

Positive-Definite Programming

Lieven Vandenberghe

Electrical Engineering Department

K. U. Leuven

Leuven, Belgium

Stephen Boyd

Electrical Engineering Department

Stanford University

Stanford, California

1 Introduction

1.1 Positive-definite programming

We consider the problem of minimizing a linear function of a variable $x \in \mathbb{R}^m$ subject to a *linear matrix inequality*:

$$\begin{aligned} & \text{minimize } c^T x, \\ & \text{subject to } F(x) \geq 0, \end{aligned} \quad (1)$$

where

$$F(x) \triangleq F_0 + \sum_{i=1}^m x_i F_i. \quad (2)$$

The problem data are the vector $c \in \mathbb{R}^m$ and $m+1$ symmetric matrices $F_0, \dots, F_m \in \mathbb{R}^{n \times n}$. The inequality sign in $F(x) \geq 0$ means that $F(x)$ is positive-semidefinite, i.e., $z^T F(x) z \geq 0$ for all $z \in \mathbb{R}^n$. This problem is called a *positive-definite program* (PDP), following Nesterov and Nemirovsky [43].

Problem (2) is a convex optimization problem since its objective and constraint are convex: if $F(x) \geq 0$ and $F(y) \geq 0$, then, for all $\lambda, 0 \leq \lambda \leq 1$,

$$F(\lambda x + (1 - \lambda)y) = \lambda F(x) + (1 - \lambda)F(y) \geq 0.$$

Although the PDP (1) may appear quite specialized, we will see that it includes many important optimization problems as special cases. For instance, consider the linear program

$$\begin{aligned} & \text{minimize } c^T x, \\ & \text{subject to } Ax \geq b, \end{aligned} \quad (3)$$

where the inequality denotes *componentwise* inequality. Since a vector $v \geq 0$ (componentwise) if and only if the matrix $\text{diag}(v)$ (i.e., the diagonal matrix with the components of v on its diagonal) is positive-semidefinite, we can express the LP (3) as a PDP with $F(x) = \text{diag}(Ax + b)$, i.e.,

$$F_0 = -\text{diag}(b), \quad F_i = \text{diag}(a_i), \quad i = 1, \dots, m,$$

where $A = [a_1 \dots a_m] \in \mathbb{R}^{n \times m}$.

Positive-definite programming can therefore be regarded as an extension of linear programming where the componentwise inequalities between vectors are replaced by matrix inequalities, or, equivalently, the first orthant is replaced by the cone of positive-semidefinite matrices. We can also view the PDP (1) as a semi-infinite linear program, since the matrix inequality $F(x) \geq 0$ is equivalent to an infinite set of linear constraints on x , i.e., $z^T F(x) z \geq 0$ for each $z \in \mathbb{R}^n$. It is therefore not surprising that the theory of positive-definite programming closely parallels linear programming theory, or that many algorithms for solving linear programs should have generalizations that handle PDPs. There are many important differences, however. For instance, the duality results are weaker for PDPs than for LPs. As another important difference, there is no simple or obvious analog of the simplex method for PDPs.

Before proceeding further we give a simple example of a *nonlinear* (convex) optimization problem that can be cast as a PDP, but not as a linear program. Consider the problem

$$\begin{aligned} & \text{minimize } \frac{(c^T x)^2}{d^T x} \\ & \text{subject to } Ax + b \geq 0, \end{aligned} \quad (4)$$

where we assume that $d^T x > 0$ whenever $Ax + b \geq 0$. We start with the standard trick of introducing an auxiliary variable t that serves as an upper bound on the objective:

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } Ax + b \geq 0 \\ & \quad \frac{(c^T x)^2}{d^T x} \leq t. \end{aligned} \quad (5)$$

In this formulation, the objective is a linear function of the variables x and t ; the nonlinear (convex) objective in (4) shows up as a nonlinear (convex) constraint in (5). These constraints, in turn, can be expressed as a linear matrix inequality in the variables x and t :

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } \begin{bmatrix} \text{diag}(Ax + b) & 0 & 0 \\ 0 & t & c^T x \\ 0 & c^T x & d^T x \end{bmatrix} \geq 0. \end{aligned} \quad (6)$$

(Here, again, $\text{diag}(v)$ represents the diagonal matrix with the elements of v on its diagonal.) Thus we have reformulated the nonlinear (convex) problem (4) as the PDP (6). More examples and applications will be given in the next section.

There are good reasons for studying positive-definite programming problems. First, positive-definiteness constraints arise directly in a number of important applications. Secondly, many convex optimization problems, e.g., linear programming and (convex) quadratically constrained quadratic programming, can be cast as PDPs. Positive-definite programming therefore offers a unified way to study the properties of and derive algorithms for a wide variety of convex optimization problems. Most importantly, however, *PDPs can be solved very efficiently, both in theory and in practice.*

Theoretical tractability follows from convexity, along with the observation that we can construct, in polynomial time, a cutting plane for the constraint set through any given infeasible point (see, e.g., [8, §2.3]). One can therefore apply the ellipsoid method of Yudin and Nemirovsky, and Shor (see [64, 56]) to solve problem (1) in polynomial time. In practice, however, the ellipsoid method is slow.

In this paper we concentrate on recently developed interior-point methods for positive-definite programming. Of course general-purpose nonlinear optimization methods (trust region methods, sequential quadratic optimization, ...) could be used, possibly after modification, to solve PDPs. Interior-point methods, however, enjoy several properties that make them especially interesting.

- It is now generally accepted that interior-point methods for linear programming are competitive with the simplex method and even faster for problems with more than 10,000 variables or constraints (see, e.g., [32]). One can expect to see the same trend more generally, especially since the very efficient simplex method has no counterpart in positive-definite programming. In our experience with positive-definite programming for control applications, we have found interior-point methods to be very efficient.
- Interior-point methods have a polynomial worst-case complexity.
- Interior-point methods are ideally suited for structured problems. We will see that every iteration of an interior-point method involves the solution of a least-squares problem with the same structure as $F(x)$ in (1). These matrices are often highly structured but not necessarily sparse. The structure can be exploited by combining interior-point methods with iterative least-squares methods such as conjugate-gradients [21] or LSQR [49]. This is not possible in the simplex method, for instance, nor in many other classical methods.

1.2 Examples and Applications

In this section we list a few examples and applications. The list is not exhaustive, and the purpose is more to give an idea of the generality of the problem. More examples are described in [43] and [8].

Quadratically Constrained Quadratic Programming

A convex quadratic constraint $(Ax + b)^T(Ax + b) - c^T x - d \leq 0$ can be written as

$$\begin{bmatrix} I & Ax + b \\ (Ax + b)^T & c^T x + d \end{bmatrix} \geq 0.$$

The left-hand side depends affinely on the vector x : it can be expressed as $F_0 + x_1 F_1 + \dots + x_m F_m \geq 0$, with

$$F_0 = \begin{bmatrix} I & b \\ b^T & d \end{bmatrix}, \quad F_i = \begin{bmatrix} 0 & a_i \\ a_i^T & c_i \end{bmatrix}, \quad i = 1, \dots, m,$$

where $A = [a_1 \dots a_m]$. Therefore, a general quadratically constrained quadratic program

$$\begin{aligned} & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i(x) \leq 0, \quad i = 1, \dots, L, \end{aligned}$$

where each f_i is a convex quadratic function $f_i(x) = (A_i x + b)^T(A_i x + b) - c_i^T x - d_i$, can be written as

$$\begin{aligned} & \text{minimize} && t \\ & \text{subject to} && \begin{bmatrix} I & A_0 x + b_0 \\ (A_0 x + b_0)^T & c_0^T x_0 + d_0 + t \end{bmatrix} \geq 0, \\ & && \begin{bmatrix} I & A_i x + b_i \\ (A_i x + b_i)^T & c_i^T x + d_i \end{bmatrix} \geq 0, \quad i = 1, \dots, L. \end{aligned}$$

This is a PDP in x and t , since one can think of the $L + 1$ matrix inequalities as diagonal blocks of one block diagonal matrix inequality $F(x, t) \geq 0$.

Matrix Norm and Maximum Eigenvalue Minimization

Suppose $A(x)$ is a (possibly rectangular) matrix that depends affinely on x : $A(x) = A_0 + x_1 A_1 + \dots + x_m A_m$. The problem of minimizing the (spectral, or maximum singular value) norm $\|A(x)\|$ over x is a PDP:

$$\begin{aligned} & \text{minimize} && t \\ & \text{subject to} && \begin{bmatrix} tI & A(x) \\ A(x)^T & tI \end{bmatrix} \geq 0. \end{aligned} \quad (7)$$

If $A(x)$ is a symmetric matrix, a related problem is to minimize the maximum eigenvalue of A :

$$\begin{aligned} & \text{minimize} && t \\ & \text{subject to} && tI - A(x) \geq 0. \end{aligned}$$

Note that both $\|A(x)\|$ and the maximum eigenvalue $\lambda_{\max}(A(x))$ are nondifferentiable functions of x .

Logarithmic Chebyshev Approximation

Suppose we wish to solve $Ax \approx b$ approximately, where $A = [a_1 \dots a_n]^T \in \mathbb{R}^{n \times m}$ and $b \in \mathbb{R}^n$. In Chebyshev approximation we minimize the infinity norm of the residual, i.e., we solve

$$\text{minimize } \max_i |a_i^T x - b_i|.$$

This can be cast as a linear program, with x and an auxiliary variable t as variables:

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } -t \leq a_i^T x - b_i \leq t, \quad i = 1, \dots, n. \end{aligned}$$

In some applications b_i has the dimension of a power or intensity, and is typically expressed on a logarithmic scale. In such cases the more natural optimization problem is

$$\text{minimize } \max_i |\log(a_i^T x) - \log(b_i)| \quad (8)$$

(assuming $b_i > 0$, and interpreting $\log(a_i^T x)$ as $-\infty$ when $a_i^T x \leq 0$).

This *logarithmic Chebyshev approximation* problem can be cast as a PDP. To see this, note that

$$|\log(a_i^T x) - \log(b_i)| = \log \max(a_i^T x/b_i, b_i/a_i^T x)$$

(assuming $a_i^T x > 0$). Problem (8) is therefore equivalent to

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } 1/t \leq a_i^T x/b_i \leq t, \quad i = 1, \dots, n, \end{aligned}$$

or,

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } \begin{bmatrix} t - a_i^T x/b_i & 0 & 0 \\ 0 & a_i^T x/b_i & 1 \\ 0 & 1 & t \end{bmatrix} \geq 0, \quad i = 1, \dots, n, \end{aligned}$$

which is a PDP. This example illustrates two important points. It shows that positive-definite programming includes many optimization problems that do not look like (1) at first sight. And secondly, it shows that the problem is much more general than linear programming, despite the close analogy.

Control and System Theory

Positive-definite programming problems arise frequently in control and system theory. Boyd, El Ghaoui, Feron and Balakrishnan catalog many examples in [8]. We will describe one simple example here.

Consider the *differential inclusion*

$$\frac{dx}{dt} \in \text{Co}\{A_1, \dots, A_L\}x(t), \quad (9)$$

where $x(t) \in \mathbb{R}^n$ and the matrices A_1, \dots, A_L are given, and $\text{Co}\{A_1, \dots, A_L\}$ denotes the *convex hull* of A_1, \dots, A_L . We seek an ellipsoidal invariant set, i.e., an ellipsoid \mathcal{E} such that for any x that satisfies (9), $x(T) \in \mathcal{E}$ implies $x(t) \in \mathcal{E}$ for all $t \geq T$. The existence of such an ellipsoid implies, for example, that all solutions of the differential inclusion (9) are bounded.

The ellipsoid $\mathcal{E} = \{x \mid x^T P x \leq 1\}$, where $P = P^T > 0$, is invariant if and only if the function $V(t) = x(t)^T P x(t)$ is nonincreasing for any solution x of (9). (In this case

we say that V is a quadratic Lyapunov function that proves stability of the differential inclusion (9).) Thus, \mathcal{E} is invariant if and only if

$$\frac{d}{dt} V(x(t)) = x(t)^T (A(t)^T P + P A(t)) x(t) \leq 0,$$

for any $x(t) \in \mathbb{R}^n$ and $A(t) \in \text{Co}\{A_1, \dots, A_L\}$. This is equivalent to $A^T P + P A \leq 0$ for all $A \in \text{Co}\{A_1, \dots, A_L\}$, which in turn is equivalent to the condition

$$A_k^T P + P A_k \leq 0, \quad k = 1, \dots, L.$$

This is a linear matrix inequality constraint in the matrix P , considered as the variable.

To find an invariant ellipsoid for the differential inclusion (9) (or verify that none exists), we need to solve the feasibility problem

$$P > 0, \quad A_k^T P + P A_k \leq 0, \quad k = 1, \dots, L \quad (10)$$

for the (matrix) variable P . Several standard methods can be used to convert this feasibility problem into a PDP that has an obvious initial feasible point. For instance, we can solve the PDP with variables $P = P^T \in \mathbb{R}^{n \times n}$ and $t \in \mathbb{R}$,

$$\begin{aligned} & \text{minimize } t \\ & \text{subject to } A_k^T P + P A_k \leq 0, \quad k = 1, \dots, L, \\ & \quad P \geq -tI, \\ & \quad P \leq I. \end{aligned}$$

(The last constraint is added, without loss of generality, to normalize the otherwise homogeneous problem.) This PDP can be initialized with $P = 0$, $t = 1$ and then solved; the optimum value of t is negative if and only if (10) is feasible.

Structural Optimization

Ben-Tal and Bendsoe in [11] consider the following problem from structural optimization. A structure of L linear elastic bars connect a set of N nodes. The geometry (topology and lengths of the bars) and the material (Young's modulus) are fixed; the task is to size the bars, i.e., determine appropriate cross-sectional areas for the bars. In the simplest version of the problem we consider one fixed set of externally applied nodal forces f_i , $i = 1, \dots, N$. (More complicated versions consider multiple loading scenarios.) The vector of (small) node displacements resulting from the load forces f will be denoted d . The objective is the elastic stored energy $\frac{1}{2} f^T d$, which is a measure of the inverse of the stiffness of the structure. We also need to take into account constraints on the total volume (or equivalently, weight), and upper and lower bounds on the cross-sectional area of each bar.

The design variables are the cross-sectional areas x_i . The relation between f and d is linear: $A(x)d = f$, where

$$A(x) \triangleq \sum_{i=1}^N x_i A_i$$

is called the stiffness matrix. The matrices A_i are all symmetric positive-semidefinite and depend only on fixed parameters (Young's modulus, length of the bars, and geometry). The optimization problem then becomes (see [11])

$$\begin{aligned} & \text{minimize} && f^T d \\ & \text{subject to} && A(x)d = f, \\ & && \sum_{i=1}^L l_i x_i \leq v, \\ & && \bar{x}_i \leq x_i \leq \bar{x}_i, \quad i = 1, \dots, L, \end{aligned}$$

where d and x are the variables, v is maximum volume, l_i are the lengths of the bars, and \bar{x}_i , \underline{x}_i the upper and lower bounds on the cross-sectional areas. For simplicity, we assume that $\bar{x}_i > 0$, and that $A(x)$ is positive-definite for all positive values of x_i . We can then eliminate d and write

$$\begin{aligned} & \text{minimize} && f^T A(x)^{-1} f \\ & \text{subject to} && \sum_{i=1}^L l_i x_i \leq v, \\ & && \bar{x}_i \leq x_i \leq \bar{x}_i, \quad i = 1, \dots, L, \end{aligned}$$

or,

$$\begin{aligned} & \text{minimize} && t \\ & \text{subject to} && \begin{bmatrix} t & f^T \\ f & A(x) \end{bmatrix} \geq 0, \\ & && \sum_{i=1}^L l_i x_i \leq v, \\ & && \bar{x}_i \leq x_i \leq \bar{x}_i, \quad i = 1, \dots, L, \end{aligned}$$

which is a PDP in x and t .

Pattern Separation by Ellipsoids

The simplest classifiers in pattern recognition use hyperplanes to separate two sets. A hyperplane $a^T x + b = 0$ separates two sets of points $\{x_i\}$ and $\{y_j\}$ if

$$\begin{aligned} a^T x_i + b &< 0 && \text{for all } i, \\ a^T y_j + b &> 0 && \text{for all } j. \end{aligned}$$

This is a set of linear inequalities in $a \in \mathbf{R}^n$ and $b \in \mathbf{R}_1$ and a solution can be found by linear programming. If the two sets cannot be separated by a hyperplane, we can try to separate them by a quadratic surface. In other words we seek a quadratic function $f(x) = x^T P x + b^T x + c$ that satisfies

$$\begin{aligned} (x_i)^T P x_i + b^T x_i + c &< 0 && \text{for all } i, && (11) \\ (y_j)^T P y_j + b^T y_j + c &> 0 && \text{for all } j. && (12) \end{aligned}$$

These inequalities are a set of linear inequalities in the variables $P = P^T \in \mathbf{R}^n$, $b \in \mathbf{R}^n$, and $c \in \mathbf{R}$, and again can be solved using linear programming.

We can put further restrictions on the quadratic surface separating the two sets. For instance, for cluster analysis we might try to find an ellipsoid that contains all the points x_i and none of the y_j (see [53]). This constraint imposes the condition $P > 0$ in addition to the linear inequalities (11) and (12) on the variables P , b , and c . Thus finding an ellipsoid that contains all the x_i variables but none of the y_j variables (or determining that no such ellipsoid exists) can be done by solving a linear matrix inequality feasibility problem.

We can optimize the shape and the size of the ellipsoid by adding an objective function and other constraints. For instance, the ratio of the largest to the smallest semi-axis length is the square root of the condition number of P . In order to make the ellipsoid as spherical as possible, one can introduce an additional variable γ , add the constraint

$$I \leq P \leq \gamma I, \tag{13}$$

and minimize γ over (11), (12) and (13). This is a PDP in the variables γ , P , b and c . This PDP will be feasible if and only if there is an ellipsoid that contains all the x_i and none of the y_j ; its optimum value is one if and only there is a sphere that separates the two sets of points.

Geometrical Problems Involving Quadratic Forms

Many geometrical problems involving quadratic functions can be expressed as PDPs. We will give one simple example. Suppose we are given m ellipsoids $\mathcal{E}_1, \dots, \mathcal{E}_m$ described as the sublevel sets of the quadratic functions

$$f_i(x) = x^T A_i x + 2b_i^T x + c_i, \quad i = 1, \dots, m,$$

i.e., $\mathcal{E}_i = \{x | f_i(x) \leq 0\}$. The goal is to find the smallest sphere that contains all m of these ellipsoids (or equivalently, contains the convex hull of their union).

The condition that one ellipsoid contain another can be expressed in terms of a matrix inequality. Suppose that the ellipsoids $\mathcal{E} = \{x | f(x) \leq 0\}$ and $\tilde{\mathcal{E}} = \{x | \tilde{f}(x) \leq 0\}$, with

$$f(x) = x^T A x + 2b^T x + c, \quad \tilde{f}(x) = x^T \tilde{A} x + 2\tilde{b}^T x + \tilde{c},$$

have nonempty interior. Then it can be shown that \mathcal{E} contains $\tilde{\mathcal{E}}$ if and only if there is a $\tau \geq 0$ such that

$$\begin{bmatrix} A & b \\ b^T & c \end{bmatrix} \leq \tau \begin{bmatrix} \tilde{A} & \tilde{b} \\ \tilde{b}^T & \tilde{c} \end{bmatrix}.$$

(The 'if' part is trivial; the 'only if' part is less trivial. See [8, 59].)

Returning to our problem, consider the sphere \mathcal{S} represented by $f(x) = x^T x - 2x_c^T x + \gamma \leq 0$. \mathcal{S} contains the ellipsoids $\mathcal{E}_1, \dots, \mathcal{E}_m$ if and only if there are nonnegative τ_1, \dots, τ_m such that

$$\begin{bmatrix} I & -x_c \\ -x_c^T & \gamma \end{bmatrix} \leq \tau_i \begin{bmatrix} A_i & b_i \\ b_i^T & c_i \end{bmatrix}.$$

Note that these conditions can be considered one large linear matrix inequality in the variables x_c, γ , and τ_1, \dots, τ_m .

Our goal is to minimize the radius of the sphere S_t , which is $r = \sqrt{x_c^T x_c - \gamma}$. To do this we express the condition $r^2 \leq t$ as the matrix inequality

$$\begin{bmatrix} I & x_c \\ x_c^T & t + \gamma \end{bmatrix} \geq 0$$

and minimize the variable t .

Putting it all together we see that we can find the smallest sphere containing the ellipsoids $\mathcal{E}_1, \dots, \mathcal{E}_m$ by solving the PDP

$$\begin{aligned} & \text{minimize} && t \\ & \text{subject to} && \begin{bmatrix} I & -x_c \\ -x_c^T & \gamma \end{bmatrix} \leq \tau_i \begin{bmatrix} A_i & b_i \\ b_i^T & c_i \end{bmatrix}, \quad i = 1, \dots, m, \\ & && \tau_i \geq 0, \quad i = 1, \dots, m, \\ & && \begin{bmatrix} I & x_c \\ x_c^T & t + \gamma \end{bmatrix} \geq 0. \end{aligned}$$

The variables here are $x_c, \tau_1, \dots, \tau_m, \gamma$, and t .

This example demonstrates once again the breadth of problems that can be reformulated as PDPs. It also demonstrates that the task of this reformulation can be nontrivial.

Other Fields

- PDPs occur in statistics, in minimum trace factor analysis (see Watson (61)), as the *educational testing problem* (see [17, 18]), and in optimum experiment design (see Pukelsheim [50]).
- PDPs have been used to compute upper or lower bounds for combinatorial optimization problems. Examples are Lovász's famous upper bound on the Shannon capacity of a graph [33], Shor's bounds for integer programming, and Alizadeh's work [4, 2, 3].
- The problem of minimizing the maximum eigenvalue of a matrix has been studied extensively by Overton; see [45] for a list of applications.

1.3 Historical Overview

A very early paper on the theoretical properties of PDPs is Bellman and Fan [9]. Other references discussing optimality conditions are Craven and Mond [13], Shapiro [55], Fletcher [18], and Allwright [5].

Many researchers have worked on the problem of minimizing the maximum eigenvalue of a symmetric matrix. See, for instance, Cullum, Donath and Wolfe [12],

Goh and Teo [24], Panier [48], Allwright [6], Overton [44, 45], Overton and Womersley [47, 46], Ringertz [51], Fan and Nekoie [20], Fan [16], and Hiriart-Urruty and Ye [26].

The history of interior-point methods is relatively young. Interior-point methods for linear programming were introduced by Karmarkar in 1984 [28], although many of the underlying principles are older (see, e.g., Fiacco and McCormick [19], Liu and Huard [31], and Dikin [15]). Karmarkar's algorithm, and the interior-point methods developed afterwards, combine a very low, polynomial, worst-case complexity with excellent behavior in practice.

Karmarkar's paper has had an enormous impact, and several variants of his method have been developed (see, e.g., the survey by Gonzaga [23]). Interior-point methods have also been extended and generalized to convex quadratic programming, and to certain linear complementarity problems (see Kojima, Megiddo, Noma and Yoshise [29]).

An important breakthrough was achieved by Nesterov and Nemirovsky in 1988 [38, 40, 39, 41, 41]. They showed that the interior-point methods for linear programming can be generalized to all convex optimization problems. The key element is the knowledge of a barrier function with certain properties (*self-concordance*). Unfortunately, although Nesterov and Nemirovsky prove that a self-concordant barrier function exists for every convex set, it is not always known how to compute it in practice.

PDPs are an important class of convex optimization problems for which self-concordant barrier functions are known, and, therefore, interior-point methods are applicable. At the same time, they offer a simple conceptual framework and make possible a self-contained treatment of interior-point methods for many convex optimization problems.

Independently of Nesterov and Nemirovsky, Alizadeh [4] has generalized interior-point methods from linear programming to positive-definite programming. Other recent articles are Jarre [27], Vandenberghe and Boyd [60], Rendl, Vanderbei and Wolkowicz [54], Yoshise [65], and Alizadeh, Haerberly and Overton [1]. An excellent reference on interior-point methods for general convex problems is Den Hertog [14].

1.4 Outline

This paper gives a survey of interior-point methods for positive-definite programming. We start with a section on duality theory. In Section 3 we introduce the barrier function for PDPs, and the concepts of central points and central path. The notion of central path is heavily used in Section 4, which discusses primal-dual methods.

This survey is not meant to be exhaustive and emphasizes primal-dual methods. The most important omissions are the projective methods of Karmarkar, and of Nesterov and Nemirovsky [43]. Our motivation for the restriction to primal-dual methods is twofold. Primal-dual methods are commonly held to be more efficient in practice, and, secondly, their behavior is often easier to analyze. Moreover, all interior-point methods are based on similar principles, and we hope that the material discussed here is sufficient as a tutorial introduction to the entire field.

2 Duality

2.1 The Dual PDP

The dual problem associated with the PDP (1) is

$$\begin{aligned} & \text{maximize} && -\text{Tr}F_0Z \\ & \text{subject to} && \text{Tr}F_iZ = c_i, \quad i = 1, \dots, m, \\ & && Z \geq 0. \end{aligned} \quad (14)$$

Here the variable is the symmetric matrix Z , which is subject to m equality constraints and the nonnegativity condition. We write $\text{Tr}X$ for the trace of a symmetric matrix, i.e., $\text{Tr}X = X_{11} + \dots + X_{nn}$. Note that the objective function in (14) is a linear function of Z .

The dual problem (14) is also a PDP, i.e., it can be put in the same form as the primal problem (1). Let us assume for simplicity that the matrices F_1, \dots, F_m are linearly independent. Then we can express

$$\{ Z \mid Z = Z^T \in \mathbf{R}^{n \times n}, \text{Tr}F_iZ = c_i \}$$

$$\{ G(y) = G_0 + y_1G_1 + \dots + y_pG_p \mid y \in \mathbf{R}^p \}$$

in the form

where $p = \frac{n(n+1)}{2} - m$ and the G_i are appropriate matrices. We define $d \in \mathbf{R}^p$ by $d_i = \text{Tr}F_0G_i$, so that $d^T y = \text{Tr}F_0(G(y) - G_0)$. Then the dual problem becomes (except for a constant term in the objective and a change of sign to transform maximization into minimization)

$$\begin{aligned} & \text{minimize} && d^T y \\ & \text{subject to} && G(y) \geq 0, \end{aligned}$$

which is a PDP. It is possible to use notation that, unlike ours, emphasizes the complete symmetry between the primal and dual problems (see, e.g., Nesterov and Nemirovsky). Our notation was designed to make the primal problem as "explicit" as possible, with x denoting a "free" variable.

As an example of the dual PDP, let us apply the definition to the linear program (3), i.e., take $F_0 = -\text{diag}(b)$ and $F_i = \text{diag}(a_i)$. In this case, the diagonal structure makes it possible to simplify the dual problem. The objective function and the equality constraints only involve the diagonal elements of Z , and, obviously, replacing the off-diagonal of a positive-definite matrix by zeros does not alter its positive-definiteness either. Instead of optimizing over all symmetric $n \times n$ matrices Z , we can therefore limit ourselves to diagonal matrices $Z = \text{diag}(z)$. Problem (14) then reduces to

$$\begin{aligned} & \text{maximize} && b^T z \\ & \text{subject to} && z \geq 0, \\ & && a_i^T z = c_i, \quad i = 1, \dots, m, \end{aligned} \quad (15)$$

which is the familiar dual of the LP (3).

This example demonstrates an important point. In general, it is often the case that the dual problem can be simplified when the matrices F_i are structured. For

example, if the matrix $F(x)$ is block diagonal, the dual variables Z can be assumed to have the same block diagonal structure.

Linear programming duality is very strong owing to the polyhedral character of the feasible set: The optimum values of (3) and (15) are always equal, except in the pathological case where both problems are infeasible. (We adopt the standard convention that the optimum value of (3) is $+\infty$ if the problem is infeasible, and the optimum value of (15) is $-\infty$ if the dual problem is infeasible.) Duality results for general PDPs are weaker, as we will see below.

Let us return to our discussion of the dual PDP. The key property of the dual PDP is that it yields bounds on the optimal value of the primal PDP, and vice versa. Suppose that Z is dual feasible, and x is primal feasible. Then we have:

$$-\text{Tr}F_0Z \leq c^T x. \quad (16)$$

The inequality follows from the simple calculation

$$c^T x + \text{Tr}ZF_0 = \sum_{i=1}^m \text{Tr}ZF_i x_i + \text{Tr}ZF_0 = \text{Tr}ZF(x) \geq 0.$$

(We used the fact that $\text{Tr}AB \geq 0$ when $A = A^T \geq 0$ and $B = B^T \geq 0$.)

Since (16) holds for any feasible x , we conclude that $p^* \geq -\text{Tr}ZF_0$, where p^* is the optimal value of the PDP (1),

$$p^* \triangleq \inf \{ c^T x \mid F(x) \geq 0 \}. \quad (17)$$

In other words: *Dual feasible matrices yield lower bounds for the primal problem.* We can interpret x as a *suboptimal point* which gives the upper bound $p^* \leq c^T x$ and Z as a *certificate* that proves the lower bound $p^* \geq -\text{Tr}F_0Z$.

If x is primal feasible and Z is dual feasible, we refer to the quantity $\eta \triangleq c^T x + \text{Tr}F_0Z$ as the *duality gap* associated with x and Z . The duality gap is the difference between the upper and lower bound; it is the width of the interval in which we have localized p^* . If we define d^* to be the optimal value in the dual problem,

$$d^* \triangleq \sup \{ -\text{Tr}F_0Z \mid Z = Z^T \geq 0, \text{Tr}F_iZ = c_i, \quad i = 1, \dots, m \}, \quad (18)$$

then we can restate the result (16) as $p^* \geq d^*$, i.e., the optimal value of the dual problem is less than or equal to the optimal value of the primal problem. Note that when the primal (or dual) problem is infeasible, the right-hand (left-hand) side becomes ∞ ($-\infty$) so the inequality trivially holds. In fact equality usually obtains, as stated in the following theorem (see Nesterov and Nemirovsky [43], or Rockafellar [52]).

Theorem 1 We have $p^* = d^*$ if any of the following conditions holds.

1. The primal problem (1) is strictly feasible.
2. The dual problem (14) is strictly feasible.

3. The primal solution set

$$X_{opt} \triangleq \{x \mid F(x) \geq 0 \text{ and } c^T x = p^*\}$$

is nonempty and bounded.

4. The dual solution set

$$Z_{opt} \triangleq \{Z \mid Z \geq 0, \text{Tr} F_i Z = c_i, \text{ and } -\text{Tr} F_0 Z = d^*\}$$

is nonempty and bounded.

Example

Consider the PDP

$$\begin{aligned} &\text{minimize } x_1 \\ &\text{subject to } \begin{bmatrix} 0 & x_1 & 0 \\ x_1 & x_2 & 0 \\ 0 & 0 & x_1 + 1 \end{bmatrix} \geq 0. \end{aligned}$$

The feasible set is $\{(x_1, x_2) \mid x_1 = 0, x_2 \geq 0\}$, and therefore $p^* = 0$. The dual problem can be simplified as

$$\begin{aligned} &\text{maximize } -z_2 \\ &\text{subject to } \begin{bmatrix} z_1 & (1-z_2)/2 & 0 \\ (1-z_2)/2 & 0 & 0 \\ 0 & 0 & z_2 \end{bmatrix} \geq 0, \end{aligned}$$

and the feasible set is $\{(z_1, z_2) \mid z_1 \geq 0, z_2 = 1\}$. The dual problem therefore has optimal value $d^* = -1$. This PDP violates all four conditions mentioned in the theorem. Both problems are feasible, but not strictly feasible, and the optimal sets X_{opt} and Z_{opt} are both unbounded. Note also the contrast with linear programming, where it is impossible to have a finite nonzero duality gap at the optimum.

Example

We take the matrix norm minimization problem mentioned in Section 1.2:

$$\begin{aligned} &\text{minimize } \|A(x)\| \\ &x \in \mathbb{R}^m \end{aligned} \tag{19}$$

where $A(x) = A_0 + x_1 A_1 + \dots + x_m A_m$, and we remind the reader that $\|A(x)\|$ is the maximum singular value of $A(x)$.

The problem (19) is a basic problem in the theory of Banach spaces; its optimum value is the norm of (the image of) A_0 in the quotient space of all $p \times q$ matrices modulo the span of A_1, \dots, A_m . In this theory we encounter the following dual of (19):

$$\begin{aligned} &\text{maximize } \text{Tr} A_0^T Q \\ &\text{subject to } \text{Tr} A_i^T Q = 0, \quad i = 1, \dots, m, \\ &\|Q\|_* \leq 1, \end{aligned} \tag{20}$$

where $\|Q\|_* = \sum \sigma_i(Q)$ is the nuclear norm of Q , which is the norm dual to the maximum singular value. It is also known that the optimal values of (19) and (20) are always equal.

Let us verify that this (Banach space) notion of duality coincides with PDP duality. The dual PDP of problem (7) is

$$\begin{aligned} &\text{maximize } -2\text{Tr} A_0^T Z_{12} \\ &\text{subject to } \text{Tr} A_i^T Z_{12} = 0, \quad i = 1, \dots, m, \\ &\text{Tr} Z_{11} + \text{Tr} Z_{22} = 1, \\ &\begin{bmatrix} Z_{11} & Z_{12} \\ Z_{12}^T & Z_{22} \end{bmatrix} \geq 0. \end{aligned} \tag{21}$$

This can be simplified. The positive-definite constraint can be rewritten as

$$\begin{bmatrix} Z_{11} & 0 \\ 0 & Z_{22} \end{bmatrix} \geq \begin{bmatrix} 0 & -Z_{12} \\ -Z_{12}^T & 0 \end{bmatrix}. \tag{22}$$

The eigenvalues of the matrix on the right are the singular values of Z_{12} , each singular value appearing twice. It is well known that if A and B are two symmetric matrices, then $A \geq B$ implies $\lambda_k(A) \geq \lambda_k(B)$, assuming the eigenvalues of A and B are taken in the same order. As a consequence, inequality (22) implies

$$2 \sum \sigma_i(Z_{12}) \leq \text{Tr} Z_{11} + \text{Tr} Z_{22} \leq 1.$$

Since the matrices Z_{11} and Z_{22} do not appear in any other constraint, nor in the objective, we see that problem (21) reduces to

$$\begin{aligned} &\text{maximize } -2\text{Tr} A_0^T Z_{12} \\ &\text{subject to } \text{Tr} A_i^T Z_{12} = 0, \quad i = 1, \dots, m, \\ &2 \sum \sigma_i(Z_{12}) \leq 1, \end{aligned}$$

which is the same as (20) with $Q = 2Z_{12}$.

Problem (19) is always strictly feasible; it suffices to choose $x = 0$ and $t > \|A_0\|$. Applying Theorem 1, we conclude that the optimal duality gap is always zero.

2.2 The Primal-Dual Formulation

Theorem 1 has important consequences for PDP algorithms. It gives conditions under which the primal-dual optimization problem

$$\begin{aligned} &\text{minimize } c^T x + \text{Tr} F_0 Z \\ &\text{subject to } F(x) \geq 0, \\ &Z \geq 0, \\ &\text{Tr} F_i Z = c_i, \quad i = 1, \dots, m, \end{aligned} \tag{23}$$

has optimum value zero. Here we minimize the duality gap $c^T x + \text{Tr} F_0 Z$ over all primal and dual feasible points. The duality gap is a linear function of x and Z , and therefore problem (23) is a PDP in x and Z .

Primal-dual methods for PDPs solve (23), assuming the primal and dual problems are strictly feasible. They generate a sequence of feasible points $x^{(k)}$ and $Z^{(k)}$, and in each step use the dual information in $Z^{(k)}$ to find good updates for $x^{(k)}$ and vice-versa.

This means that at every stage of the algorithm, we have available suboptimal primal and dual solutions x, Z . The primal solution x proves an upper bound $c^T x > p^*$ on the optimal value; the dual solution proves a lower bound $p^* > -T^T r F_0 Z$. The iteration continues until the duality gap is less than a given tolerance ϵ .

3 The Barrier Function

In this section, we introduce a barrier function for linear matrix inequality constraints and discuss its properties. This leads us to the fundamental concept of *centrality*, and the definition of *central points* and *central path*. From now on we will assume that the matrices F_i are independent.

3.1 Definition

The function

$$\phi(x) \triangleq \begin{cases} \log \det F(x)^{-1} & \text{if } F(x) > 0 \\ +\infty & \text{otherwise} \end{cases} \tag{24}$$

is a *barrier function* for $X \triangleq \{x \mid F(x) \geq 0\}$, i.e., $\phi(x)$ is finite if and only if $F(x) > 0$, and becomes infinite as x approaches the boundary of X . There are many other barrier functions for X (for example, trace can be substituted for determinant in (24)), but this one enjoys many special properties (see [43]). In particular, when $F(x) > 0$, it is analytic and strictly convex.

In the case of a set of linear inequalities $Ax \geq b$, where $A = [a_1 \dots a_n]^T$, we have $F(x) = \text{diag}(Ax - b)$, and the definition reduces to the familiar logarithmic barrier function

$$\phi(x) = \begin{cases} -\sum_{i=1}^n \log(a_i^T x - b_i) & \text{if } Ax \geq b, \\ +\infty & \text{otherwise.} \end{cases}$$

We first give formulas for the gradient $g(x)$ and Hessian $H(x)$ of ϕ . Recall that $\text{Tr} AB$ is the standard inner product of two symmetric matrices A and B ; the corresponding norm is the Frobenius norm, $\|A\|_F = (\text{Tr} A^2)^{1/2}$.

The gradient and the Hessian of ϕ can be readily derived from the following second order approximation of the function $-\log \det X$. If $X > 0$ is $n \times n$ and symmetric, then

$$\log \det (X + Y)^{-1} = \log \det X^{-1} - \text{Tr} X^{-1} Y + \frac{1}{2} \text{Tr} X^{-1} Y X^{-1} Y + o(\|Y\|^2). \tag{25}$$

From equation (25), one can immediately derive a second order approximation for $\phi(x)$:

$$\phi(x + v) \approx \phi(x) - \text{Tr} F(x)^{-1} \left(\sum_{i=1}^m v_i F_i \right) +$$

$$\begin{aligned} & \frac{1}{2} \text{Tr} F(x)^{-1} \left(\sum_{i=1}^m v_i F_i \right) F(x)^{-1} \left(\sum_{j=1}^m v_j F_j \right) \\ & = \phi(x) - \sum_{i=1}^m v_i \text{Tr} F(x)^{-1} F_i + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m v_i v_j \text{Tr} F(x)^{-1} F_i F(x)^{-1} F_j. \end{aligned}$$

We conclude that the gradient $g(x)$ and the Hessian $H(x)$ of $\phi(x)$ are given by

$$g_i(x) = -\text{Tr} F(x)^{-1} F_i = -\text{Tr} F(x)^{-1/2} F_i F(x)^{-1/2}, \tag{26}$$

and

$$\begin{aligned} H_{ij}(x) & = \text{Tr} F(x)^{-1} F_i F(x)^{-1} F_j \\ & = \text{Tr} (F(x)^{-1/2} F_i F(x)^{-1/2}) (F(x)^{-1/2} F_j F(x)^{-1/2}), \end{aligned} \tag{27}$$

for $i, j = 1, \dots, m$.

From expression (27) we can verify that ϕ is strictly convex for strictly feasible x . For $x, y \in \mathbf{R}^m$ with $F(x) > 0$,

$$\begin{aligned} y^T H(x) y & = \sum_{i,j=1}^m y_i y_j \text{Tr} (F(x)^{-1/2} F_i F(x)^{-1/2}) (F(x)^{-1/2} F_j F(x)^{-1/2}) \\ & = \text{Tr} \left(F(x)^{-1/2} \left(\sum_{i=1}^m y_i F_i \right) F(x)^{-1/2} \right)^2 \\ & = \left\| F(x)^{-1/2} \left(\sum_{i=1}^m y_i F_i \right) F(x)^{-1/2} \right\|_F^2 \geq 0. \end{aligned} \tag{28}$$

We see that $y^T H(x) y = 0$ if and only if $\sum_{i=1}^m y_i F_i = 0$. By independence of F_1, \dots, F_m , we conclude that $H(x) > 0$, i.e., ϕ is strictly convex.

Finally, we note that the barrier function ϕ is bounded below if and only if the feasible set X is bounded.

3.2 Analytic Center

3.2.1 Definition

We suppose now that the linear matrix inequality $F(x) \geq 0$ is strictly feasible and that its feasible set is bounded. Since ϕ is strictly convex, it has a unique minimizer, which we denote

$$x^* \triangleq \underset{x}{\text{argmin}} \phi(x). \tag{29}$$

We will refer to x^* as the *analytic center* of the linear matrix inequality $F(x) > 0$. It is important to note that the analytic center is a property of a linear matrix inequality and not of its solution set X . The same set X can be represented by different matrix inequalities, which have different analytic centers.

From (26) we see that x^* is characterized by

$$\text{Tr} F(x^*)^{-1} F_i = 0, \quad i = 1, \dots, m. \quad (30)$$

Thus, $F(x^*)^{-1}$ is orthogonal to the span of F_1, \dots, F_m .

In the case of a set of linear inequalities, the definition coincides with Sonnevend's definition [57, 58], i.e.,

$$x^* = \underset{Ax \geq b}{\text{argmax}} \prod_{i=1}^n (a_i^T x - b_i).$$

3.2.2 Computing the Analytic Center

Newton's method, with appropriate step length selection, can be used to efficiently compute the analytic center. Starting with a strictly feasible point $x^{(0)}$, the algorithm follows the iteration:

$$x^{(k+1)} := x^{(k)} - \alpha^{(k)} H(x^{(k)})^{-1} g(x^{(k)}), \quad (31)$$

where $\alpha^{(k)}$ is the damping factor at the k th iteration, and $g(x)$ and $H(x)$ are the gradient and Hessian of the barrier function in x .

Nesterov and Nemirovsky [43] give a simple step length rule appropriate for the general class of self-concordant barrier functions mentioned earlier. The damping factor depends on a quantity called the *Newton decrement* of ϕ at x :

$$\delta(x) \triangleq \left\| H(x)^{-1/2} g(x) \right\|.$$

The name comes from the observation that $\delta(x)^2/2$ is the difference between $\phi(x)$ and the minimum value of the quadratic approximation of ϕ at x . Alternatively, $\delta(x)$ is the length of the Newton step $-H(x)^{-1}g(x)$ measured in the norm induced by the Hessian $H(x)$.

The damping factor is:

$$\alpha^{(k)} := \begin{cases} 1 & \text{if } \delta(x^{(k)}) \leq 1/4, \\ 1/(1 + \delta(x^{(k)})) & \text{if } \delta(x^{(k)}) > 1/4. \end{cases} \quad (32)$$

Nesterov and Nemirovsky show that this step length always results in $F(x^{(k+1)}) > 0$. Moreover, for $\delta(x^{(k)}) < 1/4$, we have $\delta(x^{(k+1)}) \leq 2\delta(x^{(k)})^2$, i.e., the algorithm converges quadratically. They also give a complete convergence analysis. The main results can be summarized as follows.

- Until the region of quadratic convergence ($\delta(x) < 1/4$) is reached, the objective $\log \det F(x)^{-1}$ decreases at least by the absolute constant 0.3068 at each Newton step. (By absolute constant we mean it does not depend on n , the problem data, or the required accuracy of computing x^* .)

- Once the region of quadratic convergence is reached, at most a constant number c of Newton steps is required to compute x^* to a given accuracy. (The constant c does not depend on n or the problem data, but only on the required accuracy ϵ . Since the convergence is quadratic in this region, c grows as $\log \log 1/\epsilon$ if ϵ decreases.)

In other words, the number of Newton steps required to compute x^* given x can be bounded in terms of $-\log \det F(x) + \log \det F(x^*)$:

$$\# \text{Newton steps} \leq c + 3.26(-\log \det F(x) + \log \det F(x^*)), \quad (33)$$

where c depends only on the required accuracy of computing x^* and grows extremely slowly.

Therefore, the quantity $\psi(x) \triangleq -\log \det F(x) + \log \det F(x^*)$ has a very natural interpretation as the 'deviation from centrality' of a point x . In general, however, $\psi(x)$ can only be evaluated by computing the center x^* .

3.3 The Primal-Dual Central Path

We now return to the primal-dual formulation of Section 2.2.

For $\alpha > 0$ consider the set of strictly feasible pairs x, Z with duality gap α , i.e., $c^T x + \text{Tr} F_0 Z = \alpha$. The *analytic center* of this set is the minimizer of the barrier term $-\log \det F(x) - \log \det Z$. We denote the analytic center as $x^*(\alpha), Z^*(\alpha)$:

$$\begin{aligned} (x^*(\alpha), Z^*(\alpha)) = & \underset{\text{subject to}}{\text{argmin}} && -\log \det F(x) - \log \det Z \\ & && F(x) \geq 0, \\ & && Z \geq 0, \\ & && \text{Tr} F_i Z = c_i, \quad i = 1, \dots, m, \\ & && c^T x + \text{Tr} F_0 Z = \alpha. \end{aligned} \quad (34)$$

Thus, among all feasible pairs x, Z with the duality gap α , the pair x^*, Z^* maximizes $\det(F(x)Z)$. The pair (x^*, Z^*) converges to a primal and dual optimal pair as $\alpha \rightarrow 0$, and the curve given by (x^*, Z^*) for $\alpha > 0$ is called the *central path* for the problem (23). The central pair (x^*, Z^*) has many important properties. For our purposes here we need:

Theorem 2 $F(x^*(\alpha))Z^*(\alpha) = (\alpha/n)I$. Conversely, if x and Z are a feasible pair and $F(x)Z = (\alpha/n)I$ then $x = x^*(\alpha)$ and $Z = Z^*(\alpha)$.

In other words, centrality is characterized by $F(x)$ and Z being inverses of each other, up to a constant.

Now consider a feasible pair (x, Z) , and define $\alpha = \text{Tr} F(x)Z$. Then $(x^*(\alpha), Z^*(\alpha))$ is the central pair with the same duality gap as x, Z . Therefore

$$\log \det F(x)Z \geq \log \det F(x^*(\alpha))Z^*(\alpha) = n \log n - n \log \text{Tr} F(x)Z$$

with equality holding only when x, Z are central.

