

A Robust Algorithm for Semidefinite Programming ^{*}

Xuan Vinh Doan [†] Serge Kruk [‡] Henry Wolkowicz [§]

November 11, 2010

Contents

1	Introduction	2
1.1	Motivation and Central Problem	3
1.2	Outline	4
2	Gauss-Newton Method	4
2.1	Matrix-Free Formulation	4
2.2	Algorithm Initialization	7
2.3	Local Convergence	8
2.4	Presolve and Preconditioning	11
3	Lovász Theta Function Problem	14
3.1	Matrix-Free Formulation	14
3.2	Presolve and Preconditioning	16
3.2.1	Diagonal Preconditioning	16
3.2.2	Block-Diagonal Preconditioning	17
4	Numerics and Conclusion	19
4.1	Numerical Results	19
4.2	Concluding Remarks	31
	Bibliography	31
	Index	33

^{*}This report is available at URL: orion.math.uwaterloo.ca/~hwoikowi/henry/reports/ABSTRACTS.html

[†]Department of Combinatorics and Optimization, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada

[‡]Dept. of Mathematics and Statistics, Oakland University, Rochester, MI

[§]Department of Combinatorics and Optimization, University of Waterloo, 200 University Avenue West, Waterloo, ON N2L 3G1, Canada

List of Tables

4.1	Ratios of numbers of LSMR iterations of $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ to (\mathbf{SRSD}_o)	20
4.2	Ratios of total computational time of $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ to (\mathbf{SRSD}_o)	20
4.3	Accuracy measures for random sparse instances	20
4.4	Performance measures for <i>hard</i> random sparse instances	21
4.5	Performance measures for <i>easy</i> random sparse instances	21
4.6	Average measures for random instances with different numbers of constraints	22
4.7	Time ratios relative to $(\mathbf{SRSD}_{\text{BDiag}})$ for different numbers of constraints	22
4.8	Computational times for random instances with different sparsity	23
4.9	Average measures for random instances with different sparsity	23
4.10	Time ratios relative to $(\mathbf{SRSD}_{\text{Diag}})$ for different sparsity	24
4.11	Performance measures; random instances; strict complementarity fails	24
4.12	Performance measures; random instances; dual Slater's CQ almost fails	24
4.13	Accuracy measures for random Lovász theta function instances	25
4.14	Performance measures for <i>hard</i> random Lovász theta function instances	25
4.15	Performance measures for <i>hard</i> random Lovász theta function instances	26
4.16	Performance measures for <i>easy</i> random Lovász theta function instances	26
4.17	Performance measures for <i>easy</i> random Lovász theta function instances	26
4.18	Performance measures for <i>control</i> instances	27
4.19	Performance measures for <i>gpp100</i> instance	28
4.20	Performance measures for the first eight <i>hinf</i> instances	28
4.21	Performance measures for the remaining seven <i>hinf</i> instances	29
4.22	Performance measures for <i>hinf12</i> instance	29
4.23	Performance measures for <i>hinf5</i> instance	30
4.24	Performance measures for <i>gap</i> instances	30
4.25	Performance measures for <i>gap7</i> instance	30

Abstract

Current successful methods for solving semidefinite programs, SDP, are based on primal-dual interior-point approaches. These usually involve a symmetrization step to allow for application of Newton's method followed by block elimination to reduce the size of the Newton equation. Both these steps create ill-conditioning in the Newton equation and singularity of the Jacobian of the optimality conditions at the optimum.

In order to avoid the ill-conditioning, we derive and test a backwards stable primal-dual interior-point method for SDP. Relative to current public domain software, we realize both a distinct improvement in the accuracy of the optimum and a reduction in the number of iterations. This is true for random problems as well as for problems of special structure. Our algorithm is based on a Gauss-Newton approach applied to a single bilinear form of the optimality conditions. The well-conditioned Jacobian allows for a preconditioned (matrix-free) iterative method for finding the search direction at each iteration.

1 Introduction

We derive and test a backwards stable primal-dual interior-point method for semidefinite programming. Relative to current public domain software, we realize both a distinct improvement in accuracy and a reduction in the required number of iterations. This is for random problems as well

as for problems of special structure. The algorithm is based on using a Gauss-Newton approach with a preconditioned iterative method for finding the search direction.

Primal-dual interior-point methods are currently the methods of choice for solving semidefinite programming, SDP, problems. However, current primal-dual interior-point methods are quite unstable. They cannot provide high accuracy solutions in general; and, they often fail for ill-conditioned problems. The instability arises from two steps. Since the optimality conditions for SDP are an overdetermined system of nonlinear equations, a symmetrization step is applied so that Newton's method can be used. This symmetrization changes a (possibly) well-posed problem into an ill-posed one where the Jacobian at optimality is singular. Then, block elimination is used in order to reduce the size of the resulting Newton system. However, this block elimination does not use any type of partial pivoting and again singularities are introduced. Therefore, these algorithms are not backwards stable. (This is discussed in [6, 23]. Further details are also given in Section 1.1, below.)

We study the Gauss-Newton method for solving SDPs and illustrate that high accuracy solutions can be obtained dependably for medium sized problems. We follow the approach in [11, 23] and use a matrix free, inexact Gauss-Newton method to solve the perturbed optimality conditions.

1.1 Motivation and Central Problem

The primal-dual SDP pair we consider is

$$\begin{aligned} \text{(PSDP)} \quad p^* := \min \quad & \mathbf{C} \cdot \mathbf{X} \\ \text{s.t.} \quad & \mathcal{A}(\mathbf{X}) = \mathbf{b}, \\ & \mathbf{X} \succeq 0, \end{aligned} \tag{1.1}$$

and

$$\begin{aligned} \text{(DSDP)} \quad d^* := \max \quad & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} \quad & \mathcal{A}^*(\mathbf{y}) + \mathbf{Z} = \mathbf{C}, \\ & \mathbf{Z} \succeq 0, \end{aligned} \tag{1.2}$$

where: $\mathbf{C}, \mathbf{X}, \mathbf{Z} \in \mathcal{S}^n$, \mathcal{S}^n denotes the space of $n \times n$ real symmetric matrices equipped with the trace inner product, $\mathbf{C} \cdot \mathbf{D} = \text{trace}(\mathbf{C}\mathbf{D})$; and $\mathcal{A} : \mathcal{S}^n \rightarrow \mathbb{R}^m$ is a linear transformation, with \mathcal{A}^* its adjoint transformation.

Under a suitable *constraint qualification* assumption such as Slater's condition or strict feasibility, the primal-dual solution $(\mathbf{X}, \mathbf{y}, \mathbf{Z})$ with $\mathbf{X}, \mathbf{Z} \succeq 0$ is optimal for the primal-dual pair (1.1) and (1.2) if and only if

$$F(\mathbf{X}, \mathbf{y}, \mathbf{Z}) := \begin{pmatrix} \mathcal{A}^*(\mathbf{y}) + \mathbf{Z} - \mathbf{C} \\ \mathcal{A}(\mathbf{X}) - \mathbf{b} \\ \mathbf{Z}\mathbf{X} \end{pmatrix} = \mathbf{0}. \tag{1.3}$$

Primal-dual interior-point methods maintain $\mathbf{X}, \mathbf{Z} \succ 0$, and with the *barrier parameter* $\mu \downarrow 0$, they find approximate solutions to the following perturbed optimality conditions:

$$F_\mu(\mathbf{X}, \mathbf{y}, \mathbf{Z}) := \begin{pmatrix} \mathcal{A}^*(\mathbf{y}) + \mathbf{Z} - \mathbf{C} \\ \mathcal{A}(\mathbf{X}) - \mathbf{b} \\ \mathbf{Z}\mathbf{X} - \mu\mathbf{I} \end{pmatrix} = \mathbf{0}, \quad \mu > 0. \tag{1.4}$$

Thus, these methods are based on *path following*. The perturbed system in (1.4) is overdetermined. Under nondegeneracy assumptions, the Jacobian is full rank at optimality, e.g., [1]. Current

methods use Newton's method applied to various symmetrizations of (1.4). The two most popular symmetrizations are the so-called HRVW/KSM/M [8, 10, 17] and NT [19, 20] methods, see for example, [18]. But, the linearizations (Jacobian) of the symmetrized optimality conditions for both of these methods is singular at an optimum. This and the block elimination schemes used for finding the search directions both imply that one is solving an increasingly ill-conditioned linear system to find a search direction. Thus, it is extremely difficult to obtain high accuracy solutions; these algorithms are not backwards stable under finite precision arithmetic. Similarly, finding reasonable preconditioners for iterative methods is difficult if not impossible. In addition, the block eliminations make it difficult to exploit sparsity in the data. In this paper we propose a robust primal-dual interior/exterior-point method which uses an inexact Gauss-Newton approach with a matrix-free preconditioned conjugate gradient method. We do not change a well-conditioned system of optimality conditions to an ill-conditioned system. The method is able to attain high accuracy solutions as well as exploit sparsity.

1.2 Outline

The Gauss-Newton approach is described in Section 2. The infeasible starting point variation follows in Section 2.2. We include a proof of asymptotic convergence in Section 2.3, and a description of the preconditioning techniques in Section 2.4.

We apply our techniques to the Lovász Theta Function Problem in Section 3. Numerics and concluding remarks are given in Section 4.

2 Gauss-Newton Method

2.1 Matrix-Free Formulation

Let us consider the primal-dual SDP pair in (1.1) and (1.2). We assume the linear transformation \mathcal{A} has full rank m . It can be represented as

$$\mathcal{A}(\mathbf{X}) = \mathbf{b} \Leftrightarrow \mathbf{A}_i \cdot \mathbf{X} = b_i, \quad \forall i = 1, \dots, m.$$

The adjoint transformation \mathcal{A}^* is then defined as $\mathcal{A}^*(\mathbf{y}) = \sum_{i=1}^m y_i \mathbf{A}_i$, for all $\mathbf{y} \in \mathbb{R}^m$. Here $\mathbf{A}_i \in \mathcal{S}^n, \forall i$.

The mapping $\text{vec} : \mathcal{M}^n \rightarrow \mathbb{R}^{n^2}$, where \mathcal{M}^n is the set of all square $n \times n$ matrices, takes a matrix $\mathbf{M} \in \mathcal{M}^n$ and forms a vector $\mathbf{v} \in \mathbb{R}^{n^2}$ from its columns. The inverse mapping is $\text{Mat} := \text{vec}^{-1}$, which takes a vector $\mathbf{v} \in \mathbb{R}^{n^2}$ and forms a matrix $\mathbf{M} \in \mathcal{M}^n$ column by column. Indeed $\mathbf{M} = \text{Mat}(\text{vec}(\mathbf{M}))$, for all $\mathbf{M} \in \mathcal{M}^n$.

We also define the triangular number, $t(n) = n(n+1)/2$ and the mapping $\text{svec} : \mathcal{S}^n \rightarrow \mathbb{R}^{t(n)}$ that takes a symmetric matrix $\mathbf{S} \in \mathcal{S}^n$ and forms a vector $\mathbf{v} \in \mathbb{R}^{t(n)}$ by concatenating n vectors $\mathbf{s}_{j-1} = \sqrt{2}(s_{ij})_{1 \leq i < j}$ for all $j = 2, \dots, n$ and $\mathbf{s}_n = \text{diag}(\mathbf{S})$. This mapping is an isometry under the 2-norm. The inverse mapping $\text{sMat} = \text{svec}^{-1}$ maps a vector $\mathbf{v} \in \mathbb{R}^{t(n)}$ into a symmetric matrix $\mathbf{S} \in \mathcal{S}^n$. We also have $\text{sMat}^* = \text{svec}$, since

$$\langle \text{sMat}(\mathbf{v}), \mathbf{S} \rangle = \text{trace}(\text{sMat}(\mathbf{v})\mathbf{S}) = \mathbf{v}^T \text{svec}(\mathbf{S}) = \langle \mathbf{v}, \text{svec}(\mathbf{S}) \rangle.$$

From setting $\mathbf{A} \in \mathbb{R}^{m \times t(n)}$ with rows $\mathbf{A}_{i,:} = \text{svec}(\mathbf{A}_i)$, for all $i = 1, \dots, m$, we have:

$$\mathcal{A}(\mathbf{X}) = \mathbf{b} \Leftrightarrow \mathbf{A} \text{svec}(\mathbf{X}) = \mathbf{b}.$$

The nullspace of \mathbf{A} , $\text{null}(\mathbf{A})$, has dimension $t(n) - m$. Let us consider an orthonormal basis $\{\mathbf{q}_1, \dots, \mathbf{q}_{t(n)-m}\}$ of $\text{null}(\mathbf{A})$ and assume that we could find a primal feasible solution $\hat{\mathbf{X}} \succ 0$. Then

$$\mathcal{A}(\mathbf{X}) = \mathbf{b} \Leftrightarrow \text{svec}(\mathbf{X}) = \hat{\mathbf{x}} + \mathbf{Q}\mathbf{v},$$

where $\hat{\mathbf{x}} = \text{svec}(\hat{\mathbf{X}})$, $\mathbf{v} \in \mathbb{R}^{t(n)-m}$ and the columns of $\mathbf{Q} \in \mathbb{R}^{t(n) \times (t(n)-m)}$ are taken from the basis of $\text{null}(\mathbf{A})$.

Let us consider now the dual feasibility condition:

$$\mathcal{A}^*(\mathbf{y}) + \mathbf{Z} = \mathbf{C} \Leftrightarrow \text{svec}(\mathbf{Z}) = \mathbf{c} - \mathbf{A}^T \mathbf{y},$$

where $\mathbf{c} = \text{svec}(\mathbf{C})$. The optimality conditions in (1.4) are now equivalent to the following conditions:

$$G_\mu(\mathbf{v}, \mathbf{y}) := \text{sMat}(\mathbf{c} - \mathbf{A}^T \mathbf{y}) \text{sMat}(\hat{\mathbf{x}} + \mathbf{Q}\mathbf{v}) - \mu \mathbf{I} = \mathbf{0} \quad (\Leftrightarrow \text{vec}(G_\mu(\mathbf{v}, \mathbf{y})) = \mathbf{0}). \quad (2.1)$$

This is a single bilinear overdetermined system with $t(n)$ variables and n^2 equations and can be solved by the Gauss-Newton method. For each $\mu > 0$, there exists a unique primal-dual solution $(\mathbf{X}_\mu, \mathbf{y}_\mu, \mathbf{Z}_\mu)$ of (1.4), with $\mathbf{X}_\mu \succ 0, \mathbf{Z}_\mu \succ 0$, that lies on (and thus defines) the *central path*. The corresponding $\text{svec}(\mathbf{Z}_\mu) = \mathbf{c} - \mathbf{A}^T \mathbf{y}_\mu$ and $\text{svec}(\mathbf{X}_\mu) = \hat{\mathbf{x}} + \mathbf{Q}\mathbf{v}_\mu$, for appropriate \mathbf{v}_μ , uniquely solves (2.1).

Assume that we can find a dual feasible solution $\hat{\mathbf{y}}$ such that $\hat{\mathbf{Z}} = \mathbf{C} - \mathcal{A}^* \hat{\mathbf{y}} \succ 0$. Then $(\mathbf{0}, \hat{\mathbf{y}})$ can be used as the initial solution $(\mathbf{v}_0, \mathbf{y}_0)$ for the Gauss-Newton method. For each iteration, the search direction $(\Delta \mathbf{v}, \Delta \mathbf{y})$ is calculated by (approximately) finding the least-squares solution of the Gauss-Newton equation

$$-G_\mu(\mathbf{v}, \mathbf{y}) = G'_\mu(\mathbf{v}, \mathbf{y}) \begin{pmatrix} \Delta \mathbf{v} \\ \Delta \mathbf{y} \end{pmatrix}, \quad (2.2)$$

where $G'_\mu(\mathbf{v}, \mathbf{y}) : \mathbb{R}^{t(n)-m} \times \mathbb{R}^m \rightarrow \mathcal{M}^n$. We now compute the Jacobian $J := G'_\mu$ and its adjoint J^* .

For all $\mathbf{v} \in \mathbb{R}^{t(n)-m}, \mathbf{y} \in \mathbb{R}^m$, define $\mathcal{X} : \mathbb{R}^m \rightarrow \mathcal{M}^n$ and $\mathcal{Z} : \mathbb{R}^{t(n)-m} \rightarrow \mathcal{M}^n$ by:

$$\begin{aligned} \mathbf{Z} &:= \mathbf{C} - \text{sMat}(\mathbf{A}^T \mathbf{y}), & \mathbf{X} &:= \hat{\mathbf{X}} + \text{sMat}(\mathbf{Q}\mathbf{v}), \\ \mathcal{Z}(\mathbf{v}) &:= \mathbf{Z} \text{sMat}(\mathbf{Q}\mathbf{v}), & \mathcal{X}(\mathbf{y}) &:= -\text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X}. \end{aligned} \quad (2.3)$$

We then have $J = [\mathcal{Z} \mid \mathcal{X}]$ and the Gauss-Newton equation can be written as follows:

$$\begin{aligned} -G_\mu(\mathbf{v}, \mathbf{y}) &= \mathcal{Z}(\Delta \mathbf{v}) + \mathcal{X}(\Delta \mathbf{y}) \\ &= \mathbf{Z} \text{sMat}(\mathbf{Q}\Delta \mathbf{v}) - \text{sMat}(\mathbf{A}^T \Delta \mathbf{y}) \mathbf{X} \\ &= [\mathbf{C} - \text{sMat}(\mathbf{A}^T \mathbf{y})] \text{sMat}(\mathbf{Q}\Delta \mathbf{v}) - \text{sMat}(\mathbf{A}^T \Delta \mathbf{y}) [\hat{\mathbf{X}} + \text{sMat}(\mathbf{Q}\mathbf{v})]. \end{aligned} \quad (2.4)$$

This is again an overdetermined (linear) system with $t(n)$ decision variables $(\Delta \mathbf{v}, \Delta \mathbf{y})$, and with n^2 equations. In order to find the least squares solution, we need to compute $J^* = \begin{bmatrix} \mathcal{Z}^* \\ \mathcal{X}^* \end{bmatrix}$ and

$$J^* \circ J = \begin{bmatrix} \mathcal{Z}^* \circ \mathcal{Z} & \mathcal{Z}^* \circ \mathcal{X} \\ \mathcal{X}^* \circ \mathcal{Z} & \mathcal{X}^* \circ \mathcal{X} \end{bmatrix}, \quad (2.5)$$

since the final system of linear equations we (implicitly) solve is the normal equations

$$J^* \circ J \begin{pmatrix} \Delta \mathbf{v} \\ \Delta \mathbf{y} \end{pmatrix} = -J^* \circ G_\mu(\mathbf{v}, \mathbf{y}).$$

Consider $M \in \mathcal{M}^n$, we have:

$$\begin{aligned} \langle M, \mathcal{X}(\mathbf{y}) \rangle &= -\text{trace}(M^T \text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X}) \\ &= -\text{trace}(\text{sMat}(\mathbf{A}^T \mathbf{y})(\mathbf{X} M^T)) \\ &= -\text{trace} \left(\text{sMat}(\mathbf{A}^T \mathbf{y}) \frac{1}{2} (\mathbf{X} M^T + M \mathbf{X}) \right) \\ &= -\frac{1}{2} \mathbf{y}^T \mathbf{A} \text{svec}(\mathbf{X} M^T + M \mathbf{X}) \\ &= \langle \mathbf{y}, -\frac{1}{2} \mathbf{A} \text{svec}(\mathbf{X} M^T + M \mathbf{X}) \rangle, \end{aligned}$$

since $\text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X} \in \mathcal{M}^n$, $\text{sMat}(\mathbf{A}^T \mathbf{y}), \mathbf{X} \in \mathcal{S}^n$, and in addition $\mathbf{X} M^T + M \mathbf{X} \in \mathcal{S}^n$. Thus

$$\mathcal{X}^*(M) = -\frac{1}{2} \mathbf{A} \text{svec}(\mathbf{X} M^T + M \mathbf{X}). \quad (2.6)$$

Similarly, we have:

$$\begin{aligned} \langle M, \mathcal{Z}(\mathbf{v}) \rangle &= \text{trace}(M^T \mathbf{Z} \text{sMat}(\mathbf{Q} \mathbf{v})) \\ &= \text{trace}(\text{sMat}(\mathbf{Q} \mathbf{v})(M^T \mathbf{Z})) \\ &= \text{trace} \left(\text{sMat}(\mathbf{Q} \mathbf{v}) \frac{1}{2} (M^T \mathbf{Z} + \mathbf{Z} M) \right) \\ &= \frac{1}{2} \mathbf{v}^T \mathbf{Q}^T \text{svec}(M^T \mathbf{Z} + \mathbf{Z} M) \\ &= \langle \mathbf{v}, \frac{1}{2} \mathbf{Q}^T \text{svec}(M^T \mathbf{Z} + \mathbf{Z} M) \rangle, \end{aligned}$$

since $M^T \mathbf{Z} \in \mathcal{M}^n$, $\text{sMat}(\mathbf{Q} \mathbf{v}), \mathbf{Z} \in \mathcal{S}^n$, and in addition $M^T \mathbf{Z} + \mathbf{Z} M \in \mathcal{S}^n$. Thus

$$\mathcal{Z}^*(M) = \frac{1}{2} \mathbf{Q}^T \text{svec}(M^T \mathbf{Z} + \mathbf{Z} M). \quad (2.7)$$

Finally, the main operator $J^* J$ that is implicitly used for solving the Gauss-Newton equation (2.2) can be computed using the composition of \mathcal{Z}^* and \mathcal{X}^* (and \mathcal{Z} and \mathcal{X}). We let $\mathbf{Q}_j := \text{sMat}(\mathbf{Q}(:, j))$ denote the symmetric matrix corresponding to the j -th column of \mathbf{Q} , and for $\mathbf{S} \in \mathcal{S}^n$, let $|\mathbf{S}| := (\mathbf{S}^2)^{\frac{1}{2}}$. We get the weighted Gram matrices

$$\begin{aligned} \mathcal{Z}^* \circ \mathcal{Z}(\mathbf{v}) &= \frac{1}{2} \mathbf{Q}^T \text{svec}(\text{sMat}(\mathbf{Q} \mathbf{v}) \mathbf{Z}^2 + \mathbf{Z}^2 \text{sMat}(\mathbf{Q} \mathbf{v})) \\ &= (\langle |\mathbf{Z}| \mathbf{Q}_i, |\mathbf{Z}| \mathbf{Q}_j \rangle) \mathbf{v} \\ &= (\mathbf{Q}^T W_Z \mathbf{Q}) \mathbf{v} \\ \mathcal{X}^* \circ \mathcal{X}(\mathbf{y}) &= \frac{1}{2} \mathbf{A} \text{svec}(\text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X}^2 + \mathbf{X}^2 \text{sMat}(\mathbf{A}^T \mathbf{y})) \\ &= (\langle |\mathbf{X}| \mathbf{A}_i, |\mathbf{X}| \mathbf{A}_j \rangle) \mathbf{y} \\ &= (\mathbf{A} W_X \mathbf{A}^T) \mathbf{y} \\ \mathcal{Z}^* \circ \mathcal{X}(\mathbf{y}) &= \frac{1}{2} \mathbf{Q}^T \text{svec}(\mathbf{Z} \text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X} + \mathbf{X} \text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{Z}) \\ \mathcal{X}^* \circ \mathcal{Z}(\mathbf{v}) &= \frac{1}{2} \mathbf{A} \text{svec}(\mathbf{X} \text{sMat}(\mathbf{Q} \mathbf{v}) \mathbf{Z} + \mathbf{Z} \text{sMat}(\mathbf{Q} \mathbf{v}) \mathbf{X}), \end{aligned} \quad (2.8)$$

where W_Z, W_X are defined implicitly. Using a congruence with $U = \begin{bmatrix} \mathbf{Q}^T \text{svec} & 0 \\ 0 & \mathbf{A} \text{svec} \end{bmatrix}$, we get the form

$$J^* \circ J = \frac{1}{2} \begin{bmatrix} \mathbf{Q}^T \text{svec} & 0 \\ 0 & \mathbf{A} \text{svec} \end{bmatrix} \begin{bmatrix} (\cdot)\mathbf{Z}^2 + \mathbf{Z}^2(\cdot) & \mathbf{Z}(\cdot)\mathbf{X} + \mathbf{X}(\cdot)\mathbf{Z} \\ \mathbf{X}(\cdot)\mathbf{Z} + \mathbf{Z}(\cdot)\mathbf{X} & (\cdot)\mathbf{X}^2 + \mathbf{X}^2(\cdot) \end{bmatrix} \begin{bmatrix} \mathbf{Q}^T \text{svec} & 0 \\ 0 & \mathbf{A} \text{svec} \end{bmatrix}^T. \quad (2.9)$$

2.2 Algorithm Initialization

The approach described in the previous sections assumes that we have a feasible solution $(\hat{\mathbf{X}}, \hat{\mathbf{Z}})$. However, finding a feasible solution $(\hat{\mathbf{X}}, \hat{\mathbf{Z}})$ to start the algorithm is not an easy task; and, in general, it is as hard as solving the problem completely. In this section, we propose an algorithm with infeasible initial solutions. Given an initial (infeasible) solution $(\mathbf{X}_0, \mathbf{Z}_0)$ with $\mathbf{X}_0, \mathbf{Z}_0 \succ 0$, and corresponding $\mathbf{x}_0 = \text{svec}(\mathbf{X}_0)$ and $\mathbf{z}_0 = \text{svec}(\mathbf{Z}_0)$, we would like to find an optimal (and clearly feasible) solution $(\mathbf{X}^*, \mathbf{Z}^*)$. The feasibility of a solution (\mathbf{X}, \mathbf{Z}) or equivalently, (\mathbf{x}, \mathbf{z}) , is equivalent to the existence of (\mathbf{v}, \mathbf{y}) such that $\mathbf{x} = \hat{\mathbf{x}} + \mathbf{Q}\mathbf{v}$ and $\mathbf{z} = \mathbf{c} - \mathbf{A}^T\mathbf{y}$, where $\hat{\mathbf{x}} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{b}$, the least-squares solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$. Since we start the algorithm with an infeasible solution, the perturbed optimality conditions need to be written with $(\mathbf{v}, \mathbf{y}, \mathbf{x}, \mathbf{z})$ as decision variables:

$$F_\mu(\mathbf{v}, \mathbf{y}, \mathbf{x}, \mathbf{z}) := \begin{pmatrix} \mathbf{z} - \mathbf{c} + \mathbf{A}^T\mathbf{y} \\ \mathbf{x} - \hat{\mathbf{x}} - \mathbf{Q}\mathbf{v} \\ \text{sMat}(\mathbf{z}) \text{sMat}(\mathbf{x}) - \mu\mathbf{I} \end{pmatrix} =: \begin{pmatrix} \mathbf{r}_d \\ \mathbf{r}_x \\ \mathbf{R}_c \end{pmatrix} = \mathbf{0}. \quad (2.10)$$

To start the Gauss-Newton algorithm, the initial solution for \mathbf{v} and \mathbf{y} could simply be $\mathbf{v}_0 = \mathbf{0}$ and $\mathbf{y}_0 = \mathbf{0}$. In each iteration, the search direction $(\Delta\mathbf{v}, \Delta\mathbf{y}, \Delta\mathbf{x}, \Delta\mathbf{z})$ is calculated by (approximately) finding the least-squares solution of the Gauss-Newton equation

$$F'_\mu(\mathbf{v}, \mathbf{y}, \mathbf{x}, \mathbf{z}) \begin{pmatrix} \Delta\mathbf{v} \\ \Delta\mathbf{y} \\ \Delta\mathbf{x} \\ \Delta\mathbf{z} \end{pmatrix} = -F_\mu(\mathbf{v}, \mathbf{y}, \mathbf{x}, \mathbf{z}) = - \begin{pmatrix} \mathbf{r}_d \\ \mathbf{r}_x \\ \mathbf{R}_c \end{pmatrix}.$$

The above linearization yields

$$\begin{pmatrix} \Delta\mathbf{z} + \mathbf{A}^T\Delta\mathbf{y} \\ \Delta\mathbf{x} - \mathbf{Q}\Delta\mathbf{v} \\ \mathbf{Z} \text{sMat}(\Delta\mathbf{x}) + \text{sMat}(\Delta\mathbf{z})\mathbf{X} \end{pmatrix} = -F_\mu(\mathbf{v}, \mathbf{y}, \mathbf{x}, \mathbf{z}) = - \begin{pmatrix} \mathbf{r}_d \\ \mathbf{r}_x \\ \mathbf{R}_c \end{pmatrix},$$

where $\mathbf{X} = \text{sMat}(\mathbf{x})$ and $\mathbf{Z} = \text{sMat}(\mathbf{z})$.

We can now use the first two equations and substitute for $\Delta\mathbf{x} = \mathbf{Q}\Delta\mathbf{v} - \mathbf{r}_x$ and $\Delta\mathbf{z} = -\mathbf{A}^T\Delta\mathbf{y} - \mathbf{r}_d$ in the last equation to get:

$$G'_\mu(\mathbf{v}, \mathbf{y}) \begin{pmatrix} \Delta\mathbf{v} \\ \Delta\mathbf{y} \end{pmatrix} := \mathbf{Z} \text{sMat}(\mathbf{Q}\Delta\mathbf{v} - \mathbf{r}_x) - \text{sMat}(\mathbf{A}^T\Delta\mathbf{y} + \mathbf{r}_d)\mathbf{X} = -\mathbf{R}_c.$$

Moving constants to the right-hand side yields the equation

$$\mathbf{Z} \text{sMat}(\mathbf{Q}\Delta\mathbf{v}) - \text{sMat}(\mathbf{A}^T\Delta\mathbf{y})\mathbf{X} = \mathbf{Z} \text{sMat}(\mathbf{r}_x) + \text{sMat}(\mathbf{r}_d)\mathbf{X} - \mathbf{R}_c. \quad (2.11)$$

This system of equations is similar to (2.4) and we can solve it again using the Jacobian $J = [\mathcal{Z} \mid \mathcal{X}]$ and its adjoint J^* , where as above we have $\mathcal{X}(\mathbf{y}) = -\text{sMat}(\mathbf{A}^T \mathbf{y}) \mathbf{X}$, $\forall \mathbf{y} \in \mathbb{R}^m$, and $\mathcal{Z}(\mathbf{v}) = \mathbf{Z} \text{sMat}(\mathbf{Q} \mathbf{v})$, $\forall \mathbf{v} \in \mathbb{R}^{t(n)-m}$.

We can also check the change in the residual when taking a step α . We have, $\mathbf{x} + \alpha_p \Delta \mathbf{x} = \mathbf{x} + \alpha_p (\mathbf{Q} \Delta \mathbf{v} - \mathbf{r}_x)$ and $\mathbf{z} + \alpha_d \Delta \mathbf{z} = \mathbf{z} - \alpha_d (\mathbf{A}^T \Delta \mathbf{y} + \mathbf{r}_d)$. The new residuals are

$$\begin{aligned} \mathbf{r}_d &\leftarrow (\mathbf{z} - \alpha_d (\mathbf{A}^T \Delta \mathbf{y} + \mathbf{r}_d)) - \mathbf{c} + \mathbf{A}^T (\mathbf{y} + \alpha_d \Delta \mathbf{y}) = (1 - \alpha_d) \mathbf{r}_d, \\ \mathbf{r}_x &\leftarrow \mathbf{x} + \alpha_p (\mathbf{Q} \Delta \mathbf{v} - \mathbf{r}_x) - \hat{\mathbf{x}} - \mathbf{Q} (\mathbf{v} + \alpha_p \Delta \mathbf{v}) = (1 - \alpha_p) \mathbf{r}_p. \end{aligned}$$

This emphasizes the importance of taking a steplength of $\alpha_p = \alpha_d = 1$ early in the algorithm as we can then obtain exact feasibility from that point on.

2.3 Local Convergence

In this section, we would like to study the convergence of the Gauss-Newton method with infeasible initial solutions. Consider the following general problem

$$F(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} L(\mathbf{x}, \mathbf{y}) \\ R(\mathbf{x}) \end{pmatrix} = \mathbf{0},$$

where $L(\mathbf{x}, \mathbf{y}) = \mathbf{x} - \hat{\mathbf{x}} - \mathbf{T} \mathbf{y}$ is a linear function in \mathbf{x} and \mathbf{y} , \mathbf{T} is full-rank. We would like to solve the following linearized system of equations, and its equivalent system, in each iteration.

$$\left\{ F'(\mathbf{x}, \mathbf{y}) \begin{pmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{y} \end{pmatrix} = -F(\mathbf{x}, \mathbf{y}) \right\} \Leftrightarrow \left\{ \begin{array}{l} \Delta \mathbf{x} - \mathbf{T} \Delta \mathbf{y} = -L(\mathbf{x}, \mathbf{y}) \\ R'(\mathbf{x}) \Delta \mathbf{x} = -R(\mathbf{x}) \end{array} \right\}. \quad (2.12)$$

Thus $\Delta \mathbf{x} = \mathbf{T} \Delta \mathbf{y} - L(\mathbf{x}, \mathbf{y})$ and the second equation becomes:

$$R'(\mathbf{x}) \mathbf{T} \Delta \mathbf{y} = -R(\mathbf{x}) + R'(\mathbf{x}) L(\mathbf{x}, \mathbf{y}),$$

which is assumed to be an overdetermined system. We use the Gauss-Newton method to solve this equation and update (\mathbf{x}, \mathbf{y}) with $(\Delta \mathbf{x}, \Delta \mathbf{y})$. We obtain the following local convergence result.

Theorem 1. *Let R, L, T be defined as in (2.12), with $J_x(\mathbf{x}) := R'(\mathbf{x})$ and $J_y(\mathbf{x}) := J_x(\mathbf{x}) \mathbf{T}$. Let $R(\mathbf{x})$ be twice continuous differentiable in an open convex set D . Let $J_x(\mathbf{x}) \in \text{Lip}_\gamma(D)$ with $\|J_x(\mathbf{x})\| \leq \alpha$, for all $\mathbf{x} \in D$; and, suppose that there exists $\mathbf{x}_* \in D$, \mathbf{y}_* , and $\lambda, \sigma \geq 0$, such that $L(\mathbf{x}_*, \mathbf{y}_*) = \mathbf{0}$, $J_x(\mathbf{x}_*)^T R(\mathbf{x}_*) = \mathbf{0}$, λ is the smallest eigenvalue of $J_y(\mathbf{x}_*)^T J_y(\mathbf{x}_*)$, and*

$$\|(J_x(\mathbf{x}) - J_x(\mathbf{x}_*))^T R(\mathbf{x}_*)\|_2 \leq \sigma \|\mathbf{x} - \mathbf{x}_*\|_2, \quad \forall \mathbf{x} \in D.$$

Define the iteration

$$\begin{aligned} \Delta \mathbf{x}_k &= \mathbf{T} \Delta \mathbf{y}_k - L(\mathbf{x}_k, \mathbf{y}_k), \quad \forall k, & \Delta \mathbf{y}_k &= (J_y(\mathbf{x}_k)^T J_y(\mathbf{x}_k))^{-1} J_y(\mathbf{x}_k)^T (-R(\mathbf{x}_k) + J_x(\mathbf{x}_k) L(\mathbf{x}_k, \mathbf{y}_k)), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \Delta \mathbf{x}_k, & \mathbf{y}_{k+1} &= \mathbf{y}_k + \Delta \mathbf{y}_k. \end{aligned}$$

Let $\tau := \|\mathbf{T}\|_2$ and $c \in (1, \lambda/(\tau^2 \sigma))$. If $\tau^2 \sigma < \lambda$, then there exists $\epsilon > 0$ such that $\forall \mathbf{x}_0 \in N(\mathbf{x}_*, \epsilon)$ and \mathbf{y}_0 arbitrary, we have that: the sequences $\{\mathbf{x}_k\}, \{\mathbf{y}_k\}$ are well-defined; they converge to $(\mathbf{x}_*, \mathbf{y}_*)$ with $L(\mathbf{x}_k, \mathbf{y}_k) = \mathbf{0}, \forall k \geq 1$; and

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\| &\leq \frac{\tau^2 c}{\lambda} \left(\sigma \|\mathbf{x}_k - \mathbf{x}_*\| + \frac{\alpha \gamma}{2} \|\mathbf{x}_k - \mathbf{x}_*\|^2 \right), \\ \|\mathbf{x}_{k+1} - \mathbf{x}_*\| &\leq \frac{\lambda + \tau^2 c \sigma}{2\lambda} \|\mathbf{x}_k - \mathbf{x}_*\| < \|\mathbf{x}_k - \mathbf{x}_*\|. \end{aligned}$$

Proof. Using induction and the formulations $\Delta \mathbf{y}_k = (J_y(\mathbf{x}_k)^T J_y(\mathbf{x}_k))^{-1} J_y(\mathbf{x}_k)^T (-R(\mathbf{x}_k) + J_x(\mathbf{x}_k)L(\mathbf{x}_k, \mathbf{y}_k))$, $\Delta \mathbf{x}_k = \mathbf{T} \Delta \mathbf{y}_k - L(\mathbf{x}_k, \mathbf{y}_k)$, we see that $L(\mathbf{x}_k, \mathbf{y}_k) = \mathbf{0}, \forall k \geq 1$. Thus we have:

$$\mathbf{x}_k - \mathbf{x}_* = \mathbf{T}(\mathbf{y}_k - \mathbf{y}_*), \quad \forall k \geq 1.$$

We now continue using induction. We have:

$$\begin{aligned} \mathbf{y}_1 - \mathbf{y}_* &= \mathbf{y}_0 - \mathbf{y}_* + (J_y(\mathbf{x}_0)^T J_y(\mathbf{x}_0))^{-1} J_y(\mathbf{x}_0)^T [-R(\mathbf{x}_0) + J_x(\mathbf{x}_0)L(\mathbf{x}_0, \mathbf{y}_0)] \\ &= (J_y(\mathbf{x}_0)^T J_y(\mathbf{x}_0))^{-1} \mathbf{M}, \end{aligned}$$

where $\mathbf{M} = -J_y(\mathbf{x}_0)^T R(\mathbf{x}_0) + J_y(\mathbf{x}_0)^T J_x(\mathbf{x}_0)L(\mathbf{x}_0, \mathbf{y}_0) - J_y(\mathbf{x}_0)^T J_y(\mathbf{x}_0)(\mathbf{y}_* - \mathbf{y}_0)$. Rewriting, we obtain

$$\begin{aligned} \mathbf{M} &= -J_y(\mathbf{x}_0)^T R(\mathbf{x}_*) + J_y(\mathbf{x}_0)^T [R(\mathbf{x}_*) - R(\mathbf{x}_0) - J_y(\mathbf{x}_0)(\mathbf{y}_* - \mathbf{y}_0) + J_x(\mathbf{x}_0)L(\mathbf{x}_0, \mathbf{y}_0)] \\ &= -J_y(\mathbf{x}_0)^T R(\mathbf{x}_*) + \mathbf{T}^T J_x(\mathbf{x}_0)^T [R(\mathbf{x}_*) - R(\mathbf{x}_0) - J_x(\mathbf{x}_0)(\mathbf{T}\mathbf{y}_* - \mathbf{T}\mathbf{y}_0 - \mathbf{x}_0 + \hat{\mathbf{x}} + \mathbf{T}\mathbf{y}_0)] \\ &= -\mathbf{T}^T (J_x(\mathbf{x}_0) - J_x(\mathbf{x}_*))^T R(\mathbf{x}_*) + \mathbf{T}^T J_x(\mathbf{x}_0)^T [R(\mathbf{x}_*) - R(\mathbf{x}_0) - J_x(\mathbf{x}_0)(\mathbf{x}_* - \mathbf{x}_0)]. \end{aligned}$$

According to Dennis and Schnabel [3, Theorem 3.1.4], there exists $\epsilon_1 > 0$ such that $J_y(\mathbf{x})^T J_y(\mathbf{x})$ is nonsingular and

$$\|(J_y(\mathbf{x})^T J_y(\mathbf{x}))^{-1}\| \leq \frac{c}{\lambda}, \quad \forall \mathbf{x} \in N(\mathbf{x}_*, \epsilon_1),$$

where $c \in (1, \lambda/\sigma)$. We also have:

$$\|\mathbf{T}^T (J_x(\mathbf{x}_0) - J_x(\mathbf{x}_*))^T R(\mathbf{x}_*)\| \leq \|\mathbf{T}\| \|(J_x(\mathbf{x}_0) - J_x(\mathbf{x}_*))^T R(\mathbf{x}_*)\| \leq \tau\sigma \|\mathbf{x}_0 - \mathbf{x}_*\|, \quad \forall \mathbf{x} \in D.$$

Applying the Lipschitz condition of $J_x(\cdot)$, we obtain the following equation

$$\|R(\mathbf{x}_*) - R(\mathbf{x}_0) - J_x(\mathbf{x}_0)(\mathbf{x}_* - \mathbf{x}_0)\| \leq \frac{\gamma}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|^2.$$

Thus we have:

$$\|\mathbf{T}^T J_x(\mathbf{x}_0)^T [R(\mathbf{x}_*) - R(\mathbf{x}_0) - J_x(\mathbf{x}_0)(\mathbf{x}_* - \mathbf{x}_0)]\| \leq \frac{\tau\alpha\gamma}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|^2.$$

Combining these inequalities gives the following bound

$$\|\mathbf{y}_1 - \mathbf{y}_*\| \leq \frac{c}{\lambda} \left(\tau\sigma \|\mathbf{x}_0 - \mathbf{x}_*\| + \frac{\tau\alpha\gamma}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \right).$$

Or equivalently, we have:

$$\|\mathbf{x}_1 - \mathbf{x}_*\| \leq \frac{\tau^2 c}{\lambda} \left(\sigma \|\mathbf{x}_0 - \mathbf{x}_*\| + \frac{\alpha\gamma}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 \right), \quad \forall \mathbf{x}_0 \in N(\mathbf{x}_*, \epsilon_1).$$

Now let $\epsilon = \min \left\{ \epsilon_1, \frac{\lambda - \tau^2 c \sigma}{\tau^2 c \alpha \gamma} \right\}$. We have:

$$\begin{aligned} \|\mathbf{x}_1 - \mathbf{x}_*\| &\leq \|\mathbf{x}_0 - \mathbf{x}_*\| \left(\frac{\tau^2 c \sigma}{\lambda} + \frac{\tau^2 c \alpha \gamma}{2\lambda} \|\mathbf{x}_0 - \mathbf{x}_*\| \right) \\ &\leq \frac{\lambda + \tau^2 c \sigma}{2\lambda} \|\mathbf{x}_0 - \mathbf{x}_*\| \\ &< \|\mathbf{x}_0 - \mathbf{x}_*\|, \quad \forall \mathbf{x}_0 \in N(\mathbf{x}_*, \epsilon). \end{aligned}$$

From these results, clearly $\mathbf{x}_1 \in N(\mathbf{x}_*, \epsilon)$ and the induction step is exactly the same as in the case $k = 0$ above. \square

Applying Theorem 1 to the Gauss-Newton equation (2.10), we have:

$$L(\mathbf{z}, \mathbf{x}, \mathbf{y}, \mathbf{v}) = \begin{pmatrix} \mathbf{z} - \mathbf{c} + \mathbf{A}^T \mathbf{y} \\ \mathbf{x} - \hat{\mathbf{x}} - \mathbf{Q}\mathbf{v} \end{pmatrix},$$

where in the theorem (\mathbf{z}, \mathbf{x}) takes on the value \mathbf{x} and (\mathbf{y}, \mathbf{v}) takes on the value \mathbf{y} , with $\mathbf{T} = \begin{pmatrix} -\mathbf{A}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{Q} \end{pmatrix}$. We also have

$$R_\mu(\mathbf{z}, \mathbf{x}) = \text{sMat}(\mathbf{z}) \text{sMat}(\mathbf{x}) - \mu \mathbf{I}.$$

The optimal solution $(\mathbf{z}_*, \mathbf{x}_*)$ satisfies the condition $R_\mu(\mathbf{z}_*, \mathbf{x}_*) = 0$. Therefore, we can set $\sigma = 0$. We have: $\tau = \|\mathbf{T}\| = \sqrt{\|\mathbf{A}\|^2 + \|\mathbf{Q}\|^2}$. Now consider $R'(\mathbf{z}, \mathbf{x})$. We have:

$$R'(\mathbf{z}, \mathbf{x}) \begin{pmatrix} \Delta \mathbf{z} \\ \Delta \mathbf{x} \end{pmatrix} = \mathbf{Z} \text{sMat}(\Delta \mathbf{x}) + \text{sMat}(\Delta \mathbf{z}) \mathbf{X},$$

where $\mathbf{Z} = \text{sMat}(\mathbf{z})$ and $\mathbf{X} = \text{sMat}(\mathbf{x})$. For all $(\Delta \mathbf{z}, \Delta \mathbf{x})$, we have:

$$\left\| R'(\mathbf{z}, \mathbf{x}) \begin{pmatrix} \Delta \mathbf{z} \\ \Delta \mathbf{x} \end{pmatrix} \right\| \leq \|\mathbf{Z}\| \|\Delta \mathbf{x}\| + \|\Delta \mathbf{z}\| \|\mathbf{X}\| \leq \sqrt{\|\mathbf{Z}\|^2 + \|\mathbf{X}\|^2} \sqrt{\|\Delta \mathbf{z}\|^2 + \|\Delta \mathbf{x}\|^2}.$$

Thus, $\|R'(\mathbf{z}, \mathbf{x})\| \leq \sqrt{\|\mathbf{Z}\|^2 + \|\mathbf{X}\|^2} = \|(\mathbf{z}, \mathbf{x})\|$. Thus the parameter α can be calculated using the formulation $\alpha = \sup_{(\mathbf{z}, \mathbf{x}) \in D} \|(\mathbf{z}, \mathbf{x})\|$. Similarly, we have:

$$[R'(\mathbf{z}_1, \mathbf{x}_1) - R'(\mathbf{z}_2, \mathbf{x}_2)] \begin{pmatrix} \Delta \mathbf{z} \\ \Delta \mathbf{x} \end{pmatrix} = (\mathbf{Z}_1 - \mathbf{Z}_2) \text{sMat}(\Delta \mathbf{x}) + \text{sMat}(\Delta \mathbf{z})(\mathbf{X}_1 - \mathbf{X}_2).$$

Thus for all $(\Delta \mathbf{z}, \Delta \mathbf{x})$,

$$\left\| [R'(\mathbf{z}_1, \mathbf{x}_1) - R'(\mathbf{z}_2, \mathbf{x}_2)] \begin{pmatrix} \Delta \mathbf{z} \\ \Delta \mathbf{x} \end{pmatrix} \right\| \leq \|(\mathbf{z}_1, \mathbf{x}_1) - (\mathbf{z}_2, \mathbf{x}_2)\| \|(\Delta \mathbf{z}, \Delta \mathbf{x})\|,$$

which implies that $\|R'(\mathbf{z}_1, \mathbf{x}_1) - R'(\mathbf{z}_2, \mathbf{x}_2)\| \leq \|(\mathbf{z}_1, \mathbf{x}_1) - (\mathbf{z}_2, \mathbf{x}_2)\|$. We can set the parameter γ to be 1. The remaining parameter is the smallest eigenvalue λ of $J_y(\mathbf{z}_*, \mathbf{x}_*)^T J_y(\mathbf{z}_*, \mathbf{x}_*)$. Kruk et al. [11] show that under the strict complementary slackness condition of a unique optimal solution, the Jacobian J at the optimal solution is full rank ($\lambda > 0$). This implies that the algorithm is well-behaved under the strict complementary slackness condition. Theorem 1 indicates that there is a neighborhood around the optimal solution where the Gauss-Newton method converges quadratically. However, it is difficult to estimate the value of λ . We need all of these parameters to define the neighborhood $N(\mathbf{x}_*, \epsilon)$ in which one can simply use affine scaling and the Gauss-Newton method instead of *damped* Gauss-Newton method, i.e. we can set the barrier parameter $\mu = 0$ and take full step lengths equal to one. We call this the *crossover technique*. A heuristic rule could be built based on the infeasibility of the current solution (\mathbf{z}, \mathbf{x}) since the damped Gauss-Newton method with step sizes different from 1 (with infeasible initial solutions) will not guarantee the feasibility after one step. (See [23] for more details on other possible heuristic rules.)

2.4 Presolve and Preconditioning

The purpose of a presolve is to construct well-conditioned and sparse bases for both the range and nullspace of \mathbf{A} . These bases play an important role in the matrix-vector multiplications involved in the algorithm. If \mathbf{A} is sparse and $\mathbf{A} = [\mathbf{S} \mid \mathbf{E}]$, where $\mathbf{S} \in \mathbb{R}^{m \times m}$ is well-conditioned and $\mathbf{E} \in \mathbb{R}^{m \times (t(n)-m)}$, then we can use the columns of $\mathbf{Q} = \begin{bmatrix} -\mathbf{S}^{-1}\mathbf{E} \\ \mathbf{I} \end{bmatrix}$ for the basis of the nullspace of \mathbf{A} . The matrix inversion is fast and accurate if \mathbf{S} is a well-conditioned triangular matrix. In addition, we can also consider row and column permutations on \mathbf{A} ; therefore, the main task in the presolve is to find an $m \times m$ submatrix of \mathbf{A} (up to a permutation of rows and columns) that is well-conditioned and (approximately) triangular. In this case, \mathbf{S} and \mathbf{E} are considered to be the *main inputs* to the algorithm.

If \mathbf{A} is not sparse, it is better to keep \mathbf{A} and $\mathbf{P} = -\mathbf{S}^{-1}\mathbf{E}$ as the main inputs (with $\mathbf{Q} = [\mathbf{P}; \mathbf{I}]$). And the major task of the presolve is to find a well-conditioned submatrix \mathbf{S} . This can be done by randomly picking m columns from \mathbf{A} and reforming the submatrix \mathbf{S} until the condition number satisfies some given conditions.

The most expensive part of the algorithm is solving the overdetermined system (2.4), or equivalently, (2.11). Preconditioning is essential to reduce the computational time when using iterative methods for this part. Since we implicitly solve the normal equations, we want to find a nonsingular transformation $T : \mathbb{R}^{t(n)} \rightarrow \mathbb{R}^{t(n)}$ such that $(T^{-1})^*(J^* \circ J) \circ T^{-1}$ is well-conditioned. Instead of solving the least squares system of (2.4) or (2.11), we would like to solve the systems

$$\begin{cases} (T^{-1})^* \circ J^* \circ J \circ T^{-1}(\Delta \mathbf{q}) = (T^{-1})^* \circ J^* \circ \mathbf{R} \\ T(\Delta \mathbf{v}, \Delta \mathbf{y}) = \Delta \mathbf{q}, \end{cases} \quad (2.13)$$

where \mathbf{R} is the corresponding right-hand side of the least squares problem.

Let $T^{-1} = S = \begin{bmatrix} S_v \\ S_y \end{bmatrix}$ where $S_v : \mathbb{R}^{t(n)} \rightarrow \mathbb{R}^{t(n)-m}$ and $S_y : \mathbb{R}^{t(n)} \rightarrow \mathbb{R}^m$. We then have:

$$J \circ T^{-1}(\mathbf{q}) = J(S_v(\mathbf{q}), S_y(\mathbf{q})) = \mathcal{Z}(S_v(\mathbf{q})) + \mathcal{X}(S_y(\mathbf{q})) = \mathcal{Z} \circ S_v(\mathbf{q}) + \mathcal{X} \circ S_y(\mathbf{q}). \quad (2.14)$$

To find the adjoint $(J \circ T^{-1})^*$, we note that $(T^{-1})^* = S^* = [S_v^* \mid S_y^*]$ implies:

$$(J \circ T^{-1})^*(\mathbf{M}) = (T^{-1})^* \circ J^*(\mathbf{M}) = (T^{-1})^* \left(\begin{bmatrix} \mathcal{Z}^*(\mathbf{M}) \\ \mathcal{X}^*(\mathbf{M}) \end{bmatrix} \right) = S_v^* \circ \mathcal{Z}^*(\mathbf{M}) + S_y^* \circ \mathcal{X}^*(\mathbf{M}), \quad (2.15)$$

where \mathcal{X}^* and \mathcal{Z}^* are formulated in (2.6) and (2.7), respectively.

Consider the special case in which T is a separable transformation with respect to $\Delta \mathbf{v}$ and $\Delta \mathbf{y}$: $T(\Delta \mathbf{v}, \Delta \mathbf{y}) = \begin{bmatrix} T_v(\Delta \mathbf{v}) \\ T_y(\Delta \mathbf{y}) \end{bmatrix}$, where $T_v : \mathbb{R}^{t(n)-m} \rightarrow \mathbb{R}^{t(n)-m}$ and $T_y : \mathbb{R}^m \rightarrow \mathbb{R}^m$ are both nonsingular. If $\mathbf{q} = (\mathbf{w}, \mathbf{u})$ with $\mathbf{w} \in \mathbb{R}^{t(n)-m}$ and $\mathbf{u} \in \mathbb{R}^m$, we then have:

$$J \circ T^{-1}(\mathbf{w}, \mathbf{u}) = \mathcal{Z} \circ T_v^{-1}(\mathbf{w}) + \mathcal{X} \circ T_y^{-1}(\mathbf{u}). \quad (2.16)$$

And the adjoint is written as follows:

$$(J \circ T^{-1})^*(\mathbf{M}) = \begin{bmatrix} (T_v^{-1})^* \circ \mathcal{Z}^*(\mathbf{M}) \\ (T_y^{-1})^* \circ \mathcal{X}^*(\mathbf{M}) \end{bmatrix}. \quad (2.17)$$

The first preconditioner we would like to consider is a simple diagonal scaling transformation \mathcal{S} . Since the original system of linear equations we need to solve is $J^* \circ J(\Delta \mathbf{v}, \Delta \mathbf{y}) = J^*(\mathbf{R})$, we would want to have $\mathcal{S}^* \circ \mathcal{S}$ approximate $J^* \circ J$. We define the diagonal preconditioner \mathcal{S} as follows:

$$\mathcal{S}(\mathbf{e}_i) = \|J(\mathbf{e}_i)\|_F \mathbf{e}_i, \quad \forall i = 1, \dots, t(n), \quad (2.18)$$

where \mathbf{e}_i is the i -th unit vector in $\mathbb{R}^{t(n)}$.

The second preconditioner is the separable transformation (block diagonal) \mathcal{Q} such that $\mathcal{Q}_v^* \circ \mathcal{Q}_v$ approximates $\mathcal{Z}^* \circ \mathcal{Z}$ and $\mathcal{Q}_y \circ \mathcal{Q}_y$ approximates $\mathcal{X}^* \circ \mathcal{X}$. Using the QR decomposition, we can decompose \mathcal{Z} into $\mathcal{Q}_Z \circ \mathcal{R}_Z$, where $\mathcal{Q}_Z : \mathbb{R}^{t(n)} \rightarrow \mathcal{M}^n$ is a unitary transformation, and $\mathcal{R}_Z : \mathbb{R}^{t(n)-m} \rightarrow \mathbb{R}^{t(n)}$, with an upper triangular transformation matrix \mathbf{R}_Z . Clearly, we can set $\mathcal{Q}_v = \mathcal{P}_v \circ \mathcal{R}_Z$, where \mathcal{P}_v is the projection transformation from \mathbb{R}^{n^2} to $\mathbb{R}^{t(n)-m}$, since $\mathcal{R}_Z^* \circ \mathcal{R}_Z = \mathcal{Z}^* \circ \mathcal{Z}$. Similarly, we can set $\mathcal{Q}_y = \mathcal{P}_y \circ \mathcal{R}_X$ where $\mathcal{X} = \mathcal{Q}_X \circ \mathcal{R}_X$ under the QR decomposition and \mathcal{P}_y is the projection transformation from \mathbb{R}^{n^2} to \mathbb{R}^m .

We now would like to consider the properties of these two particular preconditioners. Consider the following measure of conditioning that depends uniformly on all the eigenvalues of a positive definite matrix \mathbf{X} ,

$$\omega(\mathbf{X}) := \frac{\text{trace}(\mathbf{X})/n}{\det(\mathbf{X})^{\frac{1}{n}}}.$$

This measure is a pseudoconvex function. Note that a function is pseudoconvex if

$$(\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) \geq 0 \implies f(\mathbf{y}) \geq f(\mathbf{x}),$$

and for pseudoconvex functions, all stationary solutions are global minimizers (see for example, [15]). We now show that we have obtained the *optimal block diagonal preconditioner* with respect to the measure ω , thus extending the corresponding diagonal preconditioner result in [4].

Proposition 1. *The measure $\omega(\mathbf{A})$ satisfies*

1. $1 \leq \omega(\mathbf{A}) \leq \kappa(\mathbf{A}) < \frac{(\kappa(\mathbf{A}) + 1)^2}{\kappa(\mathbf{A})} \leq 4\omega^n(\mathbf{A})$, where $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$,

with equality in the first and second inequality if and only if \mathbf{A} is a multiple of the identity and equality in the last inequality if and only if

$$\lambda_2 = \dots = \lambda_{n-1} = \frac{\lambda_1 + \lambda_n}{2};$$

2. $\omega(\alpha \mathbf{A}) = \omega(\mathbf{A})$, for all $\alpha > 0$;

3. if $n = 2$, $\omega(\mathbf{A})$ is isotonic with $\kappa(\mathbf{A})$.

4. The measure ω is pseudoconvex on the set of s.p.d. matrices, and thus any stationary point is a global minimizer of ω .

5. Let \mathbf{V} be a full rank $m \times n$ matrix, $n \leq m$. Then the optimal column scaling that minimizes the measure ω , i.e.

$$\min \omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D})),$$

over \mathbf{D} positive, diagonal, is given by

$$D_{ii} = \frac{1}{\|\mathbf{V}_{:,i}\|}, i = 1, \dots, n,$$

where $\mathbf{V}_{:,i}$ is the i -th column of \mathbf{V} .

6. Let \mathbf{V} be a full rank $m \times n$ matrix, $n \leq m$ with block structure $\mathbf{V} = [\mathbf{V}_1 \ \mathbf{V}_2 \ \dots \ \mathbf{V}_k]$. Then the optimal corresponding block diagonal scaling

$$\mathbf{D} = \begin{bmatrix} \mathbf{D}_1 & 0 & 0 & \dots & 0 \\ 0 & \mathbf{D}_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \mathbf{D}_k \end{bmatrix}$$

that minimizes the measure ω , i.e.

$$\min \omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D})),$$

over \mathbf{D} block diagonal, is given by

$$\mathbf{D}_i = \{\mathbf{V}_i^T \mathbf{V}_i\}^{-\frac{1}{2}}, i = 1, \dots, k.$$

Proof. Items 1 to 5 are proved in [4].

To prove 6, let the blocked \mathbf{V} be given. Then the arithmetic-geometric mean inequality yields

$$\begin{aligned} \omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D})) &= \frac{\text{trace}(\mathbf{D}^T \mathbf{V}^T \mathbf{V} \mathbf{D})/n}{\det(\mathbf{D}^T \mathbf{V}^T \mathbf{V} \mathbf{D})^{\frac{1}{n}}} \\ &= \frac{\text{trace}(\mathbf{V}^T \mathbf{V} \mathbf{D} \mathbf{D}^T)/n}{\det(\mathbf{D} \mathbf{D}^T)^{\frac{1}{n}} \det(\mathbf{V}^T \mathbf{V})^{\frac{1}{n}}}. \end{aligned}$$

Let $\bar{\mathbf{D}} = \mathbf{D} \mathbf{D}^T$, we have, $\bar{\mathbf{D}}$ is also a block diagonal matrix with $\bar{\mathbf{D}}_i = \mathbf{D}_i \mathbf{D}_i^T$ for all $i = 1, \dots, k$. Thus

$$\begin{aligned} \omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D})) &= \frac{\sum_{i=1}^k (\mathbf{V}_i^T \mathbf{V}_i) \cdot \bar{\mathbf{D}}_i / n}{\det(\mathbf{V}^T \mathbf{V})^{\frac{1}{n}} \prod_{i=1}^k \det(\bar{\mathbf{D}}_i)^{\frac{1}{n}}} \\ &= \left(\frac{1}{n \det(\mathbf{V}^T \mathbf{V})^{\frac{1}{n}}} \right) \frac{\sum_{i=1}^k \bar{\mathbf{V}}_i \cdot \bar{\mathbf{D}}_i}{\prod_{i=1}^k \det(\bar{\mathbf{D}}_i)^{\frac{1}{n}}}, \end{aligned}$$

where $\bar{\mathbf{V}}_i = \mathbf{V}_i^T \mathbf{V}_i$ for all $i = 1, \dots, k$. Consider the function $f(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k) = \frac{\sum_{i=1}^k \bar{\mathbf{V}}_i \cdot \bar{\mathbf{D}}_i}{\prod_{i=1}^k \det(\bar{\mathbf{D}}_i)^{\frac{1}{n}}}$,

we have:

$$\frac{\partial f(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k)}{\partial \bar{\mathbf{D}}_i} = \prod_{j=1}^k \det(\bar{\mathbf{D}}_j)^{-\frac{1}{n}} \left(\bar{\mathbf{V}}_i - \frac{\sum_{j=1}^k \bar{\mathbf{V}}_j \cdot \bar{\mathbf{D}}_j}{n} (\bar{\mathbf{D}}_i^{-1})^T \right).$$

Let $\bar{\mathbf{D}}_i = \bar{\mathbf{V}}_i^{-1}$ for all $i = 1, \dots, k$, we have: $\bar{\mathbf{D}}_i$ is symmetric and $\sum_{i=1}^k \bar{\mathbf{V}}_i \cdot \bar{\mathbf{D}}_i = n$ for all $i = 1, \dots, k$.

Thus $\frac{f(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k)}{\partial \bar{\mathbf{D}}_i} = \mathbf{0}$ for all $i = 1, \dots, k$. The measure ω is pseudoconvex and similarly, the

function $f(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k)$ is also pseudoconvex (quotient of a convex (linear) function and a positive concave function). Thus the matrices $(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k)$ with $\bar{\mathbf{D}}_i = \bar{\mathbf{V}}_i^{-1}$ for $i = 1, \dots, k$, minimizes $f(\bar{\mathbf{D}}_1, \dots, \bar{\mathbf{D}}_k)$. This implies that the symmetric block matrix \mathbf{D} with $\mathbf{D}_i = (\mathbf{V}_i^T \mathbf{V}_i)^{-\frac{1}{2}}$ minimizes the measure $\omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D}))$. \square

Note that the optimality condition only requires $\bar{\mathbf{D}}_i = \bar{\mathbf{V}}_i^{-1}$ for $i = 1, \dots, k$, which means we can use QR decomposition to find the optimal \mathbf{D}_i instead of matrix square root. In other words, if \mathbf{V}_i has the QR decomposition $\mathbf{V}_i = \mathbf{Q}_i \mathbf{R}_i$, we can set $\mathbf{D}_i = \mathbf{R}_i^{-1}$ for all $i = 1, \dots, k$ to minimize the measure $\omega((\mathbf{V}\mathbf{D})^T(\mathbf{V}\mathbf{D}))$.

3 Lovász Theta Function Problem

In this section, we apply the Gauss-Newton method to the Lovász theta function problem. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an undirected graph; and let $n = |\mathcal{V}|$ and $m = |\mathcal{E}|$ be the number of nodes and edges, respectively. The Lovász theta number (defined in [13]) is the optimal value of the following SDP

$$\begin{aligned} \vartheta(\mathcal{G}) := p^* := & \max \quad \mathbf{E} \cdot \mathbf{X} \\ \text{(TN)} \quad & \text{s.t.} \quad \mathbf{I} \cdot \mathbf{X} = 1 \\ & \mathbf{E}_{ij} \cdot \mathbf{X} = 0, \quad \forall (i, j) \in \mathcal{E} \\ & \mathbf{X} \succeq 0 \end{aligned} \quad (3.1)$$

The dual SDP is

$$\begin{aligned} d^* := & \min \quad z \\ \text{(DTN)} \quad & \text{s.t.} \quad z\mathbf{I} + \sum_{(i,j) \in \mathcal{E}} y_{ij} \mathbf{E}_{ij} - \mathbf{Z} = \mathbf{E}, \\ & \mathbf{Z} \succeq 0, \end{aligned} \quad (3.2)$$

where $\mathbf{y} = (y_{ij})_{(i,j) \in \mathcal{E}} \in \mathbb{R}^m$; \mathbf{I} is the identity matrix; \mathbf{E} is the matrix of all ones, $\mathbf{E}_{ij} = (e_i e_j^T + e_j e_i^T) / \sqrt{2}$ is the ij -th unit matrix in \mathcal{S}^n , and e_i is the i -th unit vector.

The theta number has important properties, e.g., it is tractable (can be computed in polynomial time) and it provides bounds for the stability and chromatic numbers of the graph, see e.g., [9, 12]. We now show how to use the Gauss-Newton approach to solve this Lovász theta function problem.

3.1 Matrix-Free Formulation

We define additional linear transformations on vectors and matrices for this problem. Let $\mathcal{G}^c = (\mathcal{V}, \mathcal{E}^c)$ be the complement graph of \mathcal{G} , i.e. \mathcal{E}^c is the edge set complement to \mathcal{E} . Let $m_c := |\mathcal{E}^c| = t(n) - m$, where $t(n) := \binom{n}{2}$, and consider two bijective index mappings, $\text{ind}_{\mathcal{E}} : \mathcal{E} \rightarrow \{1, \dots, m\}$, and $\text{ind}_{\mathcal{E}^c} : \mathcal{E}^c \rightarrow \{1, \dots, m_c\}$ with their corresponding inverses, $\text{ind}_{\mathcal{E}}^{-1}$ and $\text{ind}_{\mathcal{E}^c}^{-1}$.

We now define two linear transformations, $\text{sMat}_{\mathcal{E}} : \mathbb{R}^m \rightarrow \mathcal{S}^n$ and $\text{sMat}_{\mathcal{E}^c} : \mathbb{R}^{m_c} \rightarrow \mathcal{S}^n$ as follows:

$$\text{sMat}_{\mathcal{E}}(\mathbf{y}) = \sum_{i=1}^m y_i \mathbf{E}_{\text{ind}_{\mathcal{E}}^{-1}(i)}, \quad \text{sMat}_{\mathcal{E}^c}(\mathbf{v}) = \sum_{i=1}^{m_c} v_i \mathbf{E}_{\text{ind}_{\mathcal{E}^c}^{-1}(i)}.$$

Let $\text{svec}_{\mathcal{E}} := \text{sMat}_{\mathcal{E}}^{\dagger}$ denote the *Moore-Penrose generalized inverse* mapping from \mathcal{S}^n to \mathbb{R}^m , $\text{svec}_{\mathcal{E}}(\mathbf{S}) = \mathbf{y}$, where $y_i = S_{\text{ind}_{\mathcal{E}}^{-1}(i)} \sqrt{2}$ for all $i = 1, \dots, m$. Similarly, we can also define $\text{svec}_{\mathcal{E}^c} :=$

$\text{sMat}_{\mathcal{E}^c}^\dagger$. The mapping $\text{svec}_{\mathcal{E}^c}$ (resp. $\text{svec}_{\mathcal{E}}$) is an inverse mapping if we restrict to the subspace of symmetric matrices with zero in positions corresponding to the edge set \mathcal{E}^c (resp. \mathcal{E}). The adjoint operator $\text{sMat}_{\mathcal{E}^c}^* = \text{svec}_{\mathcal{E}^c}$, since

$$\begin{aligned}\langle \text{sMat}_{\mathcal{E}^c}(\mathbf{v}), \mathbf{S} \rangle &= \text{trace}(\text{sMat}_{\mathcal{E}^c}(\mathbf{v})\mathbf{S}) \\ &= \mathbf{v}^T \text{svec}_{\mathcal{E}^c}(\mathbf{S}) = \langle \mathbf{v}, \text{svec}_{\mathcal{E}^c}(\mathbf{S}) \rangle,\end{aligned}$$

and similarly, the adjoint operator $\text{sMat}_{\mathcal{E}}^* = \text{svec}_{\mathcal{E}}$.

Using these transformations, we can represent primal and dual feasibility as follows. Let the matrix $\mathbf{V} := \begin{pmatrix} -\mathbf{e}^T \\ \mathbf{I} \end{pmatrix} \in \mathbb{R}^{n \times (n-1)}$, we have: $\mathbf{X} \succeq 0$ is primal feasible if and only there exists $\mathbf{w} \in \mathbb{R}^{n-1}$ and $\mathbf{v} \in \mathbb{R}^{m_c}$ such that

$$\mathbf{X} = \frac{1}{n}\mathbf{I} + \text{Diag}(\mathbf{V}\mathbf{w}) + \text{sMat}_{\mathcal{E}^c}(\mathbf{v}),$$

where $\text{Diag} : \mathbb{R}^n \rightarrow \mathcal{M}^n$ is the diagonal transformation with its diagonal adjoint $\text{Diag}^* = \text{Diag}^\dagger = \text{diag} : \mathcal{M}^n \rightarrow \mathbb{R}^n$. We could also use a different matrix \mathbf{V} to maintain the isometry,

$$\mathbf{V} := \begin{pmatrix} c_1 \mathbf{e}^T \\ c_2 \mathbf{E} + \mathbf{I} \end{pmatrix} \in \mathbb{R}^{n \times (n-1)},$$

where $c_1 = -1/\sqrt{n}$ and $c_2 = -1/(n + \sqrt{n})$.

Now let us consider the dual feasibility. Similarly, $\mathbf{Z} \succeq 0$ is dual feasible if and only if there exists $z \in \mathbb{R}$ and $\mathbf{y} \in \mathbb{R}^m$ such that

$$\mathbf{Z} = -\mathbf{E} + z\mathbf{I} + \text{sMat}_{\mathcal{E}}(\mathbf{y}).$$

As in the general case, we would like to solve the perturbed optimality condition $\mathbf{Z}\mathbf{X} - \mu\mathbf{I} = \mathbf{0}$ in each iteration of the primal-dual path following interior algorithm. We again use an inexact Gauss-Newton approach with a matrix-free preconditioned conjugate gradient method.

The major operation of the Gauss-Newton approach is to calculate the Jacobian and its adjoint. The Jacobian $J := G'_\mu$ can be written as $J = [\mathcal{Z} \mid \mathcal{X}]$, where \mathcal{Z} and \mathcal{X} are two transformations defined as follows:

$$\mathcal{Z}(\mathbf{w}, \mathbf{v}) = \mathbf{Z}(\text{Diag}(\mathbf{V}\mathbf{w}) + \text{sMat}_{\mathcal{E}^c}(\mathbf{v})), \quad \mathcal{X}(z, \mathbf{y}) = (z\mathbf{I} + \text{sMat}_{\mathcal{E}}(\mathbf{y}))\mathbf{X}.$$

In order to find the adjoint J^* , we need to compute \mathcal{Z}^* and \mathcal{X}^* . Consider $\mathbf{M} \in \mathcal{M}^n$, we have:

$$\begin{aligned}\langle \mathbf{M}, \mathcal{X}(z, \mathbf{y}) \rangle &= z \text{trace}(\mathbf{M}^T \mathbf{X}) + \text{trace}(\mathbf{M}^T \text{sMat}_{\mathcal{E}}(\mathbf{y})\mathbf{X}) \\ &= z \text{trace}(\mathbf{M}^T \mathbf{X}) + \text{trace} \left(\text{sMat}_{\mathcal{E}}(\mathbf{y}) \frac{1}{2} (\mathbf{X}\mathbf{M}^T + \mathbf{M}\mathbf{X}) \right) \\ &= z \text{trace}(\mathbf{M}^T \mathbf{X}) + \frac{1}{2} \mathbf{y}^T \text{svec}_{\mathcal{E}}(\mathbf{X}\mathbf{M}^T + \mathbf{M}\mathbf{X}) \\ &= z \text{trace}(\mathbf{M}^T \mathbf{X}) + \langle \mathbf{y}, \frac{1}{2} \text{svec}_{\mathcal{E}}(\mathbf{X}\mathbf{M}^T + \mathbf{M}\mathbf{X}) \rangle.\end{aligned}$$

Thus we have:

$$\mathcal{X}^*(\mathbf{M}) = \begin{pmatrix} \text{trace}(\mathbf{M}^T \mathbf{X}) \\ \frac{1}{2} \text{svec}_{\mathcal{E}}(\mathbf{X}\mathbf{M}^T + \mathbf{M}\mathbf{X}) \end{pmatrix}. \quad (3.3)$$

Similarly, we have:

$$\begin{aligned}
\langle \mathbf{M}, \mathcal{Z}(\mathbf{w}, \mathbf{v}) \rangle &= \text{trace}(\mathbf{M}^T \mathbf{Z} (\text{Diag}(\mathbf{V}\mathbf{w}) + \text{sMat}_{\mathcal{E}^c}(\mathbf{v})) \\
&= \text{trace}(\mathbf{M}^T \mathbf{Z} (\text{Diag}(\mathbf{V}\mathbf{w})) + \text{trace} \left(\text{sMat}_{\mathcal{E}^c}(\mathbf{v}) \frac{1}{2} (\mathbf{M}^T \mathbf{Z} + \mathbf{Z}\mathbf{M}) \right) \\
&= \langle \text{diag}(\mathbf{M}^T \mathbf{Z}), \mathbf{V}\mathbf{w} \rangle + \frac{1}{2} \mathbf{v}^T \text{svec}_{\mathcal{E}^c}(\mathbf{M}^T \mathbf{Z} + \mathbf{Z}\mathbf{M}) \\
&= \langle \mathbf{w}, \mathbf{V}^T \text{diag}(\mathbf{M}^T \mathbf{Z}) \rangle + \langle \mathbf{v}, \frac{1}{2} \text{svec}_{\mathcal{E}^c}(\mathbf{M}^T \mathbf{Z} + \mathbf{Z}\mathbf{M}) \rangle,
\end{aligned}$$

Thus the adjoint \mathcal{Z}^* is

$$\mathcal{Z}^*(\mathbf{M}) = \begin{pmatrix} \mathbf{V}^T \text{diag}(\mathbf{M}^T \mathbf{Z}) \\ \frac{1}{2} \text{svec}_{\mathcal{E}^c}(\mathbf{M}^T \mathbf{Z} + \mathbf{Z}\mathbf{M}) \end{pmatrix}. \quad (3.4)$$

3.2 Presolve and Preconditioning

In this Lovász theta number problem, the bases are defined by the index mappings, $\text{ind}_{\mathcal{E}}$ and $\text{ind}_{\mathcal{E}^c}$. There are many ways to index the edge set \mathcal{E} and \mathcal{E}^c . However, it is advantageous to index the edge sets according to a *good graph partitioning*. In addition, we can also change the order of nodes since the row permutation of the matrix \mathbf{V} affects the bases. All of these issues will become clear when we discuss the diagonal and block preconditioners in the following sections.

3.2.1 Diagonal Preconditioning

In order to construct the diagonal preconditioner, we need to calculate $\|J(\mathbf{e}_i)\|_F$ for all unit vector \mathbf{e}_i in $\mathbb{R}^{t(n)}$, $i = 1, \dots, t(n)$. Since the variables are $(\mathbf{w}, \mathbf{v}, z, \mathbf{y})$, we consider four different cases, each of which corresponds to each variable.

- (i) $1 \leq i \leq n-1$: $J(\mathbf{e}_i) = \mathbf{Z} \text{Diag}(\mathbf{V}\mathbf{e}_i^w)$, where \mathbf{e}_i^w is the i -th unit vector in \mathbb{R}^{n-1} . We have: $\mathbf{V}\mathbf{e}_i^w = [-1; \mathbf{e}_i^w]$. Thus

$$\|J(\mathbf{e}_i)\|_F = \sqrt{\|\mathbf{Z}_{:,1}\|^2 + \|\mathbf{Z}_{:,i+1}\|^2}, \quad \forall i = 1, \dots, n-1.$$

- (ii) $i = n-1+j$, where $1 \leq j \leq m_c$: $J(\mathbf{e}_i) = \mathbf{Z} \text{sMat}_{\mathcal{E}^c}(\mathbf{e}_j^v)$, where \mathbf{e}_j^v is the j -th unit vector in \mathbb{R}^{m_c} . We have: $\text{sMat}_{\mathcal{E}^c}(\mathbf{e}_j^v) = \mathbf{E}_{\text{ind}_{\mathcal{E}^c}^{-1}(j)}$. Let $\text{ind}_{\mathcal{E}^c}^{-1}(j) = (k_{\mathcal{E}^c}(j), l_{\mathcal{E}^c}(j))$, we have:

$$\|J(\mathbf{e}_i)\|_F = \frac{1}{\sqrt{2}} \sqrt{\|\mathbf{Z}_{:,k_{\mathcal{E}^c}(j)}\|^2 + \|\mathbf{Z}_{:,l_{\mathcal{E}^c}(j)}\|^2}, \quad \forall i = n-1+j, j = 1, \dots, m_c.$$

- (iii) $i = n+m_c$: $J(\mathbf{e}_i) = \mathbf{X}$. Thus $\|J(\mathbf{e}_i)\|_F = \|\mathbf{X}\|_F$.

- (iv) $i = n+m_c+j$, where $1 \leq j \leq m$: $J(\mathbf{e}_i) = \text{sMat}_{\mathcal{E}}(\mathbf{e}_j^y)\mathbf{X}$, where \mathbf{e}_j^y is the j -th unit vector in \mathbb{R}^m . We have: $\text{sMat}_{\mathcal{E}}(\mathbf{e}_j^y) = \mathbf{E}_{\text{ind}_{\mathcal{E}}^{-1}(j)}$. Similar to the previous case, we obtain the following formulation:

$$\|J(\mathbf{e}_i)\|_F = \frac{1}{\sqrt{2}} \sqrt{\|\mathbf{X}_{k_{\mathcal{E}}(j),:}\|^2 + \|\mathbf{X}_{l_{\mathcal{E}}(j),:}\|^2}, \quad \forall i = n+m_c+j, j = 1, \dots, m,$$

where $(k_{\mathcal{E}}(j), l_{\mathcal{E}}(j)) = \text{ind}_{\mathcal{E}}^{-1}(j)$.

3.2.2 Block-Diagonal Preconditioning

For block preconditioner, we need to find the structure of the matrix representations of $\mathcal{Z}^* \mathcal{Z}$ and $\mathcal{X}^* \mathcal{X}$. We first calculate the columns of $\mathcal{Z}^* \mathcal{Z}$. There are two cases:

1. Column i of $\mathcal{Z}^* \mathcal{Z}$, $1 \leq i \leq n - 1$:

$\mathcal{Z}(e_i) = \mathbf{Z} \text{Diag}(\mathbf{V}e_i^w)$. Let $\mathbf{M} = \mathcal{Z}(e_i)$, we have: \mathbf{M} has only two non-zero columns, the first one, $\mathbf{M}_{:,1} = -\mathbf{Z}_{:,1}$, and the $(i + 1)$ -th one, $\mathbf{M}_{:,i+1} = \mathbf{Z}_{:,i+1}$.

Now consider $\mathcal{Z}^*(\mathbf{M})$, the first $n - 1$ elements is $\mathbf{V}^T \text{diag}(\mathbf{M}^T \mathbf{Z})$. Let $\mathbf{P} = \mathbf{M}^T \mathbf{Z}$, we have: \mathbf{P} again has only two non-zero rows, the first and $(i + 1)$ -th row. The elements of the first row of \mathbf{P} is

$$P_{1,k} = -\mathbf{Z}_{:,1}^T \mathbf{Z}_{:,k}, \quad \forall k = 1, \dots, n,$$

and the elements of the $(i + 1)$ -th row is

$$P_{i+1,k} = \mathbf{Z}_{:,i+1}^T \mathbf{Z}_{:,k}, \quad \forall k = 1, \dots, n.$$

Thus $\text{diag}(\mathbf{P})$ has two non-zero elements, the first, $-\mathbf{Z}_{:,1}^T \mathbf{Z}_{:,1}$, and the $(i + 1)$ -th element, $\mathbf{Z}_{:,i+1}^T \mathbf{Z}_{:,i+1}$. We then have $\mathbf{V}^T \text{diag}(\mathbf{P}) = \mathbf{Z}_{:,1}^T \mathbf{Z}_{:,1} \mathbf{e} + \mathbf{Z}_{:,i+1}^T \mathbf{Z}_{:,i+1} e_i^w$. The first $(n - 1) \times (n - 1)$ block of $\mathcal{Z}^* \mathcal{Z}$ can be written as $\text{Diag}(\|\mathbf{Z}_{:,2}\|_F^2, \dots, \|\mathbf{Z}_{:,n}\|_F^2) + \|\mathbf{Z}_{:,1}\|_F^2 \mathbf{E}$. The Cholesky decomposition of this block can be obtained using the faster rank-one update algorithm since the Cholesky decomposition of diagonal matrices are easy to compute and $\mathbf{E} = \mathbf{e}\mathbf{e}^T$. Another approach is to analytically find \mathbf{R}^{-1} with a special structure such that $\text{Diag}(\|\mathbf{Z}_{:,2}\|_F^2, \dots, \|\mathbf{Z}_{:,n}\|_F^2) + \|\mathbf{Z}_{:,1}\|_F^2 \mathbf{E} = \mathbf{R}^T \mathbf{R}$ and $\mathbf{R}^{-1} \mathbf{x}$ is easy to compute for any vector \mathbf{x} . We apply the following result:

Lemma 1. Consider matrix $\mathbf{A} = \mathbf{D} + \mathbf{u}\mathbf{u}^T$ where \mathbf{D} is a positive definite diagonal matrix, then \mathbf{A} can be decomposed as $\mathbf{R}^T \mathbf{R}$ with

$$\mathbf{R}^{-1} = \mathbf{D}^{-\frac{1}{2}} \left(\mathbf{I} + \frac{1}{\lambda} \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{p}\mathbf{p}^T \right),$$

where $\mathbf{p} = \mathbf{D}^{-\frac{1}{2}} \mathbf{u}$ and $\lambda = \|\mathbf{p}\|^2$.

Proof. We have: $\mathbf{A} = \mathbf{D}^{\frac{1}{2}} \mathbf{D}^{\frac{1}{2}} + \mathbf{D}^{\frac{1}{2}} \mathbf{D}^{-\frac{1}{2}} \mathbf{u}\mathbf{u}^T \mathbf{D}^{-\frac{1}{2}} \mathbf{D}^{\frac{1}{2}} = \mathbf{D}^{\frac{1}{2}} (\mathbf{I} + \mathbf{p}\mathbf{p}^T) \mathbf{D}^{\frac{1}{2}}$.

Assume $\mathbf{p}\mathbf{p}^T$ has the eigenvalue decomposition $\mathbf{V}(\lambda \mathbf{E}_{11}) \mathbf{V}^T$, we have: $\mathbf{V}\mathbf{V}^T = \mathbf{I}$, thus

$$\mathbf{A} = \mathbf{D}^{\frac{1}{2}} \mathbf{V} (\mathbf{I} + \lambda \mathbf{E}_{11}) \mathbf{V}^T \mathbf{D}^{\frac{1}{2}}.$$

$\mathbf{I} + \lambda \mathbf{E}_{11}$ is a diagonal matrix; therefore, we can compute \mathbf{A}^{-1} as follows:

$$\mathbf{A}^{-1} = \mathbf{D}^{-\frac{1}{2}} \mathbf{V} \left(\mathbf{I} + \left(\frac{1}{\lambda+1} - 1 \right) \mathbf{E}_{11} \right) \mathbf{V}^T \mathbf{D}^{-\frac{1}{2}}.$$

Applying the formulation for square roots of diagonal matrices, we have:

$$\mathbf{A}^{-1} = \mathbf{D}^{-\frac{1}{2}} \mathbf{V} \left(\mathbf{I} + \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{E}_{11} \right) \mathbf{V}^T \mathbf{V} \left(\mathbf{I} + \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{E}_{11} \right) \mathbf{V}^T \mathbf{D}^{-\frac{1}{2}}.$$

We also have $\lambda \mathbf{V} \mathbf{E}_{11} \mathbf{V}^T = \mathbf{p} \mathbf{p}^T$, thus

$$\mathbf{A}^{-1} = \mathbf{D}^{-\frac{1}{2}} \left(\mathbf{I} + \frac{1}{\lambda} \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{p} \mathbf{p}^T \right) \left(\mathbf{I} + \frac{1}{\lambda} \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{p} \mathbf{p}^T \right) \mathbf{D}^{-\frac{1}{2}} = \mathbf{S} \mathbf{S}^T,$$

where $\mathbf{S} = \mathbf{D}^{-\frac{1}{2}} \left(\mathbf{I} + \frac{1}{\lambda} \left(\frac{1}{\sqrt{\lambda+1}} - 1 \right) \mathbf{p} \mathbf{p}^T \right)$.

Let $\mathbf{R} = \mathbf{S}^{-1}$, clearly we have $\mathbf{A} = \mathbf{R}^T \mathbf{R}$ and $\mathbf{R}^{-1} = \mathbf{S}$. \square

Clearly, in order to calculate $\mathbf{R}^{-1} \mathbf{x}$, we only need to calculate $\mathbf{D}^{-\frac{1}{2}} \mathbf{x}$, $\mathbf{D}^{-\frac{1}{2}} \mathbf{u}$. For this particular matrix, we can set $\mathbf{D} = \text{Diag} \left(\|\mathbf{Z}_{:,2}\|_F^2, \dots, \|\mathbf{Z}_{:,n}\|_F^2 \right)$ and $\mathbf{u} = \|\mathbf{Z}_{:,1}\|_F \mathbf{e}$.

We now look at the remaining m_c elements of the i -th columns, $\mathbf{q} = \frac{1}{2} \text{svec}_{\mathcal{E}^c}(\mathbf{P} + \mathbf{P}^T)$. \mathbf{P} has two non-zero rows; therefore, $\mathbf{P} + \mathbf{P}^T$ has two non-zero rows and two non-zero columns. The elements of these rows and columns are the same as the elements of corresponding rows in \mathbf{P} except for four elements, $(1, 1)$, $(1, i+1)$, $(i+1, 1)$, and $(i+1, i+1)$. Note that the elements $(1, i+1)$ and $(i+1, 1)$ are zeros as $-\mathbf{Z}_{:,1}^T \mathbf{Z}_{:,i} - \mathbf{Z}_{:,i}^T \mathbf{Z}_{:,1} = 0$. Let \mathcal{E}_i^c be the set of all edges in \mathcal{E}^c starting from node i , we have: for all $j \in \text{ind}_{\mathcal{E}^c}(\mathcal{E}_i^c)$ such that $\text{ind}_{\mathcal{E}^c}^{-1}(j) = (1, l_{\mathcal{E}^c}(j)) \neq (1, i+1)$, $q_j = -\frac{1}{\sqrt{2}} \mathbf{Z}_{:,1}^T \mathbf{Z}_{:,l_{\mathcal{E}^c}(j)}$. Similarly, for all $j \in \text{ind}_{\mathcal{E}^c}(\mathcal{E}_{i+1}^c)$ such that $\text{ind}_{\mathcal{E}^c}^{-1}(j) = (i+1, l_{\mathcal{E}^c}(j))$, $q_j = \frac{1}{\sqrt{2}} \mathbf{Z}_{:,i+1}^T \mathbf{Z}_{:,l_{\mathcal{E}^c}(j)}$. Other elements of \mathbf{q} are zeros. Note that for all $n-1$ first columns, the number of non-zeros depends on $|\mathcal{E}_1^c|$. Thus in order to increase the sparsity of the resulting matrix, we should select the node 1 with the *largest* degree (in the original graph). This is one of the operations that we can consider for the presolve.

2. Column i of $\mathcal{Z}^* \mathcal{Z}$, $i = n-1+j$, where $1 \leq j \leq m_c$:

$\mathcal{Z}(\mathbf{e}_i) = \mathcal{Z} \text{sMat}_{\mathcal{E}^c}(\mathbf{e}_j^c)$. Let $\mathbf{M} = \mathcal{Z}(\mathbf{e}_i)$ and $\text{ind}_{\mathcal{E}^c}^{-1}(j) = (k, l)$, we again have that \mathbf{M} has only two non-zero columns, the k -th column, $\mathbf{M}_{:,k} = \frac{1}{\sqrt{2}} \mathbf{Z}_{:,l}$, and the l -th column, $\mathbf{M}_{:,l} = \frac{1}{\sqrt{2}} \mathbf{Z}_{:,k}$. Due to the symmetry of $\mathcal{Z}^* \mathcal{Z}$, we just need to look at the last m_c elements of these columns. Let $\mathbf{q} = \frac{1}{2} \text{svec}_{\mathcal{E}^c}(\mathbf{P} + \mathbf{P}^T)$, where $\mathbf{P} = \mathbf{M}^T \mathbf{Z}$. We have \mathbf{P} has two non-zero rows, the k -th and l -th ones. The elements are $P_{k,p} = \frac{1}{\sqrt{2}} \mathbf{Z}_{:,l}^T \mathbf{Z}_{:,p}$ and $P_{l,p} = \frac{1}{\sqrt{2}} \mathbf{Z}_{:,k}^T \mathbf{Z}_{:,p}$ for all $p = 1, \dots, n$. Similar to the previous case, we have that elements of $\mathbf{P} + \mathbf{P}^T$ are the corresponding elements of \mathbf{P} except for four elements (k, k) , (k, l) , (l, k) , and (l, l) . The element (k, l) (and (l, k)) has the value $\frac{1}{\sqrt{2}} \left(\|\mathbf{Z}_{:,k}\|_F^2 + \|\mathbf{Z}_{:,l}\|_F^2 \right)$. Thus we have: $q_j = \frac{1}{2} \left(\|\mathbf{Z}_{:,k}\|_F^2 + \|\mathbf{Z}_{:,l}\|_F^2 \right)$. Other non-zero elements of \mathbf{q} are the elements with indices in the set $\text{ind}_{\mathcal{E}^c}(\mathcal{E}_k^c \cup \mathcal{E}_l^c)$ and the computation of these elements is similar to the previous case. With this structure of non-zero elements for each column, we can see that the index mapping $\text{ind}_{\mathcal{E}^c}$ with a *good graph partitioning* results in a block structure for this $m_c \times m_c$ block of $\mathcal{Z}^* \mathcal{Z}$. More specifically, there are two cases in which the block structure can be form. In the first case, three edges (i, j) , (j, k) , and (i, k) (edges of a clique of size 3) are indexed consecutively and results in 3×3 block. However, there is no special structure of this block to be exploited. We now focus on the second case in which the set of edges from a node i , (i, j) with $j \in \mathcal{S}_i$, $\mathcal{S}_i \subset \mathcal{E}_i^c$, are indexed consecutively.

Let $\mathbf{Z}_{\mathcal{S}_i}$ be the $n \times |\mathcal{S}_i|$ submatrix of \mathbf{Z} that consists of $|\mathcal{S}_i|$ columns $\mathbf{Z}_{:,j}$, $j \in \mathcal{S}_i$, the block we obtain is $\frac{1}{2} \left(\|\mathbf{Z}_{:,i}\|_F^2 \mathbf{I} + \mathbf{Z}_{\mathcal{S}_i}^T \mathbf{Z}_{\mathcal{S}_i} \right)$. And this is another operation we need to consider for the presolve.

We continue with columns of $\mathcal{X}^* \mathcal{X}$. We also have two cases:

1. $i = 1$: $\mathcal{X}(\mathbf{e}_i) = \mathbf{X}$. Thus the first column is $\mathcal{X}^*(\mathbf{X})$ with the first element is $\text{trace}(\mathbf{X}^2)$ and the remaining m elements are $\text{svec}_{\mathcal{E}}(\mathbf{X}^2)$.
2. $i = j + 1$, where $1 \leq j \leq m$: $\mathcal{X}(\mathbf{e}_i) = \text{sMat}_{\mathcal{E}}(\mathbf{e}_j^y) \mathbf{X}$. Let $\mathbf{M} = \mathcal{X}(\mathbf{e}_i)$ and $\text{ind}_{\mathcal{E}}(j) = (k, l)$, we have: \mathbf{M} has only two non-zero rows, the k -th row, $\mathbf{M}_{k,:} = \frac{1}{\sqrt{2}} \mathbf{X}_{l,:}$; and the l -th row, $\mathbf{M}_{l,:} = \frac{1}{\sqrt{2}} \mathbf{X}_{k,:}$. Due to the symmetry of $\mathcal{X}^* \mathcal{X}$, we do not need to reconsider the first elements of these columns. Let $\mathbf{P} = \mathbf{M} \mathbf{X}$, we have, the remaining m elements are $\mathbf{q} = \frac{1}{2} \text{svec}_{\mathcal{E}}(\mathbf{P} + \mathbf{P}^T)$. Similarly, \mathbf{P} has only two non-zeros rows, the k -th and the l -th ones with the elements $P_{k,p} = \frac{1}{\sqrt{2}} \mathbf{X}_{l,:} \mathbf{X}_{:,p}$ and $P_{l,p} = \frac{1}{\sqrt{2}} \mathbf{X}_{k,:} \mathbf{X}_{:,p}$ for all $p = 1, \dots, n$. The elements of $\mathbf{P} + \mathbf{P}^T$ can be derived from elements of \mathbf{P} in a similar way shown before. The element (k, l) is $\frac{1}{\sqrt{2}} \left(\|\mathbf{X}_{k,:}\|_F^2 + \|\mathbf{X}_{l,:}\|_F^2 \right)$. Thus we have: $q_j = \frac{1}{2} \left(\|\mathbf{X}_{k,:}\|_F^2 + \|\mathbf{X}_{l,:}\|_F^2 \right)$. Other non-zero elements of \mathbf{q} are the elements with indices in the set $\text{ind}_{\mathcal{E}}(\mathcal{E}_k \cup \mathcal{E}_l)$, where \mathcal{E}_i is the set of all edges in \mathcal{E} starting from node i . With this structure of non-zero elements for each column, we can again see that the index mapping $\text{ind}_{\mathcal{E}}$ with a *good graph partitioning* results in a block structure for this $m \times m$ block of $\mathcal{X}^* \mathcal{X}$, especially when edges in \mathcal{E} from a single node are indexed consecutively. Similarly, the presolve can help us obtain this *good graph partitioning*.

4 Numerics and Conclusion

4.1 Numerical Results

We test three versions of the Gauss-Newton algorithm and three different preconditioners. The algorithm versions are: (i) a general version for full matrices; (ii) a sparse version for sparse matrices; and (iii) a specialized version for the Lovász theta function problem. The preconditioners are: (i) diagonal; (ii) two block diagonal; and (iii) multiple block diagonal. We start with the sparse version of the code. Data inputs are $\mathbf{C} \in \mathcal{S}^n$, $\mathbf{A}_i \in \mathcal{S}^n$, $i = 1, \dots, m$, and $\mathbf{b} \in \mathbb{R}^m$. The algorithm is coded in MATLAB; and, the Gauss-Newton directions are computed using LSMR developed by Fong and Saunders [5]. First, we compare different versions of the code with different preconditioners before comparing it with different solvers. We generate data inputs randomly with n from 10 to 100, $m = \lceil n(n+1)/4 \rceil$; and, the sparseness density is set to be $1/(n+m)$. We run the code without preconditioner (**SRSD**_o), with diagonal preconditioner (**SRSD**_{Diag}), and block preconditioner (**SRSD**_{BDiag}). The tolerance is set to be $\epsilon = 10^{-12}$ and crossover is used. As expected, there is a reduction in the number of LSMR iterations when preconditioners are used. This is due to the fact that the condition number of the overdetermined system that we need to solve is decreased. Table 4.1 shows the reduction ratios in number of LSMR iterations of (**SRSD**_{Diag}) and (**SRSD**_{BDiag}) as compared to that of (**SRSD**_o). However, there is a trade-off between the

time to compute the Gauss-Newton directions and the time to compute the preconditioners. Table 4.2 shows the ratios of total computational time of $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ as compared to that of $(\mathbf{SRSD}_{\text{o}})$; and, we can see that $(\mathbf{SRSD}_{\text{Diag}})$ is the best version of the code in terms of computational time for these sparse instances. Note that for these instances, the accuracy of the optimal solutions is approximately the same. However, for harder instances that we consider in other tests, $(\mathbf{SRSD}_{\text{BDiag}})$ will be the choice over $(\mathbf{SRSD}_{\text{Diag}})$ to get better accuracy since the block preconditioner produces better conditioned systems for LSMR to solve, so that we get more accurate search directions in addition to a reduction in the number of iterations of LSMR. We also test our code with LSQR, which is another code to solve overdetermined systems developed by Paige and Saunders [21]. For these instances, the differences in number of iterations and accuracy are not noticeable with only a slightly higher number of iterations for LSQR resulting in slightly higher accuracy. Similar to the choice between $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$, if we need to solve very hard instances, LSQR appears to be preferable.

n	10	20	30	40	50	60	70	80	90	100
$(\mathbf{SRSD}_{\text{Diag}})$	0.67	0.55	0.55	0.39	0.42	0.42	0.47	0.36	0.34	0.32
$(\mathbf{SRSD}_{\text{BDiag}})$	0.30	0.16	0.13	0.09	0.09	0.08	0.08	0.07	0.07	0.06

Table 4.1: Ratios of numbers of LSMR iterations of $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ to $(\mathbf{SRSD}_{\text{o}})$

n	10	20	30	40	50	60	70	80	90	100
$(\mathbf{SRSD}_{\text{Diag}})$	0.63	0.58	0.59	0.45	0.49	0.51	0.62	0.45	0.48	0.44
$(\mathbf{SRSD}_{\text{BDiag}})$	0.33	0.27	0.46	0.57	1.05	1.44	2.20	2.26	3.50	3.26

Table 4.2: Ratios of total computational time of $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ to $(\mathbf{SRSD}_{\text{o}})$

We now focus on solution accuracy for larger instances and then compare the performance of our code with other solvers. For this test, we set $n = 100$ and $m = 100$ and again generate data inputs randomly. We run $N = 1000$ instances and record the (average) of: the number of iterations, the relative norm of $\mathbf{Z}\mathbf{X}$, and the relative minimum eigenvalues of \mathbf{X} and \mathbf{Z} , see Table 4.3. We apply the crossover technique and use the diagonal preconditioner as the setting for our code. The tolerance is again set to 10^{-12} . The average shown for the two accuracy measures is *log-average*. The two measures are

$$\text{RelZXnorm} = \frac{\|\mathbf{Z}\mathbf{X}\|_F}{|\mathbf{C} \cdot \mathbf{X}| + 1}, \quad \text{Relmineig} = \frac{\min\{\lambda_{\min}(\mathbf{X}), \lambda_{\min}(\mathbf{Z})\}}{|\mathbf{C} \cdot \mathbf{X}| + 1}.$$

	Iteration	RelZXnorm	Relmineig
Average	18.66	2.62×10^{-15}	-8.79×10^{-16}
Best	14.00	2.78×10^{-16}	-3.86×10^{-17}
Worst	26.00	8.36×10^{-13}	-1.97×10^{-13}

Table 4.3: Accuracy measures for random sparse instances

We now select 20 hard instances in terms of number of iterations taken using our ($\mathbf{SRSD}_{\text{Diag}}$) code. Using these 20 instances, we compare our solver and the four SDP solvers. (See e.g., the URL www-user.tu-chemnitz.de/~helmberg/semidef.html.)

SeDuMi 1.3, CSDP 6.1.0, SDPA 7.3.1, SDPT3 4.0.

The tolerances (for the relative trace of \mathbf{ZX}) for these four solvers are set to be as small as possible without running into numerical problems. For SDPT3, it could be 10^{-13} or 10^{-14} . For the remaining solvers, 10^{-12} is sufficiently small and there are some cases, where we had to reduce the tolerance for SeDuMi or CSDP. The two accuracy measures are compared and in addition, we also record all six DIMACS error measures [16]. (The last measure is only comparable if both \mathbf{X} and \mathbf{Z} are positive semidefinite.) We also compare the computational times of all solvers for these 20 instances. Similarly, we also select 20 easy instances in terms of number of iterations taken using our ($\mathbf{SRSD}_{\text{Diag}}$) code and compare with other solvers. The numerical results are shown in Table 4.4 and 4.5, respectively.

	($\mathbf{SRSD}_{\text{Diag}}$)	SeDuMi	CSDP	SDPA	SDPT3
Iteration	23.45	19.15	16.50	16.30	24.50
RelZXnorm	7.18×10^{-15}	1.50×10^{-07}	8.92×10^{-08}	5.86×10^{-08}	1.01×10^{-10}
Relmineig	-1.12×10^{-15}	-8.72×10^{-15}	3.28×10^{-13}	1.46×10^{-13}	5.99×10^{-17}
DIMACS1	1.28×10^{-15}	7.93×10^{-11}	1.24×10^{-13}	2.48×10^{-13}	3.19×10^{-13}
DIMACS2	2.28×10^{-14}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS3	9.39×10^{-16}	3.53×10^{-16}	1.64×10^{-08}	1.44×10^{-15}	8.72×10^{-14}
DIMACS4	3.22×10^{-15}	3.42×10^{-13}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS5	1.19×10^{-15}	9.72×10^{-13}	1.54×10^{-09}	4.96×10^{-10}	1.10×10^{-13}
DIMACS6	4.65×10^{-16}	2.84×10^{-14}	1.16×10^{-09}	4.96×10^{-10}	2.01×10^{-13}
Time	61.89	1.23	0.60	0.74	1.13

Table 4.4: Performance measures for *hard* random sparse instances

	($\mathbf{SRSD}_{\text{Diag}}$)	SeDuMi	CSDP	SDPA	SDPT3
Iteration	15.20	17.50	15.10	15.55	22.40
RelZXnorm	1.64×10^{-14}	1.66×10^{-07}	1.04×10^{-06}	2.17×10^{-08}	1.77×10^{-10}
Relmineig	-5.56×10^{-15}	-1.15×10^{-14}	2.27×10^{-13}	2.20×10^{-13}	1.21×10^{-16}
DIMACS1	1.24×10^{-15}	1.34×10^{-10}	1.38×10^{-13}	5.20×10^{-14}	2.74×10^{-13}
DIMACS2	1.15×10^{-13}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS3	8.35×10^{-16}	3.57×10^{-16}	1.04×10^{-08}	1.44×10^{-15}	8.95×10^{-14}
DIMACS4	1.32×10^{-14}	4.34×10^{-13}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS5	2.19×10^{-15}	1.98×10^{-12}	1.18×10^{-09}	4.85×10^{-10}	9.34×10^{-14}
DIMACS6	1.90×10^{-15}	2.95×10^{-14}	5.47×10^{-10}	4.85×10^{-10}	1.86×10^{-13}
Time	37.00	1.16	0.53	0.59	1.05

Table 4.5: Performance measures for *easy* random sparse instances

The computational time for $(\mathbf{SRSD}_{\text{Diag}})$ is not competitive with the other solvers. One reason is that the code right now is all MATLAB code with no pre-compiled subroutines. However, the main reason is that the Gauss-Newton approach needs to solve a system of equations with $n(n+1)/2$ variables regardless of the number of constraints. We now keep $n = 100$ fixed and vary the number of constraints from $m = 500$ to $m = 5000$. The instances are generated randomly with the same sparseness density $1/(4n)$. For larger instances, we need more accuracy for the LSMR subroutine that iteratively solves for the Gauss-Newton direction. Therefore, for these instances, we use the block triangular preconditioner instead of the diagonal one. The average number of iterations and two accuracy measures for these instances are shown in Table 4.6. The corresponding computational time ratios of $(\mathbf{SRSD}_{\text{BDiag}})$ to those of other solvers are then reported in Table 4.7.

	$(\mathbf{SRSD}_{\text{BDiag}})$	SeDuMi	CSDP	SDPA	SDPT3
Iteration	14.82	15.73	15.36	16.00	21.09
RelZXnorm	6.93×10^{-15}	1.99×10^{-07}	1.91×10^{-07}	1.77×10^{-08}	2.27×10^{-09}
Relmineig	-3.69×10^{-16}	-4.94×10^{-15}	4.06×10^{-13}	9.06×10^{-14}	2.31×10^{-16}

Table 4.6: Average measures for random instances with different numbers of constraints

m	500	1000	1500	2000	2500	3000	3500	4000	4500	5000
SeDuMi	168.74	49.99	21.47	4.18	1.92	1.89	1.09	0.66	0.68	0.32
CSDP	144.02	39.62	21.64	7.58	5.89	5.56	3.53	2.40	2.88	1.62
SDPA	488.28	203.06	134.49	50.99	37.09	33.37	20.98	15.55	17.83	9.41
SDPT3	87.49	27.96	16.10	6.00	3.80	4.42	2.11	2.30	2.89	1.81

Table 4.7: Time ratios relative to $(\mathbf{SRSD}_{\text{BDiag}})$ for different numbers of constraints

In the previous test, we kept the sparseness density constant while changing the number of constraints. We now test our code with instances created with different densities. We fix $n = 100$ and $m = 2500$ while varying the density with $1/(4sn)$, $s = 1, \dots, 10$. All parameter settings are the same as in the previous test. In addition, we also run our code with the diagonal preconditioner when the data is sparse enough. It turns out in this test, for $s \geq 3$, $(\mathbf{SRSD}_{\text{Diag}})$ maintains the same level of accuracy but is more efficient than $(\mathbf{SRSD}_{\text{BDiag}})$ with respect to time. This can be explained as follows. $(\mathbf{SRSD}_{\text{BDiag}})$ in general requires fewer iterations of the LSMR subroutine than $(\mathbf{SRSD}_{\text{Diag}})$ and, if the data is not sparse enough, this saving in time can compensate for the more expensive construction of the preconditioner. On the other hand, if the data is sparser, e.g., $s \geq 3$, the reduction in the number of LSMR iterations of $(\mathbf{SRSD}_{\text{Diag}})$ gradually becomes more significant than that of $(\mathbf{SRSD}_{\text{BDiag}})$, which makes $(\mathbf{SRSD}_{\text{Diag}})$ more efficient. The time comparison between $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ is shown in Table 4.8. Note that since with $s = 2$, $(\mathbf{SRSD}_{\text{BDiag}})$ is already more efficient than $(\mathbf{SRSD}_{\text{Diag}})$, we do not run $(\mathbf{SRSD}_{\text{Diag}})$ for the instance with $s = 1$. To compare with other solvers, we select the more efficient code between $(\mathbf{SRSD}_{\text{Diag}})$ and $(\mathbf{SRSD}_{\text{BDiag}})$ for each instance. Basically, the first two instances are run with $(\mathbf{SRSD}_{\text{BDiag}})$ while the remaining ones are with $(\mathbf{SRSD}_{\text{Diag}})$. The number of iterations and two accuracy measures are reported in Table 4.9. Computational time ratios of $(\mathbf{SRSD}_{\text{Diag}})/2$ to other solvers are shown in Table 4.10.

s	2	3	4	5	6	7	8	9	10
($\mathbf{SRSD}_{\text{Diag}}$)	1080.72	201.44	182.98	230.35	102.31	76.47	80.29	57.92	52.29
($\mathbf{SRSD}_{\text{BDiag}}$)	447.55	362.80	314.22	485.58	485.93	312.37	410.73	359.74	352.83
Time ratio	0.41	1.80	1.72	2.11	4.75	4.08	5.12	6.21	6.75

Table 4.8: Computational times for random instances with different sparsity

	($\mathbf{SRSD}_{\text{Diag}}$)/2	SeDuMi	CSDP	SDPA	SDPT3
Iteration	13.80	16.80	16.20	16.20	24.30
RelZXnorm	4.47×10^{-15}	5.59×10^{-08}	1.30×10^{-07}	3.80×10^{-09}	3.65×10^{-11}
Relmineig	-7.59×10^{-16}	-2.24×10^{-16}	4.96×10^{-13}	8.71×10^{-14}	8.51×10^{-18}

Table 4.9: Average measures for random instances with different sparsity

We have seen that for well-conditioned instances, our code can achieve solutions with very high accuracy. We now move on to test the code with ill-conditioned instances, namely, the instances without strict complementary slackness or instances with which Slater’s condition almost fails. To generate instances without strict complementary slackness, we use the code developed by Wei and Wolkowicz [22]. The instances are generated with $n = 50$ and $m = 1000$. The general version of the code will be used since all matrices are dense. For these hard instances, we apply the block preconditioner without crossover. The tolerance is set to be 10^{-14} . The results are shown in Table 4.11. We note that the GN-method has significantly smaller values for $\|\mathbf{Z}\mathbf{X}\|$, RelZXnorm. This increased accuracy is most probably a result of the fact that the GN-method is derived by minimizing this measure.

The accuracy is indeed less for these hard instances, even for our code. And the stopping criterion for all tested instances but one is when the number of iterations of the LSMR subroutine reaches its maximum limit (set at $5n(n + 1)$). SeDuMi and SDPA, and SDPT3 have various kinds of numerical problems. CSDP controls the tolerance limit (around 10^{-08}); therefore, it does not incur any numerical problem even though the tolerance is set to be as small as 10^{-14} .

For instances with which Slater’s condition almost fails, we generate the instances randomly using the alternative theorem for Slater’s condition with respect to the dual problem. The settings are the same as for the previous test, with tolerance 10^{-14} , block preconditioner, and no crossover. The results are in Table 4.12.

It turns out that for other solvers, instances with which Slater’s condition almost fails are difficult to solve but not for our code. Even though no crossover is used, half of the instances have negative relative minimum eigenvalue of \mathbf{X} and \mathbf{Z} with average DIMACS2 and DIMACS4 measures being 1.05×10^{-17} and 2.53×10^{-16} , respectively. We also solve these instances with the crossover technique; the results are better in terms of computational time (number of iterations) with the trade-off being negative relative minimum eigenvalues for both \mathbf{X} and \mathbf{Z} . The average performance measures of ($\mathbf{GRSD}_{\text{BDiag}}$) with crossover is (13.70, 3.40×10^{-16} , -1.41×10^{-13} , 2.78×10^{-13} , 5.17×10^{-15} , 1.61×10^{-12} , 9.97×10^{-16} , 1.72×10^{-14} , 4.71×10^{-17} , 22.50). This is probably due to the fact that the problem is still well-conditioned even though Slater’s condition almost fails for the dual problem. (The Jacobian is non-singular at the optimal solution.) In addition, for these instances, the optimal solutions are not extremely large in magnitude, which helps the code maintain high

s	1	2	3	4	5	6	7	8	9	10
SeDuMi	2.84	2.74	1.26	1.15	1.49	0.62	0.53	0.49	0.34	0.30
CSDP	7.33	12.17	6.07	6.77	8.06	4.14	3.21	3.39	2.49	2.17
SDPA	43.35	48.63	22.13	20.95	25.47	12.33	9.25	9.82	7.08	6.04
SDPT3	4.60	8.82	5.79	5.74	7.50	3.94	2.86	2.75	2.27	2.09

Table 4.10: Time ratios relative to ($\mathbf{SRSD}_{\text{Diag}}$) for different sparsity

	($\mathbf{GRSD}_{\text{BDiag}}$)	SeDuMi	CSDP	SDPA	SDPT3
Iteration	29.20	13.60	12.10	12.00	22.70
RelZXnorm	1.32×10^{-12}	4.12×10^{-06}	6.53×10^{-05}	1.91×10^{-07}	2.97×10^{-08}
Relmineig	5.96×10^{-20}	6.20×10^{-14}	1.44×10^{-15}	3.25×10^{-12}	1.21×10^{-18}
DIMACS1	2.16×10^{-12}	2.01×10^{-07}	3.50×10^{-10}	3.79×10^{-07}	4.62×10^{-10}
DIMACS2	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS3	1.81×10^{-11}	6.34×10^{-15}	2.10×10^{-08}	2.09×10^{-14}	1.54×10^{-11}
DIMACS4	0.00×10^{00}	1.48×10^{-14}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS5	2.07×10^{-12}	9.23×10^{-09}	1.23×10^{-08}	1.58×10^{-07}	5.70×10^{-11}
DIMACS6	1.86×10^{-12}	4.99×10^{-09}	4.25×10^{-08}	3.50×10^{-07}	1.06×10^{-10}
Time	281.45	64.02	69.50	27.31	60.25

Table 4.11: Performance measures; random instances; strict complementarity fails

	($\mathbf{GRSD}_{\text{BDiag}}$)	SeDuMi	CSDP	SDPA	SDPT3
Iteration	26.00	17.10	13.80	15.20	20.40
RelZXnorm	9.90×10^{-16}	6.14×10^{-07}	4.61×10^{-07}	2.05×10^{-07}	2.13×10^{-08}
Relmineig	1.05×10^{-17}	1.06×10^{-15}	1.60×10^{-12}	1.06×10^{-10}	1.47×10^{-15}
DIMACS1	2.78×10^{-13}	1.07×10^{-10}	1.10×10^{-12}	1.17×10^{-13}	8.41×10^{-11}
DIMACS2	0.00×10^{00} (*)	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS3	1.42×10^{-12}	3.56×10^{-14}	5.21×10^{-09}	6.83×10^{-08}	2.59×10^{-13}
DIMACS4	0.00×10^{00} (*)	1.48×10^{-14}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}
DIMACS5	1.09×10^{-14}	2.77×10^{-12}	1.38×10^{-09}	2.02×10^{-14}	1.26×10^{-11}
DIMACS6	3.44×10^{-15}	3.16×10^{-12}	8.54×10^{-10}	4.43×10^{-08}	1.82×10^{-12}
Time	58.04	75.02	76.42	32.32	48.98

Table 4.12: Performance measures; random instances; dual Slater's CQ almost fails

accuracy.

The next test is for the Lovász theta function problem. We generate a random graph with $n = 100$ nodes and the number of edges is approximately $n(n-1)/4$. According to [7, 14], these are the most difficult instances to solve. We first test the specialized version of the code and the special block preconditioner that exploits the block structure of the problem. The tolerance is set to 10^{-12} and we run the code with crossover for 100 random instances. The results for (**TRSD**_{MBDiag}) (with multiple diagonal blocks) are shown in Table 4.13.

	Iteration	RelZXnorm	Relmineig
Average	18.31	4.85×10^{-15}	2.47×10^{-13}
Best	16.00	6.42×10^{-17}	-1.28×10^{-15}
Worst	23.00	9.74×10^{-13}	-3.58×10^{-11}

Table 4.13: Accuracy measures for random Lovász theta function instances

We now select 10 hard instances in terms of number of iterations and compare the results using other solvers including (**SRSD**_{Diag}) and (**TRSD**_{Diag}), both with crossover since all these instances are well-conditioned. The results are in Tables 4.14 and 4.15. Similarly, 10 easy instances are selected and the numerical results for these instances are shown in Table 4.16 and 4.17.

	SeDuMi	CSDP	SDPA	SDPT3	(TRSD _{MBDiag})
Iteration	21.00	17.10	16.00	22.60	21.00
RelZXnorm	5.00×10^{-08}	1.56×10^{-08}	2.38×10^{-07}	1.73×10^{-11}	2.67×10^{-14}
Relmineig	3.34×10^{-14}	7.09×10^{-14}	3.49×10^{-12}	2.32×10^{-18}	6.74×10^{-13}
DIMACS1	1.84×10^{-11}	6.07×10^{-14}	5.09×10^{-14}	1.27×10^{-13}	1.25×10^{-16}
DIMACS2	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	3.11×10^{-15}
DIMACS3	9.09×10^{-16}	1.09×10^{-07}	1.52×10^{-14}	8.87×10^{-13}	8.75×10^{-15}
DIMACS4	2.39×10^{-13}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	4.82×10^{-12}
DIMACS5	9.09×10^{-12}	4.44×10^{-10}	9.61×10^{-09}	1.47×10^{-13}	5.92×10^{-16}
DIMACS6	9.93×10^{-14}	2.10×10^{-10}	9.61×10^{-09}	7.08×10^{-14}	4.19×10^{-16}
Time	90.46	13.05	6.02	12.91	158.09

Table 4.14: Performance measures for *hard* random Lovász theta function instances

The final test we consider is to solve some instances in SDPLIB (see Borchers [2]). Since these instances are considered hard instances, we use block diagonal preconditioner without crossover. In addition, LSQR will be used instead of LSMR and one of the stopping criteria is when the number of LSQR iterations reaches the maximum limit, (set to be $5n(n+1)$). We run our code for all instances with $n \leq 100$, which include the set of *control*, *hinf*, *gpp*, *mcp*, *qap*, *theta*, and *truss* instances. We report the main performance measure relZXnorm, and in addition, the number of iterations, and the total computational time. Table (4.18) shows the results for the set of *control* instances.

The hardness of the *control* instances is due to the failure of strict complementary slackness. For well-conditioned problems, the (infeasible) starting points are set to be simply $\mathbf{X}_0 = \mathbf{Z}_0 = \mathbf{I}$. This controls the initial value of μ . Since we start with an infeasible solution, we would like to

	(SRSD _{Diag})	(TRSD _{Diag})	(TRSD _{MBDiag})
Iteration	23.90	21.00	21.00
RelZXnorm	5.62×10^{-14}	3.26×10^{-14}	2.67×10^{-14}
Relmineig	2.23×10^{-12}	9.51×10^{-13}	6.74×10^{-13}
DIMACS1	9.89×10^{-17}	1.11×10^{-16}	1.25×10^{-16}
DIMACS2	5.76×10^{-15}	4.46×10^{-15}	3.11×10^{-15}
DIMACS3	1.89×10^{-14}	8.71×10^{-15}	8.75×10^{-15}
DIMACS4	1.59×10^{-11}	6.80×10^{-12}	4.82×10^{-12}
DIMACS5	1.73×10^{-15}	1.87×10^{-15}	5.92×10^{-16}
DIMACS6	1.75×10^{-15}	1.85×10^{-15}	4.19×10^{-16}
Time	214.99	126.74	158.09

Table 4.15: Performance measures for *hard* random Lovász theta function instances

	SeDuMi	CSDP	SDPA	SDPT3	(TRSD _{MBDiag})
Iteration	20.20	16.10	16.00	21.40	16.30
RelZXnorm	7.34×10^{-08}	8.33×10^{-08}	9.61×10^{-08}	6.76×10^{-11}	2.18×10^{-14}
Relmineig	3.54×10^{-14}	3.69×10^{-14}	1.97×10^{-12}	1.24×10^{-17}	5.36×10^{-13}
DIMACS1	1.78×10^{-11}	1.86×10^{-14}	2.28×10^{-15}	1.85×10^{-13}	1.23×10^{-16}
DIMACS2	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	4.02×10^{-15}
DIMACS3	8.21×10^{-16}	3.38×10^{-08}	1.55×10^{-14}	1.61×10^{-12}	8.70×10^{-15}
DIMACS4	2.51×10^{-13}	0.00×10^{00}	0.00×10^{00}	0.00×10^{00}	3.18×10^{-12}
DIMACS5	1.57×10^{-11}	1.35×10^{-10}	5.28×10^{-09}	2.92×10^{-13}	1.31×10^{-15}
DIMACS6	1.13×10^{-13}	8.91×10^{-11}	5.28×10^{-09}	1.88×10^{-13}	5.44×10^{-16}
Time	88.95	13.09	6.01	12.37	63.43

Table 4.16: Performance measures for *easy* random Lovász theta function instances

	(SRSD _{Diag})	(TRSD _{Diag})	(TRSD _{MBDiag})
Iteration	20.00	16.30	16.30
RelZXnorm	9.22×10^{-17}	3.03×10^{-14}	2.18×10^{-14}
Relmineig	3.41×10^{-15}	1.72×10^{-12}	5.36×10^{-13}
DIMACS1	1.35×10^{-16}	1.23×10^{-16}	1.23×10^{-16}
DIMACS2	2.48×10^{-17}	1.71×10^{-14}	4.02×10^{-15}
DIMACS3	1.74×10^{-14}	8.56×10^{-15}	8.70×10^{-15}
DIMACS4	2.41×10^{-14}	1.26×10^{-11}	3.18×10^{-12}
DIMACS5	1.91×10^{-16}	2.13×10^{-15}	1.31×10^{-15}
DIMACS6	2.67×10^{-17}	1.34×10^{-15}	5.44×10^{-16}
Time	96.39	47.34	63.43

Table 4.17: Performance measures for *easy* random Lovász theta function instances

	(SRSD _{B_{Diag}})	SeDuMi	CSDP	SDPA	SDPT3
<i>control1</i>	7.19×10^{-11} (36, 3.23)	1.24×10^{-04} (29, 0.57)	1.15×10^{-06} (20, 0.15)	1.54×10^{-04} (28, 0.14)	3.72×10^{-06} (22, 0.84)
<i>control2</i>	6.50×10^{-07} (37, 37.12)	9.01×10^{-04} (30, 0.48)	2.86×10^{-07} (23, 0.19)	1.13×10^{-04} (30, 0.43)	3.07×10^{-05} (21, 0.72)
<i>control3</i>	1.80×10^{-05} (37, 374.44)	4.10×10^{-03} (32, 1.43)	2.26×10^{-06} (24, 1.07)	2.53×10^{-04} (35, 2.50)	3.96×10^{-04} (21, 1.32)
<i>control4</i>	4.84×10^{-05} (35, 1370.18)	4.82×10^{-03} (33, 1.83)	2.03×10^{-06} (24, 2.83)	1.45×10^{-04} (37, 8.90)	7.67×10^{-05} (21, 2.68)
<i>control5</i>	2.73×10^{-05} (38, 5250.32)	2.08×10^{-02} (34, 4.48)	2.03×10^{-06} (25, 5.36)	5.61×10^{-04} (38, 19.58)	1.36×10^{-04} (22, 5.25)
<i>control6</i>	9.67×10^{-06} (38, 18923.17)	3.70×10^{-02} (38, 9.77)	1.75×10^{-04} (61, 30.64)	1.31×10^{-03} (41, 51.02)	3.58×10^{-04} (23, 10.10)

Table 4.18: Performance measures for *control* instances

reduce the residuals to zero quickly and this can be done if the step size is close to 1 in the first few iterations. The step size is determined by the positive definiteness of the current solution. If we start with $\mathbf{X}_0 = \mathbf{Z}_0 = \mathbf{I}$ for these *control* instances, the code does not converge since the step size becomes smaller and smaller and the iterations are very close to the boundary of the semidefinite cone. In order to remedy this situation, we start with larger starting solutions, $\mathbf{X}_0 = \mathbf{Z}_0 = \epsilon^{-\alpha_0} \mathbf{I}$, where $\alpha_0 \geq 0$ and ϵ is the main tolerance. The above results are with these new starting points with $\alpha_0 = 0.25$. We can see that (**SRSD**_{B_{Diag}}) performs well as compared to the other solvers, but it does not scale well in terms of problem size since the values of m are small compared to $n(n+1)/2$.

The next instance is *gpp100* and the results are show in Table 4.19. We show additional performance measures, DIMACS1 (primal feasibility), DIMACS3 (dual feasibility), and DIMACS5 (duality gap), and also the Frobenius norm of the dual optimal \mathbf{Z} . For this instance, the Gauss-Newton direction $\Delta \mathbf{y}$ is extremely large and so is the dual solution \mathbf{Z} . The dual feasibility performance measure therefore is reduced significantly. These issues can be explained by the fact that Slater's condition almost fails for the primal problem. According to the alternative theorem, Slater's condition for the primal problem fails when there exists a dual solution \mathbf{y} such that $\mathcal{A}^* \mathbf{y} \succeq 0$, $\mathbf{y} \neq \mathbf{0}$, and $\mathbf{b}' \mathbf{y} = 0$. If this is the case, if \mathbf{y}^* is an optimal dual solution, then $\mathbf{y}^* + \alpha \mathbf{y}$ is also an dual optimal solution for all $\alpha \geq 0$. Therefore, if Slater's condition almost fails for the primal problem, the set of (dual) optimal solutions can be very large, which can explain why \mathbf{Z} has such an extremely large magnitude. In the next version of the code, we will consider some preprocessing routines to detect and resolve the issue of Slater's condition for the primal problem.

The two following tables, Table 4.20 and 4.21, are for *hinf* instances. These instances appear to have both issues that we mentioned above: stagnation due to small step sizes, and Slater's condition almost fails for the primal problem. Since these are small instances, we can check Slater's condition for the primal problem by minimizing $|\mathbf{b}' \mathbf{y}|$ subject to $\mathcal{A}^* \mathbf{y} \succeq 0$ and the additional constraint $\mathbf{e}' \mathbf{y} = 1$ to make sure that $\mathbf{y} \neq \mathbf{0}$. The optimal values of these optimization problems are in the first columns. We can see they are indeed very small except for *hinf9* and *hinf2* instances. For these instances, we use $\alpha_0 = 0.1$ for most of them.

	($\text{SRSD}_{\text{Bdiag}}$)	SeDuMi	CSDP	SDPA	SDPT3
<i>gpp100</i>	6.97×10^{-02} (41, 2414.59)	1.94×10^{-03} (31, 2.20)	1.27×10^{-03} (19, 0.56)	9.20×10^{-04} (22, 0.76)	1.58×10^{-06} (18, 0.70)
DIMACS1	2.46×10^{-14}	1.86×10^{-06}	1.27×10^{-08}	4.80×10^{-10}	1.50×10^{-09}
DIMACS3	5.30×10^{-08}	2.99×10^{-15}	4.38×10^{-10}	1.01×10^{-11}	1.12×10^{-11}
DIMACS5	7.94×10^{-08}	-1.08×10^{-06}	-8.60×10^{-10}	2.08×10^{-09}	-2.05×10^{-08}
$\ \mathbf{Z}\ _F$	6.51×10^{08}	1.71×10^{03}	7.89×10^{05}	2.88×10^{06}	6.22×10^{04}

Table 4.19: Performance measures for *gpp100* instance

	($\text{SRSD}_{\text{Bdiag}}$)	SeDuMi	CSDP	SDPA	SDPT3
<i>hinf1</i> 1.90×10^{-10}	5.13×10^{-02} (33, 1.38)	1.70×10^{-02} (21, 0.22)	1.37×10^{-03} (19, 0.04)	6.46×10^{-03} (15, 0.05)	2.32×10^{-02} (13, 0.24)
<i>hinf2</i> 4.16×10^{-05}	1.92×10^{-04} (26, 1.10)	1.35×10^{-02} (17, 0.13)	2.74×10^{-03} (61, 0.06)	3.15×10^{-05} (15, 0.05)	1.86×10^{-04} (15, 0.30)
<i>hinf3</i> 1.29×10^{-10}	3.13×10^{-01} (29, 1.35)	2.89×10^{-01} (18, 0.12)	7.22×10^{-04} (61, 0.08)	1.41×10^{-03} (13, 0.05)	6.08×10^{-04} (20, 0.36)
<i>hinf4</i> 2.47×10^{-09}	7.67×10^{-03} (45, 1.69)	1.10×10^{-02} (31, 0.34)	1.79×10^{-04} (17, 0.04)	8.95×10^{-04} (15, 0.05)	2.06×10^{-04} (21, 0.37)
<i>hinf5</i> 5.39×10^{-10}	1.17×10^{-01} (50, 1.55)	7.64×10^{-01} (18, 0.15)	1.34×10^{-02} (61, 0.09)	4.89×10^{-03} (14, 0.05)	7.57×10^{-03} (22, 0.37)
<i>hinf6</i> 2.28×10^{-09}	7.18×10^{-01} (50, 1.32)	3.78×10^{-01} (22, 0.25)	1.09×10^{-01} (61, 0.09)	4.82×10^{-01} (21, 0.07)	5.13×10^{-02} (21, 0.36)
<i>hinf7</i> 4.08×10^{-09}	6.78×10^{-02} (26, 1.96)	1.90×10^{00} (19, 0.14)	1.95×10^{-03} (61, 0.08)	2.74×10^{00} (9, 0.04)	2.23×10^{-02} (19, 0.33)
<i>hinf8</i> 8.92×10^{-10}	2.25×10^{-01} (35, 2.84)	3.74×10^{-01} (21, 0.29)	5.01×10^{-03} (61, 0.08)	1.29×10^{-01} (14, 0.05)	4.15×10^{-03} (21, 0.37)

Table 4.20: Performance measures for the first eight *hinf* instances

	(SRSD _{B_{Diag}})	SeDuMi	CSDP	SDPA	SDPT3
<i>hinf9</i> 1.53×10^{-02}	9.46×10^{-09} (29, 1.42)	5.28×10^{-03} (21, 0.21)	1.23×10^{-05} (61, 0.06)	2.40×10^{-09} (23, 0.08)	6.99×10^{-04} (22, 0.38)
<i>hinf10</i> 4.26×10^{-11}	3.06×10^{-01} (50, 9.25)	2.96×10^{-01} (25, 0.40)	6.27×10^{-03} (61, 0.07)	2.55×10^{-01} (28, 0.05)	8.13×10^{-01} (23, 0.45)
<i>hinf11</i> 7.10×10^{-11}	2.61×10^{-01} (50, 5.85)	3.23×10^{-01} (25, 0.40)	7.00×10^{-06} (61, 0.10)	1.87×10^{-02} (43, 0.12)	2.00×10^{-01} (24, 0.56)
<i>hinf12</i> 2.92×10^{-11}	2.29×10^{00} (50, 3.99)	5.34×10^{00} (38, 0.36)	1.02×10^{01} (33, 0.07)	9.02×10^{00} (28, 0.26)	1.32×10^{01} (55, 1.40)
<i>hinf13</i> 4.42×10^{-08}	3.88×10^{00} (27, 13.59)	1.19×10^{01} (15, 0.99)	1.15×10^{01} (61, 0.52)	5.63×10^{00} (12, 0.20)	6.56×10^{00} (30, 1.48)
<i>hinf14</i> 4.69×10^{-09}	4.93×10^{-02} (46, 32.00)	3.89×10^{-02} (23, 1.50)	1.55×10^{-02} (61, 0.75)	1.88×10^{-01} (15, 0.22)	8.18×10^{-02} (28, 1.50)
<i>hinf15</i> 3.16×10^{-07}	7.39×10^{00} (27, 47.83)	6.51×10^{00} (16, 0.53)	7.08×10^{-01} (61, 0.90)	7.25×10^{00} (13, 0.26)	2.06×10^{01} (27, 1.28)

Table 4.21: Performance measures for the remaining seven *hinf* instances

The performance measure for the \mathbf{ZX} norm is quite large for all these instances (for all solvers), except for the *hinf9* instance where Slater’s condition for the primal problem can be considered to be satisfied. In order to compare more thoroughly, we again look at DIMACS1, DIMACS3, and DIMACS5. The two hardest instances are *hinf12* and *hinf5*, which have optimal values with only one or two significant digits (see Borchers [2]). The next two tables, Table 4.22 and 4.23, show the results for these two instances.

	(SRSD _{B_{Diag}})	SeDuMi	CSDP	SDPA	SDPT3
<i>hinf12</i> 2.92×10^{-11}	2.29×10^{00} (50, 3.99)	5.34×10^{00} (38, 0.36)	1.02×10^{01} (33, 0.07)	9.02×10^{00} (28, 0.26)	1.32×10^{01} (55, 1.40)
DIMACS1	1.35×10^{-11}	4.40×10^{-12}	7.33×10^{-09}	1.43×10^{-08}	2.30×10^{-11}
DIMACS3	2.67×10^{-07}	1.26×10^{-07}	2.66×10^{-09}	4.28×10^{-09}	6.24×10^{-06}
DIMACS5	-9.96×10^{-05}	-3.31×10^{-04}	-3.08×10^{-02}	-5.56×10^{-02}	-1.34×10^{-05}

Table 4.22: Performance measures for *hinf12* instance

For the *hinf12* instance, the large magnitude of \mathbf{Z} makes all the solvers have negative duality gap. This can be explained by the fact that the dual feasibility error measure is quite large for all solvers due to the large magnitude of the dual solutions. If we accept the primal feasibility error in the order of 10^{-11} , SeDuMi and SDPT3 give the primal objective values of -3.14×10^{-03} and -5.45×10^{-05} , respectively. Our code (**SRSD**_{B_{Diag}}) gives the best objective value of -2.65×10^{-03} . For *hinf5*, it is clear that our code (**SRSD**_{B_{Diag}}) performs better than other solvers based on primal and dual feasibility and duality gap. With both primal and dual feasibility in the order of 10^{-10} , the primal and dual objective value are -3.62208×10^{02} and -3.62218×10^{02} , respectively, which implies the optimal value is -3.6221×10^{02} with five significant digits.

For *mcp100*, *theta1*, *theta2*, *truss1*, *truss3*, and *truss4* instances, our code (**SRSD**_{B_{Diag}}) can get

	(SRSD _{BDiag})	SeDuMi	CSDP	SDPA	SDPT3
<i>hinf5</i> 5.39×10^{-10}	1.17×10^{-01} (50, 1.55)	7.64×10^{-01} (18, 0.15)	1.34×10^{-02} (61, 0.09)	4.89×10^{-03} (14, 0.05)	7.57×10^{-03} (22, 0.37)
DIMACS1	3.66×10^{-10}	2.87×10^{-05}	1.35×10^{-06}	8.65×10^{-05}	1.37×10^{-04}
DIMACS3	4.55×10^{-10}	1.66×10^{-12}	1.11×10^{-07}	5.30×10^{-09}	4.00×10^{-09}
DIMACS5	1.31×10^{-05}	-4.07×10^{-04}	-8.55×10^{-05}	5.20×10^{-04}	-4.26×10^{-04}

Table 4.23: Performance measures for *hinf5* instance

to 10^{-14} for the relative \mathbf{ZX} norm while all other solvers reach the order of 10^{-09} . The remaining test is for the set of *gap* instances. We again observe the issue of Slater's condition for the primal problem which causes the Gauss-Newton direction $\Delta \mathbf{y}$ (and the dual solution \mathbf{Z}) to have extremely large magnitude. Table 4.24 shows the results for these instances. Except for *gap5*, the relative \mathbf{ZX} norm error measures obtained from our code (**SRSD**_{BDiag}) are quite large. However, if we consider primal and dual feasibility and duality gap, it performs as well as other solvers if not better. For example, Table 4.25 shows these results for *gap7* instance. If we allow primal and dual feasibility to be in the order of 10^{-14} and 10^{-08} , we obtain the primal and dual objective values are 4.248199×10^{02} and 4.248196×10^{02} for our code (**SRSD**_{BDiag}), which has only one significant digit less than reported in Borchers [2].

	(SRSD _{BDiag})	SeDuMi	CSDP	SDPA	SDPT3
<i>gap5</i> 4.43×10^{-11}	2.72×10^{-08} (45, 15.64)	2.58×10^{-05} (12, 0.31)	7.92×10^{-07} (14, 0.08)	2.84×10^{-04} (12, 0.05)	4.51×10^{-10} (12, 0.56)
<i>gap6</i> 6.74×10^{-11}	8.70×10^{-03} (34, 51.21)	5.00×10^{-03} (25, 0.97)	3.27×10^{-05} (16, 0.19)	2.04×10^{-02} (17, 0.15)	9.07×10^{-06} (18, 1.47)
<i>gap7</i> 2.12×10^{-10}	7.37×10^{-03} (36, 322.15)	1.69×10^{-03} (23, 1.24)	1.11×10^{-05} (16, 0.43)	1.18×10^{-02} (16, 0.46)	1.26×10^{-05} (18, 2.56)
<i>gap8</i> 1.17×10^{-09}	9.87×10^{-03} (40, 2859.53)	2.53×10^{-03} (27, 3.09)	9.89×10^{-06} (16, 0.98)	8.93×10^{-04} (21, 1.19)	8.06×10^{-05} (17, 1.32)
<i>gap9</i> 1.79×10^{-09}	2.01×10^{-03} (40, 7116.94)	1.10×10^{-03} (26, 6.09)	3.21×10^{-06} (17, 1.82)	2.56×10^{-03} (13, 1.44)	5.62×10^{-05} (18, 4.13)

Table 4.24: Performance measures for *gap* instances

	(SRSD _{BDiag})	SeDuMi	CSDP	SDPA	SDPT3
<i>gap7</i> 2.12×10^{-10}	7.37×10^{-03} (36, 322.15)	1.69×10^{-03} (23, 1.24)	1.11×10^{-05} (16, 0.43)	1.18×10^{-02} (16, 0.46)	1.26×10^{-05} (18, 2.56)
DIMACS1	7.73×10^{-14}	1.62×10^{-07}	3.36×10^{-10}	1.29×10^{-07}	5.21×10^{-07}
DIMACS3	4.91×10^{-08}	2.36×10^{-13}	1.70×10^{-08}	1.72×10^{-11}	3.21×10^{-09}
DIMACS5	4.56×10^{-07}	-2.86×10^{-07}	-3.37×10^{-06}	1.90×10^{-04}	-1.88×10^{-05}

Table 4.25: Performance measures for *gap7* instance

4.2 Concluding Remarks

We have presented a robust algorithm for SDP that is based on a matrix free, inexact Gauss-Newton method. The method takes advantage of well posedness as well as sparsity. Our numerical tests indicate that we get a reduction in the number of iterations (though each is generally more expensive) and an improvement in the accuracy of solutions, compared to current public domain software.

Though our algorithm is currently not competitive with regard to total solution time, it can be used in comparison testing of other algorithms since it provides high accuracy solutions. Further work on efficient preconditioning, scaling of the initial starting points, and parallelization is needed to make the algorithm more competitive.

Acknowledgment

We would like to thank Makoto Yamashita for helping with the installation of SDPA 7.3.1.

References

- [1] F. Alizadeh, J-P.A. Haeberly, and M.L. Overton. Complementarity and nondegeneracy in semidefinite programming. *Math. Programming*, 77:111–128, 1997. 3
- [2] B. Borchers. CSDP, a C library for semidefinite programming. *Optim. Methods Softw.*, 11/12(1-4):613–623, 1999. projects.coin-or.org/Csdp. 25, 29, 30
- [3] J.E. Dennis Jr. and R.B. Schnabel. *Numerical methods for unconstrained optimization and nonlinear equations*, volume 16 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996. Corrected reprint of the 1983 original. 9
- [4] J.E. Dennis Jr. and H. Wolkowicz. Sizing and least-change secant methods. *SIAM J. Numer. Anal.*, 30(5):1291–1314, 1993. 12, 13
- [5] D. C.-L. Fong and M. A. Saunders. Lsmr: An iterative algorithm for sparse least-squares problems. Report SOL 2010-2, Systems Optimization Laboratory, Stanford University, 2010. to appear, *SIAM J. Sci. Comp.* 19
- [6] M. Gonzalez-Lima, H. Wei, and H. Wolkowicz. A stable primal-dual approach for linear programming under nondegeneracy assumptions. *Comput. Optim. Appl.*, 44(2):213–247, 2009. 3
- [7] G. Gruber and F. Rendl. Computational experience with ill-posed problems in semidefinite programming. *Comput. Optim. Appl.*, 21(2):201–212, 2002. 25
- [8] C. Helmberg, F. Rendl, R.J. Vanderbei, and H. Wolkowicz. An interior-point method for semidefinite programming. *SIAM J. Optim.*, 6(2):342–361, 1996. 4
- [9] D.E. Knuth. The sandwich theorem. *Electronic J. Combinatorics*, 1:48pp, 1994. 14
- [10] M. Kojima, S. Shindoh, and S. HARA. Interior-point methods for the monotone semidefinite linear complementarity problem in symmetric matrices. *SIAM J. Optim.*, 7(1):86–125, 1997. 4
- [11] S. Kruk, M. Muramatsu, F. Rendl, R.J. Vanderbei, and H. Wolkowicz. The Gauss-Newton direction in semidefinite programming. *Optim. Methods Softw.*, 15(1):1–28, 2001. 3, 10
- [12] M. Laurent and F. Rendl. Semidefinite programming and integer programming. Technical Report PNA-R0210, CWI, Amsterdam, 2002. To appear as chapter of the Handbook on Discrete Optimization edited by K.Aardal, G. Nemhauser and R. Weismantel. 14

- [13] L. Lovász. On the Shannon capacity of a graph. *IEEE Trans. Inform. Theory*, 25(1):1–7, 1979. 14
- [14] J. Malick, J. Povh, F. Rendl, and A. Wiegele. Regularization methods for semidefinite programming. *SIAM Journal on Optimization*, 20(1):336–356, 2009. 25
- [15] O.L. Mangasarian. *Nonlinear Programming*. McGraw-Hill, New York, NY, 1969. 12
- [16] H.D. Mittelmann. An independent benchmarking of SDP and SOCP solvers. *Math. Program.*, 95(2, Ser. B):407–430, 2003. Computational semidefinite and second order cone programming: the state of the art. 21
- [17] R.D.C. Monteiro. Primal-dual path-following algorithms for semidefinite programming. *SIAM J. Optim.*, 7(3):663–678, 1997. 4
- [18] R.D.C. Monteiro and M.J. Todd. Path-following methods. In *Handbook of Semidefinite Programming*, pages 267–306. Kluwer Acad. Publ., Boston, MA, 2000. 4
- [19] Y.E. Nesterov and M.J. Todd. Self-scaled barriers and interior-point methods for convex programming. *Math. Oper. Res.*, 22(1):1–42, 1997. 4
- [20] Y.E. Nesterov and M.J. Todd. Primal-dual interior-point methods for self-scaled cones. *SIAM J. Optim.*, 8:324–364, 1998. 4
- [21] C.C. Paige and M.A. Saunders. LSQR: an algorithm for sparse linear equations and sparse least squares. *ACM Trans. Math. Software*, 8(1):43–71, 1982. 20
- [22] H. Wei and H. Wolkowicz. Generating and solving hard instances in semidefinite programming. *Math. Programming*, 125(1):31–45, 2010. 23
- [23] H. Wolkowicz. Solving semidefinite programs using preconditioned conjugate gradients. *Optim. Methods Softw.*, 19(6):653–672, 2004. 3, 10

Index

Mat, 4
 \mathcal{X} , 5
 \mathcal{Z} , 5
vec, 4
 $\mathbf{Q}_j := \text{sMat}(\mathbf{Q}(:,j))$, 6
 ω condition number, 12
sMat, 4
svec, 4
 $|\mathbf{S}| := (\mathbf{S}^2)^{\frac{1}{2}}$, 6
GRSD_{BDiag}, block diagonal preconditioner, 23
SRSD_{BDiag}, block diagonal preconditioner, 20
SRSD_{Diag}, diagonal preconditioner, 19
SRSD_o, no preconditioner, 19
TRSD_{Diag}, diagonal preconditioner, 25
TRSD_{MBDiag}, multiple diagonal block preconditioner, 25

adjoint linear transformation, 3

barrier parameter, 3
barrier parameter, μ , 3

central path, 5
CQ, constraint qualification, 3
crossover technique, 10, 20

DIMACS performance measures, 21
DSDP, 3

Jacobian, J , 5

local convergence theorem, 8
Lovász theta number, TN, 14

optimal block diagonal preconditioner, 12

path following, 3
primal-dual optimality conditions, 3
primal-dual SDP, 3
PSDP, 3

Slater constraint qualification, 3

TN, Lovász theta number, 14
triangular number, $t(n) = n(n+1)/2$, 4