

UTRECHT UNIVERSITY

MASTER THESIS
MATHEMATICAL SCIENCES

Methods for eigenvalue optimization with
application to semidefinite optimization

Author:
Emile BROEDERS

Supervisor:
Gerard SLEIJPEN

September, 2015

Abstract

This thesis deals with eigenvalue optimization problems and some methods which are applicable in solving those kind of problems. The goal is to develop a novel algorithm which can solve large instances of eigenvalue optimization problems.

In the second chapter is shown how semidefinite optimization problems can be turned into eigenvalue optimization problems through a rotation of basis. Although the relation between semidefinite optimization and eigenvalue optimization is not new, this approach where a semidefinite optimization problem is turned into a full eigenvalue optimization problem is.

In the third chapter SASPA (Subspace Accelerated Sensitive Pole Algorithm) is used to compute the intersection points of eigencurves. This is a novel application which allows to compute such intersection points with less computation time than through bisection methods. Moreover, SASPA also works for non-simple eigenvalues.

The fourth chapter deals with developing a novel algorithm to solve eigenvalue optimization problems based on gradients. This algorithm is able to handle large sparse problems within reasonable computation time.

Acknowledgement

I am thankful to my supervisor, Gerard Sleijpen, for encouraging and guiding me in the field of numerical mathematics. Due to him I was able to complete this thesis and conduct all research and development of methods included.

Lastly, I offer my regards to all those supported me in any way during the completion of my thesis.

Contents

1	Introduction	9
1.1	Semidefinite optimization	10
1.2	Properties of Frobenius symmetric matrices space	13
1.2.1	Matrix decomposition's	13
1.2.2	Trace product	15
1.2.3	Basic operations	16
1.3	Duality and optimality conditions	17
1.4	Some examples	20
1.4.1	Max-cut problem	20
1.4.2	Lovász theta number	23
1.4.3	Ideal GMRES and Arnoldi polynomes	25
1.4.4	Phase retrieval	27
2	Eigenvalue optimization	31
2.1	Equivalence of semidefinite problems	31
2.1.1	Perturbing semidefinite optimization problems	34
2.1.2	Some examples	38
2.2	Analysis on eigenvalues	40
2.2.1	Continuity and differentiability for the simple case	41
2.2.2	The non-simple case	45
3	Sensitive Pole Algorithm	53
3.1	Rayleigh Quotient Iteration	54
3.2	Dominant Pole Algorithm	54
3.3	Subspace accelerating with deflation	57
3.3.1	Selecting poles	59
3.3.2	Deflating	60
3.3.3	Restart search space	61
3.4	Sensitive Pole Algorithm	61
3.5	Computing cross points of eigencurves	63
3.5.1	Numerical results	66

4 Algorithms for eigenvalue optimization	71
4.1 Calculating a descent direction	72
4.1.1 Descent direction in the simple case	72
4.1.2 Descent direction in the non-simple case	74
4.1.3 Cut infeasible descent directions	76
4.2 Determine max step size	78
4.2.1 Analysis of the problem	78
4.2.2 Line search methods based on intersections	80
4.3 Combining everything	84
4.3.1 Comparison to other methods	85
5 Numerical results	89
5.1 Max-cut	89
5.2 Lovász theta number	91
5.3 Ideal GMRES and Arnoldi polynomes	94
6 Conclusion	99
A Implementation of some algorithms	101
A.1 Rotating SDP to Spectral Problems	101
A.2 SASPA to compute intersections	102
A.3 Spectral Descent	106
Bibliography	115

Chapter 1

Introduction

Applied mathematics deals with all kind of optimization problems originating from different fields. For these kind of problems the objective is to minimize or maximize a function. For instance, the expected profit of a shop or the harvest of a farm. By developing algorithms that can solve these problems we allow companies, government institutions and other organizations to optimize their utilization of scarce resources.

Convex optimization problems are problems with a specific structure. They are such that the convex sum of two possible solutions is also a solution to the problem. As an example: Say we want to maximize the profit of the harvest on one acre of land. One possible solution is to fill the land with wheat and another solution is to fill the land with corn. A convex sum of those solutions is for instance to plant wheat on one half of the field and corn on the other, which is another feasible solution to the problem.

The application of large convex optimization problems started in the second half of the twentieth century with the invention of modern computer architectures and development of the simplex algorithm by Dantzig in 1947. The simplex algorithm was developed by Dantzig as part of the United States Air Force team SCOOP (Scientific Computation of Optimum Programs) for optimization of various plans and schedules. The simplex algorithm is an algorithm to solve linear problems (to save words, when linear problems is short for linear convex optimization problems), which are a form of convex optimization problems, in an efficient manner.

Linear problems are such that if the objective value for a solution is x and we can double the solution, then the objective value for the double solution is $2x$. To illustrate this property: Say we expect to have a profit of 1000 euro on the harvest of one acre of corn field, then we expect a profit of 2000 euro on two acres. The simplex algorithm is an algorithm that returns the best solution to such a linear problem, where best means the solution which minimizes or maximizes the objective value of the problem.

In the years following Dantzig's publication of the simplex method, interest in convex optimization problems spurred. More efficient algorithms were

developed to solve linear problems and the framework of linear problems was generalized to cone optimization problems. The class of cone optimization problems consists of linear optimization, quadratic optimization, semidefinite optimization and many more exotic variants. The applications range from classic optimization, polynome fitting, estimation of statistical parameters to approximate NP-hard problems in polynomial time.

I will focus in this thesis on semidefinite optimization. These are optimization problems which can be seen as the unification of linear and quadratic problems. (Quadratic problems are optimization problems in which the objective is quadratic function.) Later will be shown how semidefinite problems are equivalent to eigenvalue optimization of a matrix. This will be used to develop a Newton-like method to solve large instances of semidefinite problems iterative.

In the next sections of this chapter all terms coined so far will be formalized. In section 1.1 a more mathematical description will be given of convex optimization, linear, quadratic and semidefinite problems. Then in section 1.2 will the most important properties of semidefinite optimization be given. In section 1.3 will optimality conditions (also known as Karush-Kuhn-Tucker conditions) and duality be explained.

1.1 Semidefinite optimization

Convex optimization problems are such that the objective is to minimize or maximize a convex function subjected to constraints. In the most general form, all these kind of problems can be written as:

Definition 1.1.1. *For a given convex space \mathcal{F} and convex functions $f, g_1, \dots, g_m : \mathcal{F} \rightarrow \mathbb{R}$, the following equation is known as a convex optimization problem (CP) with only equality constraints:*

$$\sup_{x \in \mathcal{F}} f(x) \text{ s. t. } g_i(x) = 0 \quad (\forall i = 1 : m) \quad (1.1)$$

The function f is referred to as the objective and the functions g_1, \dots, g_m are denoted as the constraints.

Moreover, there can also be additional inequality constraints such that:

$$\begin{aligned} & \sup_{x \in \mathcal{F}} f(x) \\ & \text{s. t. } g_i(x) = 0 \quad (\forall i = 1 : m) \\ & \quad h_i(x) \geq 0 \quad (\forall i = 1 : n) \end{aligned} \quad (1.2)$$

Similar, h_1, \dots, h_k denote convex functions.

Definition 1.1.2. *The space $\{y \in \mathcal{F} \text{ s. t. } g_i(y) = 0 (\forall i)\}$ is called the feasible set, which is again convex and closed. The interior elements of the feasible set are said to be strict feasible elements.*

There are many variations possible on CP 1.1.1 which all vary in properties on f , g_i and \mathcal{F} . For example, probably the best known set of convex optimization problems are known as linear optimization problems (also known as linear programs). That is f and g_i linear and $\mathcal{F} = \mathbb{R}^n$ for some n :

Definition 1.1.3. A linear program (LP) is an equation of the following form:

$$\max_{x \in \mathbb{R}^n} c^T x \text{ s. t. } Ax = b, x \geq 0 \quad (1.3)$$

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

In the definition 1.1.3, $x \geq 0$ should be read as each component of x is larger than or equals 0. Further in this thesis, each time a \leq or \geq is encountered, it should be read as a component wise (in-) equality.

The feasible set of a LP is the intersection of the positive quadrant of \mathbb{R}^n and the affine subspace $Ax = b$.

Sometimes the restriction that $x \geq 0$ is lifted or modified to e.g. $-1 \leq x \leq 1$. Obvious the feasible set is then such that it incorporates the modified constraints.

It might be observed that in the definition of LP the supremum is changed into a maximum. That is because for linear programs the maximum is always attained. This can be quickly recognized: The feasible set is closed and embedded in \mathbb{R}^n , which is a complete space. Hence all limit solutions are feasible thus the supremum can be changed without loose in a maximum.

Linear problems can efficiently be solved by the simplex or interior point algorithms and through their general nature many real life optimization problems can be captured into a linear program. For instance in logistics and scheduling problems. Linear programs will be encountered later in this thesis and form part of the backbone of the final algorithm in chapter 4.

As two more examples of convex optimization problems:

Definition 1.1.4. A quadratic program (QP) is of the form:

$$\max_{x \in \mathbb{R}^n} x^T Q x \text{ s. t. } Ax = b \quad (1.4)$$

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

Definition 1.1.5. A quadratically constrained quadratic program (QCQP) is of the form:

$$\sup_{x \in \mathbb{R}^n} x^T Q x \text{ s. t. } x^T P_i x = b_i \quad (\forall i = 1, \dots, m) \quad (1.5)$$

With $Q \in \mathbb{R}^{n \times n}$, and $P_i \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^m$ for all $i = 1, \dots, m$.

Problems like 1.1.4 and 1.1.5 arise naturally when dealing with eigenvalue problems. For instance: Finding the largest eigenvalue of a matrix, say M , is equivalent to solving:

$$\sup_{x \in \mathbb{R}^n} x^T M x \text{ s. t. } x^T x = 1 \quad (1.6)$$

This QCQP is better known as the Rayleigh quotient, which will be formalized in corollary 2.2.3.

As another example of CP, and the one which is the central theme in this thesis, is:

Definition 1.1.6. *A semidefinite optimization problem (SDP) is of the form:*

$$\sup_{X \in \mathcal{S}^n} \langle C, X \rangle \text{ s. t. } \langle A_i, X \rangle = b_i \ (\forall i = 1, \dots, m), \ X \succeq 0 \quad (1.7)$$

In which \mathcal{S}^n is the space of symmetric real matrices with trace inproduct, $C, A_1, \dots, A_m \in \mathcal{S}^n$ and $b_i \in \mathbb{R}$ for all $i = 1, \dots, m$. Furthermore $X \succeq 0$ implies that X is constrained to be positive semidefinite.

The trace inproduct is defined as:

$$\langle A, B \rangle := \text{Tr}(A^T B) = \text{Tr}(A \cdot B) = \text{Tr}(B \cdot A) \quad (1.8)$$

The space \mathcal{S}^n with trace inproduct will be more elaborated introduced in section 1.2. For now it suffices to have an intuitive understanding to see the importance of SDP.

As shorthand, the SDP is often written with linear operator $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$:

$$\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m, X \mapsto \begin{pmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{pmatrix} \quad (1.9)$$

(As an interesting side note, it follows from Riesz representation theorem that any linear operator $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$ is of the form 1.9. Hence SDPs can be written as:

$$\sup_{X \in \mathcal{S}^n} \langle C, X \rangle \text{ s. t. } \mathcal{A}(X) = b, \ X \succeq 0 \quad (1.10)$$

With $C \in \mathcal{S}^n$, $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$ and $b \in \mathbb{R}^m$.)

To demonstrate a feature that makes SDPs powerful: Both LPs and QCQP can easily be cast into SDP. First for a given linear problem: Let $\tilde{C} = \text{diag}(c)$, that is $\tilde{C} \in \mathcal{S}^n$ with on the diagonal vector c and everywhere else 0. Let $\tilde{A}_i = \text{diag}(a_i)$, symmetric matrix with on the diagonal row i of matrix A . Then the SDP version for a given LP is:

$$\begin{aligned} & \max_{X \in \mathcal{S}^n} \langle \tilde{C}, X \rangle \\ & \text{s. t. } \langle \tilde{A}_i, X \rangle = b_i \ (\forall i = 1, \dots, m) \\ & \quad X \succeq 0 \end{aligned} \quad (1.11)$$

Stated without proof: The optimum $X^* \in \mathcal{S}^n$ for the SDP is the optimum $x^* \in \mathbb{R}^n$ of the LP on the diagonal: $X^* = \text{diag}(x^*)$.

Second, also QCQP can easily be converted to SDP. Without variable substitutions:

$$\begin{aligned} & \sup_{X \in \mathcal{S}^n} \langle Q, X \rangle \\ & \text{s. t. } \langle P_i, X \rangle = b_i \quad (\forall i = 1, \dots, m) \\ & \quad X \succeq 0 \end{aligned} \tag{1.12}$$

Again, without proof: The optimum $X^* \in \mathcal{S}^n$ for the SDP is the outer product of the optimum $x^* \in \mathbb{R}^n$ of the QCQP: $X^* = x^* x^{*\text{T}}$.

As an example, consider quadratic problem 1.6. The SDP formulation is:

$$\sup_{X \in \mathcal{S}^n} \langle M, X \rangle \text{ s. t. } \langle I_n, X \rangle = 1, X \succeq 0 \tag{1.13}$$

With I_n the n by n identity matrix.

Due that the semidefinite formulation can be used easily to combine linear and quadratic aspects among other properties which will be explained in the next section, SDP's are applicable for formulating many problems. Problems where SDP are used range from graph optimization, combinatorial optimization, polynome fitting and many more. Some selected examples are discussed in section 1.4.

1.2 Properties of Frobenius symmetric matrices space

Throughout this thesis many properties of semidefinite matrices are used. It will be shown how many of these properties are used to apply semidefinite optimization in a wide variety of applications in a natural manner.

First will a few matrix decomposition's be stated. After that will a few properties of the trace inproduct be shown. Finally a few basic operators on semidefinite matrices will be shown.

If you feel confident in your knowledge of linear algebra, then you can opt for skipping this section. In this section will mostly elementary results be discussed for easing the explanation of the more advanced results in later sections.

1.2.1 Matrix decomposition's

For a start, some equivalent definitions for positive semidefinite matrices:

Definition 1.2.1. *A matrix $X \in \mathcal{S}^n$ is positive semidefinite if either (hence both) of the following properties hold:*

1. *For all $u \in \mathbb{R}^n$ this relation holds: $u^{\text{T}} X u \geq 0$.*
2. *The smallest eigenvalue of X is non-negative.*

If $X \in \mathcal{S}^n$ is positive semidefinite, then it is said that the inequality $X \succeq 0$ holds and $X \in \mathcal{S}_{\succeq 0}^n$.

The spectral decomposition theorem is very important. It states that any real symmetric matrix can be decomposed in terms of its eigenvalues and eigenvectors:

Theorem 1.2.2 (Spectral decomposition theorem). *Any real symmetric matrix $X \in \mathcal{S}^n$ can be uniquely decomposed as:*

$$X = \sum_{i=1}^n \lambda_i u_i u_i^T \quad (1.14)$$

Where $\lambda_1, \dots, \lambda_n \in \mathbb{R}$ are the eigenvalues of X counted according to multiplicity and $u_1, \dots, u_n \in \mathbb{R}^n$ the corresponding normalized eigenvectors.

This is equivalent to:

$$X = U D U^T \quad (1.15)$$

Where U is orthonormal with the eigenvectors on the columns ($U = (u_1 | u_2 | \dots | u_n)$) and D diagonal with on the diagonal the corresponding eigenvalues.

Corollary 1.2.3 (Gram decomposition). *Any matrix $X \in \mathcal{S}_{\succeq 0}^n$ can be decomposed as the sum of at most n orthogonal vectors:*

$$X = \sum_{i=1}^n x_i x_i^T \quad (1.16)$$

Where $x_i = \sqrt{\lambda_i} \cdot u_i$, the product of the square root of the i th eigenvalue and i th eigenvector of X .

Also the Cholesky for any real symmetric positive semidefinite matrix is important. With abuse of language, it gives a sense of the square root of a real symmetric positive semidefinite matrix.

Theorem 1.2.4 (Cholesky decomposition). *Any real symmetric positive semidefinite matrix $X \in \mathcal{S}_{\succeq 0}^n$ can be decomposed as the product of real lower triangle matrix L and L^T :*

$$X = L L^T \quad (1.17)$$

If $X \succ 0$ (that is positive definite), then, and only then, is that decomposition unique.

Closely related is the *LDL*-decomposition:

$$X = L D L^T \quad (1.18)$$

With L again lower triangular but with all ones on the main diagonal and D a diagonal matrix.

The reason the *LDL*-decomposition is sometimes preferred over the *LL*-decomposition is that for the *LDL*-decomposition no square roots have to be calculated. Hence its computation is slightly more efficient.

1.2.2 Trace product

Definition 1.2.5 (Trace of a matrix). *The trace is defined as the sum of the elements on the main diagonal of a matrix:*

$$\text{Tr}(X) := X_{1,1} + X_{2,2} + \cdots + X_{n,n} = \sum_{i=1}^n X_{i,i} \quad (1.19)$$

Where $(X_{i,j})_{i,j} \in \mathcal{S}^n$.

From the definition given above, it should be clear that the trace operator on \mathcal{S}^n is linear:

$$\text{Tr}(aA) = a \text{Tr}(A) \quad (a \in \mathbb{R}, A \in \mathcal{S}^n) \quad (1.20)$$

And:

$$\text{Tr}(A + B) = \text{Tr}(A) + \text{Tr}(B) \quad (A, B \in \mathcal{S}^n) \quad (1.21)$$

Another important property is that the trace is invariant under cyclic permutations:

$$\text{Tr}(ABC) = \text{Tr}(CAB) = \text{Tr}(BCA) \quad (A, B, C \in \mathcal{S}^n) \quad (1.22)$$

It can be easily deduced from the cyclic property that the trace of a symmetric real matrix equals the sum of its eigenvalues:

Lemma 1.2.6. *The trace of matrix $X \in \mathcal{S}^n$ equals the sum of its eigenvalues counted according to its multiplicity:*

$$\text{Tr}(X) = \sum_{i=1}^n \lambda_i \quad (1.23)$$

Proof. Using that the eigendecomposition of X is UDU^T , with U orthogonal and on the diagonal of D the eigenvalues of X :

$$\text{Tr}(X) = \text{Tr}(UDU^T) = \text{Tr}(U^TUD) = \text{Tr}(D) = \sum_{i=1}^n \lambda_i \quad (1.24)$$

□

Definition 1.2.7. *The trace inproduct on \mathcal{S}^n is defined as:*

$$\langle A, B \rangle = \text{Tr}(AB) \quad (1.25)$$

It follows from the linearity and cyclic properties of the trace operator on \mathcal{S}^n , that the trace inproduct is bilinear and symmetric.

An important property of the trace inproduct on $\mathcal{S}_{\geq 0}^n$ is that it is always non-negative:

Lemma 1.2.8. For all $A, B \in \mathcal{S}_{\geq 0}^n$ (symmetric positive semidefinite matrix) holds the relation:

$$\langle A, B \rangle \geq 0 \quad (1.26)$$

Proof. Let $LL^T = A$ and $MM^T = B$ both be Cholesky decomposition's. Then:

$$\langle LL^T, MM^T \rangle = \text{Tr}(LL^T MM^T) = \text{Tr}(M^T LL^T M) = \text{Tr}(C) \quad (1.27)$$

In which $C = (M^T L)(M^T L)^T$. Clearly $C \succeq 0$ since its Cholesky decomposition exists. \square

1.2.3 Basic operations

The semidefinite property is invariant under permutations:

Lemma 1.2.9. Let $X \in \mathcal{S}$ and $P \in \mathbb{R}^{n \times n}$ be non-singular. Then:

$$X \succeq 0 \Leftrightarrow P^T X P \succeq 0 \quad (1.28)$$

Proof. Assume $X \succeq 0$. Let $X = UDU^T$ be the eigendecomposition of X . Then $P^{-T}U$ forms a basis for the eigenvectors of $P^T X P$, hence the eigenvalues corresponding to the i column of $P^T X P$ is $D_{i,i}$ plus the square root of the norm of that column. Thus definitely non-negative and both sides are positive semidefinite.

For the other way around: observe that $X = P^{-T}(P^T X P)P^{-1}$. \square

The following definition is widely used:

Definition 1.2.10 (Schur complement). Consider matrix $X \in \mathcal{S}^n$ in block form:

$$X = \begin{pmatrix} A & B \\ B^T & C \end{pmatrix} \quad (1.29)$$

With $A \in \mathcal{S}^n$ non-singular. Then the matrix $C - B^T A^{-1} B$ is called the Schur complement of A in X .

Theorem 1.2.11. Let $X \in \mathcal{S}^n$ be in block form as in equation 1.2.10. X is non-singular if and only if A and $C - B^T A^{-1} B$ are both non-singular.

Proof. Define:

$$P = \begin{pmatrix} I & A^{-1}B \\ 0 & I \end{pmatrix} \quad (1.30)$$

Then $X = P^T \begin{pmatrix} A & 0 \\ 0 & C - B^T A^{-1} B \end{pmatrix} P$. Clearly P is non-singular, thus X is only singular iff $\begin{pmatrix} A & 0 \\ 0 & C - B^T A^{-1} B \end{pmatrix}$ is singular, which is only if either A is singular or $C - B^T A^{-1} B$ is singular. \square

Corollary 1.2.12. *Let $X \in \mathcal{S}^n$ be in block form as in equation 1.2.10. Then:*

$$X \succeq 0 \Leftrightarrow A \succeq 0 \text{ and } C - B^T A^{-1} B \succeq 0 \quad (1.31)$$

Proof. Note that as in theorem 1.2.11: $X = P^T \begin{pmatrix} A & 0 \\ 0 & C - B^T A^{-1} B \end{pmatrix} P$. The result follows direct from lemma 1.2.9. \square

1.3 Duality and optimality conditions

Duality theory allows us to derive optimality conditions for problems such as 1.1.3, 1.1.4, and 1.1.6. Moreover, duality theory allows to recast the problem into a corresponding dual problem which might be easier to compute than the original ‘primal’ problem.

To derive the dual problem, consider first the Lagrangian:

Definition 1.3.1 (Lagrangian). *The Lagrangian of a convex problem (i.e. 1.1.1), with only equality constraints is defined as:*

$$\mathcal{L} : \mathcal{F} \times \mathbb{R}^m \rightarrow \mathbb{R}, (x, y) \mapsto f(x) + \sum_i^m y_i g_i(x) \quad (1.32)$$

If there are in addition inequality constraints:

$$\mathcal{L} : \mathcal{F} \times \mathbb{R}^m \times \mathbb{R}_+^k \rightarrow \mathbb{R}, (x, y, z) \mapsto f(x) + \sum_i^m y_i g_i(x) + \sum_i^k z_i h_i(x) \quad (1.33)$$

The variables (y, z) are known as Lagrange multipliers.

Note that the Lagrangians allow to write convex problems as an one liner:

$$p^* = \max_{x \in \mathcal{F}} \min_{y \in \mathbb{R}^m, z \in \mathbb{R}_+^k} \mathcal{L}(x, y, z) \quad (1.34)$$

Equation 1.34 is known as the primal problem of a convex optimization problem. The dual is obtained by interchanging the max–min:

$$d^* = \min_{y \in \mathbb{R}^m, z \in \mathbb{R}_+^k} \max_{x \in \mathcal{F}} \mathcal{L}(x, y, z) \quad (1.35)$$

From now on p^* and d^* will denote the optimum objective value of the primal and the dual problem. Similar x^* and y^* are used to denote the argument for which the optimum is attained.

Example 1.3.2. *The Lagrangian of 1.1.3 is:*

$$\mathcal{L}(x, y, z) = c^T x + y^T (b - Ax) + z^T x \quad (1.36)$$

This can be rewritten as:

$$\mathcal{L}(x, y, z) = y^T b + (c^T - y^T A + z^T)x \quad (1.37)$$

Hence the dual of 1.1.3 is:

$$\begin{aligned} \min_{y \in \mathbb{R}^m, z \in \mathbb{R}^n} \quad & b^T y \\ \text{s. t.} \quad & c - A^T y + z = 0 \\ & z \geq 0 \end{aligned} \quad (1.38)$$

The more familiar dual form is obtained by removing z from the problem:

$$\begin{aligned} \min_{y \in \mathbb{R}^m} \quad & b^T y \\ \text{s. t.} \quad & c - A^T y \leq 0 \end{aligned} \quad (1.39)$$

It is possible to derive a dual program for a semidefinite program as well through a similar procedure as in the previous example:

Example 1.3.3. *The Lagrangian of 1.1.6 is:*

$$\mathcal{L}(X, y, Z) = \langle C, X \rangle + y^T (b - \mathcal{A}(X)) + \langle Z, X \rangle \quad (1.40)$$

This can be rewritten as:

$$\mathcal{L}(x, y, z) = y^T b + \langle C - \mathcal{A}^T(y) + Z, X \rangle \quad (1.41)$$

In which $\mathcal{A}^T(y) = \sum_i y_i A_i$. Note that by this definition:

$$\langle \mathcal{A}(X), y \rangle = \langle X, \mathcal{A}^T(y) \rangle \quad (1.42)$$

The dual of 1.1.6 is:

$$\begin{aligned} \inf_{y \in \mathbb{R}^m, Z \in \mathcal{S}^n} \quad & b^T y \\ \text{s. t.} \quad & C - \mathcal{A}^T(y) + Z = 0 \\ & Z \succeq 0 \end{aligned} \quad (1.43)$$

And by removing Z :

$$\begin{aligned} \inf_{y \in \mathbb{R}^m} \quad & b^T y \\ \text{s. t.} \quad & C - \mathcal{A}^T(y) \preceq 0 \end{aligned} \quad (1.44)$$

Already from observing that in the primal problem the objective is maximized, and in the dual the objective is minimized, it follows that the primal objective is less than the dual objective. This is proven for semidefinite optimization problems:

Theorem 1.3.4 (Weak duality for SDP). *For any feasible primal solution $X \in \mathcal{S}^n$ (not necessary optimal) and any feasible dual solution $y \in \mathbb{R}^m$ (also not necessary optimal) holds the following inequality:*

$$\langle C, X \rangle \leq b^T y \quad (1.45)$$

Proof. Note that any primal feasible solution $X \in \mathcal{S}^n$ implies $X \succeq 0$ and any dual feasible solution $y \in \mathbb{R}^m$ implies $C - \mathcal{A}^T(y) \preceq 0$. Hence:

$$\langle C, X \rangle - y^T b = \langle C, X \rangle - y^T \mathcal{A}(X) = \langle C - \mathcal{A}^T(y), X \rangle \leq 0 \quad (1.46)$$

□

The quantity $|\langle C - \mathcal{A}^T(y), X \rangle|$ at optimum $X \in \mathcal{S}^n$ and $y \in \mathbb{R}^m$ is known as the duality gap. When it is said that there is no duality gap, then it is implied that the duality gap equals zero. Consequently, if a given primal and dual solution have no duality gap, then the solutions are optimal.

As a side note, if the primal problem is unbounded, i.e. $p^* = \infty$, then by weak duality $d^* \geq \infty$, hence infeasible. Similar if the dual is unbounded, then the primal is infeasible.

In many cases strong duality holds, i.e. there is no duality gap. This stated here without proof:

Theorem 1.3.5 (Strong duality). *If the primal problem is bounded, that is the optimum is finite i.e. $p^* < \infty$, and strictly feasible, then the dual problem attains its infimum and there is no duality gap.*

Proof. See [5].

□

Corollary 1.3.6. *If the dual problem is bounded i.e. $d^* > -\infty$, and strictly feasible, then the primal problem attains its supremum and there is no duality gap.*

By strong duality it follows that for optimum (X^*, y^*) the following equality holds:

$$\langle C - \mathcal{A}^T(y^*), X^* \rangle = 0 \quad (1.47)$$

This equation is known as the complementary slackness condition.

If the dual solution y^* is known, and strong duality holds, then calculating X^* is equivalent to solving a lower dimensional problem than the original primal by noting that the space spanned by $C - \mathcal{A}^T(y)$ should be in the null-space of X^* :

$C - \mathcal{A}^T(y) \preceq 0$, a Cholesky decomposition $LL^T = \mathcal{A}^T(y) - C$ exists. By complementary slackness:

$$\begin{aligned} 0 &= \langle LL^T, X^* \rangle \\ &= \langle I, L^T X^* L \rangle \end{aligned} \quad (1.48)$$

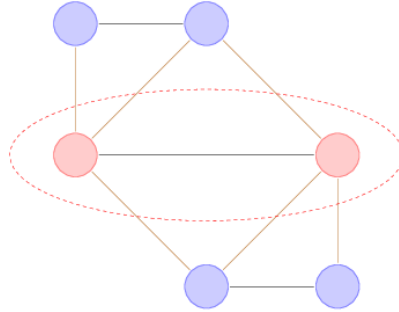


Figure 1.1: An example of a non-unique max-cut

And also $X^* \succeq 0$, hence the columns of L have to be in the null-space of X^* . Let Q has as columns an orthogonal basis perpendicular to the columns of L . Obvious Q has less columns than n ; say it has k columns which form a basis. Then by the variable substitution $X = QUQ^T$, problem 1.1.6 is reduced to the space \mathcal{S}^k .

In practice k is a small number and thus problem 1.1.6 can be solved quite fast.

1.4 Some examples

In this section four examples of semidefinite optimization problems are given. These include from graph theory: The famous Goemans-Williamson cut for solving max-cut instances [12]. The Lovász number, which gives an upper bound to the clique number and a lower bound to the chromatic number [17]. One from linear algebra: calculating ideal GMRES and Arnoldi polynomials for symmetric matrices (based on [13]). And one based on phase retrieval via matrix completion (based on [6]).

1.4.1 Max-cut problem

The max-cut problem for graphs is defined as the problem of finding the largest possible cut for a graph. That is, one wants to separate all nodes of a graph into two disjoint sets such that the number of edges between the two sets is as large as possible.

The corresponding decision problem, that is whether for a given graph G there exists a cut of at least k with k arbitrary large, is NP-complete. This problem is one of Karp's original NP-complete problems [15]. Hence it is probably impossible to derive an exact algorithm to solve the problem which runs within a reasonable time. However Goemans and Williamson accomplished to derive an 0.87856-approximation algorithm (that is, the algorithm always delivers an answer that at least is 0.87856 times the optimal value) which runs

in polynomial time using semidefinite optimization. The resulting algorithm is known as the Goemans-Williamson max-cut algorithm.

In the original paper in which Goemans and Williamson describe their approximation algorithm based on semidefinite optimization, they considered G directed and weighted [12]. However because the problem is stated here only as an example of semidefinite optimization and those details make it unnecessary more complex, those details are omitted. The algorithm described in the rest of this subsection is the specific case of the Goemans-Williamson max-cut for undirected uniform weighted graphs.

In the remainder of this subsection it is assumed that $G = (V, E)$ (graph G with nodes / vertexes V and edges E) is undirected. That is if $(i, j) \in E$ then also $(j, i) \in E$.

Definition 1.4.1. *A cut on $G = (V, E)$ is a partition of V into two sets S, \bar{S} such that $S \cap \bar{S} = \phi$ and $S \cup \bar{S} = V$. The size of the cut is defined as the number of edges starting in S and ending in \bar{S} , that is: $C(G, S) := |\{(i, j) \in E \text{ s. t. } i \in S, j \in \bar{S}\}|$.*

An example of a max-cut is given in figure 1.1.

The main idea of the Goemans-Williamson algorithm is to associate with every vertex $i \in V$ a binary variable $x_i \in \{1, -1\}$ which denotes whether $i \in S$ or $i \in \bar{S}$. Then for two vertices $i, j \in V$:

$$\frac{1 - x_i x_j}{2} = \begin{cases} 0 & \text{if } i, j \in S \text{ or } i, j \in \bar{S} \\ 1 & \text{if } i \in S, j \in \bar{S} \text{ or } i \in \bar{S}, j \in S \end{cases} \quad (1.49)$$

Hence the value for a given cut is:

$$C(G, S) = \frac{1}{2} \sum_{(i,j) \in E} \frac{1 - x_i x_j}{2} \quad (1.50)$$

(The fraction $\frac{1}{2}$ at the beginning arises because $(i, j) \in E$ implies $(j, i) \in E$, hence all edges are counted twice.)

And the max cut is thus:

$$C_{\max}(G) = \max_{x \in \mathbb{R}^{|V|}} \frac{1}{2} \sum_{(i,j) \in E} \frac{1 - x_i x_j}{2} \quad \text{s. t. } x_i \in \{-1, +1\} \quad (\forall i = 1, \dots, |V|) \quad (1.51)$$

Therefore equation 1.51 has to be solved, however the constraint that all x_i is a binary variable makes it difficult. Hence consider the relaxation to write it

as a semidefinite program:

$$\begin{aligned}
C_{\text{GW}}(G) = \max_{X \in \mathcal{S}^{|V|}} & \frac{1}{4} \sum_{(i,j) \in E} (1 - X_{i,j}) \\
\text{s. t. } & X_{i,i} = 1 \quad (\forall i \in V) \\
& X \succeq 0
\end{aligned} \tag{1.52}$$

In the ideal case the $X^* \in \mathcal{S}^n$ that maximizes problem 1.52 can be decomposed as $X^* = xx^T$ with $x_i = \pm 1$ for all i . However, because C_{GW} is a relaxation it is to be expected to find a better solution than X^* which does not permit such a decomposition. The great accomplishment of Goemans-Williamson is that they showed:

Theorem 1.4.2. *Given a undirected graph G with uniform weights, the following inequalities hold:*

$$C_{\text{GW}}(G) \geq C_{\text{max}}(G) \geq 0.87856 C_{\text{GW}}(G) \tag{1.53}$$

The full proof can be found in [12]. The proof is more algorithmic and based on the following steps:

1. Solve X for equation 1.52.
2. Perform a Cholesky decomposition $X = VV^T$ such that $X_{i,j} = v_i^T v_j$ (columns of V^T)
3. Select a random unit vector r and define a cut $S = \{i \text{ s. t. } r^T v_i \geq 0\}$
4. Check if $C(G, S) \geq 0.87856 C_{\text{GW}}(G)$, if not, go back to step 3

The number 0.87856 appears in the construction of S (step 3) as the result of some sort of randomized rounding the 'fractional' solution to fit the original problem statement.

In order to derive the dual of program 1.52 first rewrite it to a standard semidefinite problem. Note that the the sum in the objective can be rewritten as: $\sum_{(i,j) \in E} (1 - X_{i,j}) = |E| - \langle X, A_G \rangle = \langle X, L_G \rangle$; that is A_G as the adjacency matrix and L_G the Laplacian matrix which are defined as:

$$\begin{aligned}
(A_G)_{i,j} &= \begin{cases} 1 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases} \\
(L_G)_{i,j} &= \begin{cases} -1 & \text{if } (i,j) \in E \\ \text{deg}(i) & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}
\end{aligned} \tag{1.54}$$

(Here $\text{deg}(i)$ is the degree of vertex i , i.e. $\text{deg}(i) = |\{(i,j) \in E\}|$.)

And the constraint $X_{i,i} = 1$ can be reworded using $E_{i,i} = e_i e_i^T$: $X_{i,i} = \langle X, E_{i,i} \rangle = 1$.

This makes it easy to write the primal and dual pair:

$$\begin{aligned} C_{\text{GW}}(G) &= \max_{X \in \mathbb{S}^{|V|}} \frac{1}{4} \langle L_G, X \rangle \\ \text{s. t. } &\langle X, E_{i,i} \rangle = 1 \quad (\forall i \in V) \\ &X \succeq 0 \end{aligned} \tag{1.55}$$

$$\begin{aligned} C_{\text{GW}}^*(G) &= \min_{y \in \mathbb{R}^{|V|}} \sum_i y_i \\ \text{s. t. } &\frac{1}{4} L_G - \sum_i y_i E_{i,i} \preceq 0 \end{aligned} \tag{1.56}$$

1.4.2 Lovász theta number

The Lovász theta number was introduced in 1979 as an upper bound on the Shannon capacity of a graph [17]. Lovász theta number is lower bounded by the clique number of the complement graph, and upper bounded by the chromatic number of its complement [16]. Hence it is an important number for the analysis of graphs.

The clique number and the chromatic number for undirected graphs are NP-hard to compute, hence the Lovász number is interesting because it can be computed in polynomial time by using semidefinite optimization. Moreover, at the time of writing there is no algorithm yet known to compute the Shannon capacity of a graph and even for small graphs remains the Shannon capacity unknown. However, it can be easily be seen as lower bounded by the clique number. Thus if the clique number and the Lovász theta number coincide then the Shannon capacity is equal to those numbers. Lovász used this to calculate the Shannon capacity of the pentagon graph (cyclic, 5-vertices) to be $\sqrt{5}$, which was still unknown at the time [17].

In this section the Lovász theta number is presented as a semidefinite program, in the original article this is just one of the representations presented. In this section the many other representations will be omitted and stick to the semidefinite program because that is the one actually used later in chapter 5.

Definition 1.4.3. *In an undirected graph $G = (V, E)$, a clique $C \subset V$ is a set of vertices that are all adjacent to each other. A maximal clique is a clique that cannot be extended by adding another vertex.*

The clique number of a graph is the size of its largest clique, which is an upper bound to the size of all possible maximal cliques. The clique number for graph G is denoted as $\omega(G)$

Definition 1.4.4. *A coloring of a graph $G = (V, E)$ is a labeling of the vertices V such that no two adjacent vertices are labeled the same. A k -coloring is a*

partition of V into k sets, C_1, \dots, C_k , such that each set is independent, i.e. no two vertices in C_i are adjacent.

The chromatic number of a graph is a lower bound on the labels of a coloring, that is the smallest number of labels necessary to color the graph. The chromatic number for graph G is denoted as $\chi(G)$.

Definition 1.4.5. The complement of a graph $G = (V, E)$ is $\bar{G} = (V, \bar{E})$. That is, in \bar{G} two vertices are adjacent if and only if those two vertices are not-adjacent in G .

Definition 1.4.6. The Lovász number for a graph $G = (V, E)$ is defined as [17]:

$$\begin{aligned} \theta(G) = \max_{B \in \mathcal{S}^{|V|}} \langle J, B \rangle \\ \text{s. t. } \langle B, I \rangle = 1 \\ \langle B, E_{ij} \rangle = 0 \quad (\forall i, j \in V \text{ adjacent}) \\ B \succeq 0 \end{aligned} \tag{1.57}$$

In which J is the all ones matrix and E_{ij} is the matrix: $e_i e_j^T + e_j e_i^T$.

The famous sandwich theorem:

Theorem 1.4.7. For any undirected graph $G = (V, E)$ holds:

$$\omega(G) \leq \theta(\bar{G}) \leq \chi(G) \tag{1.58}$$

Proof. First $\omega(G) \leq \theta(\bar{G})$:

Let $x \in \mathbb{R}^{|V|}$ be the characteristic vector for a maximal clique C in G (i.e. if $i \in C$, $x_i = 1$ else $x_i = 0$) of size $\omega(G)$. Then a feasible solution to 1.57 for \bar{G} is:

$$B' = \frac{1}{\omega(G)} x x^T \tag{1.59}$$

Moreover, $\langle J, B' \rangle = \omega(G)$, hence $\theta(\bar{G}) \geq \omega(G)$.

Second $\theta(\bar{G}) \leq \chi(G)$:

Let B be any feasible matrix for equation 1.57; then $\langle J, B \rangle \leq \chi(G)$: Let C_1, \dots, C_k be a k -coloring such that $k = \chi(G)$. Let y_1, \dots, y_k be the characteristic vectors of C_1, \dots, C_k . Moreover, define:

$$Y = \sum_i^k (k y_i - e)(k y_i - e)^T = k^2 \sum_i^k y_i y_i^T - kJ \tag{1.60}$$

(e is the all ones vector.)

Note that $Y \succeq 0$, hence $\langle B, Y \rangle \geq 0$. Also:

$$\begin{aligned} \langle B, Y \rangle &= \langle B, k^2 \sum_i y_i y_i^T - kJ \rangle \\ &= k^2 \langle B, I \rangle - k \langle B, J \rangle \\ &= k^2 - k \langle B, J \rangle \geq 0 \end{aligned} \tag{1.61}$$

Hence for any feasible B of equation 1.57: $\langle J, B \rangle \leq \chi(G)$. \square

The dual of 1.57 is given by:

$$\begin{aligned} \min_{t \in \mathbb{R}, y \in \mathbb{R}^{|E|}} \quad & t \\ \text{s. t.} \quad & J - \mathcal{E}^T(y) - tI \preceq 0 \end{aligned} \tag{1.62}$$

In which $\mathcal{E}^T(y) = \sum_{(i,j) \in E} E_{ij} y_{ij}$.

Clearly this is already an eigenvalue optimization problem because t is minimal if and only if its equal to the maximum eigenvalue:

$$\min_{y \in \mathbb{R}^{|E|}} \lambda_{\max}(J - \mathcal{E}^T(y)) \tag{1.63}$$

1.4.3 Ideal GMRES and Arnoldi polynomes

The generalized minimal residual method (GMRES) is an iterative Krylov subspace method for numerical solving $Ax = b$ introduced by Saad and Schultz in 1986 [29]. Since its development there has been active research to convergence bounds of the algorithm and is still today considered a difficult open problem [31].

During GMRES the solution to $Ax = b$ is solved through a least squares restricted to a Krylov subspace. In each iteration the subspace is extended with one orthogonal vector, until the solution of the least squares is an exact solution. The Krylov subspace which is constructed during the algorithm up to the k th iteration is:

$$\mathcal{K}_k(A, Ab) = \text{span}\{Ab, A^2b, \dots, A^k b\} \tag{1.64}$$

Let x_k be the k iterative step of GMRES and $r_k := Ax_k - b$. Then in GMRES:

$$\|r_k\|_2 = \min_{q \in P_k} \|q(A)r_0\|_2 \tag{1.65}$$

In which P_k is the set of all polynomials of degree at most k with for all $q \in P_k$, $q(0) = 1$.

Note that the equation above is a least squares problem in disguise, $\{Aq(A)b : q \in P_k\}$ spans exactly the Krylov subspace $\mathcal{K}_k(A, Ab)$.

An easy upper bound on convergence of GMRES is:

$$\frac{\|r_k\|_2}{\|r_0\|_2} \leq \min_{q \in P_k} \|q(A)\|_2 \quad (1.66)$$

For a more thorough description of the algorithm and implementation details, see [29].

A similar iterative algorithm is Arnoldi iteration (developed by Arnoldi in 1951 [2]) which is used to compute eigenvalues of matrix A .

During Arnoldi iteration the matrix A is projected on the Krylov subspace $\mathcal{K}_k(A, b)$ (b a random unit vector). That is, let Q_k be a basis for $\mathcal{K}_k(A, b)$ and define $H_k = Q_k^T A Q_k$. Then as k increases, and the Krylov subspace keeps expanding, the matrix H_k will become more similar to A . Then the Ritz values of H_k converge to the eigenvalues of A .

It can be shown that the characteristic polynome of H_k is a monic polynome of degree k that minimizes $\|q(A)b\|_2$:

$$\min_{q \in P^k} \|q(A)b\|_2 \quad (1.67)$$

In which P^k is the set of all monic polynomials of degree k [13]. For a more inclusive description see [2].

As with GMRES, an easy upper bound on convergence is given by:

$$\min_{q \in P^k} \|q(A)\|_2 \quad (1.68)$$

The $q \in P_k$ which minimizes equation 1.66 is known as the ideal GMRES polynome and the $q \in P^k$ which minimizes equation 1.68 is known as the ideal Arnoldi polynome.

The motivation of calculating the ideal GMRES and Arnoldi polynomes is that they give an upper bound on convergence and can also be applied to testing different preconditioning strategies. The computation of the ideal GMRES/Arnoldi polynome can be formulated as a semidefinite optimization program through the following identity:

Lemma 1.4.8. *The matrix 2-norm for matrix $A \in \mathcal{S}^n$ is the largest absolute eigenvalue of A . That is:*

$$\|A\|_2 = \max_i |\lambda_i| \quad (1.69)$$

Where λ_i is a distinct labeled eigenvalue of A . Furthermore the following equality holds:

$$\|A\|_2 = \lambda_{\max} \left(\begin{pmatrix} 0 & A \\ A^T & 0 \end{pmatrix} \right) \quad (1.70)$$

Proof. Let $Z = \begin{pmatrix} 0 & A \\ A & 0 \end{pmatrix}$, denote the set of eigenvalues of Z as $\lambda(Z)$ and the eigenvalues of A as $\lambda(A)$. Then $\lambda(Z) = \pm\lambda(A)$. To see that: Label all eigenvectors of A as q_i such that $q_i^T A q_i = \lambda_i$ and $q_i^T q_j = \delta_{i,j}$ (i.e. 0 if $i \neq j$ and 1 if $i = j$). Such a decomposition is possible by the spectral decomposition theorem (theorem 1.2.2).

Then $(q_i, \pm q_i)^T$ forms two eigenvectors of Z with eigenvalues $\pm\lambda_i$. Hence because A has n eigenvalues and this way $2n$ eigenvalues for Z (which is of size $2n$) are created, it follows that Z cannot have any other eigenvalue.

Therefore the largest eigenvalue of Z is the largest absolute eigenvalue of A . \square

The computation of the ideal GMRES polynomial of degree k is:

$$\begin{aligned} & \min_{t \in \mathbb{R}, y \in \mathbb{R}^k} t \\ & \text{s. t. } \begin{pmatrix} 0 & I \\ I & 0 \end{pmatrix} + \sum_{1 \leq i \leq k} y_i \begin{pmatrix} 0 & A^i \\ (A^i)^T & 0 \end{pmatrix} - t \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} \preceq 0 \end{aligned} \quad (1.71)$$

And similar the ideal Arnoldi polynomial of degree k :

$$\begin{aligned} & \min_{t \in \mathbb{R}, y \in \mathbb{R}^k} t \\ & \text{s. t. } \begin{pmatrix} 0 & A^k \\ (A^k)^T & 0 \end{pmatrix} + \sum_{1 \leq i \leq k} y_i \begin{pmatrix} 0 & A^{i-1} \\ (A^{i-1})^T & 0 \end{pmatrix} - t \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} \preceq 0 \end{aligned} \quad (1.72)$$

1.4.4 Phase retrieval

This subsection is largely based on [6].

Phase retrieval is the problem of reconstructing a signal based on the magnitude of its Fourier transform. This problem arises in many applications due to that detectors often can only detect the squared signal.

Historically one of the first applications of phase retrieval is X-ray crystallography. In X-ray crystallography are X-rays aimed at a crystal and the angle of dispersion and the intensity is measured. From those measurements a 3D-image can be reconstructed with the density of the electrons within the crystal.

Moreover, different applications of Phase Retrieval include MRI, astronomical imaging and microscopy.

To fix the problem statement, consider a discrete one dimensional signal of size n : $x_0, \dots, x_{n-1} \in \mathbb{R}$ (notice the indexing starts with 0). Its Fourier transform is given by:

$$\tilde{x}(\omega) = \frac{1}{\sqrt{n}} \sum_t x_t \exp(-i2\pi\omega t/n), \quad \omega \in \Omega \quad (1.73)$$

The set Ω is a grid of sampled frequencies. The usual sampled frequencies are $\{0, \dots, n-1\}$.

The problem of Phase Retrieval is reconstructing x from the magnitude measurements $|\tilde{x}(\omega)|$, $\omega \in \Omega$. Note that the problem is ill-posed because there is no unique solution. If $x \in \mathbb{R}^n$ is a solution the problem, then as an example, $-x$ is also a solution. In general there are infinite many solutions. However from a practical point, all these solutions are acceptable ambiguities.

Observe that $\tilde{x}(\omega)$ can be written as $a_\omega^* x$ in which $a_\omega = \frac{1}{\sqrt{n}} \times (1, \exp(-i2\pi\omega 1/n), \dots)^*$. Hence:

$$\begin{aligned} |\tilde{x}(\omega)|^2 &= |\langle a_\omega, x \rangle|^2 \\ &= \langle a_\omega a_\omega^*, x x^* \rangle \end{aligned} \tag{1.74}$$

Denote $A_\omega = a_\omega a_\omega^*$.

As for the Max-Cut, the quadratic constraint $\langle A_\omega, x x^* \rangle$ can be relaxed by optimizing over the space of positive semidefinite symmetric matrices and substitute $X = x x^*$.

The SDP problem is a feasibility problem:

$$\begin{aligned} \text{find}_{X \in \mathcal{S}^n} \quad & X \\ \text{s. t.} \quad & \langle A_\omega, X \rangle = b_\omega \quad (\omega \in \Omega) \\ & X \succeq 0 \\ & \text{rank}(X) = 1 \end{aligned} \tag{1.75}$$

$$(b_\omega = |\tilde{x}(\omega)|^2)$$

The constraint on the rank of X is a non-convex constraint. A possibility to lift the constraint is by replacing it by a trace constraint: $\langle I, X \rangle = 1$. Moreover, as a further relaxation it can be incorporated into the objective such that the problem can be solved by standard SDP solving algorithms:

$$\begin{aligned} \min_{X \in \mathcal{S}^n} \quad & \langle I, X \rangle \\ \text{s. t.} \quad & \langle A_\omega, X \rangle = b_\omega \quad (\omega \in \Omega) \\ & X \succeq 0 \end{aligned} \tag{1.76}$$

Through these relaxations the resulting SDP offers a large degree of freedom to deal with noise and other errors made by abstracting the problem. That can be for instance be done by replacing the equality constraint with an upper and lower bound, or by using some maximum likelihood principle as was done in [6].

After obtaining an optimum $X \in \mathcal{S}^n$ of sufficient low rank, the original vector $x \in \mathbb{R}^n$ can be found through performing a Cholesky decomposition.

If X is not of sufficient low rank, then it is possible to promote in equation 1.76 solutions of lower rank through solving for a sequence of weighted trace

norms. That works as follows: Choose $\varepsilon > 0$ and $W_0 = I$. Next solve for $k = 0$:

$$\begin{aligned} \min_{X_k \in \mathcal{S}^n} \langle W_k, X_k \rangle \\ \text{s. t. } \langle A_\omega, X_k \rangle = b_\omega \quad (\omega \in \Omega) \\ X_k \succeq 0 \end{aligned} \tag{1.77}$$

Next, update $W_{k+1} = (X_k + \varepsilon I)^{-1}$ and solve 1.77 for $k = k + 1$. Repeat until X_k is of sufficient low rank or k reaches a threshold k_{\max} .

As described in [10], this scheme can be viewed as an iterative algorithm to solve:

$$\begin{aligned} \min_{X \in \mathcal{S}^n} \log(\det(X + \varepsilon I)) \\ \text{s. t. } \langle A_\omega, X \rangle = b_\omega \quad (\omega \in \Omega) \\ X \succeq 0 \end{aligned} \tag{1.78}$$

Hence, the iterative method is known as the log det heuristic.

The dual of 1.76 is:

$$\begin{aligned} - \min_{y \in \mathbb{R}^{|\Omega|}} b^\top y \\ \text{s. t. } -I - \mathcal{A}^\top(y) \preceq 0 \end{aligned} \tag{1.79}$$

Or by substituting $y' = -y$:

$$\begin{aligned} \max_{y \in \mathbb{R}^{|\Omega|}} b^\top y \\ \text{s. t. } I - \mathcal{A}^\top(y) \succeq 0 \end{aligned} \tag{1.80}$$

Note that $y = 0$ is a strict feasible dual solution. Hence if the problem is bounded and the optimum attained, then strong duality holds.

Solving the dual might not seem worthwhile because of the necessity of the primal solution for recovering the signal. However, as remarked at the end in section 1.3, if strong duality holds and the dual solution y^* is known, then the primal problem can be shown to be equivalent to a low dimensional problem.

Moreover, by solving the primal as well as the dual in the log det heuristic, a near optimum feasible solution for the next iteration is already obtained. This will be elaborated on in section 2.1.2.

Chapter 2

Eigenvalue optimization

In this chapter it will be shown how many semidefinite optimization problems can be recast into an eigenvalue optimization problem. That are problems in which the objective is to maximize or minimize the eigenvalue of a matrix through linear perturbations. The section after derives some important properties of eigenvalue problems which will be used in solving the problems.

2.1 Equivalence of semidefinite problems

It is possible to derive for most semidefinite problems an eigenvalue optimization problem. That is a problem in which the eigenvalues of a matrix operator are minimized or maximized. To do so, first perform a variable rotation, second apply a Cholesky decomposition and third conclude that the corresponding problem can be interpreted as an eigenvalue optimization problem.

Through the rest of this section the following problem will be recast into an eigenvalue optimization problem:

$$\begin{aligned} \inf_{y \in \mathbb{R}^m} b^T y \\ \text{s. t. } C - \mathcal{A}^T(y) \preceq 0 \\ Py \geq 0 \end{aligned} \tag{2.1}$$

In which Py is the vector with components of y that correspond to inequalities. This equation is similar to equation 1.44.

We want to recast problem 2.1 into a problem of the following form:

$$\begin{aligned} \inf_{z \in \mathbb{R}^k} \lambda_{\max}(\tilde{C} - \tilde{\mathcal{A}}^T(z)) \\ \text{s. t. } \tilde{P}^T z \geq 0 \end{aligned} \tag{2.2}$$

The reason being is that eigenvalue problems such as 2.2 can efficiently be solved by an eigenvalue optimization algorithm, which will be the subject for the rest of this thesis.

To do so apply a variable transformation: Define orthogonal square matrix Q of full rank such that $z = Q^T y$ and $b^T y = \|b\|_2 z_1$ holds for all y . This implies that:

$$\begin{aligned} Q &= (q_1 \mid q_2 \mid q_3 \mid \dots \mid q_m) \\ q_1 &= \hat{b} \end{aligned} \tag{2.3}$$

In which $\hat{b} = \frac{b}{\|b\|_2}$. Note that $QQ^T = Q^T Q = I$, the identity.

By using the transformation $y = Qz$, the SDP problem is modified in the following manner:

$$\begin{aligned} \inf_{z \in \mathbb{R}^m} \quad & \|b\|_2 z_1 \\ \text{s. t.} \quad & C - \mathcal{A}^T(Qz) \preceq 0 \\ & PQz \geq 0 \end{aligned} \tag{2.4}$$

Note that $\mathcal{A}^T(Qz)$ can be simplified:

Lemma 2.1.1. *A transposed linear matrix operator under rotation of its argument is still a transposed linear matrix operator. That is:*

$$\mathcal{A}^T(Qz) = \mathcal{A}_Q^T(z) = \sum_j z_j \mathcal{A}^T(q_j) \tag{2.5}$$

In which $q_1, \dots, q_m \in \mathbb{R}^m$ denote the column vectors of Q .

Proof.

$$\mathcal{A}^T(Qz) = \sum_i (Qz)_i A_i = \sum_i \left(\sum_j Q_{i,j} z_j \right) A_i \tag{2.6}$$

(In here $Q_{i,j}$ means the element in Q on the i th row, j th column.)

By swapping the sums:

$$\sum_i \left(\sum_j Q_{i,j} z_j \right) A_i = \sum_j z_j \left(\sum_i Q_{i,j} A_i \right) \tag{2.7}$$

This formulation is recognizable as a linear operator in z much alike as \mathcal{A}^T , which is a motivation to denote this formula as the operator \mathcal{A}_Q^T acting on z . Additionally, for clean notation, express $\sum_i Q_{i,j} A_i$ as $\mathcal{A}^T(q_j)$. \square

Second, next comes the crucial part in which the SDP problem is recast into an eigenvalue optimization problem.

Lemma 2.1.2. Consider the SDP problem 2.4. If $\mathcal{A}^T(b)$ is positive definite, then the SDP constraint can be altered into an eigenvalue constraint on z_1 :

$$z_1 \geq \lambda_{\max}(L^{-1}CL^{-T} - L^{-1}(\sum_{j \geq 2} z_j \mathcal{A}^T(q_j))L^{-T}) \quad (2.8)$$

In which $LL^T = \mathcal{A}^T(b)$ is the Cholesky decomposition

Proof. Note that because $\mathcal{A}^T(b)$ is positive definite, L , L^T , L^{-1} and L^{-T} are all non-singular, hence:

$$C - \mathcal{A}^T(Qz) = C - z_1 \mathcal{A}^T(b) - \sum_{j \geq 2} z_j \mathcal{A}^T(q_j) \quad (2.9)$$

$$= C - z_1 LL^T - \sum_{j \geq 2} z_j \mathcal{A}^T(q_j) \quad (2.10)$$

$$= L(L^{-1}CL^{-T} - z_1 I - L^{-1}(\sum_{j \geq 2} z_j \mathcal{A}^T(q_j))L^{-T})L^T \quad (2.11)$$

Thus the constraint $C - \mathcal{A}^T(Qz) \preceq 0$ can be written as:

$$L^{-1}CL^{-T} - L^{-1}(\sum_{j \geq 2} z_j \mathcal{A}^T(q_j))L^{-T} - z_1 I \preceq 0 \quad (2.12)$$

(See theorem 1.2.9 and note that neither L^{-1} nor L^{-T} are singular).

Because both L and L^T are positive definite.

Which is equivalent to stating that $z_1 \geq \lambda_{\max}(L^{-1}CL^{-T} - L^{-1}(\sum_{j \geq 2} z_j \mathcal{A}^T(q_j))L^{-T})$. \square

Third, and finally, check if the inequality constraint in the SDP problem is indifferent to z_1 :

Lemma 2.1.3. Consider the SDP problem 2.4. If $Pb = 0$, then $PQz \geq 0$ is independent to z_1 .

Proof.

$$PQz = z_1 Pb + PQ(0, z_2, \dots, z_m)^T = PQ(0, z_2, \dots, z_m)^T \quad (2.13)$$

\square

And finally, combining all of this into one theorem:

Theorem 2.1.4. Consider the SDP problem 2.1. If $\mathcal{A}^T(b) \succ 0$ and $Pb = 0$, then the SDP problem is equivalent to the following eigenvalue optimization problem:

$$\begin{aligned} \inf_{z \in \mathbb{R}^{m-1}} \|b\|_2 \lambda_{\max}(L^{-1}CL^{-T} - \sum_j z_j L^{-1} \mathcal{A}^T(q'_j) L^{-T}) \\ \text{s. t. } PQ'z \geq 0 \end{aligned} \quad (2.14)$$

In which q'_1, \dots, q'_{m-1}, b form a full rank orthogonal basis, $Q' = (q_1 | \dots | q_{m-1})$, and $LL^T = \mathcal{A}^T(b)$ is a Cholesky decomposition.

An implementation based on theorem 2.1.4 is included in appendix A.1.

2.1.1 Perturbing semidefinite optimization problems

This subsection uses results from section 2.2. If unfamiliar with analysis on eigenvalues then it is advised to read that section first.

When trying to use theorem 2.1.4 and $\mathcal{A}^T(b) \succeq 0$ but $\mathcal{A}^T(b)$ is not of full rank, then the trick with the Cholesky decomposition in lemma 2.1.2 is no longer valid. Therefore it is tempting to add a small perturbation such that $\tilde{\mathcal{A}}^T(\tilde{b})$ is positive definite. In this section one such possible perturbation is explored.

As illustration of the perturbation being examined: Assuming that $\mathcal{A}^T(b)$ is positive semidefinite with rank $n - 1$. Let v be such that $\mathcal{A}^T(b)v = 0$, then adding ϵvv^T with $\epsilon > 0$ to $\mathcal{A}^T(b)$ is enough to cast it positive definite. The proposed perturbation is extending b with one component ϵ and \mathcal{A}^T with the matrix vv^T . This idea, and the extension to a larger dimensional null space, is examined in the next two lemmas and theorems.

Lemma 2.1.5. *For $V \in \mathcal{S}^n$ with $\|V\|_F = \langle V, V \rangle = 1$ and $\epsilon > 0$: Let $\tilde{X} = X + \epsilon V$ and $\langle V, \tilde{X} \rangle \geq \epsilon$. If $X \succeq 0$ then and only then $\tilde{X} \succeq 0$.*

Proof. Note that if $X \succeq 0$, then \tilde{X} is the sum of two positive semidefinite matrices, hence also positive semidefinite.

For the other direction, $\tilde{X} \succeq 0$ is equivalent to that for all $Y \in \mathcal{S}_{\succeq 0}^n$ the relation $\langle \tilde{X}, Y \rangle \geq 0$ holds. Let $Y_1 = \langle Y, V \rangle V$ and $Y_2 = Y - \langle Y, V \rangle V$ for any Y . Note that both Y_1 and Y_2 are positive semidefinite. Then $\langle \tilde{X}, Y \rangle \geq 0$ and:

$$\begin{aligned}
\langle X, Y \rangle &= \langle \tilde{X} - \epsilon V, Y \rangle \\
&= \langle \tilde{X}, Y \rangle + \langle -\epsilon V, Y \rangle \\
&= \langle \tilde{X}, Y_1 + Y_2 \rangle - \epsilon \langle V, Y \rangle \\
&= \langle \tilde{X}, V \rangle \langle Y, V \rangle + \langle \tilde{X}, Y_2 \rangle - \epsilon \langle V, Y \rangle \\
&\geq \epsilon \langle Y, V \rangle + \langle \tilde{X}, Y_2 \rangle - \epsilon \langle V, Y \rangle \\
&\geq \langle \tilde{X}, Y_2 \rangle \geq 0
\end{aligned} \tag{2.15}$$

Hence for all $Y \in \mathcal{S}_{\succeq 0}^n$, $\langle X, Y \rangle \geq 0$ thus $X \succeq 0$. □

Lemma 2.1.6. *Consider problem 2.1 with $V \in \mathcal{S}_{\succeq 0}^n$. Then the following per-*

turbed primal problem has the same optimal as the unperturbed problem:

$$\begin{aligned}
& \sup_{\tilde{X} \in \mathcal{S}^n} \langle C, \tilde{X} - \varepsilon V \rangle \\
& \text{s. t. } \mathcal{A}(\tilde{X} - \varepsilon V) = b \\
& \quad \langle V, \tilde{X} \rangle \geq \varepsilon \\
& \quad \tilde{X} \succeq 0
\end{aligned} \tag{2.16}$$

Proof. Note that if the original primal problem is unbounded, then the perturbed problem is also unbounded and the statement holds trivial. Hence, assume that the original primal problem is bounded.

To continue, observe that the perturbed primal problem is equivalent to transforming $\tilde{X} = X + \varepsilon V$ for optimal \tilde{X} and X ; because with this transformation the objective of the primal problem is equivalent to the objective of the perturbed primal problem, as well as the first equality constraint.

Hence the only difference is in the semidefinite constraint, but because $\langle V, \tilde{X} \rangle \geq \varepsilon$ it follows that if $\tilde{X} \succeq 0$ then and only then $X \succeq 0$, as proven in lemma 2.1.5. Therefore the constraint is in the original problem equivalent to the perturbed constraint. \square

The next theorem examines the case where the rank of $\mathcal{A}^T(b)$ is $n - 1$. That is, if the dimension of the null space is 1. This case is intuitive easier to follow because the derivative of the corresponding eigenvalue is well defined. The multidimensional case is written down in theorem 2.1.8.

Theorem 2.1.7. *Consider the perturbed primal 2.16, with $\mathcal{A}^T(b)$ of rank $n - 1$ and $v \in \mathbb{R}^n$ normalized such that $\mathcal{A}^T(b)v = 0$ and $V = vv^T$. Denote $\tilde{\mathcal{A}}$ and \tilde{b} as the corresponding perturbed values appearing in the dual:*

$$\begin{aligned}
& \inf_{y \in \mathbb{R}^m} \tilde{b}^T y \\
& \text{s. t. } C - \tilde{\mathcal{A}}^T(y) \preceq 0 \\
& \quad Py \geq 0
\end{aligned} \tag{2.17}$$

Then $\tilde{\mathcal{A}}^T(\tilde{b}) = \mathcal{A}^T(b) + \varepsilon(\mathcal{A}^T \circ \mathcal{A}(vv^T) + vv^T)$, and $\varepsilon > 0$ can be chosen such that $\tilde{\mathcal{A}}^T(\tilde{b})$ is of full rank.

Furthermore, if $\varepsilon < \delta$ for a sufficient small δ , then $\tilde{\mathcal{A}}^T(\tilde{b}) \succ 0$ (is positive definite).

Proof. Following duality theory:

$$\tilde{b} = [(b + \varepsilon \mathcal{A}(vv^T))^T \mid \varepsilon]^T \tag{2.18}$$

$$\tilde{\mathcal{A}}^T((y_1, \dots, y_m, y_{m+1})^T) = \mathcal{A}^T((y_1, \dots, y_m)^T) + y_{m+1}vv^T \tag{2.19}$$

$$P = e_{m+1} \tag{2.20}$$

Following those relations:

$$\begin{aligned}
\tilde{\mathcal{A}}^T(\tilde{b}) &= \mathcal{A}^T(b + \varepsilon \mathcal{A}(vv^T)) + \varepsilon vv^T \\
&= \mathcal{A}^T(b) + \varepsilon \mathcal{A}^T \circ \mathcal{A}(vv^T) + \varepsilon vv^T \\
&= \mathcal{A}^T(b) + \varepsilon (\mathcal{A}^T \circ \mathcal{A}(vv^T) + vv^T)
\end{aligned} \tag{2.21}$$

Next use the continuity theorem 2.2.4 and differentiability theorem 2.2.7 to show that for small $z > 0$, $f(z)$ is of full rank:

Define:

$$\begin{aligned}
f : \mathbb{R} &\rightarrow \mathcal{S}^n, \quad z \mapsto \mathcal{A}^T(b) + zE \\
E &= \mathcal{A}^T \circ \mathcal{A}(vv^T) + vv^T
\end{aligned} \tag{2.22}$$

Because $\lambda_{n-m} > 0$ for all $1 \leq m \leq n-1$, it follows by the continuity theorem 2.2.4 that $\lambda_{n-m}(f(z)) > 0$ for small $z > 0$.

Further, observe that $\frac{\mathbf{d}\lambda_n f(z)}{\mathbf{d}z} > 1$, using theorem 2.2.7:

$$\begin{aligned}
\frac{\mathbf{d}\lambda_n f(z)}{\mathbf{d}z} &= v^T E v \\
&= v^T (\mathcal{A}^T \circ \mathcal{A}(vv^T) + vv^T) v \\
&= v^T \mathcal{A}^T \circ \mathcal{A}(vv^T) v + vv^T v \\
&= \langle \mathcal{A}^T \circ \mathcal{A}(vv^T), vv^T \rangle + 1 \\
&= \langle \mathcal{A}(vv^T), \mathcal{A}(vv^T) \rangle + 1 \\
&\geq 1
\end{aligned} \tag{2.23}$$

Hence also $\lambda_n(f(z)) > 0$ for $z > 0$ small. Thus there exists a $\delta > 0$ such that for all $0 < \varepsilon < \delta$, $f(\varepsilon)$ is of full rank. \square

This is the multidimensional case:

Theorem 2.1.8. *Consider the perturbed primal 2.16, with $\mathcal{A}^T(b)$ of rank $n-k$. That is $\lambda_{n-k+1} = \dots = \lambda_{n-1} = 0$. Let $V \in \mathcal{S}_{\geq 0}^n$ such that $\|V\|_F = 1$. Denote $\tilde{\mathcal{A}}$ and \tilde{b} as the corresponding perturbed values appearing in the dual:*

$$\begin{aligned}
&\inf_{y \in \mathbb{R}^m} \tilde{b}^T y \\
&\text{s. t. } C - \tilde{\mathcal{A}}^T(y) \preceq 0 \\
&\quad Py \geq 0
\end{aligned} \tag{2.24}$$

Then $\tilde{\mathcal{A}}^T(\tilde{b}) = \mathcal{A}^T(b) + \varepsilon(\mathcal{A}^T \circ \mathcal{A}(V) + V)$. Moreover if for each $v \in \mathbb{R}^n$ with $\mathcal{A}^T(b)v = 0$, $v^T \mathcal{A}^T \circ \mathcal{A}(V)v + v^T V v > 0$, then for any $\varepsilon < \delta$ for a sufficient small δ implies $\tilde{\mathcal{A}}^T(\tilde{b}) \succ 0$ (is positive definite).

Proof. For the first part this proof is identical to the theorem 2.1.7:

$$\tilde{b} = [(b + \varepsilon \mathcal{A}(V))^T | \varepsilon]^T \quad (2.25)$$

$$\tilde{\mathcal{A}}^T((y_1, \dots, y_m, y_{m+1})^T) = \mathcal{A}^T((y_1, \dots, y_m)^T) + y_{m+1}V \quad (2.26)$$

$$P = e_{m+1} \quad (2.27)$$

And:

$$\tilde{\mathcal{A}}^T(\tilde{b}) = \mathcal{A}^T(b) + \varepsilon(\mathcal{A}^T \circ \mathcal{A}(V) + V) \quad (2.28)$$

Also define:

$$\begin{aligned} f : \mathbb{R} &\rightarrow \mathcal{S}^n, z \mapsto \mathcal{A}^T(b) + zE \\ E &= \mathcal{A}^T \circ \mathcal{A}(V) + V \end{aligned} \quad (2.29)$$

Again continuity follows from theorem 2.2.4 which implies that for small value of $z > 0$ the dimension of the null space of $f(z)$ cannot increase.

Since this time the eigenvalue corresponding to the null space is multiple it is invalid to use theorem 2.2.7. Instead apply theorem 2.2.15: Note that the supgradient of λ_n consists of the convex hull of vector outproducts that span the null space, that is:

$$\begin{aligned} \partial\lambda_{n-k+1}(f(0)) &= \{Z \in \mathcal{S}^n \text{ s. t. } \langle Z, \mathcal{A}^T(b) \rangle = 0, Z \succeq 0\} \\ &= \text{Co}\{qq^T \text{ s. t. } \mathcal{A}^T(b)q = 0\} \end{aligned} \quad (2.30)$$

(This follows from corollary 2.2.23.)

The condition that for each $v \in \mathbb{R}^n$ with $\mathcal{A}^T(b)v = 0$, $v^T \mathcal{A}^T \circ \mathcal{A}(V)v + v^T V v$ implies that for any $Z \in \partial\lambda_{n-k+1}(f(0))$, $\langle E, Z \rangle > 0$, hence E is an ascent direction for λ_n . Because λ_n is the minimum of all eigenvalues it follows that all eigenvalues λ_{n-k+1} up to λ_n ascent in the direction of E (theorem 2.2.4).

Hence also $\lambda_n(f(z)) > 0$ for $z > 0$ small. Thus there exists a $\delta > 0$ such that for all $0 < \varepsilon < \delta$, $f(\varepsilon)$ is of full rank. \square

What's the catch of this perturbation? For a start, in the perturbed case $Pb = \varepsilon$. Hence to solve it by using the eigenvalue optimization equivalence, the problem has to be relaxed a little and perturbed. But because it is expected that $\varepsilon\lambda_{\max}$ is small positive, this will not affect the optimal solution in a significant way, albeit it is a source of error.

However, another catch of the perturbation is that the system becomes more sensitive to numerical errors. The condition number of $\tilde{\mathcal{A}}^T(b)$ will be large because of the small eigenvalue, hence the Cholesky decomposition step is prone to large numerical errors.

Thus for some instances the perturbation can be made large enough for $\tilde{\mathcal{A}}^T(b)$ to be stable, and for some instances the numerical errors will be too large in magnitude to work on the system.

2.1.2 Some examples

As examples of deriving the eigenvalue optimization problem for some SDP the examples from section 1.4 will be cast into an eigenvalue problem.

Max-cut

As explained in section 1.4.1, the dual of the max-cut is given by:

$$\begin{aligned} \min_{y \in \mathbb{R}^{|V|}} \sum_i y_i \\ \text{s. t. } \frac{1}{4}L_G - \sum_i y_i E_{i,i} \preceq 0 \end{aligned} \quad (2.31)$$

(This is the same as equation 1.56.)

To start:

$$\mathcal{A}^T(\hat{b}) = \sum_i E_{i,i} = (|V|)^{-\frac{1}{2}}I \quad (2.32)$$

Which has a trivial Cholesky decomposition: $\mathcal{A}^T(\hat{b}) = LL^T$ with $L = (|V|)^{-\frac{1}{4}}I$. Furthermore, let $m = |V|$, $q_1, \dots, q_m \in \mathbb{R}^m$ with $q_m = \hat{b} \in \mathbb{R}^m$ be an orthogonal basis. Then $\mathcal{A}^T(q_1), \dots, \mathcal{A}^T(q_m) \in \mathcal{S}^n$ are all diagonal matrices with on their diagonals q_1, \dots, q_m . That all makes SDP 2.31 equivalent to:

$$\begin{aligned} \inf_{z \in \mathbb{R}^{|V|-1}} (|V|)^{\frac{1}{2}} \lambda_{\max} \left(\frac{1}{4}L^{-1}L_G L^{-T} - \sum_i z_i L^{-1} \text{diag}(q_i) L^{-T} \right) \\ \inf_{z \in \mathbb{R}^{|V|-1}} (|V|)^{\frac{1}{2}} \lambda_{\max} \left((|V|)^{\frac{1}{2}} \frac{1}{4}L_G - \sum_i z_i (|V|)^{\frac{1}{2}} \text{diag}(q_i) \right) \\ \inf_{z \in \mathbb{R}^{|V|-1}} |V| \lambda_{\max} \left(\frac{1}{4}L_G - \sum_i z_i \text{diag}(q_i) \right) \end{aligned} \quad (2.33)$$

It follows from the fact that $\text{diag}(q_1), \dots, \text{diag}(q_{|V|-1}) \perp \text{diag}(q_m) = (|V|)^{-\frac{1}{2}}I$ that there is no sum possible such that $\sum_i \gamma_i \text{diag}(q_i) \succ 0$. (If there would be, then $q_m^T(\sum_i \gamma_i q_i) > 0$, which is a contradiction because all $q_i \perp q_m$.) Therefore the eigenvalue problem is bounded.

The difficulty of this problem lies in generating a basis q_1, \dots, q_m , which can be computational hard for large m . If memory requirements are such that q_1, \dots, q_m do not have to be optimized for as few non-zeros as possible, then a fast way to compute a basis is by using Lanczos iterations.

Note that if the SDP was sparse, then the eigenvalue problem is also sparse. That will be an important property when solving the eigenvalue problem in chapter 4.

Ideal GMRES polynome

The ideal GMRES polynome of order k for matrix A can be computed using:

$$\begin{aligned} & \min_{t \in \mathbb{R}, y \in \mathbb{R}^k} t \\ & \text{s. t. } \begin{pmatrix} 0 & I \\ I & 0 \end{pmatrix} + \sum_{1 \leq i \leq k} y_i \begin{pmatrix} 0 & A^i \\ A^i & 0 \end{pmatrix} - t \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} \preceq 0 \end{aligned} \quad (2.34)$$

As was explained in section 1.4.3.

Denote $y' = (y^T, t)^T$, then $b = e_{k+1}$, and $\mathcal{A}^T(b) = I$. Hence an orthogonal basis $q_1, \dots, q_{k+1} \in \mathbb{R}^{k+1}$ is $e_1, \dots, e_{k+1} \in \mathbb{R}^{k+1}$. That makes the eigenvalue problem transformation trivial easy:

$$\inf_{z \in \mathbb{R}^k} \lambda_{\max} \left(\begin{pmatrix} 0 & I \\ I & 0 \end{pmatrix} + \sum_{1 \leq i \leq k} z_i \begin{pmatrix} 0 & A^i \\ A^i & 0 \end{pmatrix} \right) \quad (2.35)$$

Phase retrieval

The dual of 1.76 is:

$$\begin{aligned} & - \min_{y \in \mathbb{R}^{|\Omega|}} b^T y \\ & \text{s. t. } -I - \mathcal{A}^T(y) \preceq 0 \end{aligned} \quad (2.36)$$

Or similar for the log det heuristic:

$$\begin{aligned} & - \min_{y \in \mathbb{R}^{|\Omega|}} b^T y \\ & \text{s. t. } -W_k - \mathcal{A}^T(y) \preceq 0 \end{aligned} \quad (2.37)$$

Because $A_\omega = a_\omega a_\omega^*$ and $b_\omega > 0$, it is for certain than $\mathcal{A}^T(b) \succeq 0$. If in addition $\mathcal{A}^T(b) \succ 0$, then theorem 2.1.4 can be applied. Let LL^T be a Cholesky decomposition of $\mathcal{A}^T(b)$, then problem 2.37 is equivalent to:

$$- \inf_{z \in \mathbb{R}^{|\Omega|-1}} \lambda_{\max} (L^{-1} W_k L^{-T} - \sum_j z_j L^{-1} \mathcal{A}^T(q_j) L^{-T}) \quad (2.38)$$

As mentioned at the end of 1.4.4, if the optimal y_k^* for iteration k is known, then a good dual starting point for the next iteration is $z = Q^T y_k^*$. Here are some arguments:

Let $D_k = W_{k+1} - W_k$, such that:

$$- \inf_{z \in \mathbb{R}^{|\Omega|-1}} \lambda_{\max} (L^{-1} (I + \sum_i^k D_i) L^{-T} - \sum_j z_j L^{-1} \mathcal{A}^T(q_j) L^{-T}) \quad (2.39)$$

As the log det heuristic converges it is expected that $D_k \rightarrow 0$ as $k \rightarrow \infty$. Hence the dual programs will converge. Because the term $L^{-1} \mathcal{A}^T(q_j) L^{-T}$ is constant between the programs, and the term $L^{-1}(I + \sum_i^k D_i) L^{-T}$ converges, it is expected that the optimal solution z_k^* also converges.

2.2 Analysis on eigenvalues

In order to keep notation clean, in the remainder of this chapter the subscripts and similarity transformations in equation 2.14 are omitted.

Definition 2.2.1 (Eigenvalue optimization). *An eigenvalue optimization problem is of the following form:*

$$\begin{aligned} \inf_{z \in \mathbb{R}^m} \lambda_{\max}(C - \mathcal{A}^T(z)) \\ \text{s. t. } Pz \geq 0 \end{aligned} \tag{2.40}$$

With $C \in \mathcal{S}^n$, $\mathcal{A}^T : \mathbb{R}^m \rightarrow \mathcal{S}^n$, $z \mapsto \sum_i z_i A_i$ linear operator and $P \in \mathbb{R}^{k \times m}$. Also \geq is per component defined.

In this section some properties on eigenvalues as matrix perturbation will be derived which will be used to derive important properties of eigenvalue optimization problems. It will be shown that the problem is convex, although not everywhere differentiable.

Since $C \in \mathcal{S}^n$ as well as $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$, everything is about symmetric real matrices. Hence the further discussion in this section will be restricted to symmetric real matrices.

First a remark on notation, in this section some properties of eigenvalues as perturbation are examined, i.e. functions of the form $E \mapsto A + E$ with $A, E \in \mathbb{R}^{n \times n}$. In most cases are A and E restricted to \mathcal{S}^n .

A has n eigenvalues $\lambda_i(A)$ which are labeled as:

$$\lambda_1(A) \geq \lambda_2(A) \geq \dots \geq \lambda_n(A) \tag{2.41}$$

If an eigenvalue has multiplicity $t > 1$, then it is counted t times, e.g. if λ_1 has multiplicity 2, then $\lambda_1 = \lambda_2$.

Moreover, if from the context it is clear for a given eigenvalue which matrix is meant, the argument might be omitted.

In the next subsection some theorems concerning continuity and differentiability of eigenvalues are discussed. For the first case it does not matter whether an eigenvalue is simple or non-simple, for the second case such a distinction is necessary. In the second subsection differentiability of non-simple eigenvalues is discussed.

2.2.1 Continuity and differentiability for the simple case

The following characterization due to Courant, Fischer and Weyl of λ_i forms the cornerstone of all the other derivations:

Theorem 2.2.2 (Courant–Fischer–Weyl min–max principle). *Let A be a Hermitian matrix with eigenvalues labeled as described. Then*

$$\lambda_i = \max_{\substack{\mathcal{X} \subseteq \mathcal{L}^n \\ \dim(\mathcal{X})=i}} \min_{\substack{x \in \mathcal{X} \\ \|x\|_2=1}} x^T A x \quad (2.42)$$

And:

$$\lambda_i = \min_{\substack{\mathcal{X} \subseteq \mathcal{L}^n \\ \dim(\mathcal{X})=n-i+1}} \max_{\substack{x \in \mathcal{X} \\ \|x\|_2=1}} x^T A x \quad (2.43)$$

Proof. The first equality: Because A is Hermite it is diagonalizable with eigenvectors u_k corresponding to λ_k . Since $\dim(\mathcal{X}) = i$, there is an intersection between \mathcal{X} and $\text{span}\{u_i, \dots, u_n\}$. Therefore:

$$\lambda_i \leq \max_{\substack{\mathcal{X} \subseteq \mathcal{L}^n \\ \dim(\mathcal{X})=i}} \min_{\substack{x \in \mathcal{X} \\ \|x\|_2=1}} x^T A x \quad (2.44)$$

By choosing $\mathcal{X} = \text{span}\{u_1, \dots, u_i\}$ the $\lambda_i \geq \max \min x^T A x$ inequality is observed. Hence $\lambda_i = \max \min x^T A x$.

The second equality follows by replacing A by $-A$. \square

Corollary 2.2.3 (Rayleigh quotient).

$$\lambda_1 = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^T A x = \sup_{q \in \mathbb{R}^n, \|q\|_2=1} \langle A, qq^T \rangle = \sup_{\substack{Z \in \mathcal{S}_{\geq 0}^n \\ \langle Z, I \rangle = 1}} \langle A, Z \rangle \quad (2.45)$$

And:

$$\lambda_n = \min_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^T A x = \inf_{q \in \mathbb{R}^n, \|q\|_2=1} \langle A, qq^T \rangle = \inf_{\substack{Z \in \mathcal{S}_{\geq 0}^n \\ \langle Z, I \rangle = 1}} \langle A, Z \rangle \quad (2.46)$$

Proof. Proof for the first equation is shown, the second equation follows in a similar manner.

The first equality is a direct consequence of theorem 2.2.2. The second equality follows: $Z \in \mathcal{S}_{\geq 0}^n : \langle Z, I \rangle = 1$ implies that Z has an eigendecomposition $Z = \sum_i \alpha_i q_i q_i^T$ with $0 < \alpha_i \leq 1$, $\sum_i \alpha_i = 1$ and $q_i^T q_j = \delta_{ij}$ (that is 0 if $i \neq j$ and 1 otherwise). Then:

$$\begin{aligned} \langle A, Z \rangle &= \text{Tr}(A(\sum_i \alpha_i q_i q_i^T)) \\ &= \sum_i \alpha_i \text{Tr}(A q_i q_i^T) \\ &= \sum_i \alpha_i \text{Tr}(q_i^T A q_i) \end{aligned} \quad (2.47)$$

Which is maximized if all q_i are eigenvectors corresponding to λ_{\max} . Hence the supremum of $\langle A, Z \rangle = \sum_i \alpha_i \lambda_{\max} = \lambda_{\max}$. \square

It is easy by using the min-max principle to show continuity of the eigenvalues of $A(E) = A + E$ and putting an upper bound on the variation of $\lambda_i(A(E))$ as function of E :

Corollary 2.2.4. *The perturbation of $\lambda_i(A + E)$ is bounded by:*

$$\lambda_{\min}(E) \leq \lambda_i(A + E) - \lambda_i(A) \leq \lambda_{\max}(E) \quad (2.48)$$

Which is equivalent to:

$$\lambda_i(A + E) \in [\lambda_i(A) + \lambda_{\min}(E), \lambda_i(A) + \lambda_{\max}(E)] \quad (2.49)$$

Which implies $f(E) = \lambda_i(A + E)$ is continuous in E .

Proof. First $\lambda_i(A + E) - \lambda_i(A) \leq \lambda_{\max}(E)$:

$$\begin{aligned} \lambda_i(A + E) - \lambda_i(A) &= \max_{\dim(\mathcal{X})=i} \min_{\|x\|_2=1} x^\top (A + E)x - \max_{\dim(\mathcal{X})=i} \min_{\|x\|_2=1} x^\top Ax \\ &\leq \max_{\dim(\mathcal{X})=i} \min_{\|x\|_2=1} x^\top Ax + \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^\top Ex - \max_{\dim(\mathcal{X})=i} \min_{\|x\|_2=1} x^\top Ax \\ &= \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^\top Ex \\ &= \lambda_{\max}(E) \end{aligned} \quad (2.50)$$

Second $\lambda_i(A + E) - \lambda_i \geq \lambda_{\min}(E)$:

$$\begin{aligned} \lambda_i(A + E) - \lambda_i(E) &= \min_{\dim(\mathcal{X})=n-i+1} \max_{\|x\|_2=1} x^\top (A + E)x - \min_{\dim(\mathcal{X})=n-i+1} \max_{\|x\|_2=1} x^\top Ax \\ &\geq \min_{\dim(\mathcal{X})=n-i+1} \max_{\|x\|_2=1} x^\top Ax + \min_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^\top Ex - \min_{\dim(\mathcal{X})=n-i+1} \min_{\|x\|_2=1} x^\top Ax \\ &= \min_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^\top Ex \\ &= \lambda_{\min}(E) \end{aligned} \quad (2.51)$$

\square

This gives direct a way to identify whether problem 2.40 is unbounded:

Corollary 2.2.5. *Eigenvalue optimization problem 2.40 is unbounded if and only if there exists $z \in \mathbb{R}^m$ such that $Pz = 0$ and $\mathcal{A}^\top(z) \succeq 0$ (positive definite).*

Note that the power of corollary 2.2.4 is that it does not require that λ_i is simple. In most cases it is not so easy to say something when λ_i is non-simple.

When $\lambda_i(A)$ is simple, then the derivative exists in closed form (based on [18]):

Lemma 2.2.6. *Let $\lambda_i(A) \in \mathbb{R}$ simple with corresponding eigenvector $u_i(A) \in \mathbb{R}^n$. Then $\lambda_i(A + E)$ and $u_i(A)$ are ∞ times differentiable at A .*

Proof. Define the function:

$$f : \mathbb{R}^n \times \mathbb{R} \times \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n+1}, (u, \lambda; A) \mapsto \begin{pmatrix} (A - \lambda I)u \\ u^T u - 1 \end{pmatrix} \quad (2.52)$$

Then $f(u_i, \lambda_i) = 0$. Note that f is ∞ times differentiable. The Jacobian of f to u and λ is given by:

$$D_{u, \lambda} f(u, \lambda; A) = \begin{pmatrix} A - \lambda I & u \\ u^T & 0 \end{pmatrix} \quad (2.53)$$

Which is non-singular by lemma 1.2.11 (Schur complement). Hence by the implicit function theorem there exists an open interval in $N(A) \subset \mathbb{R}^{n \times n}$ containing A . Moreover, there exists functions $u_i : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$ and $\lambda_i : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ which are continuous and ∞ times differentiable on $N(A)$. \square

Theorem 2.2.7. *Let $\lambda_i(A)$ simple with corresponding normalized eigenvector $u_i(A)$. Then the derivative of $\lambda_i(A + E)$ to E is:*

$$\frac{\mathbf{d} \lambda_i(A + E)}{\mathbf{d} E} = u_i^T(A) E u_i(A) \quad (2.54)$$

Proof. By using the identity $u_i^T(A)(A - \lambda_i(A)I)u_i(A) = 0$ and the product rule:

$$\begin{aligned} & \frac{\mathbf{d}}{\mathbf{d} E} (u_i^T(A + E)(A + E - \lambda_i(A + E)I)u_i(A + E)) \\ &= \frac{\mathbf{d} u_i^T(A + E)}{\mathbf{d} E} (A - \lambda_i(A))u_i(A) + u_i^T(A) \left(E - \frac{\mathbf{d} \lambda_i(A + E)}{\mathbf{d} E} \right) E u_i(A) \\ & \quad + u_i^T(A)(A - \lambda_i(A)) \frac{\mathbf{d} u_i^T(A + E)}{\mathbf{d} E} \\ &= u_i^T(A) \left(E - \frac{\mathbf{d} \lambda_i(A + E)}{\mathbf{d} E} \right) u_i(A) = 0 \end{aligned} \quad (2.55)$$

The last equation stems from that $(A - \lambda_i(A))u_i(A) = u_i^T(A)(A - \lambda_i(A)) = 0$. The last line can be simplified:

$$\frac{\mathbf{d} \lambda_i(A + E)}{\mathbf{d} E} = u_i^T(A) E u_i(A) \quad (2.56)$$

Because $u_i^T(A)u_i(A) = I$. \square

Corollary 2.2.8. *If for a given $y \in \mathbb{R}^m$, $\lambda_i(C - \mathcal{A}^T(y))$ is simple with eigenvector u_i , then the gradient of*

$$f : \mathbb{R}^m \rightarrow \mathbb{R}, y \rightarrow \lambda_i(C - \mathcal{A}^T(y)) \quad (2.57)$$

is given by:

$$\nabla f(y) = -(u_i^T A_1 u_i, \dots, u_i^T A_m u_i)^T \quad (2.58)$$

Also the second derivative is computable:

Lemma 2.2.9. *Let $\lambda_i(A)$ simple with corresponding eigenvector $u_i(A)$. Then the first derivative of $u_i(A + E)$ to E is:*

$$\frac{\mathbf{d} u_i(A + E)}{\mathbf{d} E} = -(A - \lambda_i I)^\dagger E u_i \quad (2.59)$$

(For ease of notation: $\lambda_j = \lambda_j(A)$ and $u_j = u_j(A)$ for all $j = 1, \dots, n$. Furthermore $(A - \lambda_i I)^\dagger$ denotes the Moore-Penrose pseudoinverse.)

Proof. Starting with differentiating the following equation:

$$\frac{\mathbf{d}}{\mathbf{d} E} ((A + E - \lambda_i I)u) = (E - \frac{\mathbf{d} \lambda_i}{\mathbf{d} E} I)u + (A - \lambda_i I) \frac{\mathbf{d} u}{\mathbf{d} E} = 0 \quad (2.60)$$

Observe that $(A - \lambda_i I)^\dagger u_i = 0$ because the kernel of $A - \lambda_i I$ is exactly spanned by u_i (by assumption that λ_i is simple). Also $(A - \lambda_i I)^\dagger \frac{\mathbf{d} u_i}{\mathbf{d} E} = \frac{\mathbf{d} u_i}{\mathbf{d} E}$ because $\frac{\mathbf{d} u_i}{\mathbf{d} E} \perp u_i$ and has therefore no components in the null space of $A - \lambda_i I$. (To see that $\frac{\mathbf{d} u_i}{\mathbf{d} E} \perp u_i$, note that $u_i^T u_i = 1$ thus $\frac{\mathbf{d} u_i^T}{\mathbf{d} E} u_i = 0$.)

Hence multiply the equation left with $-(A - \lambda_i I)^\dagger$ to obtain equation 2.59. \square

Theorem 2.2.10. *Let $\lambda_i(A)$ simple with corresponding eigenvector $u_i(A)$. Then the second derivative of $\lambda_i(A + E)$ to $E \in \mathcal{S}^n$ is:*

$$\begin{aligned} \frac{\mathbf{d}^2 \lambda_i(A + E_1 + E_2)}{\mathbf{d} E_1 \mathbf{d} E_2} &= -2u_i^T E_1 (A - \lambda_i I)^\dagger E_2 u_i \\ &= 2 \sum_{r \neq i} \frac{u_i^T E_1 u_r u_r^T E_2 u_i}{\lambda_i - \lambda_r} \end{aligned} \quad (2.61)$$

Proof. Start with equation 2.54 and differentiate it by using the chain-rule and

the lemma above:

$$\begin{aligned}
\frac{\mathbf{d}}{\mathbf{d} E_2} \frac{\mathbf{d} \lambda_i(A + E_1 + E_2)}{\mathbf{d} E_1} &= u_i^\top E_2 u_i \\
&= \frac{\mathbf{d} u_i^\top}{\mathbf{d} E_1} E_2 u_i + u_i^\top E_2 \frac{\mathbf{d} u_i}{\mathbf{d} E_1} \\
&= -u_i^\top E_1 (A - \lambda_i I)^\dagger E_2 u_i - u_i^\top E_2 (A - \lambda_i I)^\dagger E_1 u_i \\
&= -2u_i^\top E_1 (A - \lambda_i I)^\dagger E_2 u_i \\
&= 2 \sum_{r \neq i} \frac{u_i^\top E_1 u_r u_r^\top E_2 u_i}{\lambda_i - \lambda_r}
\end{aligned} \tag{2.62}$$

In the last equations it is used that $(A - \lambda_i I)^\dagger$, E_1 and E_2 are all symmetric. \square

Corollary 2.2.11. *If for a given $y \in \mathbb{R}^m$, $\lambda_i(C - \mathcal{A}^\top(y))$ is simple with eigenvector u_i , then the Hessian of*

$$f : \mathbb{R}^m \rightarrow \mathbb{R}, y \rightarrow \lambda_i(C - \mathcal{A}^\top(y)) \tag{2.63}$$

is given by:

$$(Hf(y))_{k,l} = -2 \sum_{j \neq i} \frac{u_i^\top A_k u_j u_j^\top A_l u_i}{\lambda_i - \lambda_j} \tag{2.64}$$

If λ_i is non-simple, then there are results that all eigenvalues can be labeled (counted to multiplicity) λ_i such that all λ_i curves are smooth [1]. The limitation is that the eigenvectors are not necessarily smooth. Hence, λ_i as function of perturbation E is not always in closed form. This result is stated here for completeness without proof. The reason is that the proof is a deeper result based on smoothness of the roots of a polynome under perturbation. Hence to repeat it here would require another few pages describing functional analysis on polynomes.

Theorem 2.2.12. *Let $A \in \mathcal{S}^n$ and $E \in \mathcal{S}^n$, and let $\mathcal{A}(t) = A + tE$. It is possible to label all eigenvalues of $\mathcal{A}(t)$ as $\lambda_i(t)$ for $i = 1, \dots, n$, (counted to multiplicity), such that each $\lambda_i(t)$ is a smooth function.*

Proof. See [1]. \square

2.2.2 The non-simple case

For the non simple case there is no reason to expect that the derivative of $\lambda_i(A + E)$ exists. However the subgradient of the maximum eigenvalue, and the supgradient of the minimal eigenvalue does exist. Those are sets which will be quite useful:

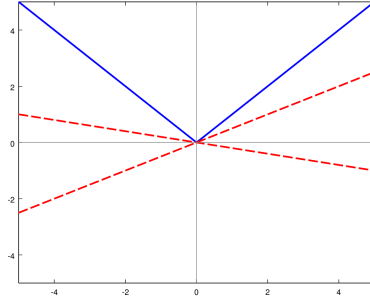


Figure 2.1: $f(x) = |x|$ is plotted with solid blue, two possible subgradients at $x = 0$ are plotted with dashed red

Definition 2.2.13 (Sub- and supgradient). *The subgradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $x_0 \in \mathbb{R}^n$ is defined as the largest closed convex set $\partial^- f(x_0)$ such that for all $v \in \partial^- f(x_0)$ holds:*

$$f(x) - f(x_0) \geq \langle v, x - x_0 \rangle \quad (\forall x \in \mathbb{R}^n) \quad (2.65)$$

Moreover, if for all $x \in \mathbb{R}^n$, $\partial^- f(x)$ is non-empty, then it is said that f is subdifferential.

Analogous at $x_0 \in \mathbb{R}^n$ is the supgradient of $g : \mathbb{R}^n \rightarrow \mathbb{R}$ the largest closed convex set $\partial^+ g(x_0)$ such that for all $v \in \partial^+ g(x_0)$ holds:

$$g(x) - g(x_0) \leq \langle v, x - x_0 \rangle \quad (\forall x \in \mathbb{R}^n) \quad (2.66)$$

Similar if for all $x \in \mathbb{R}^n$, $\partial^+ g(x)$ exists, then g is supdifferential.

If from the context it is clear whether a sub- or supgradient of f is meant, then the set will simply be denoted by ∂f , omitting the + or -.

(These definitions are based on the generalized Clark gradient defined in [8] and [7].)

The power of sub/supgradients is that they can define a descent / ascent direction for non differentiable functions. The following two are worded for the subgradient, but it should be clear that there are equivalent statements for concave supdifferential functions.

To get a bit more familiar with subgradients, consider the following example:

Example 2.2.14. *Let $f : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto |x|$, then for $x < 0$, $\partial f(x) = \{-1\}$, for $x > 0$, $\partial f(x) = \{1\}$ and for $x = 0$, $\partial f(x) = [-1, 1]$. Two possible subgradients are shown in figure 2.1.*

Lemma 2.2.15. *Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be convex subdifferential. If for some $x \in \mathbb{R}^m$ there exists $d \in \mathbb{R}^m$ such that for all $y \in \partial f(x)$, $\langle y, d \rangle < 0$, then, and only then is d a descent direction at x .*

Proof. The first implication: Assume the contradiction that such a $d \in \mathbb{R}^m$ exists but is not a descent direction. Then because f is convex either x is the global minimum for f or $-d$ is not an ascent direction. In the first case, $0 \in \partial f(x)$ and $\langle 0, d \rangle = 0$, hence d does not fit the constraints. In the second case $d \in \partial f(x)$ and $\langle d, d \rangle > 0$, thus also in this case there's a contradiction.

The second implication: Assume d is a descent direction, but there exists $y \in \partial f(x)$ such that $\langle y, d \rangle \geq 0$. Since $y \in \partial f(x)$: $f(x + td) \geq f(x) + \langle y, d \rangle$ for all $t \in \mathbb{R}$, hence d is not a descent direction, which is a contradiction. \square

Corollary 2.2.16. *Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be convex subdifferential. If for $x \in \mathbb{R}^m$, $0 \in \partial f(x)$, then and only then is x the global minimum.*

Proof. This is the previous lemma but then in other words. \square

Furthermore, subgradients are powerful tools to compute directional derivatives when the gradient does not exist (based on [11]):

Theorem 2.2.17. *Assuming that $\partial f(x)$ is bounded in a neighborhood around x : The directional derivative of f can be expressed in terms of the subgradient of f :*

$$\nabla_v f(x) = \lim_{h \downarrow 0} \frac{f(x + hv) - f(x)}{h} = \max_{g \in \partial f(x)} g^T v \quad (2.67)$$

Proof. Define $(h)_k$ as a decreasing sequence in \mathbb{R} such that $h \downarrow 0$, and define $x_k = x + h_k v$. Then:

$$\nabla_v f(x) = \lim_{k \rightarrow \infty} \frac{f(x_k) - f(x)}{h_k} \quad (2.68)$$

By the definition of the subgradient:

$$\begin{aligned} f(x_k) - f(x) &\geq h_k v^T g \quad (\forall g \in \partial f(x)) \\ f(x) - f(x_k) &\geq -h_k v^T g_k \quad (\forall g_k \in \partial f(x_k)) \end{aligned} \quad (2.69)$$

Therefore:

$$v^T g_k \geq \nabla_v f(x) \geq \max_{g \in \partial f(x)} g^T v \quad (2.70)$$

Because $\partial f(y)$ is closed and bounded for all y near x , and $\partial f(x_k)$ converges to a subset of $\partial f(x)$. It follows that there exists a sequence $g_k \in \partial f(x_k)$ that converges to $\arg \max_{g \in \partial f(x)} g^T v$. \square

Before computing the subgradient / supgradient of the maximum / minimum eigenvalue it should be marked that they are convex / concave functions. Indeed if they are not convex / concave, then subgradient / supgradient does not exist.

Lemma 2.2.18. *The function:*

$$\lambda_{\max}(A) : \mathcal{S}^n \rightarrow \mathbb{R}, A \rightarrow \lambda_{\max}(A) \quad (2.71)$$

is convex. Similar the function:

$$\lambda_{\min}(t) : \mathbb{R}^m \rightarrow \mathbb{R}, t \rightarrow \lambda_{\min}(A_0 + t_1A_1 + \cdots + t_mA_m) \quad (2.72)$$

is concave.

Proof. For the first statement a proof is given. The second statement follows in a similar manner.

The inequality for convexity follows direct from application of the Rayleigh quotient: For all $t \in [0, 1]$ and $A, B \in \mathcal{S}^n$:

$$\begin{aligned} t\lambda_{\max}(A) + (1-t)\lambda_{\max}(B) &= t \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^T Ax + (1-t) \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^T Bx \\ &\geq \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} x^T (tA + (1-t)B)x \\ &= \lambda_{\max}(tA + (1-t)B) \end{aligned} \quad (2.73)$$

□

Theorem 2.2.19 (Dubovitskii–Milyutin). *Let f be a convex function defined on \mathbb{R}^n as:*

$$f(x) = \sup_{\alpha \in \mathcal{A}} f_{\alpha}(x) \quad (2.74)$$

Where \mathcal{A} is some index-set and each f_{α} is convex and subdifferentiable. Denote the active set $I_x = \{\alpha \in \mathcal{A} \text{ s. t. } f_{\alpha}(x) = f(x)\}$.

If for all $d \in \mathbb{R}^n$:

$$\nabla_d f(x) = \sup\{\nabla_d f_{\alpha}(x) \text{ s. t. } \alpha \in I_x\} \quad (2.75)$$

holds, then:

$$\partial f(x) = \text{Co}(\cup\{\partial f_{\alpha}(x) \text{ s. t. } f(x) = f_{\alpha}(x)\}) \quad (2.76)$$

That is $\partial f(x)$ is the convex hull of the subgradients of the active functions.

Proof. Observe that $\text{Co}(\cup\{\partial f_{\alpha}(x) \text{ s. t. } f(x) = f_{\alpha}(x)\}) \subset \partial f(x)$ follows: each $g \in \partial f_{i \in I}(x)$ is in $\partial f(x)$.

The other inclusion: let $g \in \partial f(x)$. If g is not in $\text{Co}(\partial f_{i \in I_x}(x))$ then there exists a strict separating plane $d \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$ such that:

$$g^T d > \alpha > \max_{i \in I_x} \sup_{\epsilon \in \partial f_i(x)} \epsilon^T d \quad (2.77)$$

Furthermore for all $i \in I_x$:

$$\sup_{\epsilon \in \partial f_i(x)} \epsilon^T d = \lim_{t \downarrow 0} \frac{f_i(x+td) - f_i(x)}{t} \quad (2.78)$$

Hence: $\max_{i \in I_x} \lim_{t \downarrow 0} \frac{f_i(x+td) - f_i(x)}{t} < g^T d$. Because J is non-empty and finite around x , the max and lim can be interchanged (define the limit in terms of decreasing sequence $t_k \geq 0$, then there is a sub-sequence such that the statement holds). That gives:

$$\max_{i \in I_x} \sup_{\epsilon \in \partial f_i(x)} = \lim_{t \downarrow 0} \max_{i \in I_x} \frac{f_i(x+td) - f_i(x)}{t} = \lim_{t \downarrow 0} \frac{f(x+td) - f(x)}{t} \quad (2.79)$$

However, note that $g \in \partial f(x)$ implies that $f(x+td) - f(x) \geq tg^T d$ for each $t \geq 0$, hence there is a contradiction. Therefore by contradiction $g \in \partial f_{i \in I}(x)$. \square

Corollary 2.2.20. *Let f be a concave function defined on \mathbb{R}^n as:*

$$f(x) = \inf_{\alpha \in \mathcal{A}} f_\alpha(x) \quad (2.80)$$

Where each f_α is concave and supdifferentiable. Suppose the supremum is for each x attained. In that case:

$$\partial^+ f(x) = \text{Co}(\cup \{ \partial^+ f_\alpha(x) \text{ s. t. } f(x) = f_\alpha(x) \}) \quad (2.81)$$

That is $\partial^+ f(x)$ is the convex hull of the supgradients of the active functions.

As a simple example of theorem 2.2.19:

Example 2.2.21. *Let $f_1 : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto x^2$ and $f_2 : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto (x-2)^2$ and $f(x) = \max_i f_i(x)$. The for $x < 1$, $\partial f(x) = 2(x-2)$, for $x > 1$, $\partial f(x) = 2x$, and for $x = 1$, $\partial f(x) = [-2, 2]$. This example is shown in figure 2.2.*

Another example of theorem 2.2.19, which shows that there should be a neighborhood at which the supremum is attained by a finite set of functions:

Example 2.2.22. *Let for $\zeta \in \mathbb{R}$:*

$$\begin{aligned} f_\zeta(x) &= (2\zeta + \frac{\zeta}{|\zeta|})x - \zeta^2 \quad \text{for } \zeta \neq 0 \\ f_0(x) &= 0 \quad \text{for } \zeta = 0 \end{aligned} \quad (2.82)$$

Then it can be checked that:

$$\begin{aligned} f(x) &= \max_{\zeta \in \mathbb{R}} f_\zeta(x) \\ &= f_x(x) \\ &= x^2 + |x| \end{aligned} \quad (2.83)$$

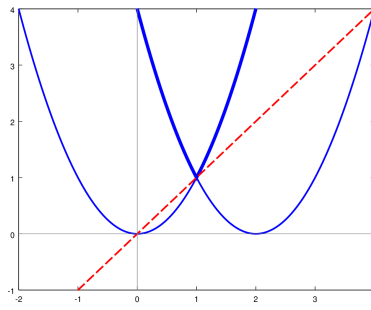


Figure 2.2: $f_1(x) = x^2$ and $f_2(x) = (x - 2)^2$ are plotted with solid blue, $f(x) = \max(f_1(x), f_2(x))$ is plotted with thick blue, and a possible subgradients at $x = 1$ is plotted with dashed red

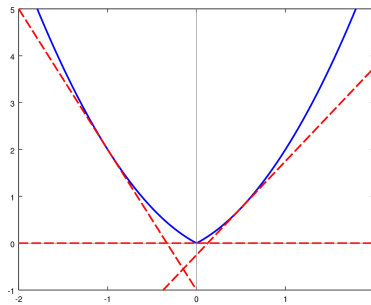


Figure 2.3: $f(x) = \max_{z \in \mathbb{R}} f_z(x)$ is plotted with thick blue, and in dashed red are plotted f_{-1} , f_0 , and $f_{0.5}$

The function $f(x)$ is shown in figure 2.3.

Clearly $f(x)$ is convex, the maximum is attained at each point $x = \zeta$ by exactly one function in the set $\{f_\zeta\}$.

Moreover, the subgradient at $x = 0$ is: $\partial f(0) = [-1, 1]$, but the gradient of the active function at $x = 0$ is $\partial f_0(0) = 0$.

Note however that for each x there exists no neighborhood $\mathcal{N}(x)$ for which the collection $\{f_y | y \in \mathcal{N}(x)\}$ is finite. Hence theorem 2.2.19 was not applicable.

The subgradient of the maximum eigenvalue follows:

Corollary 2.2.23. *The subgradient of $\lambda_{\max}(A)$ is:*

$$\begin{aligned} \partial^- \lambda_{\max}(A) &= \{Z \in \mathcal{S}_{\geq 0}^n \text{ s. t. } \langle A, Z \rangle = \lambda_{\max}, \langle Z, I \rangle = 1\} \\ &= \text{Co}\{qq^T \text{ s. t. } q \in \mathbb{R}^n, Aq = \lambda_{\max}q, \|q\|_2 = 1\} \end{aligned} \quad (2.84)$$

Similar the supgradient of λ_{\min} is:

$$\begin{aligned} \partial^+ \lambda_{\min}(A) &= \{Z \in \mathcal{S}_{\geq 0}^n \text{ s. t. } \langle A, Z \rangle = \lambda_{\min}, \langle Z, I \rangle = 1\} \\ &= \text{Co}\{qq^T \text{ s. t. } q \in \mathbb{R}^n, Aq = \lambda_{\min}q, \|q\|_2 = 1\} \end{aligned} \quad (2.85)$$

Proof. Notice that by the Rayleigh quotient $\lambda_{\max}(A) = \sup_{Z \in \mathcal{S}_{\geq 0}^n, \langle Z, I \rangle = 1} \langle A, Z \rangle$ and $\lambda_{\min}(A) = \inf_{Z \in \mathcal{S}_{\geq 0}^n, \langle Z, I \rangle = 1} \langle A, Z \rangle$ and also $\partial \langle \cdot, Z \rangle = Z$.

The corollary follows as a immediate consequence of theorem 2.2.19.

The second identity follows because the set $\{Z \in \mathcal{S}_{\geq 0}^n \text{ s. t. } \langle A, Z \rangle = \lambda_{\max}, \langle Z, I \rangle = 1\}$ can be described as the convex hull of $\{qq^T \text{ s. t. } q \in \mathbb{R}^n, Aq = \lambda_{\max}q, \|q\|_2 = 1\}$, which finishes the proof. \square

Corollary 2.2.24. *Define:*

$$f : \mathbb{R}^m \rightarrow \mathbb{R}, y \mapsto \lambda_{\max}(C - \mathcal{A}^T(y)) \quad (2.86)$$

Then the subgradient of $f(y)$ is given by:

$$\begin{aligned} \partial f(y) &= \{-\mathcal{A}(Z) \text{ s. t. } Z \in \mathcal{S}_{\geq 0}^n, \langle Z, C - \mathcal{A}^T(y) \rangle = f(y), \langle Z, I \rangle = 1\} \\ &= \text{Co}\{-\mathcal{A}(qq^T) \text{ s. t. } q \in \mathbb{R}^n, (C - \mathcal{A}^T(y))q = f(y)q, \|q\|_2 = 1\} \end{aligned} \quad (2.87)$$

Furthermore the directional derivative is given by:

$$\begin{aligned} \nabla_d f(y) &= \max_{v \in \partial f(y)} v^T d \\ &= - \min_{\substack{q \in \mathbb{R}^n \\ (C - \mathcal{A}^T(y))q = f(y)q \\ \|q\|_2 = 1}} \langle \mathcal{A}^T(d), qq^T \rangle \end{aligned} \quad (2.88)$$

Proof. The function f can be written as the composition $f = g \circ h(y)$ with $h : \mathbb{R}^m \rightarrow \mathcal{S}^n$, $y \mapsto C - \mathcal{A}^T(y)$ and $g : \mathcal{S}^n \rightarrow \mathbb{R}$, $X \mapsto \lambda_{\max}(X)$. Then by using the definition of the subgradient:

$$\begin{aligned}
f(y) - f(y_0) &= g(h(y)) - g(h(y_0)) \geq \langle W, h(y) - h(y_0) \rangle \quad (\forall W \in \partial g(h(y_0))) \\
&= \langle W, -\mathcal{A}^T(y - y_0) \rangle \\
&= \langle -\mathcal{A}(W), y - y_0 \rangle
\end{aligned} \tag{2.89}$$

Hence the elements of $\partial f(y) = \{-\mathcal{A}(W) \text{ s. t. } W \in \partial g(h(y))\}$, which equals by corollary 2.2.23 equation 2.87.

The second equation follows easily:

$$\begin{aligned}
\nabla_d f(y) &= \max_{v \in \partial f(y)} v^T d \\
&= \max_{W \in \partial g(h(y))} \langle -\mathcal{A}(W), d \rangle \\
&= - \min_{W \in \partial g(h(y))} \langle W, \mathcal{A}^T(d) \rangle \\
&= - \min_{\substack{q \in \mathbb{R}^n \\ (C - \mathcal{A}^T(y))q = f(y)q \\ \|q\|_2 = 1}} \langle \mathcal{A}^T(d), qq^T \rangle
\end{aligned} \tag{2.90}$$

□

Chapter 3

Sensitive Pole Algorithm

There are different approaches to finding the intersecting point of two eigenvalues. One approach is through a discrete line search such as a bisection method, however that is computationally expensive. Another approach is to calculate a full eigendecomposition of the initial matrix and use linear projections to estimate where the crossingpoint is. That is also computationally expensive, but for small matrices less expensive than bisection. However, this approach is unable to take non-simple eigenvalues in consideration because the gradient is not well defined.

The Sensitive Pole Algorithm is another approach based on linearization of the eigenvalues which takes in consideration eigenvalues with a large (sub)-gradient. It only computes those eigenvalues which are likely to intersect the maximum eigenvalue and thus is faster than a full eigendecomposition.

This algorithm will form an important part of a novel algorithm to solve eigenvalue optimization problems in chapter 4.

The Sensitive Pole Algorithm is closely related to the Dominant Pole Algorithm and the Rayleigh Quotient Iteration. The former was actually developed through modifying the Dominant Pole Algorithm, hence it is instructive to explain first how the Rayleigh Quotient Iteration and the Dominant Pole Algorithm works. In the first section Rayleigh Quotient Iteration is explained, second the Dominant Pole Algorithm will be explained. The third section will modify the Dominant Pole Algorithm to a subspace method for better control of convergence and faster convergence. In the fourth section the Sensitive Pole Algorithm will be explained. Finally, in the fifth section will be explained on how to use the Sensitive Pole Algorithm to compute the largest eigenvalue of a matrix as function of perturbation.

Most of this section is based on the work of Rommes and Sleijpen [28]. Rommes and Sleijpen consider the computation of simple dominant poles of a transfer function as application to model order reduction.

The Dominant Pole Algorithm as described in [28] works also on asymmetric matrices with generalized eigenvalues of matrix pair (A, E) . Instead for application in chapter 4, it is sufficient to consider symmetric matrices and restrict

to the classical eigenvalues of A .

3.1 Rayleigh Quotient Iteration

Rayleigh Quotient Iteration (RQI) allows fast computation of eigenvalues with corresponding eigenvectors if an eigenvalue and vector are roughly known. For instance near the end of other iterative eigenvalue solvers which are slower, a few RQI steps can drastically increase the precision. That can be useful for example when high precision is necessary for deflating a subspace.

Let $A \in \mathcal{S}^n$, the objective of RQI is to compute an eigenvalue with corresponding eigenvector roughly in direction $x_0 \in \mathbb{R}^n$. RQI works through repeatedly calculating the Rayleigh quotient (theorem 2.2.3) and using the result as the new shift:

Algorithm 3.1.1 Rayleigh Quotient Iteration

Require:

$A \in \mathcal{S}^n$: matrix
 $x_0 \in \mathbb{R}^n$: initial normalized vector
 $\epsilon \in \mathbb{R}_{>0}$: tolerance

Ensure:

$\lambda \in \mathbb{R}$: eigenvalue of A
 $x \in \mathbb{R}^n$: corresponding eigenvector

```
1:  $k \leftarrow 0, s_0 \leftarrow x_0^T A x_0$ 
2: while  $\|Ax_k - s_k x_k\| \geq \epsilon$  do
3:   Solve  $(A - s_k I)x_{k+1} = x_k$  for  $x_{k+1} \in \mathbb{R}^n$ 
4:   Normalize  $x_{k+1}$ :  $x_{k+1} \leftarrow \frac{x_{k+1}}{\|x_{k+1}\|}$ 
5:   Take the Rayleigh Quotient:  $s_{k+1} \leftarrow x_{k+1}^T A x_{k+1}$ 
6:    $k \leftarrow k + 1$ 
7: end while
8:  $\lambda \leftarrow s_k, x \leftarrow x_k$ 
```

It can be proven that if the Rayleigh Quotient Iteration converges, then it converges cubically. This proof will not be repeated here but can be found in [20] and [23].

3.2 Dominant Pole Algorithm

The Dominant Pole Algorithm can be understood as Rayleigh Quotient Iteration with fixed right hand to steer convergence to eigenvalues with corresponding eigenvector components in the initial vector:

Algorithm 3.2.1 Dominant Pole Algorithm

Require:

$A \in \mathcal{S}^n$: matrix

$b \in \mathbb{R}^n$: residue vector
 $s_0 \in \mathbb{R}$: initial pole estimate
 $\epsilon \in \mathbb{R}_{>0}$: tolerance

Ensure:

$\lambda \in \mathbb{R}$: dominant pole (eigenvalue of A)
 $x \in \mathbb{R}^n$: corresponding eigenvector

```

1:  $k \leftarrow 0$ 
2: while  $\|Ax_k - s_k x_k\| \geq \epsilon$  do
3:   Solve  $(A - s_k I)x_{k+1} = b$  for  $x_{k+1} \in \mathbb{R}^n$ 
4:   Normalize  $x_{k+1}$ :  $x_{k+1} \leftarrow \frac{x_{k+1}}{\|x_{k+1}\|}$ 
5:   Take the Rayleigh Quotient:  $s_{k+1} \leftarrow x_{k+1}^T A x_{k+1}$ 
6:    $k \leftarrow k + 1$ 
7: end while
8: return  $\lambda \leftarrow s_k, x \leftarrow x_k$ 
  
```

Observe the similarity with algorithm 3.1.1, only the third line is different. The effect is that Dominant Pole Algorithm can be steered towards converging to eigenvalues with large eigenvector components in $b \in \mathbb{R}^n$.

The Dominant Pole Algorithm does not converge cubically, as the Rayleigh Quotient Iteration, but quadratically:

Lemma 3.2.2. *Consider algorithm 3.2.1. Let $\lambda \in \mathbb{R}$ be an eigenvalue of A to which the algorithm converges, then $s_k \rightarrow \lambda$ quadratically.*

Proof. It can be shown that s_k follows a Newton scheme:

Define $H : \mathbb{R} \rightarrow \mathbb{R}, s \mapsto b^T (A - sI)^\dagger b$. Note that $\lim_{s \rightarrow \lambda} |H(s)| = \infty$. Also, let $G(s) = H(s)^{-1}$, then the roots of $G(s)$ are exactly the eigenvalues of A for which there exists a vector in the eigenspace $q \in \mathbb{R}^n$ such that $b^T q \neq 0$. Next it is shown that the Dominant Pole Algorithm is a Newton Method used to compute the roots of $G(s)$, which implies that $s \rightarrow \lambda$ quadratically:

The derivative of $G(s)$ to s is given by:

$$G'(s) = -\frac{H'(s)}{H^2(s)} \quad (3.1)$$

And the derivative of $H(s)$ to s is given by:

$$H'(s) = b^T (A - sI)^\dagger (A - sI)^\dagger b \quad (3.2)$$

(This can easily be checked by using the identity $(A - sI)^\dagger = \sum_i \frac{u_i u_i^T}{\lambda_i - s}$ in which λ_i is an eigenvalue of A with corresponding eigenvector u_i .)

Both those equations lead to the following Newton scheme:

$$\begin{aligned}
s_{k+1} &= s_k + \frac{G(s_k)}{G'(s_k)} \\
&= s_k + \frac{1}{H(s_k)} \frac{H^2(s_k)}{H'(s_k)} \\
&= s_k + \frac{b^\top (A - s_k I)^\dagger b}{b^\top (A - s_k I)^\dagger (A - s_k I)^\dagger b}
\end{aligned} \tag{3.3}$$

Let $x_k = (A - s_k I)^\dagger b$, then equation 3.3 can be simplified further:

$$\begin{aligned}
s_{k+1} &= s_k + \frac{b^\top (A - s_k I)^\dagger b}{b^\top (A - s_k I)^\dagger (A - s_k I)^\dagger b} \\
&= s_k + \frac{b^\top x_k}{x_k^\top x_k} \\
&= \frac{s_k x_k^\top x_k + b^\top x_k}{x_k^\top x_k} \\
&= \frac{(s_k x_k^\top + x_k^\top (A - s_k I)) x_k}{x_k^\top x_k} \\
&= \frac{x_k^\top A x_k}{x_k^\top x_k}
\end{aligned} \tag{3.4}$$

The implementation of this Newton scheme is exactly algorithm 3.2.1. Hence in the Dominant Pole Algorithm, s_k converges quadratically. \square

Theorem 3.2.3. *Consider algorithm 3.2.1. Let $\lambda \in \mathbb{R}$ be a simple eigenvalue of A to which the algorithm converges, and let $x \in \mathbb{R}^n$ be its corresponding normalized eigenvector. If $b^\top x \neq 0$, then $x_k \rightarrow x$ up to a scalar, quadratically.*

Proof. Let $\tau_k \in \mathbb{R}$ be a scalar such that:

$$\tau_k x_k = x + e_k \tag{3.5}$$

With $e_k \perp x$, the ‘error’ vector. That is, x_k is scaled such that $\tau_k x_k^\top x = 1$. It can be shown that for a sufficient large k , $\frac{1}{\tau_k} \|e_k\| \rightarrow 0$ quadratically:

$$\begin{aligned}
\|e_k\| &= \frac{1}{\tau_k} \|(I - x x^\top)(A - s_k I)^\dagger b\| \\
&= \frac{1}{\tau_k} \left\| \sum_r \frac{u_r (u_r^\top b)}{\lambda_r - s} \right\| \\
&\leq \frac{C}{\tau_k}
\end{aligned} \tag{3.6}$$

Hence $x_k = \tau_k^{-1}x + \tau_k^{-1}e_k \rightarrow \tau_k^{-1}x$ as $k \rightarrow \infty$ quadratically iff $\tau_k^{-1} \rightarrow 0$ as $k \rightarrow \infty$ quadratically. The latter statement is shown to be true:

Note that algorithm 3.2.1 computes:

$$x_k = (A - s_k I)^\dagger b = \frac{u_i u_i^\top b}{\lambda_i - s_k} \quad (3.7)$$

Hence, let $s_k \rightarrow \lambda$ ($= \lambda_j$) and x ($= u_j$) be the corresponding eigenvector, then τ_k is of the form:

$$\tau_k = \frac{x^\top b}{\lambda - s_k} \quad (3.8)$$

Moreover, as $s_k \rightarrow \lambda$ quadratically as was shown in lemma 3.2.2, it follows that $\tau_k^{-1} \rightarrow 0$ quadratically. \square

3.3 Subspace accelerating with deflation

Most of this section is based on [26]. Novel is the idea of considering different selection strategies and the inclusion of a Rayleigh Quotient Iteration step.

It is possible to accelerate algorithm 3.2.1 by modifying it into a subspace method. The main idea is that the vectors x_1, \dots, x_k are used to construct subspace X_k which spans the eigenspace associated with the eigenvalue to converge to. The matrix $S_k := X_k^\top A X_k$ is the projection of A onto X_k . The eigenvalue of S_k converges to the dominant eigenvalue of A as the span of X_k comes closer to the eigenspace associated with that eigenvalue.

S_k has different eigenvalues which are the projected eigenvalues of A . Based on some selection rule we select one of those and try to expand X_{k+1} such that in the next iteration X_k has a better cover of the corresponding eigenspace.

In standard Dominant Pole Algorithm that would be through taking the inproduct between the eigenvectors of S_k and b , then the most dominant is selected. In standard SASPA, introduced in the next section, all eigenvalues would be tested on sensitivity and the most sensitive would be selected.

The Subspace accelerated Dominant Pole Algorithm is:

Algorithm 3.3.1 Subspace Accelerated Dominant Pole Algorithm

Require:

- $A \in \mathcal{S}^n$: matrix
- $b \in \mathbb{R}^n$: residue vector
- $s_0 \in \mathbb{R}$: initial pole estimate
- $\delta \in \mathbb{R}_{>0}$: threshold from where to start Rayleigh Quotient Iteration
- $\epsilon \in \mathbb{R}_{>0}$: tolerance
- p , number of poles to compute
- k_{\min}, k_{\max} minimum and maximum size of search space for restart

Ensure:

- $\lambda_1, \dots, \lambda_p \in \mathbb{R}$: eigenvalues of A
- $q_1, \dots, q_p \in \mathbb{R}^n$: corresponding eigenvectors

```

1:  $k \leftarrow 0, i \leftarrow 1, X \leftarrow []$ 
2: while  $i \leq p$  do
3:   Solve  $(A - s_k I)x_k = b$  for  $x_k \in \mathbb{R}^n$ 
4:    $x_k \leftarrow \text{GS}(X, x_k), X \leftarrow [X, x_k]$ 
5:    $S \leftarrow X^T A X$ 
6:    $(\tilde{\lambda}, \tilde{x}) \leftarrow \text{Select}(S, X^T b)$ 
7:    $q_i \leftarrow X^T \tilde{x}$ 
8:   if  $\|Aq_i - \tilde{\lambda}q_i\| \leq \delta$  then
9:     while  $\|Aq_i - \tilde{\lambda}q_i\| \geq \epsilon$  do
10:      Solve  $(A - \tilde{\lambda}I)\hat{q}_i = q_i$  for  $\hat{q}_i \in \mathbb{R}^n$ 
11:       $q_i \leftarrow \frac{\hat{q}_i}{\|\hat{q}_i\|}$ 
12:       $\tilde{\lambda} \leftarrow q_i^T A q_i$ 
13:     end while
14:   end if
15:   while  $\|Aq_i - \tilde{\lambda}q_i\| \leq \epsilon$  do
16:      $\lambda_i \leftarrow \tilde{\lambda}$ 
17:     Deflate  $b$ , and  $X$  with  $q_i$ 
18:      $i \leftarrow i + 1$ 
19:      $\tilde{\lambda} \leftarrow \text{Select}(S, X^T b)$ 
20:      $q_i \leftarrow X^T \tilde{x}$ 
21:   end while
22:   if #columns of  $X$  is more than  $k_{\max}$  then
23:      $S \leftarrow X^T A X$ 
24:      $((\mu_1, u_1), \dots, (\mu_{k_{\min}}, u_{k_{\min}})) \leftarrow \text{Select}_{k_{\min}}(S, X^T b)$ 
25:      $\tilde{X} \leftarrow X^T [u_1, \dots, u_{k_{\min}}]$ 
26:     Orthogonalize  $\tilde{X}$ 
27:      $X \leftarrow \tilde{X}$ 
28:   end if
29:    $s_{k+1} \leftarrow \tilde{\lambda}$ 
30:    $k \leftarrow k + 1$ 
31: end while

```

A quick explanation of the inner loop of algorithm 3.3.1:

- Line 3: determine the vector x_k which will be used to expand the search spaces.
- Line 4: use a stable variant of Gram-Schmidt to orthogonalize the new vectors to the existing vectors in the search space, and add them.

- Line 5-7: select poles to converge to from the search space, for an example algorithm see 3.3.2.
- Line 8-14: elective Rayleigh Quotient Iteration for faster convergence, see algorithm 3.1.1.
- Line 15-21: check if the algorithm has converged, if so deflate and continue (deflating will be explained later on in this section).
- Line 22-28: restart of the search space.

As a side-note: The algorithm written here starts with an empty search space, however it can also be chosen to be filled already with vectors that are expected to be near optimum. For example vectors known from previous computations. Moreover, if such vectors are unknown, then it can also be chosen to fill the search space by Lanczos/Arnoldi vectors (that is a basis for the Krylov subspace $\mathcal{K}_k(A, b)$).

In simulations the fastest convergence was obtained by filling the initial search space with a few Lanczos/Arnoldi vectors.

3.3.1 Selecting poles

The matrix S is of low dimension and stems from the projected eigenvalue problem:

$$(X^T A X) \tilde{x} = \tilde{\lambda} \tilde{x} \quad (3.9)$$

Hence, for any eigenvector \tilde{x} of S , the vector $x = X^T \tilde{x}$ is an eigenvector of A projected on X^T . As the number of columns of X increase, the dimension of the projected spaces grow, which in turn implies that the approximations get closer. The idea is that S contains eigenvectors and eigenvalues which will under the right transformation increasingly precision approximate eigenvectors of A with corresponding eigenvalues.

Selecting poles for the low dimension system S is done by computing an eigendecomposition and with a dense eigenvalue method such as the QR-algorithm.

For computation of the dominant pole as in algorithm 3.2.1, that would be:

Algorithm 3.3.2 Selection strategy

Require:

- $S \in \mathbb{R}^{m \times m}$: matrix
- $b \in \mathbb{R}^n$: residue vector

Ensure:

- $\lambda \in \mathbb{R}$: dominant pole (eigenvalue of S)
- $x \in \mathbb{R}^n$: corresponding eigenvector

- 1: Compute eigenpairs of S : (λ_i, x_i)
- 2: Compute $R_i \leftarrow (x_i^T b)^2$ for all i
- 3: $j \leftarrow \arg \max_i |R_i|$

4: $\lambda \leftarrow \lambda_j, x \leftarrow x_j$

Algorithm 3.3.2 is a sample selection-rule. Other selection rules include $\arg \max_i \frac{|R_i|}{|\lambda_i|}$ and $\arg \max_i \frac{|R_i|}{|C-\lambda_i|}$.

However, note that by incorporating a selection strategy possible quadratic convergence is lost. That is because on line 15: $s_{k+1} \leftarrow \lambda$, hence s_k does not follow the Newton scheme. For the proposed selection rule (algorithm 3.3.2) this poses no problem because it can be proved that this s_{k+1} is closer to the pole than in algorithm 3.2.1.

If using a selection strategy which questionable converges, for instance which forces the algorithm to cycle between different poles without reaching one with enough precision, then it is recommended to use the Rayleigh Quotient Iteration step (line 10-16). This stops the algorithm from cycling because the poles that the algorithm cycles between are in an early stadium deflated from the search space. That allows the algorithm to converge to the other, more favorable, poles without hindrance of the other poles.

3.3.2 Deflating

Deflating means removing a vector component out of the system. For instance, deflating A with x_q means:

$$A \leftarrow (I - x_q x_q^T) A (I - x_q x_q^T) \quad (3.10)$$

That way the eigenvalue corresponding to x_q is reduced to 0. However, if A was sparse, then after applying equation 3.10, it surely is not and it is computational more difficult to solve equations of the form $(A - sI)x = b$ for $x \in \mathbb{R}^n$.

Therefore a better choice is to not deflate A but instead deflate b :

$$b \leftarrow (I - x_q x_q^T) b \quad (3.11)$$

This stops the algorithm from converging again to the already found pole because $(b^T (I - x_q x_q^T)^T x_q) (x_q^T (I - x_q x_q^T) b) = 0$.

Note that this poses no stability problem, although x_q might contain a large numerical error, because this step only happens once after convergence. After that the eigenvector component in b of the dominant pole is diminished such that the algorithm will not converge again towards that pole.

Also to make sure that the pole does not pop up again during the selection stage in the next iteration, deflate X :

$$X \leftarrow (I - x_q x_q^T) X \quad (3.12)$$

By deflating X , orthogonality of the columns will be lost. Hence they have to be reorthogonalized after deflation. That can be done through a stabilized version of Gram-Schmidt.

3.3.3 Restart search space

If the search space grows too large, then the selecting of poles on line 7-9 becomes computational more expensive. Hence by incorporating an implicit strategy for the search space it is made sure that it is always of low dimensional.

The selection criteria for the k_{\min} vectors to be included in the search space should be similar to the selection strategy on line 7-9. A suitable choice is to order them all during line 8 to some favorable ranking, return the most favorable and then on line 24 select the top k_{\min} most favorable vectors.

Note that if the system has converged and some vector is being deflated from the system, then the dimension of the search space is most likely also decreased by one. Thus the results of line 7-8 are almost surely applicable during line 23-24, which might save some computation time. However, if the Rayleigh Quotient Iteration step is incorporated, then it follows not necessarily that the search space will be decreased. Therefore if the Rayleigh Quotient Iteration is being elected, then either line 23-24 have to be computed from scratch, or the implicit restart has to be canceled for one iteration, i.e. change the ‘if’ on line 22 to an ‘else if’.

After selecting the vectors to be included in the new search space, they should be orthogonalized. Again, that can be done through applying Givens rotations or any other stable Gram-Schmidt variant on their columns. For stability it is recommended to not skip this step although it might be expected that u_1 to $u_{k_{\min}}$ are already orthogonal because they are eigenvectors.

3.4 Sensitive Pole Algorithm

This section is largely inspired by [27]. As with the previous section, novel is the inclusion of the Rayleigh Quotient Iteration step.

The Sensitive Pole Algorithm tries to compute the poles of $A \in \mathcal{S}^n$ with a large derivative under perturbation. That is, let \dot{A} denote the perturbation of A . Then we are interested in the eigenvalues of the function:

$$f(t) := A + t\dot{A} \quad (3.13)$$

The derivative of a simple eigenvalue $\lambda_i \in \mathbb{R}$ with corresponding normalized eigenvector $u_i \in \mathbb{R}^n$ of A is given by:

$$\frac{\mathbf{d}\lambda_i(f(t))}{\mathbf{d}t} = tu_i^T \dot{A} u_i \quad (3.14)$$

(Theorem 2.2.7.)

The Sensitive Pole Algorithm computes the eigenvalues $\lambda_i \in \mathbb{R}$ for which $|\frac{\mathbf{d}\lambda_i(f(t))}{\mathbf{d}t}|$ is considerable large in comparison with the other eigenvalues.

It does that by replacing $b \in \mathbb{R}^n$ in algorithm 3.2.1:

$$b \leftarrow \dot{A}x_{k-1} \quad (3.15)$$

The motivation for doing so is that ideally one wants to compute the most sensitive pole which is the eigenvalue λ_i with a large residual ($|u_i^T \dot{A}u_i|$) for some $u_i \in \mathbb{R}^n$ in the eigenspace of λ_i . If $b = \dot{A}u_i$, then the residual at step k is $(b^T x_k)(x_k^T b) = (u_i^T \dot{A}x_k)(x_k^T \dot{A}u_i)$, which converges to $(u_i^T \dot{A}u_i)^2$, with both u_i and v_i in the eigenspace corresponding to λ_i . However to compute the residual the eigenvector u_i is necessary. The best approximation of u_i at step k is x_{k-1} , hence the choice $b \leftarrow \dot{A}x_{k-1}$.

The full algorithm for Sensitive Pole Algorithm is:

Algorithm 3.4.1 Sensitive Pole Algorithm

Require:

- $A \in \mathcal{S}^n$: matrix
- $\dot{A} \in \mathcal{S}^n$: perturbation
- $s_0 \in \mathbb{R}$: initial pole estimate
- $\epsilon \in \mathbb{R}_{>0}$: tolerance

Ensure:

- $\lambda \in \mathbb{R}$: sensitive pole (eigenvalue of A)
- $x \in \mathbb{R}^n$: corresponding eigenvector

- 1: $k \leftarrow 0, x_{-1} = (1, \dots, 1)^T \in \mathbb{R}^n$
 - 2: **while** not converged **do**
 - 3: $b \leftarrow \dot{A}x_{k-1}$
 - 4: Solve $(A - s_k I)x_k = b$ for $x_k \in \mathbb{R}^n$
 - 5: $s_{k+1} \leftarrow s_k + \frac{c^T x_k}{x_k^T x_k} \quad \triangleright \left(\frac{x_k^T A x_k}{x_k^T x_k} \right)$
 - 6: **if** $\|Ax_k - s_{k+1}x_k\| \leq \epsilon$ **then**
 - 7: **return** $\lambda \leftarrow s_{k+1}, x \leftarrow x_k$
 - 8: **end if**
 - 9: $k \leftarrow k + 1$
 - 10: **end while**
-

And also based on algorithm 3.3.1, a subspace accelerated variant:

Algorithm 3.4.2 Subspace Accelerated Sensitive Pole Algorithm

Require:

- $A \in \mathcal{S}^n$: matrix
- $\dot{A} \in \mathcal{S}^n$: perturbation
- $s_0 \in \mathbb{R}$: initial pole estimate
- $\delta \in \mathbb{R}_{>0}$: threshold from where to start Rayleigh Quotient Iteration
- $\epsilon \in \mathbb{R}_{>0}$: tolerance
- p , number of poles to compute
- k_{\min}, k_{\max} minimum and maximum size of search space for restart

Ensure:

- $\lambda_1, \dots, \lambda_p \in \mathbb{R}$: dominant sensitive poles (eigenvalues of A)
- $q_1, \dots, q_p \in \mathbb{R}^n$: corresponding eigenvectors

- 1: $k \leftarrow 0, i \leftarrow 1, x_{-1} = (1, \dots, 1)^T \in \mathbb{R}^n$

```

2: while  $i \leq p$  do
                                     ▷ Extend search space
3:    $b \leftarrow Ax_{k-1}$ 
4:   if  $i \geq 2$  then
                                     ▷ Deflate  $b$  with  $q_1, \dots, q_{i-1}$ 
5:      $b \leftarrow \text{GS}([q_1, \dots, q_{i-1}], b)$ 
6:   end if
7:   Solve  $(A - s_k I)x_k = b$  for  $x_k \in \mathbb{R}^n$ 
8:    $x_k \leftarrow \text{GS}(X, x_k)$ ,  $X \leftarrow [X, x_k]$ 
                                     ▷ Select most promising
9:    $S \leftarrow X^T A X$ 
10:   $(\tilde{\lambda}, \tilde{x}) \leftarrow \text{Select}(S, X^T b)$ 
11:   $q_i \leftarrow X^T \tilde{x}$ ,  $q_i \leftarrow \frac{q_i}{\|q_i\|}$ 
                                     ▷ Elective Rayleigh Quotient Iteration
12:  if  $\|Aq_i - \tilde{\lambda}q_i\| \leq \delta$  then
13:    while  $\|Aq_i - \tilde{\lambda}q_i\| \geq \epsilon$  do
14:      Solve  $(A - \tilde{\lambda}I)\hat{q}_i = q_i$  for  $\hat{q}_i \in \mathbb{R}^n$ 
15:       $q_i \leftarrow \frac{\hat{q}_i}{\|\hat{q}_i\|}$ 
16:       $\tilde{\lambda} \leftarrow q_i^T A q_i$ 
17:    end while
18:  end if
                                     ▷ Check convergence and deflate
19:  if  $\|Aq_i - \tilde{\lambda}q_i\| \leq \epsilon$  then
20:     $\lambda_i \leftarrow \tilde{\lambda}$ 
21:    Deflate  $X$  with  $q_i$ 
22:     $i \leftarrow i + 1$ 
23:  end if
                                     ▷ Check size search space and implicit restart
24:  if #columns of  $X$  is more than  $k_{\max}$  then
25:     $S \leftarrow X^T A X$ 
26:     $((\mu_1, u_1), \dots, (\mu_{k_{\min}}, u_{k_{\min}})) \leftarrow \text{Select}_{k_{\min}}(S, X^T b)$ 
27:     $\tilde{X} \leftarrow X^T [u_1, \dots, u_{k_{\min}}]$ 
28:    Orthogonalize  $\tilde{X}$ 
29:     $X \leftarrow \tilde{X}$ 
30:  end if
31:   $s_{k+1} \leftarrow \lambda$ 
32:   $k \leftarrow k + 1$ 
33: end while

```

3.5 Computing cross points of eigencurves

One of the problems that has to be dealt with during eigenvalue optimization is to estimate when the largest eigenvalue gets overtaken by another eigenvalue during a perturbation, i.e. when two eigenvalue curves cross.

For clarity, in this section the following problem is considered: let $A \in \mathcal{S}^n$

be the unperturbed matrix, $\dot{A} \in \mathcal{S}^n$ the perturbation and $A(t) = A + t\dot{A}$ the perturbation function. By theorem 2.2.12 it is possible to label all eigenvalues of $A(t)$ as $\lambda_i(t)$ with at $t = 0$:

$$\lambda_1(0) \geq \lambda_2(0) \geq \dots \geq \lambda_n(0) \quad (3.16)$$

(If at $t = 0$, an eigenvalue has multiplicity τ , then there is a sequence $\lambda_i(0) = \dots = \lambda_{i+\tau-1}(0)$.)

The functions $\lambda_i(t)$ for $t \geq 0$ are called eigencurves.

Let further on in this section τ denote the multiplicity of λ_1 at $t = 0$. The problem is stated as estimating what the minimal $t^* > 0$ is for which there exists an index $i^* > \tau$ such that $\lambda_1(t^*) = \lambda_{i^*}(t^*)$:

If A is small with only simple eigenvalues, then the problem can be solved by performing a full eigendecomposition of A such that for each $\lambda_i(t)$ a unique eigenvector $u_i \in \mathbb{R}^n$ at $t = 0$ is known. Then all directional derivatives $v_i = u_i^T \dot{A} u_i$ (theorem 2.2.7) are calculated and let $T = \{(\frac{\lambda_1(0) - \lambda_i(0)}{v_i - v_1}, i) \text{ s. t. } i > 1\}$:

$$(t^*, i^*) = \arg \min_{(t, i) \in T} t \text{ s. t. } t > 0 \quad (3.17)$$

However, A is for many instances large and sparse with possible non-simple eigenvalues. In such cases not all derivatives of the eigencurves are known, nor is a full eigendecomposition feasible to compute.

Despite that a full eigendecomposition is infeasible, it is still possible to calculate equation 3.17: with the Sensitive Pole Algorithm it is possible to compute a subset of T which contains the most promising elements.

By using algorithm 3.4.2 with some kind of deflation and the standard selection rule it is possible to calculate the eigencurves with a large absolute derivative. By employing a selection rule based on the linear approximation of the intersections convergence is steered to more promising eigencurves.

The modified selection rule is:

Algorithm 3.5.1 Selection strategy for intersecting maximum eigencurve

Require:

- $S \in \mathbb{R}^{m \times m}$: matrix
- $\dot{S} \in \mathbb{R}^{m \times m}$: perturbation
- $\lambda_{\max} \in \mathbb{R}$: maximum eigenvalue at $t = 0$
- $v_{\max} \in \mathbb{R}$: derivative $\frac{d\lambda_{\max}}{dt}$

Ensure:

- $\lambda \in \mathbb{R}$: sensitive pole (eigenvalue of S)
- $x \in \mathbb{R}^n$: corresponding eigenvector

- 1: Compute eigenpairs of S : (λ_i, x_i)
 - 2: Compute $R_i \leftarrow \frac{\lambda_{\max} - \lambda_i}{x_i^T \dot{S} x_i - v_{\max}}$ for all i
 - 3: $j \leftarrow \arg \max_i R_i$ s. t. $R_i > 0$
 - 4: $\lambda \leftarrow \lambda_j$, $x \leftarrow x_j$
-

Moreover, as it is unpractical to run 3.4.2 for a long time it is necessary to think of a stopping criteria. It is undesirable to stop the method too soon since than it is unlikely to have found the crossing eigencurve, on the other hand it is also desired to stop as soon as possible and not waste computing cycles to non-nearest eigencurves.

In simulations it was chosen to stop the algorithm as soon as the converged to eigencurve has the nearest intersection point ever seen by the algorithm.

The full algorithm:

Algorithm 3.5.2 SASPA to compute nearest intersection

Require:

- $A \in \mathcal{S}^n$: matrix
- $\dot{A} \in \mathcal{S}^n$: perturbation
- $\lambda \in \mathbb{R}$: largest eigenvalue at $t = 0$
- $q_0, \dots, q_k \in \mathbb{R}^n$: corresponding eigenvectors
- $f \in \mathbb{R}^n$: forward derivative largest eigencurve
- $\epsilon_{RQI} \in \mathbb{R}_{>0}$: threshold from where to start Rayleigh Quotient Iteration
- $\epsilon_{tol} \in \mathbb{R}_{>0}$: tolerance
- k_{\min}, k_{\max} minimum and maximum size of search space for restart
- $iter_{\max}$, maximum number of iterations

Ensure:

- $\hat{t} \in \mathbb{R}_+$: intersection point
- $\hat{\lambda} \in \mathbb{R}$: nearest crossing eigenvalues at $t = 0$
- $\hat{q} \in \mathbb{R}^n$: corresponding eigenvector

- 1: $iter \leftarrow 1, num \leftarrow k, s \leftarrow \lambda, x \leftarrow (1, \dots, 1)^T \in \mathbb{R}^n, \hat{t} \leftarrow \infty, \tau \leftarrow \infty$
- 2: **while** $iter \leq iter_{\max}$ **do**
 - 3: $b \leftarrow \dot{A}x$ ▷ Extend search space
 - 4: $b \leftarrow \text{GS}([q_0, \dots, q_{num}], b)$ ▷ Deflate b with q_0, \dots, q_{num}
 - 5: Solve $(A - sI)x = b$ for $x \in \mathbb{R}^n$
 - 6: $x \leftarrow \text{GS}(X, x), X \leftarrow [X, x]$ ▷ Select most promising
 - 7: $S \leftarrow X^T A X$
 - 8: Compute all eigenpairs (s_i, x_i) of S
 - 9: Compute for all pairs corresponding cross points: $t_i \leftarrow \frac{\lambda - s_i}{x_i^T X^T \dot{A} X x_i - q_0^T \dot{A} q_0}$
 - 10: Reorder all triplets (t_i, s_i, x_i) such that $t_i \geq t_{num}$ for all $i \leq j$
 - 11: $(t, s, x) \leftarrow (t_1, s_1, Xx_1)$
 - 12: $x \leftarrow \frac{x}{\|x\|}$ ▷ Elective Rayleigh Quotient Iteration
- 13: **if** $\|Ax - sx\| \leq \epsilon_{RQI}$ **then**
- 14: **while** $\|Ax - sx\| \geq \epsilon_{tol}$ **do**
- 15: Solve $(A - sI)\dot{x} = x$ for $\dot{x} \in \mathbb{R}^n$
- 16: $x \leftarrow \frac{\dot{x}}{\|\dot{x}\|}$
- 17: $s \leftarrow x^T A x$

```

18:          $t \leftarrow \frac{\lambda - s}{x^T \dot{A} x_i - q_0^T \dot{A} q_0}$ 
19:     end while
20: end if
▷ Update stop criteria
21:  $\tau \leftarrow \min(\tau, t)$ 
▷ Check convergence and deflate
22: if  $\|Ax - sx\| \leq \epsilon_{tol}$  then
23:      $num \leftarrow num + 1$ 
24:      $q_{num} \leftarrow x$ 
25:     Deflate  $X$  with  $q_{num}$ 
26:     if  $t \leq \hat{t}$  then
27:          $\hat{\lambda} \leftarrow s, \hat{q} \leftarrow q_{num}, \hat{t} \leftarrow t$ 
28:     end if
29:     if  $t \leq \tau$  then
30:         return
31:     end if
32:      $(t, s, x) \leftarrow (t_2, s_2, Xx_2)$ 
33:      $x \leftarrow \frac{x}{\|x\|}$ 
34: end if
▷ Check size search space and implicit restart
35: if #columns of  $X$  is more than  $k_{max}$  then
36:      $X \leftarrow [Xx_1, Xx_2, \dots, Xx_{k_{min}}]$ 
37:     Orthogonalize and normalize  $X$ 
38: end if
39:      $iter \leftarrow iter + 1$ 
40: end while

```

The Matlab implementation being used is included in appendix A.2.

3.5.1 Numerical results

As a demonstration of algorithm 3.5.2, some numerical results are presented.

The described experiment is with A taken from the 10th DIMACS challenge [3]. The matrix A is chosen to be the Laplacian matrix from a graph in the Random Geometric Graphs set. That are random geometric graphs labeled `rgg_n_2_X_s0.graph`, each with 2^X vertices with X ranging from 15 up to 24. Each vertex is a random point in the unit square and edges connect vertices whose Euclidean distance is below $0.55 \frac{\log(n)}{n}$.

The non-zero elements of the Laplacian matrix of `rgg_n_2_15_s0.graph` are shown in figure 3.1.

As perturbation \dot{A} , a random sparse diagonal matrix is chosen with elements uniform distributed between -1 and 1 . If the largest eigenvalue has a positive gradient, then \dot{A} is flipped, that way as t increases the largest eigenvalue will always decrease.

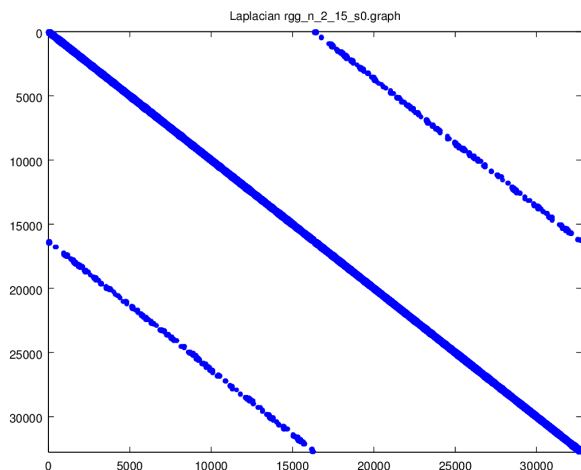


Figure 3.1: Non-zero elements of the Laplacian matrix for `rgg_n_2_15_s0.graph`

t^* through bisection	t^* through SASPA	iterations SASPA
2.0887	2.1283	21
0.5380	0.5662	13
1.1191	1.3674	33
0.6787	0.7236	19

Table 3.1: Numerical results for `rgg_n_2_15_s0.graph` (size of Laplacian is 32768×32768 , spectrum is $[0, 25.896]$)

An example of such a perturbation is shown in 3.2, the matrix A is the Laplacian matrix of `rgg_n_2_15_s0.graph`. It can be seen that the nearest intersection of the largest eigencurve is around 0.6. The algorithm computes the intersection to be at 0.56.

As parameters for the algorithm, in all cases is $\epsilon = 10^{-12}$, $k_{\min} = 4$ and $k_{\max} = 20$. The initial λ_{\max} and its corresponding eigenvector are computed at the initial stage with the built in Matlabmethod `eigs`.

The results for random diagonal matrices are shown in tables 3.1, 3.2 and 3.3. Each row represents a different test case. The first column is the intersection calculated by using a golden section search method, which is the most accurate value for t^* known, the second column through SASPA with the described rule and the third column the number of iterations necessary for SASPA.

The exact cross-point of the eigencurves was computed by using a golden section search, which is a variant of bisection. Per computation of the intersection point up to four decimals accurate, 19 computations of the largest eigenvalue were necessary. Thus even if computing the largest eigenvalue can be done within a few iterations of Jacobi–Davidson, Rayleigh Quotient Iteration

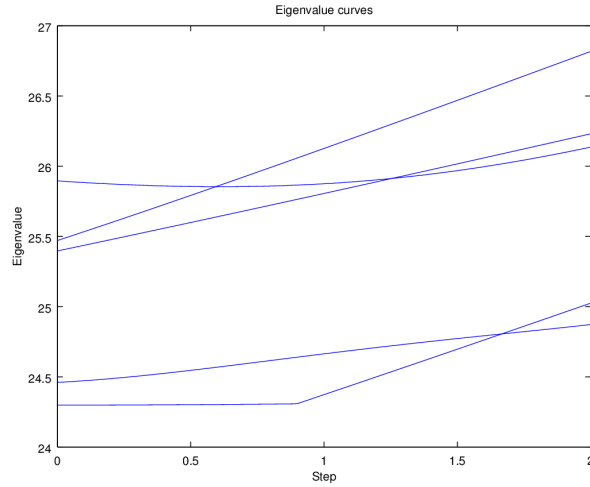


Figure 3.2: The largest five eigenvalue of $A + t\dot{A}$ with A the Laplacian matrix for `rgg_n_2_15_s0.graph` and \dot{A} as described

t^* through bisection	t^* through SASPA	iterations SASPA
0.7537	0.9641	11
1.6514	1.7891	21
1.0576	1.3721	17
2.2711	2.4224	9

Table 3.2: Numerical results for `rgg_n_2_16_s0.graph` (size of Laplacian is 65536×65536 , spectrum is $[0, 28.580]$)

t^* through bisection	t^* through SASPA	iterations SASPA
2.1570	2.0325	21
0.6842	1.7685	20
0.1985	0.2199	5
2.2538	3.0589	31

Table 3.3: Numerical results for `rgg_n_2_17_s0.graph` (size of Laplacian is 131072×131072 , spectrum is $[0, 29.655]$)

or the power method, it still takes a lot more time to find the intersection point through the bisection method than with SASPA.

Although the reported t^* is less accurate than the t^* determined by the bisection method, SASPA with the proposed selection rule is able to compute the intersection point for different matrix perturbations within a few iterations.

For larger matrices, the error increases. That can be attributed to the fact that second order effects increase in magnitude for the larger instances. That is because the spectrum is more dense. Hence by theorem 2.2.11 it can be expected that the second order derivatives of the eigenvalues are larger for dense spectrum's than for less dense spectrum's.

To get a sense of the computation time difference: By using the `eigs` method from `Matlab`, the instances from `rgg_n.2.15.s0.graph` take 30–35 seconds to compute; by using SASPA with selection rule it takes 24–30 seconds for 50 iterations (that includes initial computation of λ_{\max} and the corresponding eigenvector).

For `rgg_n.2.17.s0.graph` the computation time by bisection rose to 125–175 seconds; by SASPA to 140–150 seconds for 50 iterations. The instances in which the bisection method was faster were those with the intersection point close to 0. That is because in those instances the eigencurves diverge fast, hence the computation for the largest eigenvalue takes less time. However, in those instances SASPA converges within a few iterations, thus it is wasteful to run actually 50 iterations.

If for `rgg_n.2.17.s0.graph` the maximum number of iterations is reduced to 25, then computation by SASPA takes 100–115 seconds to converge. And in most cases, the method converges to the nearest intersection. However, it should be noted that the method does not converge always to the nearest within 25 iterations.

Hence SASPA is in most cases much faster than the bisection method. If the required accuracy is not too large then SASPA is a suitable choice for these kind of computations. If accuracy is very important, then the bisection method is better suited.

Note that these computations time have a sense of arbitrariness. They depend on the current load of the computer on which they are run, the underlying method used by `eigs`, the availability of multicore processors, etc. Most of those parameters cannot be controlled on the machine that these methods are tested on. Hence the computation times that are reported should not be taken as absolute values, but only used to form a general sense of expected computation times of both methods relative to each other.

Following those computational results, the proposed method is faster than a bisection method based on `eigs` from `Matlab`, if the maximum number of iterations is bounded correct. Therefore implementations should focus to bound the number of iterations based on the best estimate so far of the nearest intersection.

Chapter 4

Algorithms for eigenvalue optimization

In this chapter it will be described how to solve eigenvalue optimization problems of the following form:

$$\begin{aligned} \inf_{z \in \mathbb{R}^m} \lambda_{\max}(C - \mathcal{A}^T(z)) \\ \text{s. t. } Pz \geq 0 \end{aligned} \tag{4.1}$$

With $C \in \mathcal{S}^n$, $\mathcal{A}^T : \mathbb{R}^m \rightarrow \mathcal{S}^n$, $z \mapsto \sum_i z_i A_i$, $P \in \mathbb{R}^{k \times m}$ with normalized rows and \geq component wise.

These problems arise in the study of semidefinite optimization problems as was explained in chapter 2. The emphasis in this chapter will be on sparse instances of eigenvalue problems. Also it will be marked when the distinction between sparse and non-sparse instances becomes important.

The basic idea of solving problem 4.1 in an iterative manner is as follows:

1. Start with a feasible solution $z \in \mathbb{R}^m$, i.e. $Pz \geq 0$, e.g. $z = 0$
2.
 - (a) Calculate a descent direction
 - (b) Calculate global optimum in that direction
 - (c) Update z
 - (d) Repeat until steps are too small for steady improvement

Except for the first step, which is self-evident, each of these steps will be explained with more rigor in the next sections.

On notation: in the next section the formula $\lambda_{\max}(C - \mathcal{A}^T(z))$ might be denoted as $\lambda_{\max}(z)$. The reason being that in chapter 2 already many results were obtained for functions such as $\lambda_{\max}(C - \mathcal{A}^T(z))$ which are better resembled by writing it as a function $\lambda_{\max}(z)$. The optimum (i.e. the global minimum)

is denoted as $z^* \in \mathbb{R}^m$. Also, let z_i be the optimum point being calculated in the i th iteration, and d_i and t_i the search direction and the step size in the i th iteration. If from the context is clear which iteration is meant, then the subscript is omitted.

As a quick recap on some results derived in chapter 2: The function $\lambda_{\max}(z)$ is convex, and if for $z \in \mathbb{R}^m$, $\lambda_{\max}(C - \mathcal{A}^T(z))$ is simple, then and only then does the gradient $\nabla\lambda_{\max}(z)$ exists. For all $z \in \mathbb{R}^m$ with $\lambda_{\max}(C - \mathcal{A}^T(z))$ multiple is the subgradient defined. Only the subgradient at $z^* \in \mathbb{R}^m$ contains the all zero vector.

4.1 Calculating a descent direction

This section is concerned with finding a feasible descent direction $d \in \mathbb{R}^m$ such that there exists a step size $t > 0$ such that $\lambda_{\max}(z + td) < \lambda_{\max}(z)$ and $P(z + td) \geq 0$.

The first section gives some examples and explains the steepest descent for when $\lambda_{\max}(z)$ is simple. The second section generalizes the first section to include the non-simple case. The third section is concerned with cutting infeasible directions from d .

Although the first section might seem easy it is noteworthy that the examples presented there provide for a better insight for the non-simple case, which is considerable harder. That is because if $\lambda_{\max}(z)$ is no longer simple, than the problem is not smooth and the only tools available are the subgradients introduced in section 2.2.

4.1.1 Descent direction in the simple case

In the case that for a given $z \in \mathbb{R}^m$, $\lambda_{\max}(C - \mathcal{A}^T(z))$ is simple, then it is easy to compute the steepest descent direction. For the steepest direction only the gradient is necessary which only depends on the largest eigenvalue with corresponding eigenvector. Hence the steepest descent direction is preferable if the instance of problem 4.1 is large and sparse.

If for a given $z \in \mathbb{R}^m$, $\lambda_{\max}(C - \mathcal{A}^T(z))$ is simple with normalized eigenvector q , then the gradient $\nabla\lambda_{\max}(z)$ is well defined (corollary 2.2.8). Thus an obvious choice for descent direction would be $d := -\nabla\lambda_{\max}(z)^T$:

Definition 4.1.1. *If $\lambda_{\max}(C - \mathcal{A}^T(z))$ is simple, then the steepest descent direction is:*

$$d := -\nabla\lambda_{\max}(z) = (q^T A_1 q, \dots, q^T A_m q)^T \quad (4.2)$$

The next three examples show how the steepest direction works in practice for unconstrained problems:

Example 4.1.2. *Consider problem 4.1 with $m = 2$ and no inequalities and:*

$$C = \begin{pmatrix} 2 & 1 \\ 1 & 0 \end{pmatrix}, A_1 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, A_2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad (4.3)$$

Analytically the eigenvalues of $C - \mathcal{A}^T(z)$ are $\lambda_1 = 1 + \sqrt{(z_1 - 1)^2 + (z_2 - 1)^2}$ and $\lambda_2 = 1 - \sqrt{(z_1 - 1)^2 + (z_2 - 1)^2}$.

Thus the eigenvalues at $z = 0$ are $\lambda_1 = 1 + \sqrt{2}$ and $\lambda_2 = 1 - \sqrt{2}$; hence λ_1 is simple. The eigenvector q corresponding to λ_1 is $q = (1 + \sqrt{2}, 1)^T$. Then:

$$\nabla f(0) = -(q^T A_1 q, q^T A_2 q) \approx -(4.83, 4.83) \quad (4.4)$$

Hence $d = -\nabla f(0)^T$ is a descent direction. And indeed, the optimum is at $z^* = (1, 1)^T$ at which $C - \mathcal{A}^T(z)$ is the identity and $\lambda_1 = \lambda_2 = 1$.

The next example shows that for simple instances with $m = 2$ the steepest direction is not necessarily the best direction:

Example 4.1.3. Consider problem 4.1 with $m = 2$, and no inequalities, and:

$$C = \begin{pmatrix} -2 & 1 & -2 \\ 1 & 2 & 0 \\ -2 & 0 & 2 \end{pmatrix}, A_1 = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}, A_2 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (4.5)$$

Analytically the eigenvalues are:

$$\begin{aligned} \lambda_1 &= \frac{1}{2} \left(\sqrt{4(z_2 - 1)^2 + \frac{1}{5}((5z_1 + 12)^2 + 16)} - z_1 \right) \\ \lambda_2 &= 2 \\ \lambda_3 &= \frac{1}{2} \left(-\sqrt{4(z_2 - 1)^2 + \frac{1}{5}((5z_1 + 12)^2 + 16)} - z_1 \right) \end{aligned} \quad (4.6)$$

λ_1 is minimized for $z^* = (-2, 1)^T$, at which $\lambda_1 = 2$.

The eigenvalues of C (that is at $z = 0$) are $\lambda_1 = 3$, $\lambda_2 = 2$ and $\lambda_3 = -3$. Hence λ_1 is simple. The eigenvector q corresponding to λ_1 is $q = (-1, -1, 2)^T$. Then:

$$\nabla f(0) = -(q^T A_1 q, q^T A_2 q) = (3, -2) \quad (4.7)$$

And $d = \frac{-\nabla f(0)^T}{\|\nabla f(0)^T\|_2} \approx (-0.83, 0.55)^T$. Hence the global minimizer z^* is not necessarily of the form $z_0 + td$ for $m \geq 2$.

The next example shows that steepest directions turn a simple largest eigenvalue fast in a non-simple one:

Example 4.1.4. Consider problem 4.1 with $m = 2$, no inequalities and:

$$C = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}, A_1 = \begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}, A_2 = \begin{pmatrix} 0 & 0 \\ 0 & -2 \end{pmatrix} \quad (4.8)$$

Analytically the eigenvalues are:

$$\begin{aligned} \lambda_1 &= 2 + z_1 \\ \lambda_2 &= 3 - z_1 + 2z_2 \end{aligned} \quad (4.9)$$

Which is clearly unbounded for $(z_1, z_2) \rightarrow (-\infty, -\infty)$.

Hence at $z = 0$ the largest eigenvalue is $\lambda_{\max} = 3$ and corresponding eigenvector is $q = (0, 1)^T$. The corresponding gradient is

$$\nabla f(0) = -(q^T A_1 q, q^T A_2 q) = (-1, 2) \quad (4.10)$$

After one step (assuming an exact line search in direction $d = -\nabla \lambda_{\max}(z)$) $z = (\frac{1}{4}, -\frac{1}{2})^T$ and $\lambda_1 = \lambda_2$.

Hence the steepest direction is not necessarily the best direction, but in practice for unconstrained problems it works quite decent. It seems that for problems in which the largest eigenvalues stay simple that a good strategy is to use some algorithm which minimizes a quadratic function, such as conjugated gradients, e.g. the Fletcher-Reeves method. However in practice, most instances of eigenvalue optimization in which the starting eigenvalue is simple, have within a few iterations a non-simple largest eigenvalue.

Note that the steepest descent direction d is not necessarily a feasible descent direction because of the constraint $Pz \geq 0$ in equation 4.1. This will be returned on this in subsection 4.1.3.

4.1.2 Descent direction in the non-simple case

Often $\lambda_{\max}(C - \mathcal{A}^T(z))$ is non-simple. Hence the gradient $\nabla \lambda_{\max}(C - \mathcal{A}^T(z))$ does not exist. However corollary 2.2.24 gives a generalization of the gradient and directional derivative which is applicable for this situation which allows to compute a steepest direction.

Before continuing, The following terms and definitions will be encountered in the rest of this section: Denote the multiplicity of $\lambda_{\max}(z)$ with t and the normalized eigenvectors that span the whole eigenspace as q_1, \dots, q_t . Moreover, let Q be the matrix which spans the eigenspace, that is with q_1, \dots, q_t as columns.

If $\lambda_{\max}(C - \mathcal{A}^T(z))$ is non-simple, then the steepest descent direction is:

$$\begin{aligned}
d &= \arg \min_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \nabla_d \lambda_{\max}(z) \\
&= \arg \min_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} - \min_{q \in Q} \langle \mathcal{A}^T(d), qq^T \rangle \\
&= \arg \max_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \min_{q \in Q} \langle \mathcal{A}^T(d), qq^T \rangle \\
&= \arg \max_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \lambda_{\min}(Q^T \mathcal{A}^T(d)Q) \\
&= \arg \max_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \min_{\substack{m \in \mathbb{R}^t \\ \|m\|=1}} m^T Q^T \mathcal{A}^T(d)Qm \\
&= \arg \max_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \min_{\substack{m \in \mathbb{R}^t \\ \|m\|=1}} \langle Qmm^T Q^T, \mathcal{A}^T(d) \rangle \\
&= \arg \max_{\substack{d \in \mathbb{R}^m \\ \|d\|=1}} \min_{\substack{m \in \mathbb{R}^t \\ \|m\|=1}} \langle \mathcal{A}(Qmm^T Q^T), d \rangle
\end{aligned} \tag{4.11}$$

A linear program will be devised to solve equation 4.11.

This LP is based on the observation that in the last identity $Qmm^T Q$ can be written as $\sum_{i,j} M_{i,j} q_i q_j^T$ with $M = mm^T$. Also because $\|m\| = 1$, each diagonal element of M is non-negative and $\sum_i M_{i,i} = 1$. Hence the last identity in 4.11 can be written as a convex sum over $q_i q_i^T$ plus some cross terms.

In experiments it was confirmed that a good strategy is to minimize with d the elements in the convex sum and constraint d to be such that the cross terms are zero. The LP used:

$$\begin{aligned}
&\min_{\substack{\delta \in \mathbb{R} \\ d \in \mathbb{R}^m}} \delta \\
&\text{s. t. } \delta + d^T \mathcal{A}(q_i q_i^T) \geq 0 \quad (\forall i = 1, \dots, t) \\
&\quad d^T \mathcal{A}(q_i q_j^T) = 0 \quad (\forall i, j = 1, \dots, t \quad i \neq j) \\
&\quad -1 \leq d \leq 1
\end{aligned} \tag{4.12}$$

(This LP is similar to the quadratic program in the successive quadratic programming approach explained in [21] and [22] and shortly described in section 4.3. However in this LP there are less constraints on d with respect to other eigenvalues than λ_{\max} . Also the QP in [22] uses a quadratic objective to stop d in the optimum from being too large, but that is not a requirement in this case because only the direction of d is important.)

If there are active constraints, i.e. $P_a \neq 0$, then linear program 4.12 can be

augmented with the active constraints:

$$\begin{aligned}
& \min_{\substack{\delta \in \mathbb{R} \\ d \in \mathbb{R}^m}} \delta \\
& \text{s. t. } \delta + d^T \mathcal{A}(q_i q_i^T) \geq 0 \quad (\forall i = 1, \dots, t) \\
& \quad d^T \mathcal{A}(q_i q_j^T) = 0 \quad (\forall i, j = 1, \dots, t \quad i \neq j) \\
& \quad P_a d \geq 0 \\
& \quad -1 \leq d \leq 1
\end{aligned} \tag{4.13}$$

This linear program has $\frac{t(t+1)}{2}$ -constraints in m -variables, hence is actually quite small and suitable to be solved through a simplex solver.

Example 4.1.5. Consider problem 4.1 with $m = 2$, no inequalities and:

$$C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, A_1 = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, A_2 = \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix} \tag{4.14}$$

Analytically the eigenvalues are:

$$\begin{aligned}
\lambda_1 &= 1 - z_1 \\
\lambda_2 &= 1 - 2z_2
\end{aligned} \tag{4.15}$$

Clearly the problem is unbounded in the direction $d = (2, 1)^T$. Solving linear program 4.12 with $q_1 = e_1$ and $q_2 = e_2$ yields:

$$\begin{aligned}
\mathcal{A}(q_1 q_1^T) &= (1, 0)^T \\
\mathcal{A}(q_2 q_2^T) &= (0, 2)^T \\
d &= \left(1, \frac{1}{2}\right)^T \\
\delta &= \nabla_d \lambda_{\max}(z) = -1
\end{aligned} \tag{4.16}$$

4.1.3 Cut infeasible descent directions

If the problem is constrained, i.e. $Pz \geq 0$ with P not trivial, then the region in which $\lambda_{\max}(C - \mathcal{A}^T(z))$ is minimized is within a convex polytope. That is because $Pz \geq 0$ implies that each row of P describes a half space, and all those together in $Pz \geq 0$ is the intersection of all those half spaces. To continue, the following definitions give some more grasp of the problem.

Definition 4.1.6. The set $Pz \geq 0$ is called the feasible set and is denote by \mathcal{F} , the interior of \mathcal{F} , i.e. $Pz > 0$, is called the strict feasible set and is denoted as $\text{int}(\mathcal{F})$ and the boundary of \mathcal{F} is denoted as $\partial \mathcal{F}$.

Definition 4.1.7. The active constraints on $Pz \geq 0$ is the index set \mathcal{I} , such that $p_i^T z = 0$, with p_i the i th row of P as column vector, for all $i \in \mathcal{I}$. The active constraints matrix P_a has as rows exactly the active rows of P . That is $P_a^T = (p_{I_1} \mid \dots \mid p_{I_l})$.

If there are no active constraints at $z_i \in \mathcal{F}$, then $z_i \in \text{int}(\mathcal{F})$.

Constraint p_i (or short i) is said to be violated if $p_i^T z < 0$. Hence if there are violated constraints at z , then $z \notin \mathcal{F}$.

Clearly if $z \in \text{int}(\mathcal{F})$, then there is no constraint on search direction d , however there is a constraint on the maximum step size. This point is something that will be discussed further in section 4.2.1.

If for a given $z \in \partial \mathcal{F}$ with active constraints, then those constraints also put a constraint on descent direction d . Specially, if P_a is the active constraints matrix at z , then $P_a d \geq 0$ for any feasible descent direction d .

Hence, for a given descent direction d , vector $\tilde{d} \in \mathbb{R}^m$ is sought, which is a feasible descent direction, which minimizes $\|d - \tilde{d}\|_2$. If all rows of P_a are orthogonal, then that is easy by using some kind of Gram-Schmidt decomposition and discard the violating directions in P . However, the difficulty of this problem is that the rows of P_a are not necessarily orthogonal.

A strategy is to calculate d as described in the previous subsections and then try to find a vector $r \in \mathbb{R}^m$ of small norm such that $\tilde{d} = d + r$ is feasible. This can be solved by the following quadratic program:

$$\begin{aligned} \min_{r \in \mathbb{R}^m} r^T r \\ \text{s. t. } P_a(d + r) \geq 0 \end{aligned} \tag{4.17}$$

Quadratic program 4.17 has l constraints in m variables, which is not too large and can be solved fast by a simplex solver. (See for instance [33]; the idea is akin to linear simplex by rotating through vertexes of the polytope described by $P(d + r) = 0$ while also checking the edges for optimality.)

Algorithm 4.1.8 Algorithm to cut infeasible directions from d

Require:

$z \in \mathbb{R}^m$: current location

$P \in \mathbb{R}^{l \times m}$: constraint matrix, i.e. $Pz \geq 0$

$d \in \mathbb{R}^m$: descent direction

Ensure:

$\tilde{d} \in \mathbb{R}^m$: feasible descent direction

- 1: **function** CUT INFEASIBLE(z, P, d)
 - 2: $P_a \leftarrow \{p_i \text{ s. t. } (Pz)_i = 0\}$, i.e. the set of active constraints
 - 3: Solve QP 4.17 for $r \in \mathbb{R}^m$
 - 4: $\tilde{d} \leftarrow d + r$
 - 5: **return** \tilde{d}
 - 6: **end function**
-

4.2 Determine max step size

After determining a descent direction $d \in \mathbb{R}^m$ in section 4.1, the next concern is computing a optimal step size $t^* \in \mathbb{R}^+$ that minimizes:

$$f(t) : \mathbb{R}^+ \rightarrow \mathbb{R}, t \mapsto \lambda_{max}(z + td) = \lambda_{max}(C - \mathcal{A}^T(z + td)) \quad (4.18)$$

Such that $P(z + t^*d) \geq 0$. That is, the step size should not be too large such that $z + td$ becomes infeasible. Obvious if there are no inequalities, then the problem is unbounded.

This is a non-smooth convex problem, where non-smoothness arises if eigenvalue curves cross. For large sparse problems, f is hard to compute directly because the eigenvalues tend to cluster. That is, the gap between two distinct eigenvalues becomes very small. Hence the power method will converge very slow and neither the Lanczos method works efficiently because of the large multiplicity of the eigenvalues. Moreover, the matrix $C - \mathcal{A}^T(z + td)$ is near the optimum far from diagonal dominant, hence the Davidson method does not work as well.

The first subsection gives a short analysis of the problem and states some notation which will be used further in this section. In the second subsection an algorithm is developed for large problems which cannot be tackled by classic line (secant or Newton method) searches.

4.2.1 Analysis of the problem

Some words on notation: by theorem 2.2.12 it is possible to label all eigencurves of $C - \mathcal{A}^T(z + td)$ at $t = 0$ as:

$$\lambda_1(z) \geq \dots \geq \lambda_n(z) \quad (4.19)$$

With $t \rightarrow \lambda_i(z + td)$ smooth. If $\lambda_{max}(z)$ has multiplicity τ , then $\lambda_1(z) = \dots = \lambda_\tau(z) > \lambda_{\tau+1}(z)$. Note that this labeling of the eigencurves does not imply that $\lambda_2(z + td) \leq \lambda_1(z + td)$ for all $t \in \mathbb{R}^+$. In fact minimizing equation 4.18 involves the computation of t^* such that $\lambda_1(z + t^*d) \leq \lambda_{max}(z + t^*d)$.

Moreover, in order to handle the problem some kind of derivatives are necessary:

Definition 4.2.1. *Observe that the forward derivative of $f(t)$ to t exists and is defined as:*

$$D_t f(t) = \lim_{h \downarrow 0} \frac{\lambda_{max}(z + td + hd) - \lambda_{max}(z + td)}{h} \quad (4.20)$$

Similar the backward derivative is defined as:

$$D'_t f(t) = \lim_{h \uparrow 0} \frac{\lambda_{max}(z + td + hd) - \lambda_{max}(z + td)}{h} \quad (4.21)$$

From these observations the following observation can be deduced:

Theorem 4.2.2. *Consider the function f from equation 4.18. If for $t' \in \mathbb{R}^+$, function f is non-smooth at t' , then $D_t f(t')$ is left discontinuous and at least two eigencurves $\lambda_i(t)$, $\lambda_j(t)$ with $i \neq j$ cross each other at t' .*

Proof. This follows direct from theorem 2.2.12: All eigencurves can be labeled such that they are smooth. Hence if f is non-smooth at t' , then $\lambda_{\max}(z + t'd)$ is non-simple and thus at least two distinct eigencurves cross each other. Moreover, there is a neighborhood N of t' such that for all $t \in N$, $t \geq t'$, $D_t f(t)$ follows one eigencurve other than the distinct curve followed $D_t f(t)$ for $t < t'$. Thus for $t \in N \cap (-\infty, t')$, $D_t f(t)$ is continuous and for $t \in N \cap [t', \infty)$, $D_t f(t)$ is continuous, \square

Also:

Theorem 4.2.3. *Consider the function f from equation 4.18. Denote the minimizer with t^* , i.e. $f(t^*) \leq f(t)$ for all $t \in \mathbb{R}^+$. Then the following optimality condition holds at t^* , given that $z + t^*d \in \text{int}(\mathcal{F})$:*

$$\begin{aligned} D_t f(t) &\leq 0 \quad (\forall t \in [0, t^*]) \\ D_t f(t) &\geq 0 \quad (\forall t \in [t^*, \infty)) \end{aligned} \tag{4.22}$$

Or equivalent:

$$\begin{aligned} D'_t f(t) &\leq 0 \quad (\forall t \in [0, t^*]) \\ D'_t f(t) &\geq 0 \quad (\forall t \in (t^*, \infty)) \end{aligned} \tag{4.23}$$

Proof. All these conditions follow from noting that f is convex and the identity:

$$f(t) = \int_0^t D_t f(\tau) \mathbf{d}\tau + f(0) \tag{4.24}$$

\square

By theorem 2.2.17 it follows that:

$$D_t f(t) = \max_{\substack{u \in \mathbb{R}^m \\ (C - \mathcal{A}^T(z+td))u = \lambda_{\max}(z+td)u \\ \|u\|_2=1}} -u^T \mathcal{A}^T(d)u \tag{4.25}$$

And similar:

$$D'_t f(t) = - \max_{\substack{u \in \mathbb{R}^m \\ (C - \mathcal{A}^T(z+td))u = \lambda_{\max}(z+td)u \\ \|u\|_2=1}} u^T \mathcal{A}^T(d)u \tag{4.26}$$

Hence the computation of $D_t f(t)$ can be seen as an eigenvalue problem: Let $q_1(t), \dots, q_{l(t)}(t)$ be a basis for the eigenspace to the largest eigenvalue, then any $u \in \mathbb{R}^m$ can be expressed as:

$$u = Q(t)z \quad \text{with} \quad Q(t) = (q_1(t), \dots, q_{l(t)}(t)) \quad \text{and} \quad \|z\| = 1 \quad (4.27)$$

Equation 4.25 is then equivalent to:

$$D_t f(t) = -\lambda_{\min}(Q(t)^T \mathcal{A}^T(d)Q(t)) \quad (4.28)$$

Similar the backward derivative:

$$D'_t f(t) = -\lambda_{\max}(Q(t)^T \mathcal{A}^T(d)Q(t)) \quad (4.29)$$

And checking the optimality conditions 4.2.3:

$$\begin{aligned} D_t f(t^*) &= -\lambda_{\min}(Q(t^*)^T \mathcal{A}^T(d)Q(t^*)) \geq 0 \\ D'_t f(t^*) &= -\lambda_{\max}(Q(t^*)^T \mathcal{A}^T(d)Q(t^*)) \geq 0 \end{aligned} \quad (4.30)$$

The bounded case

If there exists a $\hat{t} > 0$ such that $z + td \notin \mathcal{F}$ ($P(z + td) \not\geq 0$) for all $t > \hat{t}$, then the problem is bounded. Hence optimality conditions 4.30 do not hold necessary in the optimum.

It can be that $D_t f(\hat{t}) < 0$, thus that in the unbounded case $t^* > \hat{t}$, but that would imply $z + t^*d \notin \mathcal{F}$. Therefore t^* should be bounded by \hat{t} .

To compute the minimal \hat{t} (further denoted as t_{\max} : Let \mathcal{I} be the index set such that $P_{\mathcal{I}}d < 0$, that is the row index set of all constraints which will be violated if $t \rightarrow \infty$. Then t_{\max} is:

$$t_{\max} = \min_{i \in \mathcal{I}} -\frac{P_i z}{P_i d} \quad (4.31)$$

4.2.2 Line search methods based on intersections

If all eigencurves are assumed to be linear then the minimum for f is obtained for certain at an intersection between two eigencurves. Hence a good start for minimizing f is by using the algorithm developed in section 3.4 to find the nearest intersection point for the largest eigencurve.

After a nearby intersection of two eigencurves is found at \hat{t} , the following step is to check optimality conditions from theorem 4.2.3. That can be done through checking the intersecting curve, if it has a positive slope then because of smoothness (theorem 2.2.12) it follows that $D_t f(t) \geq 0$ for $t \geq \hat{t}$.

Therefore the main idea is to iterative compute the nearest intersection \hat{t} , if the slope of the crossing eigencurve is positive, then stop, else assume that the crossing eigencurve is the new maximum and search for the next intersection.

The algorithm based on section 3.4 to compute the nearest intersection:

Algorithm 4.2.4 SASPA to compute nearest intersection

Require:

- $A \in \mathcal{S}^n$: matrix
- $\dot{A} \in \mathcal{S}^n$: perturbation
- $\lambda \in \mathbb{R}$: largest eigenvalue at $t = 0$
- $q_0 \in \mathbb{R}^n$: corresponding eigenvector
- $\epsilon_{RQI} \in \mathbb{R}_{>0}$: threshold from where to start Rayleigh Quotient Iteration
- $\epsilon_{tol} \in \mathbb{R}_{>0}$: tolerance
- $\epsilon_{stop} \in \mathbb{R}$: additional stop criteria
- k_{\min}, k_{\max} minimum and maximum size of search space for restart
- $iter_{\max}$, maximum number of iterations

Ensure:

- $\hat{t} \in \mathbb{R}_+$: intersection point
- $\hat{\lambda} \in \mathbb{R}$: nearest crossing eigenvalues at $t = 0$
- $\hat{q} \in \mathbb{R}^n$: corresponding eigenvector

- 1: $iter \leftarrow 1, num \leftarrow 0, s \leftarrow \lambda, x = (1, \dots, 1)^T \in \mathbb{R}^n, \hat{t} \leftarrow \infty$
- 2: **while** $iter \leq iter_{\max}$ **do**
 - 3: $b \leftarrow \dot{A}x$ ▷ Extend search space
 - 4: $b \leftarrow \text{GS}([q_0, \dots, q_{num}], b)$ ▷ Deflate b with q_0, \dots, q_{num}
 - 5: Solve $(A - sI)x = b$ for $x \in \mathbb{R}^n$
 - 6: $x \leftarrow \text{GS}(X, x), X \leftarrow [X, x]$ ▷ Select most promising
 - 7: $S \leftarrow X^T A X$
 - 8: Compute all eigenpairs (s_i, x_i) of S
 - 9: Compute for all pairs corresponding cross points: $t_i \leftarrow \frac{\lambda - s_i}{x_i^T X^T \dot{A} X x_i - q_0^T \dot{A} q_0}$
 - 10: Reorder all triplets (t_i, s_i, x_i) such that $t_i \geq t_{num}$ for all $i \leq j$
 - 11: $(t, s, x) \leftarrow (t_1, s_1, X x_1)$
 - 12: $x \leftarrow \frac{x}{\|x\|}$ ▷ Elective Rayleigh Quotient Iteration
 - 13: **if** $\|Ax - sx\| \leq \epsilon_{RQI}$ **then**
 - 14: **while** $\|Ax - sx\| \geq \epsilon_{tol}$ **do**
 - 15: Solve $(A - sI)\dot{x} = x$ for $\dot{x} \in \mathbb{R}^n$
 - 16: $x \leftarrow \frac{\dot{x}}{\|\dot{x}\|}$
 - 17: $s \leftarrow x^T A x$
 - 18: **end while**
 - 19: **end if** ▷ Check convergence and deflate
 - 20: **if** $\|Ax - sx\| \leq \epsilon_{tol}$ **then**
 - 21: $num \leftarrow num + 1$
 - 22: $q_{num} \leftarrow x$
 - 23: Deflate X with q_{num}
 - 24: **if** $t \leq \hat{t}$ **then**

```

25:          $\hat{\lambda} \leftarrow s, \hat{q} \leftarrow q_{num}, \hat{t} \leftarrow t$ 
26:     end if
27:     if  $t \leq \epsilon_{stop}$  then
28:         return  $(\hat{t}, \hat{\lambda}, \hat{q})$ 
29:     end if
30:      $(t, s, x) \leftarrow (t_2, s_2, Xx_2)$ 
31:      $x \leftarrow \frac{x}{\|x\|}$ 
32: end if
▷ Check size search space and implicit restart
33: if #columns of  $X$  is more than  $k_{max}$  then
34:      $X \leftarrow [Xx_1, Xx_2, \dots, Xx_{k_{min}}]$ 
35:     Orthogonalize and normalize  $X$ 
36: end if
37:      $iter \leftarrow iter + 1$ 
38: end while

```

Note that it is very similar to algorithm 3.5.2; the only difference is the additional stopping criteria based on the discussion of the numerical results in section 3.4.

Moreover, if the maximum eigencurve is a parabola, which can happen near the optimum, then there is not a nearby intersection. Hence at the nearest intersection the optimality conditions 4.2.3 are not full filled.

Note that the parabola which interpolates $\lambda_{max}(t)$ is convex because the Hessian (2.2.11) of the largest eigenvalue is almost everywhere equal or greater than 0. Therefore a viable tactic is to compute the optimum through a parabola fit.

Denote the nearest intersection point by $\hat{t} \in \mathbb{R}$. The parabola interpolation for $\lambda_{max}(0)$, $\lambda_{max}(\hat{t})$ and forward derivative $f(0)$ at $t = 0$ is given by:

$$\lambda_{max}(t) \approx t^2 \left(\frac{\lambda_{max}(\hat{t}) - \lambda_{max}(0)}{\hat{t}^2} - \frac{f(0)}{\hat{t}} \right) + tf(0) + \lambda_{max}(0) \quad (4.32)$$

If $\frac{\lambda_{max}(\hat{t}) - \lambda_{max}(0)}{\hat{t}^2} - \frac{f(0)}{\hat{t}} > 0$, then the minimum of that parabola is attained at:

$$t^* = \frac{-f(0)}{2 \left(\frac{\lambda_{max}(\hat{t}) - \lambda_{max}(0)}{\hat{t}^2} - \frac{f(0)}{\hat{t}} \right)} \quad (4.33)$$

If $\frac{\lambda_{max}(\hat{t}) - \lambda_{max}(0)}{\hat{t}^2} - \frac{f(0)}{\hat{t}} < 0$, then the parabola has a maximum, however this is not something bound to happen because the parabola tends to be convex and not concave. If this happens then the interpolating interval $[0, \hat{t}]$ is very small and the interpolation parabola is a bad fit. In that case a better method would be a line search through bisection.

The full algorithm to determine the step size for unbounded problems:

Algorithm 4.2.5 Method to compute optimal step size

Require: $A = C - \mathcal{A}^T(z)$: current location, matrix form $\mathcal{A}^T(d) \in \mathcal{S}^n$: search direction, matrix form $f(0) \in \mathbb{R}$: forward derivative at $t = 0$ $t_{\max} \in \mathbb{R}$: see equation 4.31**Ensure:** $t^* \in \mathbb{R}$: minimizer of $\lambda_{\max}(A - t\mathcal{A}^T(d))$ $Q(t^*) \in \mathbb{R}^{n \times l}$: orthogonal normalized basis of the eigenspace to λ_{\max}

- 1: Compute nearest intersection (t, s, q) with algorithm 4.2.4
 - 2: $t \leftarrow \min(t, t_{\max})$
 - 3: Compute orthonormal eigenbasis $Q(t)$ for $A - t\mathcal{A}^T(d)$ corresponding to $\lambda_{\max}(t)$
 - 4: Compute optimality conditions 4.30:
 $f \leftarrow -\lambda_{\min}(Q(t)^T \mathcal{A}^T(d)Q(t))$
 $b \leftarrow \lambda_{\max}(Q(t)^T \mathcal{A}^T(d)Q(t))$
 - 5: **while** $f > 0$ and $b < 0$ **do** ▷ Proceed by interpolation
 - 6: $t' \leftarrow \frac{-f(0)}{2(\frac{\lambda_{\max}(t) - \lambda_{\max}(0)}{t} - f(0))}$
 - 7: Compute orthonormal eigenbasis $Q(t')$ for $A - t'\mathcal{A}^T(d)$ corresponding to $\lambda_{\max}(t')$
 - 8: Compute optimality conditions 4.30:
 $f \leftarrow -\lambda_{\min}(Q(t')^T \mathcal{A}^T(d)Q(t'))$
 $b \leftarrow \lambda_{\max}(Q(t')^T \mathcal{A}^T(d)Q(t'))$
 - 9: $t \leftarrow t'$
 - 10: **end while**
 - 11: **if** $t < t_{\max}$ and $f < 0$ **then** ▷ Further improvement is possible
 - 12: $A' \leftarrow A - t\mathcal{A}^T(d)$
 - 13: Repeat algorithm 4.2.5 with: $A \leftarrow A'$ and $f(0) \leftarrow f$; returns t^* and $Q(t^*)$
 - 14: $t \leftarrow t + t^*$
 - 15: $Q(t) \leftarrow Q(t^*)$
 - 16: **end if**
 - 17: $t^* \leftarrow t$
 - 18: $Q(t^*) \leftarrow Q(t)$
-

A practical note on implementation: Often a basis for the eigenspace corresponding to the largest eigenvalue has to be computed. Luckily the matrix is permuted in such a way that it is known which eigencurves meet after the permutation. Hence it is expected that the eigenvectors which compromise the basis, are somewhat near the eigenvectors of the eigencurves at the beginning. Therefore the basis can be constructed by updating the old eigenvectors to the new matrix.

There are a number of iterative methods available to compute a basis includ-

ing: QR Method, Lanczos, Subspace Accelerated Rayleigh Quotient Iteration (SARQI) and Jacobi–Davidson [30].

The first method works well for small matrices and the second method for which no prior information on the eigenvectors are known. However, as described we have already quite a good idea on how the basis should look like. Therefore the choice is limited to the last three methods.

During simulation were the best results obtained by using Jacobi–Davidson in combination with `BiCGstab`.

4.3 Combining everything

The full method to solve eigenvalue equation 4.1 based on the previous sections is iterative: Compute steepest descent direction $d \in \mathbb{R}^m$, compute step size $t^* \in \mathbb{R}$, update $z \leftarrow z + t^*d$ and repeat.

The full algorithm written out:

Algorithm 4.3.1 Algorithm to solve eigenvalue problem 4.1

Require:

$C \in \mathcal{S}^n$: initial matrix

$A_1, \dots, A_m \in \mathcal{S}^n$: matrix component of \mathcal{A}^T

$P \in \mathbb{R}^{l \times m}$: constraints

Ensure:

$z \in \mathbb{R}^m$: argument that minimizes equation 4.1

$A \in \mathcal{S}^n$: the matrix $C - \mathcal{A}^T(z)$

```

1:  $z \leftarrow 0 \in \mathbb{R}^m$ ,  $A \leftarrow C$ 
2: Compute a basis  $Q$  for the eigenspace corresponding to  $\lambda_{\max}(C)$ 
3: while not converged do
4:   if  $\text{rank}(Q) = 1$  then ▷  $\lambda_{\max}$  is simple
5:     Compute  $d \in \mathbb{R}^m$  as  $-\nabla \lambda_{\max}(z) = \mathcal{A}(QQ^T)$ 
6:   else
7:     Compute  $\delta \in \mathbb{R}$  and  $d \in \mathbb{R}^m$  from LP 4.12
8:   end if
9:    $d \leftarrow \text{CUT INFEASIBLE}(z, P, d)$  ▷ See algorithm 4.1.8
10:   $f \leftarrow -\lambda_{\min}(Q^T \mathcal{A}^T(d)Q)$ 
11:  if  $f \geq 0$  then ▷ Check convergence
12:    return  $(z, A)$ 
13:  end if
14:  Compute  $t_{\max}$  by equation 4.31
15:  Compute  $t^*$  and  $Q$  by algorithm 4.2.5
16:   $z \leftarrow z + t^*d$ 
17:   $A \leftarrow A - t^* \mathcal{A}^T(d)$ 
18: end while

```

An implementation is included in appendix A.3.

On implementation: When the first basis Q that has to be computed, no prior information is yet known. Thus a suitable iterative method is Lanczos Method. During simulations the best results were obtained by combining Lanczos Method with Jacobi–Davidson.

When the algorithm is in the main loop the eigenspace is updated each time during algorithm 4.2.5.

4.3.1 Comparison to other methods

Algorithm 4.3.1 is not the only known method to solve eigenvalue problems. Not all known algorithms will be discussed in full detail; but will give a brief overview of other known algorithms and highlight the main differences.

Successive Quadratic Programming Algorithm

Overton’s method, also known as Successive Quadratic Programming Algorithm, developed around 1988, is an iterative method to solve eigenvalue problems of the form 4.1 [21] [22].

It works by solving each iteration a QP similar to 4.12, to compute a step direction d and step t in one go. Denote y_i as the argument in the i th iteration. Then to compute y_{i+1} : let λ_1 the largest eigenvalue of $C - \mathcal{A}^T(y_i)$, q_1, \dots, q_k be eigenvectors corresponding to the largest k eigenvalues, and q_{k+1}, \dots, q_m the other eigenvectors, $W \in \mathcal{S}_{\geq 0}^n$:

$$\begin{aligned}
 & \min_{\substack{\delta \in \mathbb{R} \\ d \in \mathbb{R}^m}} \delta + d^T W d \\
 & \text{s. t. } \delta + d^T \mathcal{A}(q_i q_i^T) = \lambda_i - \lambda_1 \quad (\forall i = 1, \dots, k) \\
 & \quad \delta + d^T \mathcal{A}(q_i q_j^T) = 0 \quad (\forall i, j = 1, \dots, k, \quad i \neq j) \\
 & \quad \delta + d^T \mathcal{A}(q_i q_i^T) \geq \lambda_i - \lambda_1 \quad (\forall i = k + 1, \dots, m) \\
 & \quad -\rho \leq d \leq \rho
 \end{aligned} \tag{4.34}$$

Then $y_{i+1} = y_i + d$.

Note the similarity with the earlier devised LP: If taken k to be the multiplicity of λ_1 , then the first two equality’s are similar to the LP 4.12.

The argument for taking k larger is that near the optimum those eigencurves most likely coalesce. Hence it seems a good idea to let the linearized eigencurves join when taking step d .

The third inequality is to limit the size of d to be such that the linearized largest eigencurve will not cross any of the other eigencurves. Also the fourth inequality is to limit d from growing too large.

Moreover, in the objective the positive semidefinite matrix W also prevents d from growing too large. If W is chosen to be the Hessian matrix, then the method can be shown to be a Newton method.

Overton’s method has been used for many different applications and is seen as a robust method. However one of its drawbacks is the necessity of a full

eigendecomposition of $C - \mathcal{A}^T(y_1)$ which limits the method to small problems. Furthermore, solving a QP instead of a LP is considerable harder, and when y_i is far from optimum, then it can be argued that those extra computation cycles are unnecessary.

Spectral Bundle Method

The Spectral Bundle Method has been developed by Helmberg and Randl in 1997 [14] and is also a dual solver algorithm for semidefinite optimization problems. That is, during each iteration a full feasible dual solution is used as starting point from which a descent direction is computed, similar as algorithm 4.3.1.

Starting from a semidefinite optimization problem, first they assume that the rank of X is constant such that the additional constraint:

$$\langle I, X \rangle = a \tag{4.35}$$

can be added. This implies that the dual is extended:

$$\begin{aligned} \min_{y \in \mathbb{R}^m, \lambda \in \mathbb{R}} \quad & a\lambda + b^T y \\ \text{s. t.} \quad & C - \mathcal{A}^T(y) - \lambda I \preceq 0 \end{aligned} \tag{4.36}$$

Which has for optimum (λ^*, z^*) , $\lambda^* = \lambda_{\max}(C - \mathcal{A}^T(y))$. Hence the dual problem can be written as:

$$\min_{y \in \mathbb{R}^m} a\lambda_{\max}(C - \mathcal{A}^T(y)) + b^T y \tag{4.37}$$

Note that this dual is a mixed eigenvalue optimization and linear optimization problem.

Also note that the subgradient of:

$$f(y) := a\lambda_{\max}(C - \mathcal{A}^T(y)) + b^T y \tag{4.38}$$

can be shown to be:

$$\partial f(y) = \{b - \mathcal{A}(W) \text{ s. t. } \langle W, C - \mathcal{A}^T(y) \rangle = \lambda_{\max}(C - \mathcal{A}^T(y)), \text{Tr}(W) = 1, W \succeq 0\} \tag{4.39}$$

Next they proceed solving this optimization problem by projecting the eigenvalue optimization problem on some subspace, which consists of vectors computed in earlier iterations such as eigenvectors to the then largest eigenvalue. Then on that projected space they try to minimize $f(y + d) + \frac{u}{2} \langle d, d \rangle$, which is similar to minimizing $f(y + d)$ with the constraint that $\|d\|_2$ becomes not too large. The parameter u is a check on how large $\|d\|$ may grow. The reason that u should not be too small and let $\|d\|$ become too large is because then the difference between f and the projection of f becomes too large.

For a full discussion, see [14] or [4].

The main ingredients of the Spectral Bundle method are the proximity idea, the projection of f onto a subspace which is local closely related to the original, and the use of subgradients to minimize f local.

The similarities with 4.3.1 are that it works by using subgradients and trying to follow the largest eigenvalues. It also works by remembering information from earlier iterations on the eigenbasis, but instead of algorithm 4.3.1 it works by keeping the change local. Algorithm 4.3.1 has no such limit and uses multiple Jacobi-Davidson steps to correct its eigenbasis.

Chapter 5

Numerical results

In this chapter some numerical results of algorithm 4.3.1 are presented.

The test cases are those described in section 1.4 and section 2.1.2.

5.1 Max-cut

This section follows section 1.4.1.

In section 2.1.2 was the max-cut eigenvalue version presented as:

$$\inf_{z \in \mathbb{R}^{|V|-1}} |V| \lambda_{\max} \left(\frac{1}{4} LG - \sum_i z_i \text{diag}(q_i) \right) \quad (5.1)$$

The size of the instance only depends on the number of edges, hence this problem is well suited for benchmarking. In this section some convergence results for the graphs `toruspm3--8--50` and `toruspm3--15--50` are presented.

The graphs `toruspm3--8--50` and `toruspm3--15--50` were taken from the seventh DIMACS Implementation Challenge [24]. They arise from the study of spin glass model through the Ising model.

In the Ising model there are n particles with each spin ± 1 . Particles will align their spin in such a way to minimize the total energy. The total energy depends on, with $s_i = \pm 1$ the spin of particle i and $w_{i,j}$ the interaction between particle i and j :

$$H = \sum_{\substack{i,j \\ i \neq j}}^n w_{i,j} s_i s_j \quad (5.2)$$

In typical simulations all particles are placed on a grid, and for any two disconnected particles $w_{i,j} = 0$, and for the connected particles $w_{i,j}$ is taken to be random according to a given distribution. For instance, $w_{i,j} = \pm 1$ with 50% probability.

Table 5.1: Results for solving max-cut instances

Name:	Max-cut:	Iterations to 10%:	To 1%:	To 0.1%:
<code>toruspm3--8--50</code>	527.8	19	37	184
<code>toruspm3--15--50</code>	3474.4	62	—	—

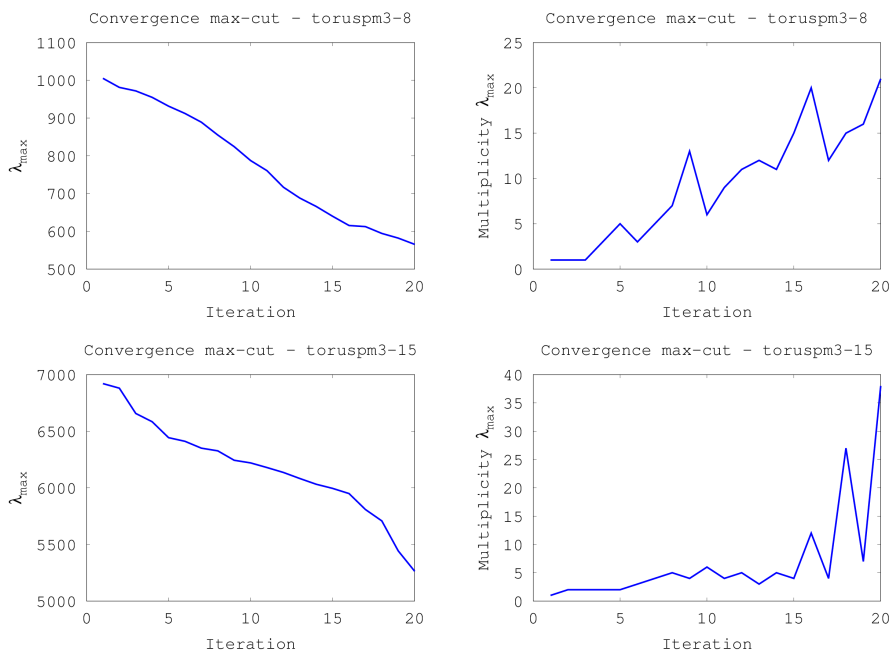


Figure 5.1: Convergence results on maxcut instances

The maximum energy of the particles on such a grid is a maximal cut, hence these kind of models can be studied computational by solving the max-cut problem. For more information, see [9].

The graphs `toruspm3--8--50` and `toruspm3--15--50` are from such models. They represent a grid on a three dimensional torus of size 8 and 15 (that is, `toruspm3--8--50` has $|V| = 8^3 = 512$ vertices and `toruspm3--15--50` has $|V| = 15^3 = 3375$ vertices) with edge weight $w_{i,j} = \pm 1$ with 50% probability.

The results of solving these max-cut instances with algorithm 4.3.1 are presented in table 5.1:

In figure 5.1 convergence of the algorithm is shown. It can be seen that the initial convergence is fast and that the rank of the largest eigenvalue is increasing as the algorithm gets closer to the optimum.

These simulations show that algorithm 4.3.1 has initial fast convergence and is able to solve large eigenvalue optimization problems.

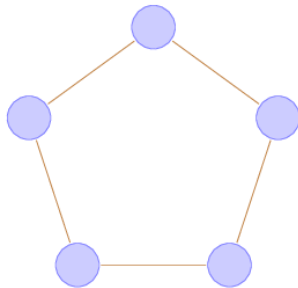


Figure 5.2: A cycle graph of order 5

5.2 Lovász theta number

This section follows section 1.4.2.

The theta number for a graph $G = (V, E)$ can be calculated through the following eigenvalue optimization problem:

$$\min_{y \in \mathbb{R}^{|E|}} \lambda_{\max}(J - \mathcal{E}^T(y)) \quad (5.3)$$

In which $\mathcal{E}^T(y) = \sum_{(i,j) \in E} y_{i,j} E_{i,j}$ with $E_{i,j} = e_i e_j^T + e_j e_i^T$. Note that this is equation is the same as 1.63.

The theta number for a few different types of graphs is known, hence those graphs are perfect suited to test algorithm 4.3.1 on. In this section some numerical results on cycle graphs and Kneser $(2, k)$ graphs, for which theta is known.

Those two set of graphs are interesting not because the theta number is hard to compute for those, it is not, but because they highlight two features of the method: Graphs with n nodes are represented by large matrices of order n , hence it is expected that for large graphs the computation of the eigenbasis forms a bottleneck. Graphs with m edges result in a semidefinite problem with m variables, which implies that the linear program in the method for computing a descent direction has m variables. Hence it is expected that for graphs with many edges solving the linear program forms a bottleneck. The set of cycle graphs contains graphs which have as many vertices as edges, it is expected that computing the eigenbasis will form the most time consuming step for large values of n . On the other hand the set of Kneser $(2, k)$ graphs contains graphs for which the number of edges is exponential larger than the number of vertices, therefore it is expected that for those graphs solving the linear program will be the most time consuming.

Table 5.2: Results for computation of $\theta(C_n)$

Order:	Vertices:	Edges:	θ :	Iterations:	Multiplicity of λ_{\max} :	Time (s):
10	10	10	5.0000	2	2	0.04
100	100	100	50.000	2	2	0.06
101	101	101	50.488	2	3	0.04
500	500	500	250.00	2	2	0.86
501	501	501	250.50	2	3	0.62
1000	1000	1000	500.00	2	2	6.52
1001	1001	1001	500.50	2	3	5.17
5000	5000	5000	2500.0	2	5	1122
5001	5001	5001	2500.5	2	4	1035

Cycle graphs

A cycle graph of order n has n vertices and $n + 1$ edges such that each vertex is connected with two neighbours and the graph forms a cycle. For an example, see figure 5.2 which depicts a cycle graph of order 5.

The theta number for cycle graphs of order n can be shown to be [17]:

$$\theta(C_n) = \begin{cases} \frac{n}{2} & \text{if } n \text{ is even} \\ \frac{n \cos(\pi/n)}{1 + \cos(\pi/n)} & \text{else} \end{cases} \quad (5.4)$$

Results of computation of the theta number for cycle graphs of varying order are presented in table 5.2:

Note that the computation time presented in table 5.2 are not a good representation on the quality of the algorithm because the computer on which these instances were run was also idling with other tasks. The column is only shown to give a qualitative idea on how well the algorithm scales and typical run times on an ordinary laptop.

The computation time for the large instances surge drastic. That is because during the algorithm an eigenbasis to the largest eigenvalue has to be computed. For the smaller instances that is done by using Lanczos method. For the larger instances that is no longer feasible hence the basis is computed by using an implementation of Jacobi–Davidson based on [30].

Kneser graphs

A Kneser graph (n, k) is such that its vertices correspond to k -element subsets of a set of n elements. Two vertices are connected if the their subsets are disjoint. For an example, see figure 5.3.

The theta number for Kneser graph (n, k) is:

$$\theta(K_{n,k}) = \binom{n-1}{k-1} \quad (5.5)$$

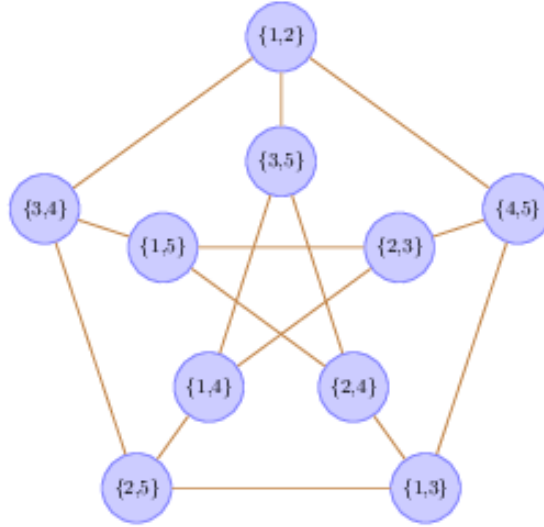


Figure 5.3: Kneser graph $(5, 2)$, also known as the Petersen graph

Table 5.3: Results for computation of $\theta(K_{2,k})$

k :	Vertices:	Edges:	θ :	Iterations:	Multiplicity of λ_{\max} :	Time (s):
5	10	15	4	2	5	0.46
10	45	630	9	3	10	0.29
15	105	4095	14	2	15	20.59
20	190	14535	19	2	20	567.15

Some results of computation of the theta number for different Kneser graphs $(2, k)$ are presented in table 5.3:

Note that as was the case with table 5.2, the column showing the computation time is not very accurate. It only serves to indicate that the computation time increases exponential in k .

The computation time rises exponential in k because the linear program for computing a descent direction in the algorithm grows in size linear to the number of edges, which increase exponential. Hence most of the computation time is spend on solving the linear program.

Moreover, note that for the Kneser graph $(2, 10)$ one additional iteration was necessary. That is due to small numerical errors made in the first iteration which are being corrected in the second iteration.

Table 5.4: Results for computation of ideal Arnoldi polynome of degree 8

Name:	$\ q(A)\ _2$	Iterations to 10%:	To 1%:	To 0.1%:	To 0.01%:	To 0.001%:
Grcar	1766.31	16	52	305	809	1409
Ellipse	7710.27	22	40	100	177	243
Lemniscate 1	1	7	9	10	11	14

when only small steps are taken (such as near the optimum), it is expected that the Spectral Bundle method has faster convergence. Moreover, for large and sparse problems both algorithm 4.3.1 and the Spectral Bundle method are expected to be faster than the Successive Quadratic Programming method.

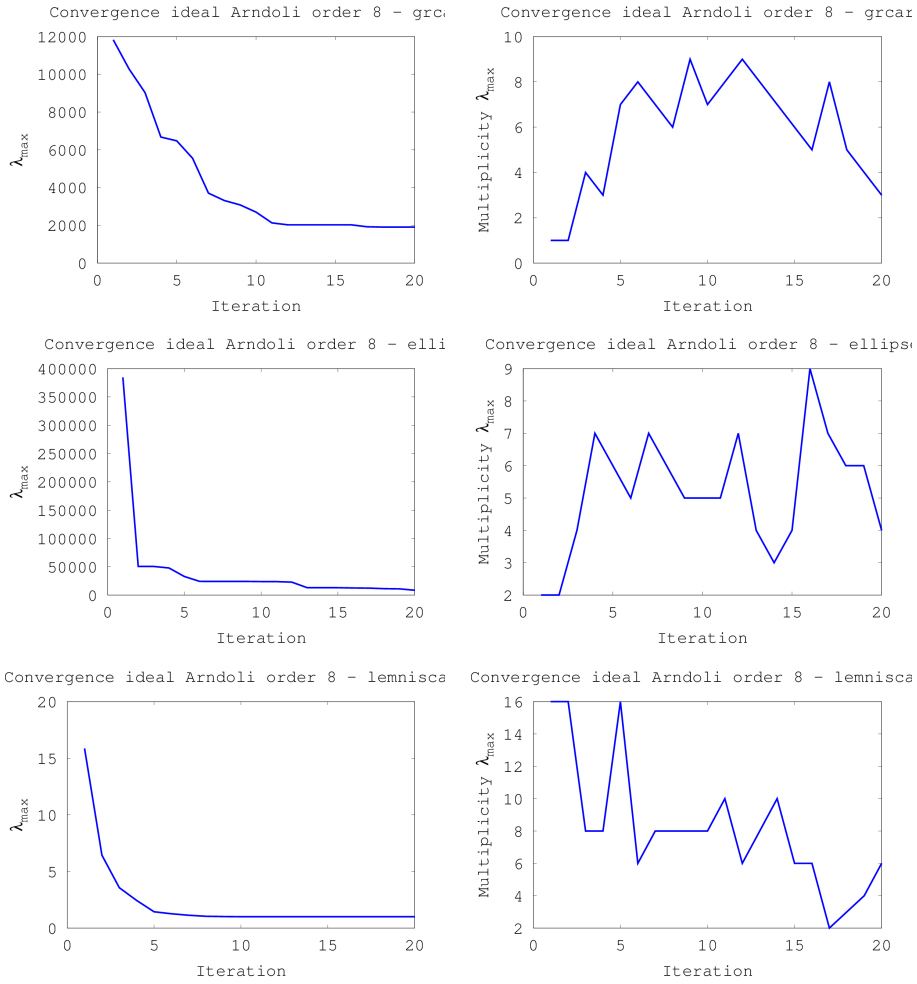


Figure 5.4: Convergence results for computation of ideal Arnoldi polynome of degree 8

Chapter 6

Conclusion

In this thesis I have presented my research on eigenvalue optimization problems and how those problems are linked to semidefinite optimization problems, how SASPA can be used to compute intersecting eigencurves and developed a novel algorithm to solve eigenvalue optimization problems based on those insights.

That semidefinite optimization problems and eigenvalue optimization problems are closely linked is well-known. However the relation developed in chapter 2, section 2.1 based on rotating and perturbing is original. This relation is useful in simplifying large sparse semidefinite optimization problems such as max-cut. However, a sever drawback is that for some instances more memory is required to store the eigenvalue optimization problem than the semidefinite problem.

Also was shown how SASPA can be used to compute intersecting eigencurves. SASPA was recently developed by Rommes and Martins in 2008 [27]. The original application was computing sensitive poles in the analysis of large scale system matrices. The idea that SASPA can also be used to compute the intersection of eigencurves presented in chapter 3 is novel. The adapted method works efficient for large matrices for which eigenvalues have a large derivative. This method has been used successfully to develop a novel algorithm for solving eigenvalue optimization problems.

The new algorithm developed in chapter 4 for solving eigenvalue optimization problems is largely inspired by [21], [14] and [4]. Each presented an algorithm to tackle these kind of problems based on gradients, Newton steps and subgradients. Novel is that algorithm 4.3.1 is based on only partial information each step (that is that no full eigendecomposition is necessary) and computes for each step a new descent direction. This algorithm has been shown to be efficient for computing large and sparse instances of eigenvalue optimization problems in chapter 5.

Future research possibilities include modifying SASPA and algorithm 4.3.1 for parallelization, adapt algorithm 4.3.1 for use in matrix rank minimization problems and incorporating an adapted Conjugated Gradients method such as the Fletcher-Reeves method.

Paralysing SASPA will benefit the method immensely. There are many

matrix-vector products which can be calculated much faster on a multi-core machine than on a single-core. It is exciting to see how much SASPA can benefit by going parallel.

As a beneficiary side-effect, if SASPA can be paralysed in combination with a paralysed version of Jacobi-Davidson (such as in [19]), then algorithm 4.3.1 can also be paralysed. It is expected that a paralysed version of this algorithm is able to solve considerable larger eigenvalue optimization problems in shorter time.

Many current applications which are closely related to semidefinite optimization problems are matrix rank minimization problems. That are problems in which in addition to regular semidefinite constraints there is also a constraint on the rank of X . Two important examples were shown in section 1.4: max-cut and phase retrieval. Phase retrieval is very important because it arises in many applications including MRI, astronomical imaging and microscopy.

Algorithm 4.3.1 can be used to solve the dual of relaxations of such problems to semidefinite optimization problems. It is expected that this algorithm in combination with other methods, such as the log det heuristic [10], is able to solve these kind of problems. Solving these matrix-rank optimization problems through semidefinite optimization is not new [25], but the considering of a dual solving algorithm through eigenvalue optimization is novel. Due to the limitation of primal solvers on the instance size (primal solvers use much more memory), rank minimization problems can currently only be solved for small instances. It is interesting to see if with a (paralysed) dual solving algorithm larger instances can be solved.

Also, algorithm 4.3.1 has a drawback near the optimum in which computation of a descent direction takes considerable more time. It is suspected that many descent directions that are computed at this stage are not orthogonal and form a zig-zag path around the optimum. A Conjugated Gradient version such Fletcher-Reeves method will improve the convergence rate.

In conclusion, through research into eigenvalue optimization problems and SASPA a novel algorithm was developed which is shown to be able to tackle large eigenvalue optimization instances. However, its current implementation is severely limited by the fact that it runs single-core. With more research this method can be turned into a very fast method for solving eigenvalue optimization problems with a possible application to phase retrieval.

Appendix A

Implementation of some algorithms

A.1 Rotating SDP to Spectral Problems

This is the code based on theorem 2.1.4:

```
scripts/rotateSdpToSpectral.m
1 function [rotC, rotA, scale] = rotateSdpToSpectral (b, C, A)
2
3     % First, determine rotation matrix that transforms b into e1
4     [Q, R] = qr(b);
5
6     if R(1) < 0
7         Q = -Q;
8         R = -R;
9     end
10    scale = R(1);
11
12    % Second, determine matrices A^*Q_{1..n}
13    Aq = cell(size(A, 1), 1);
14
15    for i = 1:size(A, 1)
16
17        Aq{i} = sparse(size(A{1}, 1), size(A{1}, i));
18
19        for j = 1:size(A, 1)
20            Aq{i} = Aq{i} + Q(j, i) .* A{j};
21        end
22
23    end
24
25    % Third, make a decomposition of Aq{1}
26    if Aq{1} == speye(size(A{1}, 1))
27
28        rotA = Aq{2:size(A, 1)};
29        rotC = C;
```

```

30
31     else
32
33         smallestEig = jd(Aq{1}, 'sa');
34
35         %     if smallestEig < 0
36         %         perturb = smallestEig + 0.01;
37         %         disp(sprintf('Warning: instance is permuted by %e',
perturb))
38         %         b = [b; perturb];
39         %         A{end} = speye(size(A{1}));
40         %         [rotC, rotA, scale] = rotateSdpToSpectral(b, C, A);
41         %     else
42         %         Rchol = chol(Aq{1});
43         %         Rchol_inv = inv(Rchol);
44         %         rotC = Rchol_inv' * C * Rchol_inv;
45
46         %         rotA = cell(size(A, 1) - 1, 1);
47         %         for i = 1:size(A, 1) - 1
48         %             rotA{i} = Rchol_inv' * Aq{i+1} * Rchol_inv;
49         %         end
50         %     end
51     end
52 end

```

A.2 SASPA to compute intersections

This is the code based on algorithm 3.5.2:

```

scripts/intersectMaxEigencurve.m
1 function [varargout] = intersectMaxEigencurve(varargin)
2
3     [A, dA, lambda_max, U, v_max, tol, itmax, kmin, kmax] = ...
4         readOptions(varargin{:});
5
6     x = ones(size(A, 1), 1) + 0.05.*(2.*rand(size(A, 1), 1) -
7         1);
8     X = lanczos(A, x, kmin);
9
10    deflate = orth(U);
11
12    X = X - U*(U'*X);
13    X = orth(X);
14
15    S = X' * A * X;
16    [V, D] = eig(S);
17    [s, i] = max(diag(D));
18    x = X*V(:, i);
19
20    t = Inf; tm = Inf;
21
22    u = [];
23    lambda = 0;

```



```

78     res = norm(A*x - s.*x);
79
80     printf( '%d \t\t %f \t\t %f \t\t %e\n', ...
81             iteration, s, res, T(1))
82
83     tm = min(tm, T(1));
84
85     if res > tol && res < 1
86         [x, s] = RQL(A, x, s, tol);
87     end
88
89     if norm(A*x - s.*x) < tol
90
91         printf( '_____
92                 converged
93                 _____\n' )
94         printf( '\tt: %e \tlambda: %f\tres: %e\n',
95                 ...
96                 T(1), s, norm(A*x-s.*x))
97
98         if T(1) < t
99             u = x;
100             lambda = s;
101             t = T(1);
102             s = lambda_max;
103
104             if ~isinf(t) && t <= tm
105                 converged = true;
106             end
107         elseif ~isinf(t)
108             converged = true;
109         end
110
111         printf( '\tt*: %e\tlambda*: %f\n', t, lambda
112                )
113         printf( '
114                _____
115                n' )
116
117         deflate = orth([deflate x]);
118
119         X = X - deflate * (deflate'*X);
120         X(:, 1) = X(:, 1) ./ norm(X(:, 1));
121         X = orth(X);
122         i = 2; while i < size(X, 2)
123             x = X(:, i) - X(:, 1:i-1) * (X(:,
124                 1:i-1)'* X(:, i));
125             if norm(x) > 1e-6
126                 X(:, i) = x ./ norm(x);
127                 i = i + 1;
128             else
129                 X = [X(:, 1:i-1), X(:, i+1:
130                     size(X, 2))];
131                 end
132             end
133         end
134     end
135     X = orth(X);
136     x = X * (2 .* rand(size(X, 2), 1) - 1);

```



```

127
128         elseif size(X, 2) > kmax
129
130             X = X * V(:, 1:kmin);
131             X = orth(X);
132
133         end
134
135         if T(1) > t
136             s = lambda_max;
137         end
138
139         iteration = iteration + 1;
140     end
141
142     if nargout == 1
143         varargout{1} = u;
144     elseif nargout == 2
145         varargout{1} = u;
146         varargout{2} = lambda;
147     elseif nargout == 3
148         varargout{1} = u;
149         varargout{2} = lambda;
150         varargout{3} = full(t);
151     elseif nargout >= 4
152         varargout{1} = u;
153         varargout{2} = lambda;
154         varargout{3} = full(t);
155         varargout{4} = full(tm);
156     end
157 end
158
159 function [A, dA, lambda_max, U, v_max, tol, itmax, kmin, kmax] =
160     ...
161     readOptions(varargin)
162
163     A = varargin{1};
164     dA = varargin{2};
165
166     i = 3;
167
168     if length(varargin) >= i+1
169         lambda_max = varargin{i}; i=i+1;
170         U = varargin{i}; i=i+1;
171     else
172         [V, D] = eigs(A, 1, 'la');
173         lambda_max = D(1,1);
174         U = V(:, 1);
175     end
176
177     if length(varargin) >= i
178         v_max = varargin{i}; i=i+1;
179     else
180         v_max = max(U' * dA * U);
181     end
182

```

```

183     if length(varargin) < i, tol = 1e-8;     else tol = varargin{i};
184         end, i=i+1;
185     if length(varargin) < i, itmax = 200;    else itmax = varargin{i}
186         }; end, i=i+1;
187     if length(varargin) < i, kmin = 4;      else kmin = varargin{i};
188         end, i=i+1;
189     if length(varargin) < i, kmax = 8;      else kmax = varargin{i};
190         end, i=i+1;
191 end
192 function V = lanczos(A, v0, n)
193     V = zeros(size(A, 1), n);
194     V(:, 1) = v0;
195     for i = 1:n-1
196         v = A*V(:, i);
197         a = v'*V(:, i);
198         if i > 1, v = v - a.*V(:, i) - b.*V(:, i-1);
199         else, v = v - a.*V(:, i);
200         end
201         b = norm(v);
202         V(:, i+1) = v ./ b;
203     end
204 end
205 end
206 end
207 end

```

A.3 Spectral Descent

:

This is the code based on algorithm 4.3.1.

```

scripts/spectralDescent.m
1 function [lambda_max, z, A, conv_his, rank_his, f_his] =
2 spectralDescent (C, opA, P, z0, outz, output)
3
4     tau = 0.1; dtol = 1e-10; tstart = tic;
5
6     n = size(C, 1); m = size(opA, 1);
7     if size(z0, 1) == m && size(z0, 2) == 1
8         z = z0; A = C;
9         for i = 1:m, A = A - z(i).*opA{i}; end
10    else
11        z = zeros(m, 1); A = C;
12    end
13    [Q, lambda_max] = eigenspace(A, 'la', [], tau);
14    lambda_max = max(diag(lambda_max));
15    conv_his = [lambda_max];
16    rank_his = [];
17    f_his = [];

```

```

18
19     not_converged = true;
20     while not_converged && tau > 1e-8
21
22         if size(output, 2) > 1
23             fid = fopen(output, 'a');
24             fprintf(fid, ...
25                 '%f : Calculating steepest
26                 direction, rank Q: %d\tlambda:
27                 %d\n', ...
28                 toc(tstart), size(Q, 2), lambda_max
29                 );
30             fclose(fid);
31         end
32
33         d = zeros(m, 1);
34         if size(Q, 2) == 1
35             for i = 1:m
36                 d(i) = Q' * opA{i} * Q;
37             end
38             d = d ./ norm(d);
39         else
40             t = size(Q, 2);
41
42             f = zeros(m+1, 1); f(1) = 1;
43             Al = zeros(t*(t+1)/2, m+1);
44             Al(:, 1) = ones(t*(t+1)/2, 1);
45             temp = eye(t);
46             Al(:, 1) = temp(tril(true(t)));
47             for i = 1:m
48                 temp = Q' * opA{i} * Q;
49                 Al(:, i+1) = temp(tril(true(t)));
50             end
51             b = zeros(t*(t+1)/2, 1);
52             lb = -ones(m+1, 1); lb(1) = -Inf;
53             ub = ones(m+1, 1); ub(1) = Inf;
54
55             x = linprog(f, Al, b, [], [], lb, ub);
56             T = repmat('L', (t*(t+1)/2), 1);
57             T(temp(tril(true(t))) == 0) = 'S';
58             x = glpk(f, Al, b, lb, ub, T);
59             d = x(2:m+1);
60
61             if x(1) >= -dtol || isnan(x(1))
62                 [Q, lambdas] = eigenspace(A, 'la',
63                 [], tau);
64                 lambdas = diag(lambdas); lambda_max
65                 = max(lambdas);
66                 tau = lambda_max - min(lambdas) - 1
67                 e-8;
68                 if size(output, 2) > 1
69                     fid = fopen(output, 'a');
70                     fprintf(fid, ...
71                         '%f : No descent
72                         direction found
73                         , decreasing
74                         tau to: %f\n',

```

```

66         ...
67         toc(tstart), tau);
68         fclose(fid);
69     end
70     if tau < 1e-8
71         if size(output, 2) > 1
72             fid = fopen(output,
73                 'a');
74             fprintf(fid, ...
75                 '%f :
76                 Converged
77                 , Df: %
78                 e\
79                 tLambda
80                 : %e\
81                 tTau: %
82                 f\n',
83                 ...
84                 toc(tstart)
85                 , x(1),
86                 lambda_max
87                 , tau);
88             fclose(fid);
89         end
90         not_converged = false;
91         rank_his = [rank_his; size(Q, 2)];
92         break
93     end
94     cut = abs(lambdas - lambda_max) <
95         tau;
96     Q = Q(:, cut);
97     continue;
98 end
99
100 rank_his = [rank_his; size(Q, 2)];
101
102 if size(P, 1) > 0
103     b = P*z; I = find(b <= 0);
104     r = quadprog(speye(m), [], -P(I, :), P(I,
105         :)*d);
106     d = d + r;
107 end
108
109 dA = sparse(n, n);
110 for i = 1:m
111     dA = dA + d(i) .* opA{i};
112 end
113
114 f = -jd(Q'*dA*Q, 'sa');
115 b = -jd(Q'*dA*Q, 'la');
116
117 if (f >= -dtol || b >= dtol) && size(Q, 2) > 1
118     lambdas = diag(Q'*A*Q);
119     tau = lambda_max - min(lambdas) - 1e-8;

```

```

106         if size(output, 2) > 1
107             fid = fopen(output, 'a');
108             fprintf(fid, ...
109                 '%f : No descent direction
110                 found, decreasing tau
111                 to: %f\n', ...
112                 toc(tstart), tau);
113             fclose(fid);
114         end
115         if tau < 1e-8
116             if size(output, 2) > 1
117                 fid = fopen(output, 'a');
118                 fprintf(fid, ...
119                     '%f : Converged, Df
120                     : %e\tLambda: %
121                     e\tTau: %f\n',
122                     ...
123                     toc(tstart), f,
124                     lambda_max);
125                 fclose(fid);
126             end
127             not_converged = false;
128             break
129         end
130         cut = abs(lambdas - lambda_max) < tau;
131         Q = Q(:, cut);
132         continue;
133     end
134
135     if f >= -dtol
136
137         if size(output, 2) > 1
138             fid = fopen(output, 'a');
139             fprintf(fid, ...
140                 '%f : Converged, Df: %e\tBf
141                 : %e\tLambda: %e\n',
142                 ...
143                 toc(tstart), f, b,
144                 lambda_max);
145             fclose(fid);
146         end
147         not_converged = false;
148         break
149     end
150
151     f_his = [f_his; f];
152     if size(output, 2) > 1
153         fid = fopen(output, 'a');
154         fprintf(fid, ...
155             '%f : Starting step calculation, Df
156             : %e\n', ...
157             toc(tstart), f);
158         fclose(fid);
159     end
160
161     while f < -dtol

```

```

153
154 [q, lambda, t, tm] = intersectMaxEigencurve
      (A, -dA, lambda_max, Q, f);
155
156 if t == Inf && tm == Inf, t = 10/max(max(dA
      )); end
157 if t == Inf, t = tm; end
158
159 if size(output, 2) > 1
160     fid = fopen(output, 'a');
161     fprintf(fid, ...
162             '%f : Found intersection at
              , t: %e\n', ...
              toc(tstart), t);
163     fclose(fid);
164 end
165
166 if size(P, 1) > 0
167     tmax = - min((P*z) ./ (P*d));
168     t = min(tmax, t);
169 end
170
171 [Qt, lambda_t, dist] = eigenspace(A - t.*dA
      , 'la', [Q q], tau);
172 lambda_t = max(diag(lambda_t));
173
174 if size(output, 2) > 1
175     fid = fopen(output, 'a');
176     fprintf(fid, ...
177             '%f : After step, rank Q: %
              d\tLambda: %e\n', ...
              toc(tstart), size(Qt, 2),
              lambda_t);
178     fclose(fid);
179 end
180
181 ft = -jd(Qt'*dA*Qt, 'sa');
182 bt = -jd(Qt'*dA*Qt, 'la');
183
184 if size(output, 2) > 1
185     fid = fopen(output, 'a');
186     fprintf(fid, ...
187             '%f : New Df: %e\t Bf: %e\n
              ', ...
              toc(tstart), ft, bt);
188     fclose(fid);
189 end
190
191 while ft > dtol && bt > dtol
192
193     par = -f*t/(2*((lambda_t-lambda_max
194             )/t-f));
195     int = (lambda_max-lambda_t+t*bt)/(
196             bt-f);
197
198     if int > 1e-12 && int < t
199         if par > int && par < t, nt
200

```

```

201         = par;
202     else , nt
203         = int;
204     end
205 elseif par > 1e-12 && par < t
206     nt = par;
207 else , break , end
208
209 [nQt, nlambda_t, ndist] =
210     eigenspace(A - nt.*dA, 'la', Qt
211     , tau);
212 nlambda_t = max(diag(nlamba_t));
213
214 if size(nQt, 2) == 0, break , end
215
216 t = nt; Qt = nQt; lambda_t =
217     nlambda_t; dist = ndist;
218
219 ft = -jd(Qt'*dA*Qt, 'sa');
220 bt = -jd(Qt'*dA*Qt, 'la');
221
222 if size(output, 2) > 1
223     fid = fopen(output, 'a');
224     fprintf(fid, ...
225         '%f : Through
226         parabola fit, t
227         :%e\trank Q: %d
228         \tLambda: %e\
229         tDf: %e\tBf: %e
230         \n', ...
231         toc(tstart), t,
232         size(Qt, 2),
233         lambda_t, ft,
234         bt);
235     fclose(fid);
236 end
237
238 if ft > -dtol && bt < dtol && lambda_t >
239     lambda_max
240
241     lambdas = diag(Q'*A*Q);
242     if max(lambdas) - min(lambdas) < 2e
243         -8, break , end
244     tau = max(lambdas) - min(lambdas) -
245         1e-8;
246
247     if size(output, 2) > 1
248         fid = fopen(output, 'a');
249         fprintf(fid, ...
250             '%f : Decreased tau
251             to: %f\n', ...
252             toc(tstart), tau);
253         fclose(fid);
254     end
255
256 cut = abs(lambdas - lambda_max) <

```

```

241         tau;
242         Q = Q(:, cut);
243     elseif (ft > -dtol || abs(lambda_t -
244             lambda_max) < 1e-8) && size(Qt, 2) ==
245             size(Q, 2)
246         tau = tau + dist + 1e-8;
247         [Qt, lambda_t] = eigenspace(A-t.*dA
248             , 'la', Qt, tau);
249         lambda_t = max(diag(lambda_t));
250     if size(output, 2) > 1
251         fid = fopen(output, 'a');
252         fprintf(fid, ...
253             '%f : Increased tau
254                 to: %f\n', ...
255                 toc(tstart), tau);
256         fclose(fid);
257     end
258     if lambda_t >= lambda_max
259         [Q, lambda_max] =
260             eigenspace(A, 'la', [],
261                 tau);
262         lambda_max = max(diag(
263             lambda_max));
264     end
265     end
266     if lambda_t < lambda_max
267         A = A - t.*dA;
268         z = z + t.*d;
269         Q = Qt;
270         lambda_max = lambda_t;
271         f = ft;
272         b = bt;
273     if size(outz, 2) > 1
274         save(outz, 'z', 'A', 'Q', '
275             lambda_max', ...
276                 not_converged
277                 ', '
278                 conv_his
279                 ', '
280                 rank_his
281                 ', '
282                 f_his')
283         ;
284     end
285     else, break, end
286 end

```



```

281
282     conv_his = [conv_his; lambda_max];
283     if size(output, 2) > 1
284         fid = fopen(output, 'a');
285         fprintf(fid, ...
286             '\n%f : New iteration, lambda: %e\n', ...
287                 toc(tstart), lambda_max);
288         fclose(fid);
289     end
290     if size(outh, 2) > 1
291         save(outh, 'z', 'A', 'Q', 'lambda_max', '
292             not_converged', ...
293             'conv_his', 'rank_his', '
294             f_his');
295     end
296     if size(outh, 2) > 1
297         tend = toc(tstart);
298         save(outh, 'z', 'A', 'Q', 'lambda_max', '
299             not_converged', 'tend', ...
300             'conv_his', 'rank_his', 'f_his');
301     end
end

```


Bibliography

- [1] Dmitri Alekseevsky, Andreas Kriegl, Peter W Michor, and Mark Losik. Choosing roots of polynomials smoothly. *Israel Journal of Mathematics*, 105(1):203–233, 1998.
- [2] Walter Edwin Arnoldi. The principle of minimized iterations in the solution of the matrix eigenvalue problem. *Quarterly of Applied Mathematics*, 9(1):17–29, 1951.
- [3] David A Bader, Henning Meyerhenke, Peter Sanders, and Dorothea Wagner. *Graph partitioning and graph clustering*, volume 588. American Mathematical Soc., 2013.
- [4] Mischja Alexander van Bossum. Semidefinite optimization, a spectral approach. 2002.
- [5] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [6] Emmanuel J Candes, Yonina C Eldar, Thomas Strohmer, and Vladislav Voroninski. Phase retrieval via matrix completion. *SIAM Review*, 57(2):225–251, 2015.
- [7] Francis H Clarke, Yuri S Ledyaev, Ronald J Stern, and Peter R Wolenski. *Nonsmooth analysis and control theory*. Springer Science & Business Media, 1998.
- [8] Frank H Clarke. *Optimization and nonsmooth analysis*. John Wiley & Sons, 1983.
- [9] Caterina De Simone, Martin Diehl, Michael Jünger, Petra Mutzel, Gerhard Reinelt, and Giovanni Rinaldi. Exact ground states of ising spin glasses: New experimental results with a branch-and-cut algorithm. *Journal of Statistical Physics*, 80(1-2):487–496, 1995.
- [10] Maryam Fazel, Haitham Hindi, and Stephen P Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. In *American Control Conference, 2003. Proceedings of the 2003*, volume 3, pages 2156–2162. IEEE, 2003.

- [11] Roger Fletcher. *Practical methods of optimization*. John Wiley & Sons, 1990.
- [12] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42(6):1115–1145, November 1995.
- [13] Anne Greenbaum and Lloyd N Trefethen. Gmres/cr and arnoldi/lanczos as matrix approximation problems. *SIAM Journal on Scientific Computing*, 15(2):359–368, 1994.
- [14] Christoph Helmberg and Franz Rendl. A spectral bundle method for semidefinite programming. *SIAM Journal on Optimization*, 10(3):673–696, 2000.
- [15] Richard M Karp. *Reducibility among combinatorial problems*. Springer, 1972.
- [16] Donald Ervin Knuth. *The sandwich theorem*. Stanford University, Department of Computer Science, 1993.
- [17] László Lovász. On the shannon capacity of a graph. *Information Theory, IEEE Transactions on*, 25(1):1–7, 1979.
- [18] Jan R Magnus. On differentiating eigenvalues and eigenvectors. *Econometric Theory*, 1(02):179–191, 1985.
- [19] Margreet Nool and Auke van der Ploeg. A parallel jacobi–davidson-type method for solving large generalized eigenvalue problems in magnetohydrodynamics. *SIAM Journal on Scientific Computing*, 22(1):95–112, 2000.
- [20] Alexander M Ostrowski. On the convergence of the rayleigh quotient iteration for the computation of the characteristic roots and vectors. i. *Archive for Rational Mechanics and Analysis*, 1(1):233–241, 1957.
- [21] Michael L Overton. On minimizing the maximum eigenvalue of a symmetric matrix. *SIAM Journal on Matrix Analysis and Applications*, 9(2):256–268, 1988.
- [22] Michael L Overton. Large-scale optimization of eigenvalues. *SIAM Journal on Optimization*, 2(1):88–120, 1992.
- [23] Beresford N Parlett. The rayleigh quotient iteration and some generalizations for nonnormal matrices. *Mathematics of Computation*, 28(127):679–693, 1974.
- [24] G Pataki and SH Schmieta. The dimacs library of mixed semidefinite-quadratic-linear programs. 2002.

- [25] Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [26] Joost Rommes and Nelson Martins. Efficient computation of transfer function dominant poles using subspace acceleration.
- [27] Joost Rommes and Nelson Martins. Computing large-scale system eigenvalues most sensitive to parameter changes, with applications to power system small-signal stability. *Power Systems, IEEE Transactions on*, 23(2):434–442, 2008.
- [28] Joost Rommes and Gerard LG Sleijpen. Convergence of the dominant pole algorithm and rayleigh quotient iteration. *SIAM Journal on Matrix Analysis and Applications*, 30(1):346–363, 2008.
- [29] Youcef Saad and Martin H Schultz. Gmres: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal on scientific and statistical computing*, 7(3):856–869, 1986.
- [30] Gerard LG Sleijpen and Henk A Van der Vorst. A jacobi–davidson iteration method for linear eigenvalue problems. *SIAM Review*, 42(2):267–293, 2000.
- [31] David Titley-Peloquin, Jennifer Pestana, and Andrew J Wathen. Gmres convergence bounds that depend on the right-hand-side vector. *IMA Journal of Numerical Analysis*, 34(2):462–479, 2014.
- [32] Kim-Chuan Toh and Lloyd N Trefethen. The chebyshev polynomials of a matrix. *SIAM Journal on Matrix Analysis and Applications*, 20(2):400–419, 1998.
- [33] Philip Wolfe. The simplex method for quadratic programming. *Econometrica: Journal of the Econometric Society*, pages 382–398, 1959.