellipsoids of the same shape, randomly rotated around the origin. The orthogonally invariant distribution is essentially a uniform distribution over orthogonal matrices and is invariant under left and right orthogonal multiplication, i.e., the measure of

$$QX = \{QX : X \in X\}$$

for any measurable $X \subset O^d$ and $Q \in O^d$ is equal to that of $X$. In other words, it is the Haar measure on the set $O^d$. For more information see [25]. A simple algorithm for sampling from this distribution requires generating a random matrix $U$ with independent standard normal entries and then calculating its QR decomposition, $QR = U$. The matrix $Q$ follows the orthogonally invariant distribution. A proof of this can be found in Stewart, [34]. This choice of distribution is crucial for the analysis that follows, because it allows us to use the orthogonal invariance property.

Thus, our model for the original formulation of the robust IMRT problem is the sampled problem

$$\max \quad c^T x \tag{2.8}$$

$$\text{subject to} \quad x^T W_i x \leq 1, i = 1, ..., n,$$

where $W_i = Q_i^T A Q_i$ and $n$ is the total number of constraints. Let $V_1$ be the optimal value of Problem (2.8). The feasible region of (2.8) is defined as the intersection of ellipsoids that contain the unit ball in $\mathbb{R}^d$. So, we easily get the inequality $V_1 \geq 1$.

We next introduce the approximation to (2.8), which is written as

$$\max \quad c^T x \tag{2.9}$$

$$\text{subject to} \quad x^T C_i x \leq 1, i = 1, ..., n.$$

The matrices $C_i$ are symmetric positive semidefinite low-rank approximations to $W_i$ for each $i = 1, ..., n$. We haven't yet specified the relationship between the matrices $C_i$ and $W_i$. In what follows we will assume that one has access to the matrices $W_i$ and then some algorithm is used to calculate the low-rank approximation $C_i$. We further assume that all matrices $C_i, i = 1, ...n$ are of the same rank.

The optimal values of Problems (2.8) and (2.9) are random variables defined on the same probability space. The probability distribution of $V_1$ depends on the matrix $A$ and also on the number of constraints $n$. The distribution of the optimal value of the Problem (2.9) depends additionally on the type of approximation used and on the rank of the matrices $C_i, i = 1, ..., n$. We are interested in deriving properties of these optimal values and also of their relative difference.

The stylized optimization problems that we have created are sampled optimization problems with independent constraints. Our goal is to derive properties of the optimal value of such problems. Recent research, [8, 9, 14, 21], has focused on using sampled optimization problems, which are straightforward to solve in practice, in order to approximate the much harder to solve in practice

chance-constrained optimization problems. Our method for analyzing sampled problems is based on the connection between sampled and chance-constrained problems. In the following sections we present the relevant results from chance-constrained optimization theory. For a comprehensive discussion of the theory and applications of such problems, see [33, 32].

## 2.3   Chance-constrained optimization theory

Chance-constrained optimization theory has a long history, dating back to the work of Charnes and Cooper for linear programs in 1959, [11] and has found applications in areas such as portfolio optimization under value-at-risk constraints, staffing of call centers and emergency services. The main goal of this approach is to reduce sensitivity in the solution of optimization problems with respect to some unknown parameter. It assumes that the parameters are distributed according to some known probability $\mathbb{P}$ on a set $\Xi$. The goal is not to satisfy all constraints, but to find a solution that violates a set of constraints that has small $\mathbb{P}$-probability. A chance-constrained problem with linear objective function and convex constraints is given in general form by

$$\max \quad c^T x \tag{2.10}$$

$$\text{subject to} \quad \mathbb{P}\left(f(x, \xi) > 0\right) \leq \epsilon.$$

where $x \in \mathbb{R}^d$ is the optimization variable, $f(x, \xi)$ is a convex function in $x$ for all $\xi \in \Xi$, and $\xi$ is the unknown random parameter that is assumed to lie in the set $\Xi$ The parameter $\epsilon \in (0, 1)$ controls the probability that the optimal solution of (2.10) violates the constraints.

The main drawback of chance-constrained optimization problems is that

18

they are extremely hard to solve in practice. Even if the function $f(x, \xi)$ is convex, the feasible region of (2.10) is not convex in general. Moreover, merely calculating the probability $\mathbb{P}(f(x, \xi) > 0)$ for a fixed $x$ is a nontrivial task that typically involves calculating a multidimensional integral. Another weakness of chance-constrained programs lies in the fact that knowledge of the probability measure $\mathbb{P}$ on the set of parameters $\Xi$ is assumed. In a practical setting $\mathbb{P}$ would have to be estimated, introducing another source of sensitivity in the optimal solution. For more details see [32],[35].

Very closely related to chance-constrained optimization theory is robust optimization, introduced by Ben-Tal and Nemirovski, [4, 5, 6]. Robust optimization gives a similar framework for reducing the sensitivity of the optimal solution of optimization problems with respect to uncertainty in parameter values. In this framework one seeks a solution which simultaneously satisfies all possible constraint instances. In general form, a robust optimization problem with linear objective function and convex constraints is given by

$$\max \quad c^T x \tag{2.11}$$

$$\text{subject to} \quad f(x, \xi) \leq 0, \quad \forall \xi \in \Xi,$$

where $x \in \mathbb{R}^d$ is the optimization variable, $f(x, \xi)$ is a convex function in $x$ for all $\xi \in \Xi$, and $\xi$ is the unknown parameter that is assumed to lie in the set $\Xi$. Problems of this type can include an infinite number of constraints. In special cases and under some regularity conditions they can be solved by reformulating the constraints in (2.11) as a finite collection of constraints.

Motivated by the computational complexity of chance-constrained problems, Calafiore and Campi [8, 9] and de Farias and Van Roy [14] independently proposed tractable approximations. The idea is to approximate Problem (2.10)

by sampling independent and identically distributed parameters $\xi_i$ under the distribution $\mathbb{P}$ and to solve instead the much easier sampled problem

$$\max \quad c^T x \tag{2.12}$$

$$\text{subject to} \quad f(x, \xi_i) \leq 0, \quad i = 1, ..., n.$$

De Farias and Van Roy, [14], study problems where the constraint function $f(\cdot, \xi)$ is linear. They use results from Computational Learning Theory to give a lower bound on the number of sampled constraints needed in order to guarantee that the feasible region of Problem (2.10) is included in the feasible region of (2.11) with probability at least $1 - \delta$. Calafiore and Campi [8, 9] consider general convex functions $f(x, \xi)$ and provide a similar bound on the number of constraints that need to be sampled so that the optimal solution of the sampled Problem (2.12) is feasible for (2.10) with probability at least $1 - \delta$.

## 2.4   Fundamental inequality

Before proceeding with the fundamental inequality, we review the results in [8, 9] in more detail. Let $\hat{x}$ and $V$ be the optimal solution and the optimal value of the sampled Problem (2.12). For any $\epsilon \in (0, 1)$ we denote by $\hat{x}(\epsilon)$, $G(\epsilon)$ and $X(\epsilon)$ the optimal solution, the optimal value and the feasible region of the chance-constrained Problem (2.10) respectively. Also, we define $\hat{x}_k$ to be the optimal solution to the sampled problem that is obtained if we remove the $k$-th constraint from (2.12).

**Definition 2.4.1** (Support Constraint) *The k-th constraint $f(x, \xi_k) \leq 0$ is called a support constraint for the Problem (2.12) if $c^T \hat{x}_k > c^T \hat{x}$.*

The results are based on a fundamental property of convex programs given in the following theorem.

**Theorem 2.4.1 (Theorem 2 in [8])** *A convex program in $\mathbb{R}^d$ has at most $d$ support contraints.*

Using this, they prove the following result. Since this is the basic building block of our main result, Theorem 2.4.2, we give here the proof as presented in [21].

**Proposition 2.4.1 (Theorem 1 in [9])** *Fix $\epsilon > 0$. Let $\hat{x}$ denote the optimal solution of the sampled Problem (2.12). Then*

$$P\left(\hat{x} \notin X(\epsilon)\right) \leq \binom{n}{d}(1 - \epsilon)^{n-d}. \tag{2.13}$$

**Proof:**

The sampled Problem (2.12) is a convex program in $\mathbb{R}^d$ with $n$ constraints. Let $I \subseteq \{1, 2, ..., n\}$, with $|I| = d$. Let

$$\Xi_I^n = \left\{(\xi_1, \xi_2, ..., \xi_n) : \text{ all the support constraints } \subseteq I\right\}.$$

Then Theorem 2.4.1 implies that $\Xi^n = \Xi \times \Xi \cdots \times \Xi$ can be expressed as

$$\Xi^n = \cup_{\{I \subseteq \{1,...,n\}:|I|=d\}} \Xi_I^n.$$

We define $\hat{x}_I$ to be the optimal solution of the sampled problem with only the samples $i \in I$ present, and $\mathcal{A}_I$ to be the event

$$\mathcal{A}_I = \{(\xi_i)_{i \in I} : \hat{x}_I \notin X(\epsilon)\}.$$

21

We have

$$\mathbb{P}((\xi_1, ..., \xi_n) : \hat{x} \notin \mathcal{X}(\epsilon)) \leq \sum_{\{I \subseteq \{1,...,n\}:|I|=d\}} \mathbb{P}((\xi_1, ..., \xi_n) \in \Xi_I^n : \hat{x}_I \notin \mathcal{X}(\epsilon))$$

$$= \sum_{\{I \subseteq \{1,...,n\}:|I|=d\}} \mathbb{P}(\mathcal{A}_I) \, \mathbb{P}((\xi)_{i \notin I} : f(\hat{x}_I, \xi_i) \leq 0 | \mathcal{A}_I)$$

$$= \sum_{\{I \subseteq \{1,...,n\}:|I|=d\}} \mathbb{P}(\mathcal{A}_I) \prod_{i \notin I} \mathbb{P}(\xi_i : f(\hat{x}_I, \xi_i) \leq 0 | \mathcal{A}_I),$$

where each probability in the sum can be written as a product because $\{\xi_i\}_{i=1}^n$ are i.i.d. samples. Since $\hat{x}_I \notin \mathcal{X}(\epsilon)$, it follows that for $i \notin I$,

$$\mathbb{P}(\xi_i : f(\hat{x}_I, \xi_i) \leq 0 \,|\, \mathcal{A}_I) \leq 1 - \epsilon.$$

Thus,

$$\mathbb{P}((\xi_1, ..., \xi_n) : \hat{x} \notin \mathcal{X}(\epsilon)) \leq (1 - \epsilon)^{n-d} \sum_{\{I \subseteq \{1,...,n\}:|I|=d\}} \mathbb{P}((\xi_i)_{i \in I} : \hat{x}_I \notin \mathcal{X}(\epsilon))$$

$$\leq \binom{n}{d}(1 - \epsilon)^{n-d}.$$

$\square$

Proposition 2.4.1 gives an upper bound on the probability that the optimal solution of the sampled problem is infeasible for the corresponding chance-constrained problem. By inverting inequality (2.13), one can get a lower bound on the number of constraints that need to be sampled so that the optimal solution of the sampled problem is infeasible for the chance-constrained problem with small probability.

The following result gives an upper bound on the probability that the optimal value of the sampled problem is greater than the optimal value of the chance-constrained problem. The main argument of the proof is based on looking for feasible points for the sampled problem, in the direction of the optimal solution of the chance-constrained problem.

**Proposition 2.4.2 (Theorem 2 in [9])** *Let V be the optimal value of the sampled problem and for $\epsilon > 0$ let $G(\epsilon)$ be the optimal value of the corresponding chance-constrained problem. Then we have*

$$\mathbb{P}(V \geq G(\epsilon)) \geq (1 - \epsilon)^n.$$

We can now combine Propositions 2.4.1 and 2.4.2 to get the following result that describes the behavior of the tail probability of the optimal value $V$ of the sampled problem.

**Theorem 2.4.2** *Assume that for all $v \in \mathbb{R}$ we have $\mathbb{P}(V = v) = 0$. Then, for any $\epsilon \in (0, 1)$ we have*

$$(1 - \epsilon)^n \leq \mathbb{P}(V > G(\epsilon)) \leq \binom{n}{d}(1 - \epsilon)^{n-d}. \tag{2.14}$$

*Let $\mathcal{G} = (\lim_{\epsilon \downarrow 0} G(\epsilon), \lim_{\epsilon \uparrow 1} G(\epsilon))$. Then, for any $v \in \mathcal{G}$, there exists a $G^{-1}(v) \in (0, 1)$, where $G^{-1}(\cdot)$ is a left inverse of $G(\cdot)$, such that*

$$(1 - G^{-1}(v))^n \leq \mathbb{P}(V > v) \leq \binom{n}{d}(1 - G^{-1}(v))^{n-d}. \tag{2.15}$$

*Also, if $V(n)$ is the optimal value of Problem (2.12) with n constraints, then the sequence $\{V(n)\}_{n=1}^{\infty}$ of random variables satisfies*

$$\lim_{n \to \infty} \frac{\log \mathbb{P}(V(n) > v)}{n} = \log(1 - G^{-1}(v)),$$

*for any $v$ in $\mathcal{G}$.*

**Proof:**

Since $G(\epsilon)$ is the optimal value of the chance-constrained problem, using Proposition 2.4.1 we get

$$\mathbb{P}(V > G(\epsilon)) \leq \mathbb{P}(\hat{x} \notin X(\epsilon)) \leq \binom{n}{d}(1 - \epsilon)^{n-d}.$$

From Proposition 2.4.2 and using our assumption we have

$$\mathbb{P}(V > G(\epsilon)) = \mathbb{P}(V \geq G(\epsilon)) \geq (1 - \epsilon)^n.$$

We thus get that

$$(1 - \epsilon)^n \leq \mathbb{P}(V > G(\epsilon)) \leq \binom{n}{d}(1 - \epsilon)^{n-d}.$$

Let $\mathcal{G} = (\lim_{\epsilon \downarrow 0} G(\epsilon), \lim_{\epsilon \uparrow 1} G(\epsilon))$ and consider an arbitrary $v \in \mathcal{G}$. We first prove that there exist no nontrivial intervals of constancy for $G$. If there exists an interval $[\epsilon_1, \epsilon_2]$ such that for all $\epsilon \in [\epsilon_1, \epsilon_2]$, $G(\epsilon) = v$, we get that

$$(1 - \epsilon_1)^n \leq \mathbb{P}(V > v) \leq \binom{n}{d}(1 - \epsilon_2)^{n-d}.$$

Considering large $n$ leads to a contradiction.

If there exists a unique $\epsilon(v)$ such that

$$G(\epsilon(v)) = v,$$

we get that

$$(1 - \epsilon(v))^n \leq P(V > v) \leq \binom{n}{d}(1 - \epsilon(v))^{n-d}.$$

Otherwise, we define

$$\epsilon_1(v) = \sup \{\epsilon | G(\epsilon) \leq v\}$$

and

$$\epsilon_2(v) = \inf \{\epsilon | G(\epsilon) \geq v\}.$$

From the definition and since there is no $\epsilon$ such that $G(\epsilon) = v$ we have $\epsilon_1(v) = \epsilon_2(v)$. We have for any $\omega < \epsilon_1(v)$ and $\zeta > \epsilon_2(v)$ that

$$(1 - \zeta)^n \leq \mathbb{P}(V > v) \leq \binom{n}{d}(1 - \omega)^{n-d}.$$

This implies that

$$(1 - \epsilon(v))^n \leq \mathbb{P}(V > v) \leq \binom{n}{d}(1 - \epsilon(v))^{n-d}, \tag{2.16}$$

where $\epsilon(v) = \epsilon_1(v) = \epsilon_2(v)$.

In both cases we set $G^{-1}(v) = \epsilon(v)$ and so, we get the inequality

$$\left(1 - G^{-1}(v)\right)^n \leq \mathbb{P}(V > v) \leq \binom{n}{d}\left(1 - G^{-1}(v)\right)^{n-d} \tag{2.17}$$

for any $v \in \mathcal{G}$.

Recall that we denote by $V(n)$ the optimal value of the sampled problem with $n$ constraints. From (2.17) we easily get that

$$\lim_{n \to \infty} \frac{\log \mathbb{P}(V(n) > v)}{n} = \log(1 - G^{-1}(v)).$$

$\square$

Theorem 2.4.2 gives a characterization of the tail probability of the optimal value $V$ of the sampled problem that depends on the optimal value $G(\epsilon)$ of the corresponding chance-constrained problem. In order to apply this theorem to a specific sampled optimization problem, we must have some information about $G(\epsilon)$, which, as we explained previously, is very hard to obtain in general. In the type of problems that we will analyze in subsequent chapters, we can either explicitly solve the related chance-constrained problems, or get asymptotic results for their optimal values.

# CHAPTER 3

## THE SVD APPROXIMATION

## 3.1   The approximation

In this chapter we present the first approximation that we will study in our model. It is based on the SVD of the matrices $W_i$ in the optimization Problem (2.8). We start with some motivation for using this approximation. We then present the approximate optimization problem and apply the method that we developed in Chapter 2.

Recall that the initial optimization problem in our model is given by

$$\max \ \ c^T x$$

$$\text{subject to} \ \ x^T W_i x \leq 1, i = 1, ..., n,$$

where $W_i = Q_i^T A Q_i$. The matrices $\{Q_i\}_{i=1}^n$ are independent random orthogonal matrices, uniformly distributed in the set $O^d$ of orthogonal matrices in $\mathbb{R}^{d \times d}$ and $A$ is a diagonal matrix such that

$$1 = A_{11} \geq A_{22} \geq \cdots \geq A_{dd} \geq 0.$$

We assume that the matrices $W_i, i = 1, .., n$ are approximated by the matrices

$$C_i = Q_i^T A^{(k)} Q_i, \ i = 1, ..., n.,$$

where $A^{(k)} = \text{diag}(A_{11}, ..., A_{kk}, 0, ..., 0)$.

The use of this approximation is appealing for many reasons. First, $C_i$ is a

solution to the problem

$$\min \ \|W_i - X\| \tag{3.1}$$

$$\text{subject to} \ \ \text{rank}(X) \leq k, X \in \mathbb{R}^{d \times d},$$

where $\|\cdot\|$ denotes either the Frobenius or the 2-norm. Furthermore, $C_i$ is a symmetric positive semidefinite matrix, and thus the feasible region of the approximate problem is an intersection of constraints of the same type as the constraints of the original problem. Finally, such an approximation is appealing from a geometric point of view. Each constraint is replaced by an ellipsoidal cylinder that has the same principal directions as the corresponding original constraint and is unbounded along the principal directions with the $d - k$ largest semi-axes. A simple example with 2 constraints can be seen in Figure 3.1.



Figure 3.1: Illustration of the approximation in a sampled problem in $d = 2$ dimensions and with $n = 2$ constraints.

We thus have the following sampled problems in our model, an original

27

problem of the form

$$\max \quad c^T x \tag{3.2}$$

$$\text{subject to} \quad x^T Q_i^T A Q_i x \le 1, i = 1, ..., n$$

and the approximation to (3.2) given by

$$\max \quad c^T x \tag{3.3}$$

$$\text{subject to} \quad x^T Q_i^T A^{(k)} Q_i x \le 1, i = 1, ..., n.$$

We denote by $V_1$ and $V_2$ the optimal values of Problems (3.2) and (3.3) respectively. We define the relative error of the approximation to be

$$R_2 = \frac{V_2 - V_1}{V_1}. \tag{3.4}$$

Problems (3.2) and (3.3) have constraints of the same form, the only difference being the substitution of the diagonal matrix $A$ with the matrix $A^{(k)}$. This implies that the probability distributions of the optimal values $V_1$ and $V_2$ are of the same form, but with different parameters. In what follows in this chapter, we focus on Problem (3.2), since the same analysis holds for the approximating Problem (3.3).

## 3.2 Main result

In this section we apply the method that we developed in Section 2.4 to the optimal value $V_1$ of Problem (3.2). We first formulate and solve the chance-constrained problem corresponding to the sample Problem (3.2). We then give an analytic expression for the optimal value of the chance-constrained problem and apply Theorem 2.4.2 to the optimal value $V_1$ of Problem (3.2). We conclude

with an asymptotic result for the optimal value of the chance-constrained prob-
lem, which will be used in Chapter 5.

We begin by presenting a robust optimization version of Problem (3.2) with
an uncountable number of constraints, which can be thought of as a limiting
case of the sampled problem with infinite number of constraints. This is the
problem

$$\max \ c^T x \tag{3.5}$$

$$\text{subject to} \ x^T Q^T A Q x \leq 1, \forall Q \in O^d.$$

Problem (3.5) has an infinite number of constraints which can be equivalently
written as

$$\sup \left\{ x^T Q^T A Q x | Q \in O^d \right\} \leq 1$$

$$\Leftrightarrow \|x\|_2^2 \sup \left\{ u^T A u \mid \|u\|_2 = 1 \right\} \leq 1$$

$$\Leftrightarrow \|x\|_2^2 \leq 1,$$

where the last inequality follows from the fact that $\|A\|_2 = 1$. Thus the feasible
set of (3.5) is the closed unit ball in $\mathbb{R}^d$, the optimal solution is $c$ and the optimal
value is equal to 1. Intuitively, the feasible region of the chance-constrained and
the sampled problem resemble the unit ball for small $\epsilon$ and a large number of
constraints $n$ respectively.

The chance-constrained problem corresponding to Problem (3.2) for $\epsilon \in (0, 1)$
is given by

$$\max \ c^T x \tag{3.6}$$

$$\text{subject to} \ x \in \mathcal{X}_1(\epsilon),$$

where

$$\mathcal{X}_1(\epsilon) = \left\{ x \in \mathbb{R}^d \mid \mathbb{P}(x^T Q^T A Q x > 1) \leq \epsilon \right\}.$$

29

It turns out that due to the symmetry in our problem, the feasible region $\mathcal{X}_1(\epsilon)$ is simple as well. We have the following proposition.

**Proposition 3.2.1** *Let $\mathcal{X}_1(\epsilon)$ and $G_1(\epsilon)$ be the feasible region and the optimal value of the chance-constrained Problem (3.6) respectively. Then $\mathcal{X}_1(\epsilon)$ is a closed ball centered at 0 and the optimal value $G_1(\epsilon)$ satisfies the equation*

$$P\left( \frac{\sum_{j=1}^d A_{jj} Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{G_1^2(\epsilon)} \right) = \epsilon, \tag{3.7}$$

*where $\left\{Y_j\right\}_{j=1}^d$ are independent standard normal random variables.*

**Proof:** Our first observation is that the feasible region $\mathcal{X}_1(\epsilon)$ of the chance-constrained problem is orthogonally invariant. This is a direct consequence of the orthogonal invariance property of the distribution of the orthogonal matrix $Q$. Also, for any $x \in \mathcal{X}_1(\epsilon)$ and any $\lambda \in [0, 1]$ we have $\lambda x \in \mathcal{X}_1(\epsilon)$, so the feasible region $\mathcal{X}_1(\epsilon)$ is a ball centered at 0. This allows us to easily find the optimal solution of the chance-constrained problem.

Our next step is to find the radius of $\mathcal{X}_1(\epsilon)$. We use the fact that for any unit vector $u \in \mathbb{R}^d$, the vector $s = Qu$ is uniformly distributed on the unit sphere in $\mathbb{R}^d$. Let $F(z) = P(s^T A s \leq z)$ be the distribution function of the random variable $s^T A s$. We express the radius in terms of the function $F(\cdot)$.

Let $\lambda > 0$ be such that $\lambda e_1 \in \mathcal{X}_1(\epsilon)$. We have $\lambda e_1 \in \mathcal{X}_1(\epsilon)$ if and only if

$$\mathbb{P}(\lambda^2 e_1^T Q^T A Q e_1 > 1) \leq \epsilon,$$

which holds if and only if

$$\lambda \leq \sqrt{\frac{1}{F^{-1}(1 - \epsilon)}}.$$

From this we see that

$$\mathcal{X}_1(\epsilon) = B\left(0, \sqrt{\frac{1}{F^{-1}(1-\epsilon)}}\right).$$

The optimal solution is

$$\hat{x}_1(\epsilon) = \sqrt{\frac{1}{F^{-1}(1-\epsilon)}} \; c$$

and the optimal value is

$$G_1(\epsilon) = \sqrt{\frac{1}{F^{-1}(1-\epsilon)}}.$$

For a point $s$ that is uniformly distributed on the boundary of the unit ball we have

$$s \overset{D}{=} \left(\frac{Y_1}{\|Y\|_2}, \frac{Y_2}{\|Y\|_2}, \dots, \frac{Y_d}{\|Y\|_2}\right)^T,$$

where $\{Y_j\}_{j=1}^{d}$ are independent standard normal random variables and $\overset{D}{=}$ denotes equality in distribution. We thus get that $G_1(\epsilon)$ satisfies the equation

$$P\left(\frac{\sum_{j=1}^{d} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{G_1^2(\epsilon)}\right) = \epsilon. \qquad (3.8)$$

For the function $G_1^{-1}$ we get the expression

$$P\left(\frac{\sum_{j=1}^{d} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{v^2}\right) = G_1^{-1}(v). \qquad (3.9)$$

$\square$

Proposition 3.2.1 gives us an analytic expression for the optimal value $G_1(\epsilon)$ of the $\epsilon$-chance-constrained problem that involves the diagonal entries of the matrix $A$ and the ratio of sums of squared normal random variables. The optimal value $G_1(\epsilon)$ does not depend on the objective function vector $c$, because of the symmetry in the distribution of the constraints. Figure 3.2 shows the robust problem, the chance-constrained and a sampled problem in 2 dimensions

Combining Proposition (3.2.1) and Theorem (2.4.2) we easily get the following result for the optimal value $V_1$.

Figure 3.2: Feasible regions of the robust, sampled and chance-constrained problems in 2 dimensions. The thicker lines give the constraints of the sampled problem and the dotted line represents the boundary of the feasible region of the chance-constrained problem.

**Proposition 3.2.2** *Let $V_1$ be the optimal value of Problem (3.2). We have that for any $v > 1$,*

$$(1 - G_1^{-1}(v))^n \le \mathbb{P}(V_1 > v) \le \binom{n}{d}(1 - G_1^{-1}(v))^{n-d}, \tag{3.10}$$

*where*

$$G_1^{-1}(v) = P\left(\frac{\sum_{j=1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right)$$

*and $\{Y_j\}_{j=1}^d$ are independent standard normal random variables.*

Our next goal is to get a better understanding of the behavior of the function $G_1^{-1}(\cdot)$, which is given by (3.9). We will focus on the asymptotic behavior of this function as $v$ approaches 1. As we will show in Chapter 5, the asymptotic

behavior of $G_1^{-1}(\cdot)$ near 1, characterizes the asymptotic behavior of the optimal value $V_1$ of the sampled problem when the number of constraints $n$ tends to infinity.

Before we proceed with our analysis, we recall a few definitions related to the asymptotic behavior of real functions and sequences of random variables.

**Definition 3.2.1** *Let $f, g$ be real functions. We then write $f(x) = O(g(x))$ as $x \to \infty$ if for any real function $\alpha(\cdot)$ such that $\alpha(x) \to \infty$, we have $\frac{f(x)}{g(x)} \frac{1}{\alpha(x)} \to 0$. Similarly we have $f(x) = \Omega(g(x))$, if $\frac{f(x)}{g(x)} \alpha(x) \to \infty$ and $f(x) = \Theta(g(x))$ if $f(x) = \Omega(g(x))$ and $f(x) = O(g(x))$.*

This notation can also be used when studying the behavior of the function $f$ around a point $a$, by simply replacing the limits above with limits as $x \to a$.

The following lemma gives the asymptotic behavior of $G_1^{-1}(v)$ as $v \downarrow 1$, for a special choice of the diagonal matrix $A$.

**Lemma 3.2.1** *Consider the function $G_1^{-1}(\cdot)$ corresponding to the chance-constrained Problem (3.6) with*

$$A = \begin{bmatrix} I_k & 0 \\ 0 & \theta I_{d-k} \end{bmatrix}, \theta \in [0, 1).$$

*Then*

$$\lim_{v \downarrow 1} \frac{G_1^{-1}(v)}{(v-1)^{(d-k)/2}} = \left( \frac{2}{1-\theta} \right)^{(d-k)/2} \frac{2}{(d-k)B\left( \frac{d-k}{2}, \frac{k}{2} \right)},$$

*where $B(\alpha, \beta)$ denotes the beta function.*

**Proof:** From (3.9) we have for $v < \frac{1}{\sqrt{\theta}}$ that

$$G_1^{-1}(v) = P\left(\frac{\sum_{j=1}^{k} Y_j^2 + \theta \sum_{j=k+1}^{d} Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{v^2}\right)$$

$$= P\left((1 - v^{-2}) \sum_{j=1}^{k} Y_j^2 > (v^{-2} - \theta) \sum_{j=k+1}^{d} Y_j^2\right)$$

$$= P\left(\frac{\sum_{j=k+1}^{d} Y_j^2}{\sum_{j=1}^{k} Y_j^2} < \frac{1 - v^{-2}}{v^{-2} - \theta}\right)$$

$$= P\left(\frac{\sum_{j=k+1}^{d} Y_j^2/(d-k)}{\sum_{j=1}^{k} Y_j^2/k} < \frac{k}{d-k}\frac{v^2 - 1}{1 - \theta v^2}\right).$$

The random variable

$$\frac{\sum_{j=k+1}^{d} Y_j^2/(d-k)}{\sum_{j=1}^{k} Y_j^2/k}$$

follows the $F$ distribution with $d - k, k$ degrees of freedom. Because of the connection between the distribution function of an $F$ random variable and the incomplete beta function, see [1], we get that

$$G_1^{-1}(v) = I\left(\frac{1 - v^{-2}}{1 - \theta}; \frac{d-k}{2}, \frac{k}{2}\right),$$

where

$$I(x; \alpha, \beta) = \frac{\int_0^x t^{\alpha-1}(1 - t)^{\beta-1}dt}{B(\alpha, \beta)}$$

and

$$B(\alpha, \beta) = \int_0^1 t^{\alpha-1}(1 - t)^{\beta-1}dt.$$

It is easy to see that

$$\lim_{x \downarrow 0} \frac{I(x; \alpha, \beta)}{x^\alpha} = \frac{1}{\alpha B(\alpha, \beta)}.$$

We thus get that

$$\lim_{v \downarrow 1} \frac{G_1^{-1}(v)}{(v - 1)^{(d-k)/2}} = \left(\frac{2}{1 - \theta}\right)^{(d-k)/2} \frac{2}{(d - k)B\left(\frac{d-k}{2}, \frac{k}{2}\right)}.$$

$\square$

34

Using Lemma 3.2.1 we prove the following result.

**Proposition 3.2.3** *Consider the chance-constrained Problem (3.6) and let l be the number of diagonal elements of A that are equal to 1. Then we have that as $v \downarrow 1$,*

$$G_1^{-1}(v) = \Theta\left((v-1)^{(d-l)/2}\right).$$

**Proof:**

We have that

$$G_1^{-1}(v) = P\left(\frac{\sum_{j=1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right)$$

$$= P\left(\frac{\sum_{j=1}^l Y_j^2 + \sum_{j=l+1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right).$$

It is easy to see that

$$m(v) \le G_1^{-1}(v) \le M(v),$$

where

$$m(v) = P\left(\frac{\sum_{j=1}^l Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right)$$

and

$$M(v) = P\left(\frac{\sum_{j=1}^l Y_j^2 + \sum_{j=l+1}^d A_{l+1,l+1}Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right).$$

From Lemma 3.2.1 we have that

$$m(v) = \Theta\left((v-1)^{(d-l)/2}\right)$$

and

$$M(v) = \Theta\left((v-1)^{(d-l)/2}\right).$$

We thus get that

$$G_1^{-1}(v) = \Theta\left((v-1)^{(d-l)/2}\right).$$

□

Proposition 3.2.3 gives the asymptotic behavior of the function $G_1^{-1}(\cdot)$, close to 1. The important fact is that the asymptotic rate depends only on the number $l$ of elements equal to 1 in the diagonal matrix $A$. As we will see in Chapter 5, the behavior of $G_1^{-1}(v)$ close to 1 characterizes the behavior of the optimal value of the sampled problem when the number of constraints $n$ tends to infinity.

# CHAPTER 4

## THE NYSTRÖM APPROXIMATION

In this chapter, we analyze an alternative algorithm for approximating a symmetric positive definite matrix with a low-rank one. The algorithm is based on the idea of using the subspace spanned by a few columns of the matrix to create the approximation. In Section 4.1 we present the algorithm and show the connection with algorithms of the same type. In Section 4.2 we analyze the behavior of the approximation in a special case in our model. In the last two sections of this chapter we analyze the approximation in our model, using the method that we developed in Chapter 2.

## 4.1 The approximation

The SVD-based algorithm that we analyzed in Chapter 3 is appealing from a theoretical point of view since it returns an optimal low-rank matrix approximation with respect to the Frobenius and the 2-norm. Its main drawback though, is its computational complexity. Computing the SVD of a $d \times d$ matrix requires $O(d^3)$ operations and $O(d^2)$ memory. In applications where $d$ is large, using such an algorithm may be problematical.

Recent research, [18, 36, 22, 19, 20, 2] has focused on algorithms that are sub-optimal, in the sense that they return a low-rank approximation that does not attain the smallest possible error. From a theoretical point of view, research has focused on providing upper bounds on the approximation error. The main advantage of such algorithms is that they are able to produce an approximation

using at most $O(d)$ operations and memory. In the context of our model, approximating all the constraints using such an algorithm is practical, since it requires $O(dn)$ operations, which is small compared to the $O\big((n+d)d^2\big)$ operations required for a single step of an interior point method for the original problem.

Typically, such algorithms do not preserve the structure of the original matrix, so we focus on algorithms that can approximate a symmetric positive semidefinite matrix with a low-rank symmetric positive semidefinite matrix. We will describe two related algorithms that have been proposed in the literature and then present the algorithm that we will analyze.

In what follows if $W$ is a $d \times d$ matrix and $I, J$ are two disjoint subsequences of $\{1, 2, ..., d\}$, let $W(I, J)$ denote the submatrix of $W$ that is composed of the intersection of the rows of $W$ in $I$ and the columns in $J$. We denote by $W(:, J)$ the submatrix of $W$ composed of the columns in $J$ and similarly by $W(I, :)$ the submatrix of $W$ composed of the rows in $I$. We denote by $[IJ]$ the concatenation of the two sequences $I, J$. If $x$ is a vector in $\mathbb{R}^d$, $x_I$ denotes the vector that contains only the elements with indices in $I$.

In [36], a randomized method to approximate a symmetric positive semidefinite matrix $W$ was proposed in the context of Support Vector Machines. The algorithm chooses $k$ columns from $W$ uniformly at random and without replacement, in order to construct the approximation. The algorithm works as follows.

**Algorithm 4.1.1**

- **Input:** A $d \times d$ symmetric positive semidefinite matrix $W$, integer $k$.

- **Output:** A $d \times d$ symmetric positive semidefinite matrix $C$ of rank $k$.

- Sample $k$ columns of $W$ uniformly at random and without replacement. Let $I$ be this sequence of indices.

- Set $R = W(:, I)$ and $F = W(I, I)$.

- Return $C = RF^{-1}R^T$.

In [36] there is no theoretical analysis of the behavior of this algorithm and issues such as the existence of the inverse were not addressed. In computational experiments that were performed, the procedure was shown to work well. This method has been referred to as the Nyström method, because it can be interpreted in terms of the Nyström technique for solving linear integral equations [17].

Drineas and Mahoney, [20], analyzed an algorithm similar to Algorithm 4.1.1 but more general. They use sampling of columns with respect to a general probability distribution $\{p_j\}_{j=1}^d$. They provide an upper bound on the approximation error, when columns are sampled using the judiciously chosen probabilities

$$p_j = \frac{W_{jj}^2}{\sum_{t=1}^d W_{tt}^2}, \ j = 1, ..., d.$$

The bound holds in expectation and with high probability. The algorithm works as follows.

**Algorithm 4.1.2**

- **Input:** A $d \times d$ symmetric positive semidefinite matrix $W$, $\{p_j\}_{j=1}^d$, such that $\sum_{j=1}^d p_j = 1$, integers $m \le d$ and $k \le m$.

- **Output:** A $d \times d$ symmetric positive semidefinite matrix $C$.

- Pick $m$ columns of $W$ in i.i.d. trials, with replacement and with respect to the probabilities $\{p_j\}_{j=1}^d$; let $I$ be the set of indices of the sampled columns.

- Scale each sampled column (whose index is $j \in I$) by dividing its elements by $\sqrt{mp_j}$; Let $R$ be the $d \times m$ matrix containing the sampled columns rescaled in this manner. Let $F$ be the $m \times m$ submatrix of W whose entries are $W_{ij}/(m\sqrt{p_i p_j})$, $i, j \in I$.

- Compute $F_k$, the best rank-$k$ approximation to $F$ with respect to the Frobenius norm.

- Return $C = RF_k^+ R^T$, where $F_k^+$ denotes the Moore-Penrose generalized inverse of the matrix $F_k$.

The algorithm that we will analyze in our model shares ideas from Algorithms 4.1.1 and 4.1.2. It can be thought of as a deterministic version of Algorithm 4.1.1. The choice of columns is made using the idea from Algorithm 4.1.2 of picking the columns that contain large diagonal elements. Let $I$ be the set of indices of the columns that are chosen. The algorithm then creates the unique symmetric positive semi-definite matrix such that the columns in $I$ are approximated exactly and the remaining columns are linear combinations of the ones in $I$.

**Algorithm 4.1.3**

- **Input**: A $d \times d$ symmetric positive semi-definite matrix $W$ and a positive integer $k$.

- **Output**: A $d \times d$ symmetric positive semi-definite matrix $C$ of rank at most $k$.

- Pick the $k$ columns of $W$ with the largest diagonal entries $W_{ii}$. Name $I$ the set of those $k$ indices.

- Let $R$ be the $d \times k$ matrix formed by the columns of $W$ in $I$ and $F$ be the submatrix of $W$ created by the intersection of the rows and columns in $I$.

- Return $C = RF^+R^T$.

We assume that in the case that there are two or more equal diagonal elements, we choose from them uniformly at random and without replacement.

In [3], the following result is proven which provides some useful properties of Algorithm 4.1.3.

**Proposition 4.1.1 (Proposition 1 in [3])** *Let $W$ be a $d \times d$ symmetric positive semidefinite matrix. Let $I$ be a subset of $\{1, 2, ..., n\}$ and let $J$ be its ordered complement in $\{1, 2, ..., n\}$. If the columns of $W$ with indices in $I$ are chosen during the application of Algorithm 4.1.3 to $W$, then $C$ is the unique $d \times d$ matrix $C$ such that*

- *C is symmetric*

- *The column space of C is spanned by the columns of W with indices in I.*

- *The columns of C and W with indices in I are equal.*

*The matrices $C$ and $W - C$ are positive semidefinite. Furthermore, the matrix $C$ is such that*

$$C([IJ], [IJ]) = \begin{bmatrix} W(I, I) & W(I, J) \\ W(I, J)^T & W(J, I)W(I, I)^+W(I, J) \end{bmatrix}.$$

Using this result, it is easy to see that the output of Algorithm 4.1.3 can be alternatively seen as the result of an incomplete Cholesky factorization with symmetric permutations that uses only the columns of $W$ in $I$.

We then use Algorithm 4.1.3 in our stylized optimization problem. Our model consists of the original optimization Problem (2.8)

$$\max \ c^T x$$

$$\text{subject to} \ x^T Q_i A^T Q_i x \leq 1, i = 1, ..., n,$$

with optimal value $V_1$, and the approximation to it given by

$$\max \ c^T x \tag{4.1}$$

$$\text{subject to} \ x^T C_i x \leq 1, i = 1, ..., n,$$

where $C_i$ is the output of Algorithm 4.1.3 with inputs $W_i = Q_i^T A Q_i$ and $k$, for $i = 1, ..., n$. We denote by $V_3$ the optimal value of Problem (4.1) and we define the relative error

$$R_3 = \frac{V_3 - V_1}{V_1}. \tag{4.2}$$

We write the approximating matrix as $C(Q)$ whenever we want to stress the dependence on the random orthogonal matrix $Q$. Similarly, when we want to emphasize the dependence on the number of constraints, we denote by $V_3(n)$ the optimal value of Problem (4.1). We have already studied Problem (2.8), so we now focus on Problem (4.1).

## 4.2 The near-spherical case

We begin our analysis with a simple case, $A = I_d$ that highlights some fundamental properties of the optimal value $V_3$ of Problem (4.1). In this case, the

constraints of the original problem are unit balls centered at 0. We then extend the results to the case where $A$ is close to $I_d$, i.e., when the constraints of the original problem are near-spherical.

We first consider the extreme case, $A = I_d$ and $c = e_1$. Then, Problem (2.8) becomes

$$\max\ e_1^T x \tag{4.3}$$

$$\text{subject to}\ \ \|x\|_2 \leq 1$$

and we have $V_1 = 1$. Algorithm 4.1.3 then approximates each $W_i$ with the matrix $C_i$, where $C_i$ is a diagonal matrix with $k$ diagonal elements equal to 1 and $d - k$ elements equal to 0. The set of indices such that the corresponding diagonal elements of $C_i$ are equal to 1 is chosen uniformly at random from the set $\{1, 2, ..., d\}$. The feasible regions of Problems (2.8) and (4.1) in this case can be seen in Figure 4.2.

From the form of the approximating problem, we see that we have $V_3 = 1$ if and only if the first column of at least one of the matrices $W_1, ..., W_n$ is chosen during the approximation and $V_3 = \infty$ otherwise. Let $\mathcal{A}_n$ be this event. We then have that for any $v > 1$,

$$\mathbb{P}(V_3(n) < v) = \mathbb{P}(\mathcal{A}_n).$$

Since columns are chosen uniformly and independently across matrices, it is easy to see that

$$\mathbb{P}(V_3(n) < v) = 1 - \left(1 - \frac{k}{d}\right)^n.$$

Also, in this case we have that the relative difference in the optimal values of the original and the approximate optimization problems is

$$R_3(n) = \frac{V_3(n) - V_1(n)}{V_1(n)} = V_3(n) - 1.$$

Figure 4.1: Feasible regions of Problems (2.8) and (4.1), when $A = I_2$ and $k = 1$.

So we see that

$$\mathbb{P}(R_3(n) > r) = \left(1 - \frac{k}{d}\right)^n,$$

for any $r > 0$.

We conclude that in this case, the approximate problem has optimal value $V_3$ that is equal to the optimal value $V_1 = 1$ of the original problem with high probability. In contrast, for the SVD-based approximation, when $A = I_d$ we have from Proposition 3.2.2 that for any $r > 0$,

$$P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right)^n \leq \mathbb{P}\left(R_2(n) > r\right) \leq \binom{n}{d} P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right)^{n-d},$$

where $\{Y_j\}_{j=1}^d$ are independent standard normal random variables. For $r$ close to 0, we have that

$$\left(1 - \frac{k}{d}\right) < P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right).$$

44

This implies that

$$\lim_{n\to\infty} \frac{\log \mathbb{P}(R_3(n) > r)}{n} < \lim_{n\to\infty} \frac{\log \mathbb{P}(R_2(n) > r)}{n},$$

for $r$ close to 0. We thus have that $\mathbb{P}(R_3(n) > r)$ converges to 0 with respect to $n$ faster than $\mathbb{P}(R_2(n) > r)$, i.e., for large $n$, the relative error of the approximation using Algorithm 4.1.3 is small with higher probability than the error of the SVD-based approximation.

We can derive a similar result in the case $A = I_d$ for arbitrary unit vector $c \in \mathbb{R}^d$ in the objective function. We consider the original problem

$$\max \ c^T x \tag{4.4}$$

$$\text{subject to} \ \|x\|_2 \leq 1.$$

We assume that $c$ has $l \leq d$ non-zero elements and without loss of generality we assume that $c = [c_I^T \ 0]^T$, where $c_I$ is a unit vector in $R^l$. It is easy to see that we have $V_3 = 1$ if and only if there exists an $i$ such that the columns with indices $1, 2, ..., l$ of $W_i$ are picked during the application of Algorithm 4.1.3 to the matrix $W_i$. This implies that $V_3$ can become arbitrarily close to 1 only if the rank $k$ of the approximating matrices is at least as large as the number $l$ of non-zero elements in the objective function vector $c$. We will see in the analysis that follows in this chapter that this is a general property of this approximation.

If $k \geq l$ it is easy to see by extending the argument that we used in the case $c = e_1$, that

$$\mathbb{P}(V_3(n) = 1) = 1 - \left(1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}}\right)^n.$$

We thus get that under the condition $k \geq l$, the relative error $R_3(n)$ is equal to 0 with probability that decays quickly with respect to the number of constraints

*n*. In contrast, for the SVD-based approximation, when $A = I_d$ we have from Proposition 3.2.2 that for any $r > 0$,

$$P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right)^n \leq \mathbb{P}\left(R_2(n) > r\right) \leq \binom{n}{d}P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right)^{n-d},$$

where $\{Y_j\}_{j=1}^d$ are independent standard normal random variables. For *r* close to 0, we have that

$$1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}} < P\left(\frac{\sum_{j=1}^k Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{(1+r)^2}\right).$$

This implies that

$$\lim_{n\to\infty} \frac{\log \mathbb{P}(R_3(n) > r)}{n} < \lim_{n\to\infty} \frac{\log \mathbb{P}(R_2(n) > r)}{n},$$

for *r* close to 0. We thus have that $\mathbb{P}(R_3(n) > r)$ converges to 0 with respect to *n* faster than $\mathbb{P}(R_2(n) > r)$, i.e., for large *n*, the relative error of the approximation using Algorithm 4.1.3 is small with higher probability than the error of the SVD-based approximation.

The cases that we examined so far were somewhat artificial, since $A = I_d$. Our conclusion was that the optimal value of the approximate problem is a good approximation to the optimal value of the original problem with probability that increases to 1 quickly with respect to the number of constraints. If the case $A = I_d$ is representative of the behavior of the algorithm, then we expect similar conclusions to hold when *A* is close to $I_d$ as well. In order to investigate whether this is true, we assume that the matrix *A* depends on a parameter $\theta \in \mathbb{R}$ and we write $A = A(\theta)$. We also assume that as $\theta \uparrow 1$, we have $A(\theta) \to I_d$. The choice of $\theta \uparrow 1$ is arbitrary and results in no loss of generality. We will first present the analysis of the case $c = e_1$ and then give the results for arbitrary $c \in \mathbb{R}^d$.

We write the original optimization Problem (2.8) as

$$\max \quad e_1^T x \tag{4.5}$$

$$\text{subject to} \quad x^T Q_i^T A(\theta) Q_i x \le 1 \quad , i = 1, ..., n$$

and the approximate problem as

$$\max \quad e_1^T x \tag{4.6}$$

$$\text{subject to} \quad x^T C_i(\theta) x \le 1 \quad , i = 1, ..., n,$$

where $C_i(\theta)$ is the rank-$k$ approximation to the matrix $Q_i^T A(\theta) Q_i$ given by Algorithm 4.1.3. Let $V_3(\theta)$ be the optimal value of Problem (4.6).

Recall that $A(\theta)$ is a $d \times d$ diagonal matrix and $\{Q_i\}_{i=1}^n$ are random, independent, $d \times d$ orthogonal matrices following the uniform distribution in $O^d$. We further assume that $A_{11}(\theta) = f_1(\theta) \equiv 1$ and $A_{jj}(\theta) = f_j(\theta), j = 2, ..., d$, where $f_j, j = 2, ..., d$ are functions defined in an open interval that includes 1 and which satisfy the following conditions:

- $f_j(\theta) < 1$ for $\theta < 1$.

- $f_j(1) = 1$.

- $f_j$ are increasing functions with $f_j'(1) > 0, j = 2, ..., d$ and $f_j'(1), j = 2, ..., d$ are all distinct.

In all the proofs below we work with orthogonal matrices $Q_1, ..., Q_n$ such that the matrix inversion during Algorithm 4.1.3 is possible and the diagonal elements of $W_i = Q_i^T A(\theta) Q_i$ are unequal when $\theta \ne 1$ lies in a neighbourhood of 1, so that we do not need to randomly sample columns. This can be done without affecting our results, since the set of such $Q_1, ..., Q_n$ has probability 1.

47

We first prove a lemma that describes the way that columns are picked by Algorithm 4.1.3, when $\theta$ belongs in a suitable neighbourhood of 1.

**Lemma 4.2.1** *Let $W(\theta) = Q_i^T A(\theta) Q_i$ be one of the symmetric positive semidefinite matrices in the constraints of Problem (4.5). Let $I(\theta)$ be the set of $k$ indices of the columns of $W(\theta)$ that are picked by Algorithm 4.1.3 . Also, let $C(\theta)$ be the output of Algorithm 4.1.3 . Then there exists a $\theta_0 < 1$ such that for $\theta \in (\theta_0, 1)$, $I(\theta) = J$, which does not depend on $\theta$. We also have that*

$$\lim_{\theta \uparrow 1} C_{ij}(\theta) = \begin{cases} 1 & \text{if } i = j, \, i, j \in J \\ 0 & \text{otherwise} \end{cases}.$$

**Proof:** Let $C(\theta)$ be the rank-$k$ approximation to $W = Q^T A(\theta) Q$ given by Algorithm 4.1.3 . Then $C(\theta)$ is of the form

$$C(\theta) = W(:, I(\theta)) W^{-1}(I(\theta), I(\theta)) W(I(\theta), :),$$

where $I$ is the set of the indices of the columns of $W$ that contain the $k$-largest diagonal elements.

For integers $1 \leq m, j \leq d$, $m \neq j$ the elements of matrix $W$ as a function of $\theta$ are

$$W_{jj}(\theta) = \sum_{k=1}^{d} f_k(\theta) Q_{kj}^2$$

and

$$W_{mj}(\theta) = \sum_{k=1}^{d} f_k(\theta) Q_{kj} Q_{km}.$$

We also have that

$$\lim_{\theta \uparrow 1} W_{jj}(\theta) = W_{jj}(1) = 1, \tag{4.7}$$

$$\lim_{\theta \uparrow 1} W_{mj}(\theta) = W_{mj}(1) = 0, \tag{4.8}$$

$$W_{jj}'(1) = \sum_{k=1}^{d} f_k'(1) Q_{kj}^2.$$

48

So, there exists a $\theta_0 < 1$ such that for $\theta \in (\theta_0, 1)$ the $k$ columns that will be chosen are the columns with indices $l_1, l_2, ..., l_k$ such that $\{W'_{l_j l_j}(1), j = 1, ..., k\}$ are the $k$-smallest values of the set $\{W'_{ll}(1), l = 1, ..., d\}$. Then, we have that for all $\theta \in (\theta_0, 1)$,

$$I(\theta) = \{l_1, l_2, ..., l_k\} = J.$$

Suppose that the matrix $W(\theta)$ after some symmetric permutations is in the form

$$W(\theta) = \begin{bmatrix} R_1(\theta) & R_2(\theta) \\ R_2^T(\theta) & R_3(\theta) \end{bmatrix}$$

such that $R_1(\theta)$ is $(d-k) \times (d-k)$, $R_3(\theta)$ is $k \times k$ and the last $k$ columns of $W$ are the ones that are picked for $\theta$ arbitrarily close to 1. Then

$$C(\theta) = \begin{bmatrix} R_2(\theta) \\ R_3(\theta) \end{bmatrix} R_3^{-1}(\theta) \begin{bmatrix} R_2^T(\theta) & R_3^T(\theta) \end{bmatrix} = \begin{bmatrix} R_2(\theta)R_3^{-1}(\theta)R_2^T(\theta) & R_2(\theta) \\ R_2^T(\theta) & R_3(\theta) \end{bmatrix}.$$

From (4.7) and (4.8) it is easy to see that

$$\lim_{\theta \uparrow 1} R_3(\theta) = I_k$$

and

$$\lim_{\theta \uparrow 1} R_2(\theta) = 0.$$

This gives us that

$$\lim_{\theta \uparrow 1} C_{ij}(\theta) = \begin{cases} 1 & \text{if } i = j \in J \\ 0 & \text{otherwise} \end{cases}.$$

$\square$

Lemma 4.2.1 shows that under our assumptions the choice of columns by Algorithm 4.1.3 is stable for $A(\theta)$ close to $I_d$ with probability 1.

When $A = I_d$ we saw that a necessary and sufficient condition so that $V_3 = 1$ is that during application of Algorithm 4.1.3 to the matrices $W_1, ...W_n$, the first

column of at least one of these matrices is chosen. We demonstrate in the following simple example that this condition does not suffice for

$$\lim_{\theta \uparrow 1} V_3(\theta) = 1.$$

Consider the following family of optimization problems in 2 dimensions with 2 constraints,

$$\max \ e_1^T x$$

$$\text{subject to} \ x^T Q_i^T A(\theta) Q_i x \leq 1, i = 1, 2,$$

where $A(\theta) = \text{diag}(1, \theta)$ and

$$Q_i = \begin{bmatrix} \cos(t_i) & \sin(t_i) \\ \sin(t_i) & \cos(t_i) \end{bmatrix}, i = 1, 2.$$

We assume $\theta < 1$ and we pick $t_i \in (0, \pi/4), i = 1, 2$. Under this choice of parameters, both constraints are approximated choosing the first column and it is easy to see that the optimal value of the approximate problem is such that

$$\lim_{\theta \uparrow 1} V_3(\theta) = \left| \frac{\sin(2 t_2) + \sin(2 t_1)}{\sin(2 t_2) - \sin(2 t_1)} \right| > 1.$$

We next provide a result that describes the behavior of the optimal value $V_3$ as the diagonal matrix $A$ tends to the identity matrix, $I_d$. As we show in the proof of this result, a sufficient condition for

$$\lim_{\theta \uparrow 1} V_3(\theta) = 1,$$

is that the matrices $Q_1, ..., Q_n$ are such that columns with all indices $1, ..., d$ are chosen during the application of Algorithm 4.1.3 to the matrices $W_1, ..., W_n$.

**Proposition 4.2.1** *Let $V_3(\theta)$ be the optimal value of Problem (4.6). Then*

$$1 - \mathbb{P}\left( \lim_{\theta \uparrow 1} V_3(\theta) = 1 \right) = O\left( \left( 1 - \frac{k}{d} \right)^n \right). \tag{4.9}$$

**Proof:** From Lemma 4.2.1 we have that

$$\lim_{\theta \uparrow 1} C_i(\theta) = D_i,$$

where $D_i$ is a diagonal matrix, with $k$ elements in the diagonal equal to 1, and $d - k$ elements equal to 0. We also have that the choice of columns is stable in a neighbourhood of 1 and each subset of columns of size $k$ has equal probability of being chosen.

Suppose that the orthogonal matrices $\{Q_i\}_{i=1}^n$ are such that columns with all possible indices $1, 2, ..., d$ are chosen for $\theta$ arbitrarily close to 1 during the application of Algorithm 4.1.3 to the matrices $W_1, ..., W_n$. Let $\mathcal{K} \subset \left(O^d\right)^n$ be the set of $\{Q_i\}_{i=1}^n$ that satisfy this property. We will prove that under this condition we have

$$\lim_{\theta \uparrow 1} V_3(\theta) = 1.$$

We first prove that under this condition, the feasible region of Problem (4.6) is bounded uniformly in $\theta$ in a neigbourhood of 1. Let $x \in \mathbb{R}^d$ be an arbitrary feasible point of Problem (4.6). Let $v \in \{1, 2, ..., d\}$ be such that

$$|x_v| = \max\{|x_1|, |x_2|, ..., |x_d|\}.$$

Because of our assumption, there exists a $\kappa \in \{1, 2, ..., n\}$, such that column $v$ is picked during application of Algorithm 4.1.3 to matrix $Q_\kappa^T A Q_\kappa$. Let $B = C_\kappa$. Then constraint $\kappa$ gives us

$$\sum_{l,m=1}^d B_{lm}(\theta) x_l x_m \leq 1$$

$$\Leftrightarrow B_{vv}(\theta) x_v^2 + 2 \sum_{l \neq v} B_{vl}(\theta) x_l x_v + \sum_{m,l \neq v} B_{lm}(\theta) x_l x_m \leq 1. \tag{4.10}$$

Since the matrix $B$ is positive semidefinite, we have that the last term in the left hand side of (4.10) is non-negative. Then, inequality (4.10) implies that

$$B_{vv}(\theta)x_v^2 + 2\sum_{l \neq v} B_{vl}(\theta)x_l x_v \leq 1$$

$$\Leftrightarrow B_{vv}(\theta)x_v^2 \leq 1 - 2\sum_{l \neq v} B_{vl}(\theta)x_l x_v. \tag{4.11}$$

We define $m(\theta)$ to be the maximum in absolute value of all elements of the matrices $C_i(\theta), i = 1, ..., n$ that tend to 0 as $\theta$ increases to 1 and $M(\theta)$ to be the minimum of all elements of $C_i(\theta), i = 1, ..., n$ that tend to 1. Then, we have

$$\lim_{\theta \uparrow 1} m(\theta) = 0$$

and

$$\lim_{\theta \uparrow 1} M(\theta) = 1.$$

For $\theta$ in a neighbourhood of 1, we have that (4.11) implies

$$M(\theta)x_v^2 \leq 1 + 2(d-1)m(\theta)x_v^2$$

$$\Rightarrow x_v^2 \leq \frac{1}{M(\theta) - 2(d-1)m(\theta)}. \tag{4.12}$$

Since

$$|x_v| = \max\{|x_1|, |x_2|, ..., |x_d|\},$$

we see that

$$|x_1| \leq \sqrt{\frac{1}{M(\theta) - 2(d-1)m(\theta)}}. \tag{4.13}$$

From (4.13) we have

$$V_3(\theta) \leq \sqrt{\frac{1}{M(\theta) - 2(d-1)m(\theta)}}.$$

Since $\lim_{\theta \uparrow 1} m(\theta) = 0$ and $\lim_{\theta \uparrow 1} M(\theta) = 1$, we get that

$$\lim_{\theta \uparrow 1} V_3(\theta) = 1.$$

Thus, for any $(Q_1, ..., Q_n) \in \mathcal{K}$ we have that

$$\lim_{\theta \uparrow 1} V_3(\theta) = 1.$$

From that, we easily get that

$$\mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right) \geq \mathbb{P}(\mathcal{K}).$$

Each subset of indices of size $k$ from $\{1, 2, ..., d\}$ has the same probability of being picked, and we require each index from $\{1, 2, ..., d\}$ to be picked at least once, in any of the $n$ constraints. This problem is equivalent to a version of the now classic coupon collector's problem. In this problem a population $S$ of $s$ distinct elements is sampled with replacement in groups of $k$ elements at each time. The quantity of interest is the sample size necessary for the acquisition of the set $S$. De Moivre in [15] and [16] first derived the probability that all elements of $S$ will be obtained after $n$ samples for the case $k = 1$. Laplace, in a memoir , [27], and then in [28], generalized De Moivre's result to the case $k \geq 1$. So, we have that

$$P(\mathcal{K}) = \sum_{j=k}^{d} (-1)^{d-j} \binom{d-k}{j-k} \left[\frac{\binom{j}{k}}{\binom{d}{k}}\right]^{n-1},$$

from which the result follows directly.

$\square$

Proposition 4.2.1 describes the behavior of the optimal $V_3$ when the diagonal matrix $A$ is close to the identity. It implies that when $A$ approaches $I_d$, the probability that $V_3$ does not approach 1 decays at least as fast as $(1 - k/d)^n$, as the number of constraints tends to infinity. Furthermore, comparing it with the limiting case $A = I_d$, we have that as $n$ increases to infinity, it implies that $\mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right)$

and $\mathbb{P}(V_3(1) = 1)$ tend to 1 at the same rate. Also, for $v > 1$, we get that

$$\lim_{\theta \uparrow 1} \mathbb{P}(V_3(\theta) > v) \leq \mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right) = O\left(\left(1 - \frac{k}{d}\right)^n\right).$$

A similar analysis can be performed under the same assumptions for the problem

$$\max \quad c^T x \tag{4.14}$$

$$\text{subject to} \quad x^T C_{(\theta)} x \leq 1 \quad , i = 1, ...n,$$

where $c \in \mathbb{R}^d$ is an arbitrary unit vector and for each $i$, $C_i$ is the output of Algorithm 4.1.3 with inputs $Q_i^T A Q_i$ and $k$. If $l$ is the number of non-zero elements of the vector $c$, and the rank of the approximating matrices is $k \geq l$, then the optimal value $V_3(n)$ of the sampled problem can become arbitrarily close to 1. Furthermore, we can prove the following result.

**Proposition 4.2.2** *Let $V_3(\theta)$ be the optimal value of Problem (4.14). Then*

$$1 - \mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right) = O\left(\left(1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}}\right)^n\right). \tag{4.15}$$

The proof of Proposition 4.2.2 is presented in the appendix. It is based on similar arguments as in the proof of the case $c = e_1$. As in the case $c = e_1$, it implies that $\mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right)$ and $\mathbb{P}(V_3(1) = 1)$ tend to 1 at the same rate. Also, for $v > 1$, it implies that

$$\lim_{\theta \uparrow 1} \mathbb{P}(V_3(\theta) > v) \leq \mathbb{P}\left(\lim_{\theta \uparrow 1} V_3(\theta) = 1\right) = O\left(\left(1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}}\right)^n\right).$$

## 4.3 Simple objective function

The analysis of the near-spherical case described the behavior of the approxima-
tion when the diagonal matrix $A$ is close to the identity matrix. In this section
we apply the method that we developed in Chapter 2, in order to derive results
about the optimal value $V_3$ for any diagonal matrix $A$. We first analyze the case
$c = e_1$.

Recall that Problem (4.1) is given under this assumption by

$$\max \quad e_1^T x$$

$$\text{subject to} \quad x^T C(Q_i)x \leq 1 \quad, i = 1, ...n,$$

where $C(Q_i)$ is the output of Algorithm 4.1.3 with inputs $W_i = Q_i^T A Q_i$ and $k$. The
matrices $Q_i, i = 1, .., n$ are independent random uniformly distributed orthogo-
nal matrices and $A$ is a diagonal matrix with

$$1 = A_{11} \geq A_{22} \geq \cdots \geq A_{dd} \geq 0.$$

We denote by $V_3$ the optimal value of Problem (4.1).

We first consider the robust optimization problem corresponding to the ap-
proximating Problem (4.1). This is

$$\max \quad e_1^T x \tag{4.16}$$

$$\text{subject to} \quad x^T C(Q)x \leq 1, \quad \forall Q \in O^d.$$

Since the robust problem can be thought of as a limiting case of the $\epsilon$-chance-
constrained problem as $\epsilon \downarrow 0$, we denote by $X_3(0)$ the feasible region of this
problem and by $G_3(0)$ its optimal value. For this problem, we assume that if
during application of Algorithm 4.1.3 to the matrix $Q^T A Q$, there are diagonal

elements that are equal, then all possible subsets of columns of size $k$ are chosen, i.e., instead of choosing uniformly at random which columns to pick, all possible choices are created, leading to multiple low-rank approximations and multiple constraints. This is done in order to avoid introducing any randomness to this problem. Figure 4.3 shows the feasible region of the robust problem, when $d = 2$, $k = 1$ and $A_{22} = 0.25$.



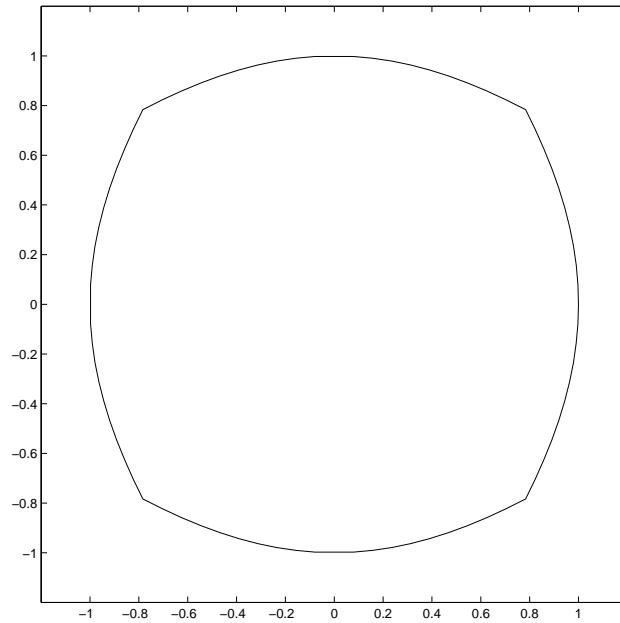Figure 4.2: Feasible region of the robust problem when $d = 2$, $k = 1$ and $A_{22} = 0.25$.

The following lemma gives us the solution and the optimal value of the robust problem.

**Lemma 4.3.1** *For the robust optimization Problem (4.16) we have* $G_3(0) = 1$ *and the optimal solution is* $\hat{x}_3(0) = e_1$.

**Proof:** We first check that $x = e_1$ is feasible. This is easy to see since for every

56

$Q \in \mathcal{O}^d$ we have

$$e_1^T C(Q) e_1 \leq 1.$$

Hence, we have $G_3(0) \geq 1$.

We consider a feasible point of the form $x = \lambda u$, where $u \neq e_1$ is a unit vector in $\mathbb{R}^d$ and $\lambda > 0$. Without loss of generality we assume that $u_1 > 0$ and $u_2 \neq 0$. We have that

$$\lambda u \in \mathcal{X}_0$$

$$\Leftrightarrow \lambda \sqrt{u^T C(Q) u} \leq 1, \forall Q \in \mathcal{O}^d. \tag{4.17}$$

Taking $Q = I_d$ in (4.17), we get that

$$\lambda \leq \frac{1}{\sqrt{u^T C(I_d) u}} \leq \frac{1}{\sqrt{\sum_{j=1}^{k} A_{jj} u_j^2}}.$$

This implies that

$$e_1^T x = \lambda e_1^T u \leq \frac{u_1}{\sqrt{\sum_{j=1}^{k} A_{jj} u_j^2}} < 1.$$

So we get that $G_3(0) = 1$ and $\hat{x}_3(0) = e_1$ is the unique optimal solution.

$\square$

Since the optimal value of the robust problem is 1, we see that as $n$ increases, the optimal value $V_3$ of the sampled problem can become arbitrarily close to 1.

We then consider the chance-constrained problem corresponding to (4.1). For $\epsilon \in (0, 1)$ it is given by

$$\max \quad e_1^T x \tag{4.18}$$

$$\text{subject to} \quad x \in \mathcal{X}_3(\epsilon).$$

We have $x \in \mathcal{X}_3(\epsilon)$ if and only if $\mathbb{P}\left(x^T C(Q) x > 1\right) \leq \epsilon$, where $Q$ is uniformly distributed in $O^d$ under $\mathbb{P}$. Figure 4.3 shows the feasible region of the chance-constrained problem for $\epsilon = 0.05$.



Figure 4.3: Feasible region of the chance-constrained problem when $d = 2$, $A_{22} = 0.25$ and $\epsilon = 0.05$.

The feasible region is not convex in general, as can be seen from Figure 4.4.

We can provide a useful characterization of the feasible region $\mathcal{X}_3(\epsilon)$ as follows. If $x = \lambda u$, where $u$ is a unit vector in $\mathbb{R}^d$, we have

$$x = \lambda u \in \mathcal{X}_3(\epsilon)$$

if and only if

$$\mathbb{P}\left(u^T C(Q) u > \frac{1}{\lambda^2}\right) \leq \epsilon.$$

Let $F(\cdot; u)$ be the distribution function of the random variable $u^T C(Q) u$. We then
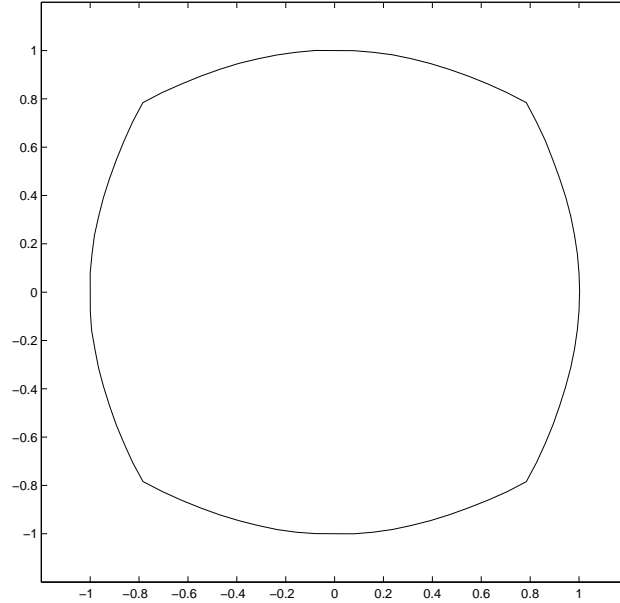
Figure 4.4: Feasible region of the chance-constrained problem when $d = 2$, $A_{22} = 0.25$ and $\epsilon = 0.85$.

have that

$$1 - F\left(\frac{1}{\lambda^2}; u\right) \leq \epsilon$$

$$\Leftrightarrow |\lambda| \leq \frac{1}{\sqrt{F^{-1}(1 - \epsilon; u)}}.$$

We first present our main result about the optimal value of the sampled problem. The result is an application of Theorem 2.4.2 to Problem (4.1).

**Proposition 4.3.1** *Let $V_3$ be the optimal value of the sampled Problem (4.1) with $c = e_1$ and $G_3(\epsilon)$ denote the optimal value of the corresponding chance-constrained problem. Then, for any $v > 1$, there exists a $G_3^{-1}(v) \in (0, 1]$, where $G_3^{-1}(\cdot)$ is a left inverse function of $G_3(\cdot)$, such that*

$$(1 - G_3^{-1}(v))^n \leq \mathbb{P}(V_3 > v) \leq \binom{n}{d}(1 - G_3^{-1}(v))^{n-d}. \tag{4.19}$$

**Proof:** For the sampled problem we have for $v > 1$ that

$$\mathbb{P}(V_3 = v) = 0.$$

59

This is a consequence of the fact that for $A \neq I_d$, no specific constraint can be sampled with positive probability. Since $\mathbb{P}(V_3 < v) > 0$ for all $v > 1$, we get that

$$\lim_{\epsilon \downarrow 0} G_3(\epsilon) = 1.$$

We thus get from Theorem 2.4.2 that for any $v > 1$, we have

$$(1 - G_3^{-1}(v))^n \leq \mathbb{P}(V_3 > v) \leq \binom{n}{d}(1 - G_3^{-1}(v))^{n-d}.$$

$\square$

Proposition 4.3.1 gives us an important theoretical result about the optimal value of the sampled problem (4.1) but it does not provide a useful representation for $G_3^{-1}(v)$. This would require knowledge of the optimal value of the chance-constrained problem. Since solving the chance-constrained problem is hard, we will instead provide a bound for $G_3^{-1}(v)$. We start with some preliminary results about the feasible region of the chance-constrained problem. For the remainder of this section, we focus on rank-1 approximations. From the analysis that follows it is straightforward to see that the same bound holds for higher rank approximations.

The following lemma shows that the chance-constrained problem has a feasible region that is symmetric with respect to the axes.

**Lemma 4.3.2** *For any $\epsilon > 0$, the feasible region $X_3(\epsilon)$ of Problem (4.18) with $k = 1$ is symmetric with respect to the axes.*

**Proof:** Without loss of generality we will prove symmetry with respect to the $x_1$ axis. Let $x \in X_3(\epsilon)$. We define the set of constraints that are violated by $x$,

$$C_x = \left\{ Q \in O^d \ \Big| \ \frac{|W(p, :)x|}{\sqrt{W_{pp}}} > 1 \right\},$$

60

where $p$ depends on $Q$ and is the index of the column of the matrix $Q^T A Q$ that was picked during the approximation. Let

$$U = \text{diag}(1, -1, ..., -1) \in \mathbb{R}^{d \times d}.$$

Then we have that $y = Ux$ is symmetric to $x$ with respect to the $x_1$ axis. The $ij$ element of $W = Q^T A Q$ is

$$W_{ij} = \sum_{k=1}^{d} Q_{ki} Q_{kj} A_{kk}.$$

The $ij$ element of the matrix $U^T W U = (QU)^T A (QU)$ is

$$(U^T W U)_{ij} = \sum_{k=1}^{d} Q_{ki} Q_{kj} A_{kk} U_{ii} U_{jj}.$$

We thus get that the diagonal elements of $U^T W U$ and $W$ are the same. The set of constraints violated by $y$ is

$$
\begin{aligned}
C_y &= \left\{ Q \in O^d \mid \frac{|W(p, :)y|}{\sqrt{W_{pp}}} > 1 \right\} \\
&= \left\{ Q \in O^d \mid \frac{|W(p, :)Ux|}{\sqrt{W_{pp}}} > 1 \right\} \\
&= \left\{ Q \in O^d \mid \frac{|e_p^T U U^T Q^T A Q U x|}{\sqrt{e_p^T Q^T A Q e_p}} > 1 \right\} \\
&= \left\{ Q \in O^d \mid \frac{|e_p^T U (QU)^T A (QU) x|}{\sqrt{e_p^T (QU)^T A Q U e_p}} > 1 \right\} \\
&= \left\{ Q \in O^d \mid \frac{|e_p^T (QU)^T A (QU) x|}{\sqrt{e_p^T (QU)^T A Q U e_p}} > 1 \right\} \\
&= C_x U^T.
\end{aligned}
$$

But then, since $Q$ is uniformly distributed in $O^d$, we get that

$$\mathbb{P}(C_y) = \mathbb{P}(C_x U^T) = \mathbb{P}(C_x) \leq \epsilon,$$

and therefore, $y \in \mathcal{X}_3(\epsilon)$.

61

$\square$

In the following lemma, we prove that the midpoint of two feasible points for the $\epsilon$-chance-constrained problem is feasible for the $2\epsilon$-chance-constrained problem.

**Lemma 4.3.3** *Let* $x, y \in \mathcal{X}_3(\epsilon)$, *where* $\epsilon < 1/2$. *Then*

$$z = \frac{1}{2}x + \frac{1}{2}y \in \mathcal{X}_3(2\epsilon).$$

**Proof:** We define the sets

$$C_x^c = \left\{ Q \in \mathcal{O}^d \mid \frac{|W(p, :)x|}{\sqrt{W_{pp}}} \leq 1 \right\}$$

and

$$C_y^c = \left\{ Q \in \mathcal{O}^d \mid \frac{|W(p, :)y|}{\sqrt{W_{pp}}} \leq 1 \right\}.$$

Since $x, y \in \mathcal{X}_3(\epsilon)$, we have

$$\mathbb{P}\left(C_x^c\right), \mathbb{P}\left(C_y^c\right) > 1 - \epsilon.$$

But for any $Q \in C_x^c \cap C_y^c$ we have that

$$\frac{|W(p, :)z|}{\sqrt{W_{pp}}} = \frac{|W(p, :)(x + y)|}{2\sqrt{W_{pp}}} \leq \frac{1}{2}\frac{|W(p, :)x|}{\sqrt{W_{pp}}} + \frac{1}{2}\frac{|W(p, :)y|}{\sqrt{W_{pp}}} \leq 1.$$

So $C_z^c \supseteq C_x^c \cap C_y^c$, whence $\Rightarrow \mathbb{P}(C_z^c) > 1 - 2\epsilon$ and so $z \in \mathcal{X}_3(2\epsilon)$.

$\square$

In the following lemma, we give an upper bound on the norm of feasible points of the chance-constrained problem of the form $z = \lambda e_1$.

**Lemma 4.3.4** *Let $z = \lambda e_1$ for some $1 < \lambda < \sqrt{2}$. Then we have that*

$$z \in \mathcal{X}_3(\epsilon) \Leftrightarrow P\left(\frac{\sum_{j=1}^{d} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{\lambda^2}\right) \leq \epsilon,$$

*where $\{Y_j\}_{j=1}^{d}$ are independent standard normal random variables.*

**Proof:** Let $p$ be the index of the column of $W = Q^T A Q$ that contains the largest diagonal element. We will prove that for $p \neq 1$, $\frac{|W_{1p}|}{\sqrt{W_{pp}}} \leq \frac{\sqrt{2}}{2}$. Without loss of generality we assume that $W_{1p} > 0$. Assume that

$$\frac{W_{1p}}{\sqrt{W_{pp}}} > \frac{\sqrt{2}}{2}.$$

Then, since we have $W_{1p} \leq \sqrt{W_{11} W_{pp}}$, we get that

$$W_{pp} \geq W_{11} \geq W_{1p} > \frac{1}{2}.$$

Let $g = \frac{\sqrt{2}}{2}(e_1 + e_p)$. We have that

$$g^T W g = \frac{1}{2} W_{11} + \frac{1}{2} W_{pp} + W_{1p} > 1.$$

But this is a contradiction since we have $x^T W x \leq 1$ for any $x \in \mathbb{R}^d$ with $\|x\|_2 = 1$.

We have that $z \in \mathcal{X}_3(\epsilon)$ if and only if

$$\Leftrightarrow \mathbb{P}\left(\frac{|W_{1p}|}{\sqrt{W_{pp}}} > \frac{1}{\lambda}\right) \leq \epsilon. \tag{4.20}$$

For $\lambda < \sqrt{2}$, (4.20) becomes

$$\mathbb{P}\left(W_{11} > \frac{1}{\lambda^2}\right) \leq \epsilon$$

$$\Leftrightarrow P\left(\frac{\sum_{j=1}^{d} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{\lambda^2}\right) \leq \epsilon.$$

$\square$

63

We can now combine these results in order to prove our main result.

**Proposition 4.3.2** *Let $G_3(\epsilon)$ be the optimal value of the $\epsilon$-chance-constrained Problem (4.18). Then for $v > 1$ that is sufficiently close to 1, we have that*

$$G_3^{-1}(v) \geq \frac{G_1^{-1}(v)}{2},$$

*where*

$$G_1^{-1}(v) = P\left(\frac{\sum_{j=1}^{d} A_{jj}Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{v^2}\right)$$

*and $Y_1, ..., Y_d$ are independent standard normal random variables.*

**Proof:** Assume that there exists an optimal solution $x = \hat{x}_3(\epsilon)$ for the $\epsilon$-chance constrained problem. Then we have that $e_1^T x = G_3(\epsilon)$. Also, let $y$ be the point that is symmetric to $x$ with respect to the $x_1$ axis. Then we have from Lemma 4.3.2 that $y \in \mathcal{X}_3(\epsilon)$ and $e_1^T y = G_3(\epsilon)$ as well. Furthermore, from Lemma 4.3.3 we have that

$$z = \frac{1}{2}x + \frac{1}{2}y = G_3(\epsilon)e_1 \in \mathcal{X}_3(2\epsilon).$$

This means that for $G_3(\epsilon) < \sqrt{2}$ we have

$$P\left(\frac{\sum_{j=1}^{d} A_{jj}Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{G_3(\epsilon)^2}\right) \leq 2\epsilon.$$

Consider $G_1(2\epsilon)$, the optimal value of the $2\epsilon$-chance-constrained problem corresponding to the original sampled Problem (2.8), that satisfies the equation

$$P\left(\frac{\sum_{j=1}^{d} A_{jj}Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{G_1^2(2\epsilon)}\right) = 2\epsilon.$$

For $\epsilon$ such that $G_3(\epsilon) < \sqrt{2}$ we then have

$$G_3(\epsilon) \leq G_1(2\epsilon).$$

If an optimal solution does not exist, there exists a sequence $\{x_\nu\}_{\nu=1}^\infty \subset \mathcal{X}_3(\epsilon)$ such that

$$\lim_{\nu \to \infty} e_1^T x_\nu = G_3(\epsilon).$$

By using the sequence $\{y_\nu\}_{\nu=1}^\infty \subset \mathcal{X}_3(\epsilon)$ of points symmetric to $\{x_\nu\}_{\nu=1}^\infty$ with respect to the $x_1$ axis, we can get following similar arguments that

$$G_3(\epsilon) \leq G_1(2\epsilon).$$

$\square$

Proposition 4.3.2 gives a bound on the quantity $G_3^{-1}(v)$ for $v$ close to 1 with respect to $G_1^{-1}(v)$. Since we have

$$\frac{G_1^{-1}(v)}{2} \leq G_3^{-1}(v) \leq G_1^{-1}(v),$$

we get the inequality

$$(1 - G_1^{-1}(v))^n \leq \mathbb{P}(V_3 > v) \leq \binom{n}{d}\left(1 - \frac{G_1^{-1}(v)}{2}\right)^{n-d},$$

that relates the optimal value of the approximate sampled Problem (4.1) with the optimal value of the chance-constrained Problem (3.6). We will discuss the implications that this result has on the asymptotic behavior of $V_3(n)$ as $n$ tends to infinity in Chapter 5.

In order to connect the results of this section with the previous results for the near-spherical case, we consider the setting that we introduced in Section 4.2. For $v > 1$, let

$$G_1^{-1}(v; \theta) = P\left(\frac{\sum_{j=1}^d A_{jj}(\theta)Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right).$$

We have that

$$\lim_{\theta \uparrow 1} \mathbb{P}(V_3(\theta) > v) \leq \lim_{\theta \uparrow 1} \binom{n}{d}\left(1 - G_1^{-1}(v; \theta)/2\right)^{n-d} = \binom{n}{d}(1 - 1/2)^{n-d}.$$

We thus get that

$$\lim_{n \to \infty} \frac{\log \left( \lim_{\theta \uparrow 1} \mathbb{P}(V_3(\theta) > v) \right)}{n} \leq \log \left( 1 - \frac{1}{2} \right).$$

From the analysis in Section 4.2 we have that

$$\lim_{n \to \infty} \frac{\log \left( \lim_{\theta \uparrow 1} \mathbb{P}(V_3(\theta) > v) \right)}{n} \leq \log \left( 1 - \frac{k}{d} \right),$$

a result that is weaker for small $k$, but stronger for large $k$. This suggests that our bound

$$G_3^{-1}(v) \leq \frac{G_1^{-1}(v)}{2}$$

is not tight and that the results of Section 4.2 are complementary to the results given in this section.

## 4.4   General case

In this section, we analyze Problem (4.1) without the assumption $c = e_1$. We denote by $l$ the number of non-zero elements of $c$ and by $k$ the rank of the approximating matrices. For simplicity and without loss of generality we assume that the non-zero elements of $c$ are the first $l$ ones, so that $c = [c_I \; 0]$, where $I = \{1, 2, ..., l\}$. Furthermore, we assume throughout this section that the diagonal matrix $A$ has positive diagonal entries, i.e.,

$$1 = A_{11} \geq A_{22} \geq \cdots \geq A_{dd} > 0.$$

The approximate problem is given by

$$\max \; c^T x \tag{4.21}$$

$$\text{subject to} \; x^T C(Q_i) x \leq 1 \; , i = 1, ...n,$$

where $C(Q_i)$ is the output of the approximation Algorithm 4.1.3 with inputs $W_i = Q_i^T A Q_i$ and $k < d$. We denote by $V_1$ and $V_3$ the optimal values of Problems (3.2) and (4.21) respectively.

We first present the robust optimization problem corresponding to (4.21). This is

$$\max \ c^T x \tag{4.22}$$

$$\text{subject to } \ x^T C(Q) x \leq 1 \quad \text{for every } Q \in O^d.$$

We let $X_3(0)$ be the feasible region of this problem and $G_3(0)$ its optimal value. As in the previous section, we assume that when during application of Algorithm 4.1.3 there are diagonal elements of the matrix $W_i$ that are equal, the approximation returns multiple matrices created with all possible choices of columns.

We start with the following lemma about Problem (4.22).

**Lemma 4.4.1** *Consider an orthogonal matrix $Q \in O^d$ such that the index set of the columns that are picked by Algorithm 4.1.3 is $J$, where $J$ does not contain $I = \{1, 2, ..., l\}$. Then we have that there exists a constant $h < 1$ that does not depend on $Q$, such that*

$$c^T C(Q) c \leq h.$$

**Proof:** Let $H = I \cap J^c$ be the non-empty set of the indices that are in $I$ but not in $J$ and $K = I \cap J$. For notational simplicity we write $C = C(Q)$. We then have that

$$c^T C c = c_K^T C(K, K) c_K + 2 c_K^T C(K, H) c_H + c_H^T C(H, H) c_H.$$

But we have that the rows and columns of $C$ that belong in $K$ are equal to the

respective rows and columns of $W = Q^T A Q$. So, we get that

$$c^T C c = c_K^T W(K, K) c_K + 2 c_K^T W(K, H) c_H + c_H^T C(H, H) c_H.$$

By adding and subtracting the quantity $c_H^T W(H, H) c_H$, we get that

$$c^T c = c^T W c + c_H^T (C(H, H) - W(H, H)) c_H.$$

But since we have that $W = Q^T A Q$ is of full rank, the Schur complement

$$W(H, H) - C(H, H)$$

is a positive definite matrix. If

$$\sup \left\{ c_H^T (C(H, H) - W(H, H)) c_H \mid Q \in O^d \right\} = 0,$$

there would exist an orthogonal matrix $Q$ such that $W(H, H) - C(H, H)$ is positive semidefinite. This implies that

$$\sup \left\{ c_H^T (C(H, H) - W(H, H)) c_H \mid Q \in O^d \right\} < 0.$$

Since we have $c^T W c \leq 1$, there exists an $h < 1$ that depends only on $A$ and $c$ such that

$$c^T C c \leq h.$$

$\square$

Lemma 4.4.1 states that if at least one of the columns $1, 2, ..., l$ is not picked during the application of Algorithm 4.1.3 to $W = Q^T A Q$, then $c^T C(Q) c$ is bounded away from 1. It is clear from the proof of the lemma, that the constant $h$ depends on $c$ and $A$. For a fixed $A$, $h$ approaches 1 only if the smallest in absolute value non-zero element of $c$ tends to 0, otherwise it stays bounded away from 1.

Our next result, states that if $k < l$, then the optimal value of the robust Problem (4.22) is greater than 1.

**Proposition 4.4.1** *For the robust optimization Problem (4.16), if we have $k < l$, then $G_3(0) > 1$.*

**Proof:** Since we have $k < l$, there is no choice of columns that contains all the indices of the non-zero elements of $c$. From Lemma 4.4.1, there exists a $\lambda_0 > 1$ such that $\lambda_0^2 c^T C(Q)c \leq 1$ for all $Q \in O^d$. Then $x = \lambda_0 c$ is feasible, and we have that

$$G_3(0) \geq \lambda_0 c^T c > 1.$$

$\square$

This result implies that if $c$ has a large number of non-zero elements, Algorithm 4.1.3 results in an approximation with error that does not become arbitrarily small, even if the number of constraints tends to infinity.

We next prove that if the number of columns $k$ used in the approximation is greater or equal to $l$, then the optimal value of the robust Problem (4.22) is equal to 1.

**Proposition 4.4.2** *For the robust optimization Problem (4.16), if we have $k \geq l$, then $G_3(0) = 1$ and an optimal solution is $\hat{x}_3(0) = c$.*

**Proof:** Since the feasible region of Problem (4.22) contains the unit ball in $\mathbb{R}^d$, we have $G_3(0) \geq 1$. Let $x$ be a feasible vector for Problem (4.22). Then we have that for every $Q \in O^d$,

$$x^T C(Q)x \leq 1.$$

Let $I = \{1, 2, ..., k\}$ and $J = \{k + 1, ..., d\}$. For arbitrary orthogonal matrices $Q_1 \in O^k$

and $Q_2 \in O^{d-k}$ we consider the $d \times d$ orthogonal matrix

$$Q = \begin{bmatrix} Q_1 & 0 \\ 0 & Q_2 \end{bmatrix}.$$

We have that the first $l$ columns of $W = Q^T A Q$ are picked and also

$$C(I, I) = Q_1^T A(I, I) Q_1$$

and

$$C(I, J) = 0, C(J, J) = 0.$$

So we get that

$$x^T C(Q) x \leq 1$$

$$\Leftrightarrow x_I^T W(I, I) x_I \leq 1.$$

Since this holds for any $Q_1 \in O^d$, we get that

$$\|x_I\|_2 \leq 1.$$

This implies that

$$c^T x = c_I^T x_I \leq \|c_I\|_2 \|x_I\|_2 \leq 1. \tag{4.23}$$

We thus get $G_3(0) = 1$ and $x = c$ is an optimal solution.

$\square$

We next present the chance-constrained analog of Problem (4.22). We assume in what follows that the number of columns $k$ chosen during the approximation is such that $k \geq l$. This is necessary in order to guarantee that the approximate sampled problem will be a good approximation to the original problem for a

large number of constraints. For $\epsilon \in (0, 1)$ the chance-constrained problem is given by

$$\max \ c^T x \tag{4.24}$$

$$\text{subject to} \ x \in \mathcal{X}_3(\epsilon).$$

We have $x \in \mathcal{X}_3(\epsilon)$ if and only if $\mathbb{P}\left(x^T C(Q)x > 1\right) \le \epsilon$.

Applying Theorem 2.4.2 to Problem (4.21), we get a result that connects the optimal value $V_3$ of the sampled problem with the optimal value $G_3(\epsilon)$ of the chance-constrained problem.

**Proposition 4.4.3** *Let $V_3$ be the optimal value of the sampled Problem (4.21) with $k \ge l$ and $G_3(\epsilon)$ denote the optimal value of the corresponding chance-constrained problem. Then we have that for any $v > 1$ there exists a $G_3^{-1}(v) \in (0, 1]$, where $G_3^{-1}(\cdot)$ is a left inverse of $G_3(\cdot)$, such that*

$$(1 - G_3^{-1}(v))^n \le \mathbb{P}(V_3 > v) \le \binom{n}{d}(1 - G_3^{-1}(v))^{n-d}. \tag{4.25}$$

**Proof:**

Inequality (4.25) holds for $v$ in the image of $G(\cdot)$ because of Theorem 2.4.2. Since $\mathbb{P}(V_3 < v) > 0$ for all $v > 1$, we get that

$$\lim_{\epsilon \downarrow 0} G_3(\epsilon) = 1.$$

Using Theorem 2.4.2 we get that for any $v > 1$, there exists a $G_3^{-1}(v)$ and we have

$$(1 - G_3^{-1}(v))^n \le \mathbb{P}(V_3 > v) \le \binom{n}{d}(1 - G_3^{-1}(v))^{n-d}.$$

$\square$

Proposition 4.3.1 gives us an important theoretical result about the optimal value of the sampled problem (4.21) but it does not provide a useful representation for $G_3^{-1}(v)$. This would require solving the chance-constrained problem, or providing a bound on its optimal value.

In the following chapter, we capitalize on the results that we have for the optimal values of the problems that we have studied so far in order to derive asymptotic results. More specifically, we capitalize on Inequality (2.15), in order to describe the asymptotic behavior of the optimal values $V_i(n), i = 1, 2, 3$, as $n$ tends to infinity. This leads to asymptotic results for the relative errors $R_2(n)$ and $R_3(n)$.

# CHAPTER 5

## ASYMPTOTICS AND DISCUSSION

## 5.1 Preliminaries

In this section we focus on the behavior of the optimal values of the optimization problems in our model when the number $n$ of constraints tends to infinity. This is important, since we are interested in problems with a large number of constraints. We will derive asymptotic results for the optimal values and the relative errors. The basic tool that we will use for the analysis that follows is the fundamental inequality (2.15). We present the analysis through the original sampled Problem (3.2), but the analysis holds for the approximating Problems (3.3) and (4.1) as well.

Consider the original optimization Problem (3.2). We denote its optimal value by $V_1(n)$ in order to stress the dependence on the number of constraints $n$. The main result that we have for $V_1(n)$ states that for any $v > 1$, we have

$$(1 - G_1^{-1}(v))^n \le \mathbb{P}(V_1(n) > v) \le \binom{n}{d}(1 - G_1^{-1}(v))^{n-d},$$

where

$$G_1^{-1}(v) = P\left(\frac{\sum_{j=1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} > \frac{1}{v^2}\right)$$

and $\{Y_j\}_{j=1}^d$ are independent standard normal random variables. The idea is to find decreasing sequences of real numbers $\{Z_1(n)\}_{n=1}^\infty$ and $\{z_1(n)\}_{n=1}^\infty$ such that

$$\lim_{n\to\infty} \mathbb{P}(V_1(n) < Z_1(n)) = 1 \tag{5.1}$$

and

$$\lim_{n\to\infty} \mathbb{P}(V_1(n) > z_1(n)) = 1. \tag{5.2}$$

This implies that

$$\lim_{n\to\infty} \mathbb{P}(z_1(n) < V_1(n) < Z_1(n)) = 1, \tag{5.3}$$

giving us a description of how $V_1(n)$ behaves for large $n$.

Since for any $v > 1$,

$$\lim_{n\to\infty} \mathbb{P}\left(V_1(n) > v\right) = 0,$$

we get that as $n \to \infty$, $V_1(n)$ converges to 1 in $\mathbb{P}$-probability. Our plan is to use the results in Theorem 2.4.2 and Proposition 3.2.3 in order to get sufficient conditions on $Z_1(n)$ and $z_1(n)$ so that (5.3) holds. This will allow us to describe, in some sense that will be made precise later in this section, how fast the optimal value $V_1(n)$ converges to 1 in $\mathbb{P}$-probability. Applying (5.3) to the sampled Problems (3.3) and (4.1) in our model, leads to similar results for their respective optimal values $V_2(n)$ and $V_3(n)$. Furthermore, we describe in the same sense, how fast the relative errors

$$R_2(n) = \frac{V_2(n) - V_1(n)}{V_1(n)}$$

and

$$R_3(n) = \frac{V_3(n) - V_1(n)}{V_1(n)}$$

converge to 0.

The following proposition gives sufficient conditions for (5.1) and (5.2) to hold.

**Proposition 5.1.1** *Let $V_1(n)$ be the optimal solution to the sampled Problem (3.2) and let $G_1(\epsilon)$ denote the optimal value of the corresponding chance-constrained problem. Then for sequences $\{Z_1(n), z_1(n)\}_{n=1}^{\infty}$, we have that as $n \to \infty$,*

$$\lim_{n\to\infty} \frac{n}{\log n} G_1^{-1}(Z_1(n)) = \infty \implies \lim_{n\to\infty} \mathbb{P}(V_1(n) < Z_1(n)) = 1,$$

$$\lim_{n\to\infty} nG_1^{-1}(z_1(n)) = 0 \Rightarrow \lim_{n\to\infty} \mathbb{P}(V_1(n) > z_1(n)) = 1.$$

**Proof:** From (3.10), we have that a sufficient condition for

$$\lim_{n\to\infty} \mathbb{P}(V_1(n) > z_1(n)) = 1$$

is $\lim_{n\to\infty} \left(1 - G^{-1}(z_1(n))\right)^n = 1$, which is equivalent to

$$\lim_{n\to\infty} n \log(1 - G_1^{-1}(z_1(n))) = 0. \tag{5.4}$$

In order for (5.4) to hold we must have $\lim_{n\to\infty} G_1^{-1}(z_1(n)) = 0$. Since

$$\lim_{x\to 0} \frac{\log(1-x)}{x} = -1,$$

a condition equivalent to (5.4) is

$$\lim_{n\to\infty} nG_1^{-1}(z_1(n)) = 0. \tag{5.5}$$

From (3.10) a sufficient condition for

$$\lim_{n\to\infty} \mathbb{P}(V_1(n) < Z_1(n)) = 1$$

is $\lim_{n\to\infty} \binom{n}{d}\left[1 - G_1^{-1}(Z_1(n))\right]^{n-d} = 0$, which is equivalent to

$$\Leftrightarrow \lim_{n\to\infty} \log\binom{n}{d} + (n-d)\log\left[1 - G_1^{-1}(Z_1(n))\right] = -\infty. \tag{5.6}$$

Equation (5.6) holds if and only if

$$\lim_{n\to\infty} d \log n - nG_1^{-1}(Z_1(n)) = -\infty. \tag{5.7}$$

A sufficient condition for (5.7) is

$$\lim_{n\to\infty} \frac{n}{\log n} G^{-1}(Z_1(n)) = \infty. \tag{5.8}$$

$\square$

It is easy to see that the result of Proposition 5.1.1 holds for the sampled Problems (3.3) and (4.1) as well. We thus have that for $i = 1, 2, 3$,

$$\lim_{n\to\infty} \frac{n}{\log n} G_i^{-1}(Z_i(n)) = \infty \Rightarrow \lim_{n\to\infty} \mathbb{P}(V_i(n) < Z_i(n)) = 1$$

and

$$\lim_{n\to\infty} nG_i^{-1}(z_i(n)) = 0 \Rightarrow \lim_{n\to\infty} \mathbb{P}(V_i(n) > z_i(n)) = 1.$$

Before we proceed with our analysis, we recall a few definitions related to the asymptotic behavior of sequences of random variables.

**Definition 5.1.1** *If $\{X(n)\}_{n=1}^{\infty}$ is a sequence of random variables defined on some probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and $g$ is a real function, we have $X(n) = O_{\mathbb{P}}(g(n))$, if for any real sequence $\alpha(n)$ such that $\alpha(n) \to \infty$, we have*

$$\frac{X(n)}{g(n)} \frac{1}{\alpha(n)} \xrightarrow{\mathbb{P}} 0 \text{ as } n \to \infty$$

*Similarly, we have $X(n) = \Omega_{\mathbb{P}}(g(n))$, if for any real sequence $\alpha(n)$ such that $\alpha(n) \to \infty$, we have*

$$\frac{X(n)}{g(n)} \alpha(n) \xrightarrow{\mathbb{P}} \infty \text{ as } n \to \infty.$$

*We write $X(n) = \Theta_{\mathbb{P}}(g(n))$ if $X(n) = \Omega_{\mathbb{P}}(g(n))$ and $X(n) = O_{\mathbb{P}}(g(n))$.*

## 5.2 The SVD problem

From Proposition 3.2.3 we know the asymptotic behavior of $G_1^{-1}(\cdot)$ close to 1. We can combine this result with Proposition 5.1.1 and get the following asymptotic result for the optimal value $V_1(n)$.

**Proposition 5.2.1** *For the optimal value $V_1(n)$ of the sampled Problem (3.2) we have*

$$V_1(n) - 1 = O_{\mathbb{P}}\left(\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}\right)$$

*and*

$$V_1(n) - 1 = \Omega_{\mathbb{P}}\left(n^{-2/(d-l_1)}\right),$$

*where $l_1$ is the number of diagonal elements of the matrix A that are equal to 1.*

**Proof:** From Proposition 3.2.3 we have that as $v \downarrow 1$,

$$G_1^{-1}(v) = \Theta\left((v-1)^{(d-l_1)/2}\right).$$

Using this result and Proposition 5.1.1, we see that for a sequence $\{Z_1(n)\}_{n=1}^{\infty}$ such that

$$\lim_{n \to \infty} n \left(Z_1(n) - 1\right)^{(d-l_1)/2} = 0, \tag{5.9}$$

we have

$$\lim_{n \to \infty} \mathbb{P}\left(V_1(n) > Z_1(n)\right) = 1.$$

Consider any sequence $\{\alpha(n)\}_{n=1}^{\infty}$ such that

$$\lim_{n \to \infty} \alpha(n) = \infty.$$

Then we have for any $M > 0$ that

$$\lim_{n \to \infty} \mathbb{P}\left((V_1(n) - 1)\alpha(n)n^{2/(d-l_1)} > M\right) =$$

$$\lim_{n \to \infty} \mathbb{P}\left(V_1(n) > 1 + \frac{M}{\alpha(n)}n^{-2/(d-l_1)}\right) = 1,$$

where the last limit is equal to 1 because the sequence $1 + \frac{M}{\alpha(n)}n^{-2/(d-l_1)}$ satisfies condition (5.9). Thus, we have

$$(V_1(n) - 1)\alpha(n)n^{2/(d-l_1)} \overset{\mathbb{P}}{\to} \infty.$$

77

Similarly, for a sequence $\{z_1(n)\}_{n=1}^{\infty}$ such that

$$\lim_{n\to\infty} (z_1(n) - 1)^{(d-l_1)/2} \frac{n}{\log n} = \infty, \tag{5.10}$$

we have

$$\lim_{n\to\infty} \mathbb{P}\left(V_1(n) < z_1(n)\right) = 1.$$

Consider any sequence $\{\alpha(n)\}_{n=1}^{\infty}$ such that

$$\lim_{n\to\infty} \alpha(n) = \infty.$$

Then we have for any $M > 0$ that

$$\lim_{n\to\infty} \mathbb{P}\left(\frac{V_1(n) - 1}{\alpha(n)} \left(\frac{n}{\log n}\right)^{2/(d-l_1)} < M\right) =$$
$$\lim_{n\to\infty} \mathbb{P}\left(V_1(n) < 1 + M\alpha(n) \left(\frac{n}{\log n}\right)^{-2/(d-l_1)}\right) = 1,$$

where the last limit is equal to 1 because the sequence $1 + M\alpha(n)\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}$ satisfies condition (5.10). Thus, for any sequence $\{\alpha(n)\}_{n=1}^{\infty}$ such that $\lim_{n\to\infty} \alpha(n) = \infty$, we get

$$\frac{V_1(n) - 1}{\alpha(n)} \left(\frac{n}{\log n}\right)^{2/(d-l_1)} \xrightarrow{\mathbb{P}} 0.$$

$\square$

Proposition 5.2.1 describes how fast the optimal value $V_1(n)$ converges to 1 as $n$ tends to infinity. We see that this depends only on the dimension $d$ of the problem and the number $l_1$ of singular values that are equal to 1.

Similarly, we can apply the same technique to the SVD-based approximating Problem (3.3) and get that

$$V_2(n) - 1 = O_{\mathbb{P}}\left(\left(\frac{n}{\log n}\right)^{-2/(d-l_2)}\right) \tag{5.11}$$

and

$$V_2(n) - 1 = \Omega_{\mathbb{P}}\left(n^{-2/(d-l_2)}\right), \tag{5.12}$$

where $l_2$ is the number of diagonal elements of the matrix $A^{(k)}$ that are equal to 1.

We can use Proposition 5.2.1 in order to get asymptotic results for the relative error $R_2(n)$. We distinguish two cases that we will treat separately, $l_1 > l_2$ and $l_1 = l_2$. This choice of cases is dictated by the fact that the asymptotic behavior of the optimal value of the sampled problems depends on the number of diagonal elements equal to 1.

**Proposition 5.2.2** *Let $R_2(n)$ be the relative error as defined in (3.4). Let $l_1$ and $l_2$ be the number of diagonal elements equal to 1 in $A$ and $A^{(k)}$ respectively. Then, if $l_1 > l_2$, we have*

$$R_2(n) = O_{\mathbb{P}}\left(\left(\frac{n}{\log n}\right)^{-2/(d-l_2)}\right)$$

*and*

$$R_2(n) = \Omega_{\mathbb{P}}\left(n^{-2/(d-l_2)}\right).$$

**Proof:** From Proposition 5.2.1 we have that for any sequences $\{\alpha_1(n), \alpha_2(n)\}_{n=1}^{\infty}$ with $\lim_{n\to\infty} \alpha_1(n) = \infty$ and $\lim_{n\to\infty} \alpha_2(n) = \infty$,

$$\lim_{n\to\infty} \mathbb{P}\left(V_2(n) > 1 + \frac{1}{\alpha_2(n)} n^{-2/(d-l_2)}\right) = 1$$

and

$$\lim_{n\to\infty} \mathbb{P}\left(V_1(n) < 1 + \alpha_1(n)\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}\right) = 1.$$

Thus

$$\lim_{n\to\infty} \mathbb{P}\left(R_2(n) > \frac{\frac{1}{\alpha_2(n)} n^{-2/(d-l_2)} - \alpha_1(n)\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}}{1 + \alpha_1(n)\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}}\right) = 1.$$

79

For any sequence $\{\gamma(n)\}_{n=1}^{\infty}$, satisfying $\lim\limits_{n\to\infty} \gamma(n) = \infty$, we have

$$\lim_{n\to\infty} \mathbb{P}\left( R_2(n) n^{2/(d-l_2)} \gamma(n) > \frac{\frac{\gamma(n)}{\alpha_2(n)} - \gamma(n)\alpha_1(n) \left(\frac{n}{\log n}\right)^{-2/(d-l_1)} n^{2/(d-l_2)}}{1 + \alpha_1(n) \left(\frac{n}{\log n}\right)^{-2/(d-l_1)}} \right) = 1.$$

By choosing appropriate sequences $\alpha_1(n)$ and $\alpha_2(n)$, we have that

$$\lim_{n\to\infty} \frac{\frac{\gamma(n)}{\alpha_2(n)} - \gamma(n)\alpha_1(n) \left(\frac{n}{\log n}\right)^{-2/(d-l_1)} n^{2/(d-l_2)}}{1 + \alpha_1(n) \left(\frac{n}{\log n}\right)^{-2/(d-l_1)}} = \infty.$$

We thus see that

$$R_2(n) n^{2/(d-l_2)} \gamma(n) \overset{\mathbb{P}}{\to} \infty.$$

The claim

$$R_2(n) \left(\frac{n}{\log n}\right)^{\frac{2}{d-l_2}} \frac{1}{\gamma(n)} \overset{\mathbb{P}}{\to} 0$$

follows easily from the inequality $R_2(n) \le V_2(n) - 1$ and Proposition 5.2.1.

$\square$

This result states that when $l_2 < l_1$, the relative error $R_2(n)$ converges to 0 at the same rate as $V_2(n)$. In other words, since the optimal value $V_2(n)$ of the approximating problem converges to 1 at a slower rate than the optimal value $V_1(n)$ of the original problem, the error is dominated by $V_2(n)$.

We next treat the case $l_1 = l_2$. In this case, we can only provide a lower bound on how fast $R_2(n)$ converges to 0.

**Proposition 5.2.3** *Let $R_2(n)$ be the relative error as defined in (3.4). Let $l_1$ and $l_2$ be the number of diagonal elements equal to 1 in $A$ and $A^{(k)}$ respectively. Then, if $l_1 = l_2$, we have*

$$R_2(n) = O_{\mathbb{P}}\left( \left(\frac{n}{\log n}\right)^{-2/(d-l_1)} \right).$$

**Proof:** The result is a direct implication of the inequality $R_2(n) \leq V_2(n) - 1$ and (5.12).

$\square$

Proposition 5.2.3 states that when $l_2 = l_1$, the relative error $R_2(n)$ converges to 0 at least as fast as $V_2(n)$ and $V_1(n)$. Our asymptotic result in Proposition 5.2.1 is not fine enough to give us a lower asymptotic bound on $R_2(n)$ as well. Under the condition $l_1 = l_2$ though, we have that both $V_1(n) - 1$ and $V_2(n) - 1$ roughly behave as $n^{-2/(d-l_1)}$ for large $n$, which suggests that $R_2(n)$ should also behave as $n^{-2/(d-l_1)}$ asymptotically.

## 5.3 The Nyström problem

In this section we focus on asymptotic results about the optimal value $V_3(n)$ of Problem (4.1) and the relative difference

$$R_3(n) = \frac{V_3(n) - V_1(n)}{V_1(n)},$$

between the optimal values of Problems (4.1) and (3.2). We will present the results in the case $c = e_1$. In this case we have a bound on the function $G_3^{-1}(v)$ for $v$ close to 1.

We have through Propositions 4.3.1 and 4.3.2, for $v$ close to 1, that

$$(1 - G_3^{-1}(v))^n \leq \mathbb{P}(V_3(n) > v) \leq \binom{n}{d}(1 - G_3^{-1}(v))^{n-d}$$

and

$$G_1^{-1}(v) \leq G_3^{-1}(v) \leq \frac{G_1^{-1}(v)}{2}.$$

81

We thus get the following Proposition that describes the asymptotic behavior of $G_3^{-1}(v)$ near 1.

**Proposition 5.3.1** *If $c = e_1$, when $v \downarrow 1$,*

$$G_3^{-1}(v) = \Theta\left((v-1)^{\frac{d-l_1}{2}}\right),$$

*where $l_1$ is the number of diagonal elements of $A$ that are equal to 1.*

**Proof:** The result follows directly from the inequality

$$G_1^{-1}(v) \leq G_3^{-1}(v) \leq \frac{G_1^{-1}(v)}{2}.$$

and the fact that

$$G_1^{-1}(v) = \Theta\left((v-1)^{\frac{d-l_1}{2}}\right).$$

$\square$

Based on Proposition 5.3.1, we prove the following asymptotic result for the optimal value $V_3(n)$.

**Proposition 5.3.2** *For the optimal value $V_3(n)$ of the sampled Problem (4.1) with $l = 1$, where $l$ is the number of non-zero elements of the vector $c$, we have*

$$V_3(n) - 1 = O_{\mathbb{P}}\left(\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}\right)$$

*and*

$$V_3(n) - 1 = \Omega_{\mathbb{P}}\left(n^{-2/(d-l_1)}\right),$$

*where $l_1$ is the number of diagonal elements of the matrix $A$ that are equal to 1.*

**Proof:** The proof follows the same steps as the proof of Proposition 5.2.1.

$\square$

This result states that if the number of nonzero components of $c$ is 1, then $V_3(n)$ converges to 1 roughly as fast as $n^{-2/(d-l_1)}$.

We get the following result for the relative error $R_3(n)$.

**Proposition 5.3.3** *Consider Problems (3.2) and (4.1) with $l = 1$, where $l$ is the number of non-zero elements of the vector $c$. Let $R_3(n)$ be the relative error as defined in (4.2). Let $l_1$ be the number of diagonal elements equal to 1 in A. Then, we have*

$$R_3(n) = O_{\mathbb{P}}\left(\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}\right).$$

**Proof:** The result is a direct implication of the inequality

$$R_3(n) \leq V_3(n) - 1$$

and Corollary 5.3.2.

$\square$

Proposition 5.3.3 states that when the number of nonzero components of $c$ is 1, the relative error $R_3(n)$ converges to 0 at least as fast as $\left(\frac{n}{\log n}\right)^{-2/(d-l_1)}$. Our asymptotic results in Propositions 5.2.1 and 5.3.2 are not fine enough to give us a lower asymptotic bound on $R_3(n)$ as well.

CHAPTER 6

**CONCLUSIONS**

## 6.1 Discussion

So far in this dissertation, we have derived a number of properties of the optimal values of the problems in our model and the relative errors of the two approximations. In this section we will analyze these results and discuss their implications for applying low-rank approximations to optimization problems in practice.

The fundamental result in this dissertation is Theorem 2.4.2, which characterizes the probability distributions of sampled problems with independent convex constraints and linear objective function. Through Proposition 3.2.1, we have proved that for every $v > 1$, the optimal value $V_1$ of the original Problem (3.2) satisfies the inequality

$$P\left(\frac{\sum_{j=1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{v^2}\right)^n \leq \mathbb{P}(V_1 > v) \leq \binom{n}{d}P\left(\frac{\sum_{j=1}^d A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{v^2}\right)^{n-d},$$

where $\{Y_j\}_{j=1}^d$ are standard normal random variables. Furthermore, for Problem (3.3) we have that

$$P\left(\frac{\sum_{j=1}^k A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{v^2}\right)^n \leq \mathbb{P}(V_2 > v) \leq \binom{n}{d}P\left(\frac{\sum_{j=1}^k A_{jj}Y_j^2}{\sum_{j=1}^d Y_j^2} < \frac{1}{v^2}\right)^{n-d}.$$

These inequalities quantify the connection between the optimal values $V_1$ and $V_2$ and the singular values $A_{11}, ..., A_{dd}$ of the constraint matrices $W_i$.

We can use this result to give some theoretical justification for the success of low-rank approximation in optimization problems where the constraint matrices have the same singular values and the singular values decay quickly. If

in our model the $d - k$ smallest singular values $A_{k+1,k+1}, ..., A_{dd}$ are close to 0, the probabilities

$$P\left(\frac{\sum_{j=1}^{d} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} < \frac{1}{v^2}\right) \text{ and } P\left(\frac{\sum_{j=1}^{k} A_{jj} Y_j^2}{\sum_{j=1}^{d} Y_j^2} < \frac{1}{v^2}\right)$$

are close to each other. Then, based on our analysis, we expect the optimal value $V_2(n)$ of the approximating problem with rank $k$ constraints to be a good approximation for the optimal value $V_1(n)$ of the original problem.

In Section 5.2 we have derived asymptotic results about the optimal values $V_1(n)$ and $V_2(n)$ as the number of constraints $n$ tends to infinity. Our main motivation for this was to explain the behavior of low-rank approximations in problems with a large number of constraints. In the asymptotic analysis for the optimal value $V_1(n)$, we have shown that the crucial quantity is the number $l_1$ of singular values equal to 1. We proved in Proposition 5.2.1 that $V_1(n) - 1$ essentially behaves as $n^{-2/(d-l_1)}$, for a large number of constraints $n$. Since typically in applications where a low-rank approximation is attractive, the number of variables $d$ is large, this implies that the optimal value $V_1(n)$ converges to 1 slowly.

Our main quantity of interest though, is the relative difference $R_2(n)$ in the optimal values of Problems (3.2) and (3.3). Because of the importance of the number of singular values $l_1$ and $l_2$ that are equal to 1 in the matrices $A$ and $A^{(k)}$ respectively, we have two cases: $l_1 > l_2$ and $l_1 = l_2$. We interpret the former case as an approximation of an optimization problem with constraints that have a few large singular values with a low-rank problem that does not approximate all the large singular values. The latter case corresponds to an approximation where all the large singular values of the original problem are approximated in the low-rank problem.

If $l_1 > l_2$, from Proposition 5.2.2 we have that for large $n$, $R_2(n)$ essentially

behaves as $n^{-2/(d-l_2)}$. This means that the approximation error is dominated by the optimal value $V_2(n)$ and the relative error converges to 0 with the rate depending heavily on $l_2$ the number of elements in the matrix $A^{(k)}$ that are equal to 1. Intuitively, this means that unless the number of constraints is large, a low-accuracy low-rank approximation will not result in small error in terms of optimal values. On the other hand, based on this result we see that if the number of constraints $n$ in the problem is large, even a rough low-rank approximation can give satisfactory results.

Under the assumption $l_1 = l_2$, we see from Proposition 5.2.3 that the relative error $R_2(n)$ converges to 0 roughly as fast as $n^{-2/(d-l_1)}$. We thus get a similar conclusion, i.e., a large number of constraints is required in order for the approximation error to get small.

We then consider the approximating Problem (4.1). We have proved that if $A$ is of full rank and the rank of the approximating matrices is $k \geq l$, where $l$ is the number of non-zero elements of the objective function vector $c$, then

$$(1 - G_3^{-1}(v))^n \leq \mathbb{P}(V_3 > v) \leq \binom{n}{d}(1 - G_3^{-1}(v))^{n-d},$$

for any $v > 1$. Through the analysis in Section 4.2 we have shown that when $A$ is close to the identity and $v$ close to 1, $\mathbb{P}(V_3(n) < v)$ tends to 1 faster than $\mathbb{P}(V_2(n) < v)$, suggesting that $V_3(n)$ is a better approximation for $V_1(n)$ for large $n$.

In the special case $l = 1$, i.e., when the vector $c$ is parallel to one of the axes, we have shown that for $v$ close to 1, we have

$$\frac{1}{2}P\left(\frac{\sum_{j=1}^{d} A_{jj}Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{v^2}\right) \leq G_3^{-1}(v) \leq P\left(\frac{\sum_{j=1}^{d} A_{jj}Y_j^2}{\sum_{j=1}^{d} Y_j^2} > \frac{1}{v^2}\right).$$

This suggests that under the condition $l = 1$, Problem (4.1) provides a good approximation to the original Problem (3.2), especially when the singular val-

ues $A_{11}, ..., A_{dd}$ are not close to 0. Furthermore, we have shown in this case that for large $n$, the optimal values $V_1(n)$ and $V_3(n)$ converge to 1 roughly as fast as $n^{-2/(d-l_1)}$. Thus, in problems where $l_2 < l_1$, we see that $V_3(n)$ converges to 1 faster than $V_2(n)$, making the approximation based on Algorithm 4.1.3 attractive. This is a somewhat surprising result, since the SVD results in an optimal low-rank approximation. Algorithm 4.1.3, though, provides a better approximation to the optimal value because it approximates the feasible region close to the optimal value, when the objective function vector $c$ is parallel to one of the axes.

If the rank of the approximating matrices is $k < l$, we have proved that the approximation based on Algorithm 4.1.3 results in error that cannot be arbitrarily close to 0, even for a large number of constraints. This implies that in problems where $c$ has many non-zero components, using Algorithm 4.1.3 is not practical.

In such a case, one could adapt Algorithm 4.1.3 as follows. If $U \in O^d$ is such that $Uc = e_1$, we define $y \in \mathbb{R}^d = Ux$. Then the original problem takes the form

$$\max \ e_1^T y \tag{6.1}$$

$$\text{subject to} \ \ y^T(UQ_i^T)A(Q_iU^T)y \leq 1, i = 1, ..., n.$$

In our model, due to the fact that $\{Q_i\}_{i=1}^n$ are uniformly distributed, the optimal value of this problem has the same distribution as $V_1$. Applying Algorithm 4.1.3 with $k \geq 1$ to the matrices

$$\tilde{W}_i = UQ_i^T AQ_iU^T, i = 1, ..., n,$$

gives us an approximating problem with optimal value $\tilde{V}_3$ that satisfies for $v$ close to 1 the relationship

$$\left(1 - G_1^{-1}(n)\right)^n \leq \mathbb{P}(\tilde{V}_3 > v) \leq \binom{n}{d}\left(1 - \frac{G_1^{-1}(v)}{2}\right)^{n-d}.$$

This means, that under the assumptions of our model, using such an approximation leads to an optimization problem with optimal value that is close to the optimal value of the original problem, even if the rank of the approximating matrices is $k = 1$. Furthermore, it is not clear whether such an approach should provide a good approximation for a general optimization problem with quadratic constraints

Reflecting on the motivating problem, it is not clear whether the conclusions that we have derived from our model extend to the IMRT problem. The symmetric positive definite matrices in the constraints of the IMRT problem do not have identical singular values, but according to the discussion in Section 2.2, their singular values are close to each other. It would be interesting to see if the results extend to this case or whether there is a continuity result for the optimal values of such problems with respect to the singular values.

Furthermore, according to the analysis of our model, the optimal value of the approximating problem decays to 1 at most as fast as $n^{-2/(d-k)}$. Since $d$ is in the order of 1000 and $k$ is in the order of 10, this is consistent with the empirical findings in [13] and suggests that the test that was used in order to evaluate the quality of the solution was not suitable. A more appropriate test would be to estimate the covariance matrices by sampling a large number of shifts from a continuous distribution and then check whether the optimal solution is feasible for the resulting constraints.

In summary, our main goal in this dissertation was to explore an empirical phenomenon in terms of rigorous analysis of a stylized model. What we have proved should thus serve as an indication of what might hold in general. One can try to generalize our model, incorporating, for example, more general types

of constraints, or assuming a different distribution for their random components. This would make the model more realistic, but it would probably require different mathematical techniques to analyze it.

We have also demonstrated in our model that matrices with rapidly decaying eigenvalues can be replaced by a low-rank approximation, with minor loss of accuracy. This is a common observation in applications of low-rank approximations in various fields. Thus, our work serves in the direction of extending this idea in optimization applications and unifying our view of low-rank approximations.

Finally, although our work is mainly concerned with the theoretical explanation of a phenomenon, it leads to observations that show possible future directions for the application of low-rank approximations in optimization problems in practice. We have shown that in special cases, using Algorithm 4.1.3 in our model leads to better approximations in terms of optimal values. Also, in our model, we have proposed a way of applying Algorithm 4.1.3 that is adapted to the optimization problem and results in an approximation that performs well in our model.

## 6.2    Future directions

In reflecting on the research questions addressed, we observed several interesting research topics. These would allow us to further evaluate and expand the findings of this dissertation.

An obvious first direction for future research is to analyze more sophisticated

stylized optimization problems. It is of interest to see if similar results hold in optimization problems with constraints that are of different shape, without common centers and sampled from a general probability distribution. Inequality (2.15) holds for any sampled problem with independent convex contraints, but the corresponding chance-constrained problems are hard to solve, unless the probability distributions and the form of the constraints are carefully chosen. Also, it is not clear if the results hold for problems with constraints that are not sampled independently.

The fundamental result in this dissertation is Inequality (2.15), that provides a link between sampled and chance-constrainted problems. This inequality is based on a result that holds for general convex constraints, so it could be possible to derive a better bound, that holds for the model that we are considering. Apart from being interesting from a theoretical point of view, such an improvement could lead to tighter asymptotic results for the optimal values.

A weakness of our method is that it fails to provide results about the joint distribution of the optimal values of the sampled problems that we are studying. It would be of interest to investigate ways to generalize the chance-constrained methodology, in order to provide results about the joint distribution of the optimal values of the sampled problems, when the random components of the constraints are common. This could lead to improved bounds on the number of constraints required so that the relative difference of the optimal values is below some level with certain probability.

Finally, in our model, we have seen that although the SVD-based algorithm provides an optimal matrix approximation, in some cases it fails to provide as good an approximation compared to Algorithm 4.1.3. By taking into account

the structure of the problem that is approximated, it may be the case that there exist algorithms that require small number of operations and which provide a very good approximation to the problem. As an example, we can mention the adaptation of Algorithm 4.1.3 in our model that we presented in Section 6.1. Using information about the objective function during the approximation, the resulting approximating problem has an optimal value that is close to the optimal value of the original problem, at least in our model. Thus, a natural direction for future research is to design matrix approximation algorithms that can be adapted to the optimization problem.

# APPENDIX A

## PROOFS OF RESULTS IN SECTION 4.2

We begin with a lemma, that provides an extension to the coupon collector probability that we used in the case $c = e_1$. It provides the probability that after sampling uniformly at random $n$ times with replacement subsets of size $k$ from a set of size $d$, all $d$ elements are picked at least once and a certain subset of size $l \leq k$ is included in at least one of the random samples of size $k$.

**Lemma A.0.1** *Let $S$ be a set with $d$ elements and let $L \subseteq S$, where $L$ has $l$ elements. Suppose we sample $k \geq l$ distinct elements of $S$ with replacement. Each subset of size $k$ of $S$ has the same probability of being picked. Suppose we sample independently $n$ times. We define the events*

$$\mathcal{B} = \{ \text{After sampling } n \text{ times, all elements of } S \text{ have been sampled} \},$$

$$C = \{ \text{After sampling } n \text{ times, one of the samples contains all } l \text{ elements of } L \}.$$

*Let $p_n$ be the probability of the event*

$$\mathcal{A} = \mathcal{B} \cap C.$$

*We have that*

$$1 - p_n = O\left( \left( 1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}} \right)^n \right).$$

**Proof:** Let $(X_i)_{i=1}^{\infty}$ be a stochastic process with values in

$$F = \{ k, k+1, ..., d \} \times \{ 0, 1, 2, ..., l \} \times \{ 0, 1 \},$$

such that for each $i$, $X_i = (X_i(1), X_i(2), X_i(3))$, where

- $X_i(1)$ is equal to 0 if no sample out of the first $i$ contains the set $\{1, 2, ..., l\}$ and 1 otherwise.

- $X_i(2)$ is the number of different elements of $\{1, 2, ..., l\}$ sampled in the first $i$ samples,

- $X_i(3)$ is the number of different elements of $\{1, 2, ..., d\}$ sampled in the first $i$ samples.

It is easy to see that $(X_i)_{i=1}^{\infty}$ is a Markov chain. We arrange the state space using a lexicographical ordering (considering the state vector in the form $X(1), X(2), X(3)$.) Then, the resulting one step probability transition matrix $P$ is upper triangular. We will focus on the diagonal elements, because these are equal to the eigenvalues.

Assume that the current state is $(v_1, v_2, v_3)$, where $v_2 < l$. This also implies that $v_1 \leq d - (l - v_2)$ and $v_3 = 0$. Then we have that

$$P(X_{i+1} = (v_1, v_2, v_3)|X_i = (v_1, v_2, v_3)) = \frac{\binom{v_1}{k}}{\binom{d}{k}}.$$

If $v_2 = l, v_3 = 0$, then

$$P(X_{i+1} = (v_1, v_2, v_3)|X_i = (v_1, v_2, v_3)) = 1 - \frac{\binom{v_1-l}{k-l}}{\binom{d}{k}}.$$

Finally, if $v_2 = l, v_3 = 0$, then we have that

$$P(X_{i+1} = (v_1, v_2, v_3)|X_i = (v_1, v_2, v_3)) = \frac{\binom{v_1}{k}}{\binom{d}{k}}.$$

Then the second largest eigenvalue is given by the second largest of the above values, which is

$$\lambda_2 = 1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}}.$$

We have that

$$p_n = P(X_n = (d, l, 1)) = \sum_j h_j(n-1)\lambda_j^{n-1},$$

where $h_j$ are polynomials with degree that depends on the algebraic multiplicity of $\lambda_j$. Since we need $\lim_{n\to\infty} p_n = 1$, and the second largest eigenvalue is $\lambda_2$, we get that

$$1 - p_n = O(\lambda_2^n).$$

$\square$

We use this lemma, in the following proposition.

**Proposition A.0.1** *Let $V_3(\theta)$ be the optimal value of Problem (4.14). Then*

$$1 - \mathbb{P}\left(\lim_{\theta\uparrow 1} V_3(\theta) = 1\right) = O\left(\left(1 - \frac{\binom{d-l}{k-l}}{\binom{d}{k}}\right)^n\right). \tag{A.1}$$

**Proof:**

We fix the orthogonal matrices $Q_1, ...., Q_n$ and we assume that they are such that for every $\theta$ in a neighbourhood of 1, the following hold:

1. Columns with all possible indices $\{1, 2, ..., d\}$ are picked during application of Algorithm 4.1.3 to matrices $Q_i^T A(\theta) Q_i, i = 1, ..., n$,

2. Columns with indices $\{1, 2, ..., l\}$ are picked together during application of Algorithm 4.1.3 to at least one of the matrices $Q_i^T A(\theta) Q_i, i = 1, ..., n$.

Let $\mathcal{K} \subset (O^d)^n$ be the set of all such $(Q_1, ..., Q_n)$. As in the proof of Proposition (4.2.1), we define $M(\theta)$ to be the minimum of all elements of matrices $C_i(\theta)$ that

94

tend to 1 as $\theta$ goes to 1 and $m(\theta)$ to be the maximum in absolute value of all elements of the matrices $C_i(\theta)$ that tend to 0. We then have that in a neighbourhood of 1, it holds that

$$|x_i| \leq \sqrt{\frac{1}{M(\theta) - 2(d-1)m(\theta)}} = K(\theta).$$

Let $Q_p^T A(\theta) Q_p$ be one of the matrices such that columns $\{1, ....l\}$ are picked during application of Algorithm 4.1.3 . We let $C(\theta) = C_p(\theta)$, for simplicity in notation. Then we have that the $p$-th constraint gives, for $\theta$ in a neighbourhood of 1,

$$x^T C(\theta) x \leq 1$$

$$\Leftrightarrow \sum_{j,v=1}^{l} C_{jv}(\theta) x_j x_v + \sum_{j,v=l+1}^{d} C_{jv}(\theta) x_j x_v + 2 \sum_{j=1}^{l} \sum_{v=l+1}^{d} C_{vl}(\theta) x_j x_v \leq 1$$

Since $C$ is a symmetric positive semidefinite matrix, we get

$$\sum_{j,v=1}^{l} C_{jv}(\theta) x_j x_v \leq 1 - 2 \sum_{j=1}^{l} \sum_{v=l+1}^{d} C_{vl}(\theta) x_j x_v$$

$$\Rightarrow \sum_{j,v=1}^{l} C_{jv}(\theta) x_j x_v \leq 1 + 2l(d-l)m(\theta)K(\theta)^2$$

$$\Leftrightarrow \sum_{j=1}^{l} C_{jj}(\theta) x_j^2 \leq 1 + 2l(d-l)m(\theta)K(\theta)^2 - \sum_{j,v=1, j\neq v}^{l} C_{jv}(\theta) x_j x_v$$

$$\Rightarrow \sum_{j=1}^{l} C_{jj}(\theta) x_j^2 \leq 1 + 2l(d-l)m(\theta)K(\theta)^2 + l(l-1)m(\theta)K(\theta)^2$$

$$\Leftrightarrow \sum_{j=1}^{l} x_j^2 \leq 1 + \left(2ld - l^2 - l\right)m(\theta)K(\theta)^2 + \sum_{j=1}^{l}(1 - C_{jj}(\theta))x_j^2$$

$$\Rightarrow \sum_{j=1}^{l} x_j^2 \leq 1 + \left(2ld - l^2 - l\right)m(\theta)K(\theta)^2 + l(1 - M(\theta))K(\theta)^2$$

Let $x_{1l} = \begin{bmatrix} x_1 \cdots x_l \end{bmatrix}^T$. Then we have that

$$\|x_{1l}\|_2 \leq \sqrt{1 + \alpha(\theta)},$$

95

where $\lim\limits_{\theta\uparrow 1} \alpha(\theta) = 0$. Since for every $x$ in the feasible region we have that $\left|c^T x\right| \leq$ $\|c\|_2 \|x_{1l}\|_2 = \|x_{1l}\|_2$, it follows that

$$V_3(\theta) \leq \sqrt{1 + \alpha(\theta)},$$

which implies

$$\lim_{\theta\uparrow 1} V_3(\theta) = 1.$$

So, we have that

$$\mathbb{P}\left(\lim_{\theta\uparrow 1} V_3(\theta) = 1\right) \geq \mathbb{P}(\mathcal{K}).$$

and then the result follows from Lemma A.0.1.

$\square$

# BIBLIOGRAPHY

[1] M. Abramowitz and I. A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover, New York, 1964.

[2] Dimitris Achlioptas and Frank McSherry. Fast computation of low-rank matrix approximations. *J. ACM*, 54(2), 2007.

[3] F. Bach and M. Jordan. Predictive low-rank decomposition for kernel methods. In *Proceedings of the 22nd International Conference on Machine Learning*, 2005.

[4] A. Ben-Tal and A. Nemirovski. Robust convex optimization. *Math. Oper. Res.*, 23(4):769–805, 1998.

[5] A. Ben-Tal and A. Nemirovski. Robust solutions of uncertain linear programs. *Operations Research Letters*, 25(1):1–13, 1999.

[6] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization*. SIAM, Philadelphia, PA, 2001.

[7] M.W. Berry, S.T. Dumais, and G.W. O'Brian. Using linear algebra for intelligent information retrieval. *SIAM Review*, 37(4):573–595, 1995.

[8] G. Calafiore and M.C. Campi. Uncertain convex programs: randomized solutions and confidence levels. *Math. Program.*, 102(1):25–46, 2005.

[9] G. Calafiore and M.C. Campi. The scenario approach to robust control design. *IEEE Transactions on Automatic Control*, 51(5):742–753, 2006.

[10] M. Carol, W. H. Grant, D. Pavord, P. Eddy, H. S. Targovnik, B. Butler, S. Woo, J. Figura, V. Onufrey, R. Grossman, and R. Selkar. Initial clinical experience with the Peacock intensity modulation of a 3-D conformal radiation therapy system. *Sterotactic and Functional Neurosurgery*, 66(1-3):3034, 1996.

[11] A. Charnes and W.W. Cooper. Chance-constrained programming. *Management Sci.*, 6:73–79, 1959.

[12] M. Chu, Y. Zinchenko, S.G. Henderson, and M. Sharpe. Robust optimization for intensity modulated radiation therapy treatment planning under uncertainty. *Physics in Medicine and Biology*, 50(23):5463–5477, 2005.

[13] Millie Chu. *Robust Intensity Modulated Radiation Therapy Treatment Planning.* PhD thesis, Cornell University, 2006.

[14] D.P. de Farias and B. Van Roy. On constraint sampling in the linear programming approach to approximate dynamic programming. *Math. Oper. Res.*, 29(3):462–478, 2004.

[15] A. de Moivre. De mensura sortis, seu, de probabilitate eventuum in ludis a casu fortuito pendentibus. *Philosophical Transactions of the Royal Society of London A*, 27:213–264, 1711.

[16] A. de Moivre. *The Doctrine of Chances: or, a Method of Calculating the Probabilities of Events in Play.* London, 1718. Reprint of the posthumous third edition by Chelsea, New York, 1967.

[17] L. M. Delves and J. L. Mohamed. *Computational methods for integral equations.* Cambridge University Press, New York, NY, USA, 1986.

[18] Amit Deshpande and Santosh Vempala. Adaptive Sampling and Fast Low-Rank Matrix Approximation. In *Proc. of RANDOM*, 2006.

[19] Petros Drineas, Ravi Kannan, and Michael W. Mahoney. Fast Monte Carlo algorithms for matrices II: Computing a low-rank approximation to a matrix. *SIAM Journal on Computing*, 36(1):158–183, 2006.

[20] Petros Drineas and Michael W. Mahoney. On the Nyström method for approximating a Gram matrix for improved kernel-based learning. *Journal of Machine Learning Research*, 6(12):p2153 – 2175, 2005.

[21] Emre Erdogan and Garud Iyengar. Ambiguous chance constrained problems and robust optimization. *Math. Program.*, 107(1-2):37–61, 2006.

[22] S. Fine and K. Scheinberg. Efficient SVM training using low-rank kernel representations. *Journal of Machine Learning Research*, 2:243–264, 2001.

[23] G. H. Golub and C. F. van Loan. *Matrix Computations.* Johns Hopkins University Press, 3rd edition, 1996.

[24] Intensity Modulated Radiation Therapy Collaborative Working Group. Intensity-modulated radiotherapy: current status and issues of interest. *Int. J. Radiation Oncology Biol. Phys.*, 51(4):880914, 2001.

[25] P. R. Halmos. *Measure Theory*, volume 18 of *Graduate Texts in Mathematics*. Springer, NY, 1974.

[26] Jon M. Kleinberg. Authoritative sources in a hyperlinked environment. *J. ACM*, 46(5):604–632, 1999.

[27] P. S. Laplace. Mémoir sur les suites récurro-récurrentesat leurs usages dans la théorie des hasards. *Mémoires de l' Academie des Sciences de Paris*, 6:353–371, 1774.

[28] P. S. Laplace. *Théorie Analytique des Probabilités*. Courcier, Paris, 1812.

[29] C. Clifton Ling, Chandra Burman, Chen S. Chui, Stephen A. Leibel Gerald J. Kutcher, Thomas LoSasso, Radhe Mohan, Thomas Bortfeld, Larry Reinstein, Spiridon Spirou, X. H. Wang, Quiwen Wu, Michael Zelefsky, and Zvi Fuks. Conformal radiation treatment of prostate cancer using inversely-planned intensity-modulated photon beams produced with dynamic multileaf collimation. *Int. J. Radiation Oncology Biol. Phys.*, 35(4):721–730, 1996.

[30] H. Murase and S. K. Nayar. Visual learning and recognition of 3-D objects from appearance. *Int. J. Comput. Vision*, 14(1):5–24, 1995.

[31] C. H. Papadimitriou, P. Raghavan, H. Tamaki, and S. Vempala. Latent semantic indexing: A probabilistic analysis. *Journal of Computer and System Sciences*, 61(2):217 – 235, 2000.

[32] A. Prékopa. *Stochastic Programming*. Kluwer, Norwell, MA, 1995.

[33] A. Ruszczynski and A. Shapiro editors. *Stochastic Programming*. Handbook in Operations Research and Management Science. Elsevier, 2003.

[34] G.W. Stewart. The efficient generation of random orthogonal matrices with an application to condition estimation. *SIAM Journal in Numerical Analysis*, 17(3):403–409, 1980.

[35] S. Vajda. *Probabilistic Programming*. Academic, New York, 1972.

[36] C. K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In T. Leen, T. Dietterich, and V. Tresp, editors, *Neural Information Processing Systems 13*, pages 682–688. MIT Press, 2001.