

# The Trust Region Subproblem and Semidefinite Programming

Charles Fortin <sup>\*</sup>      Henry Wolkowicz <sup>†</sup>

March 8, 2003

University of Waterloo  
Department of Combinatorics & Optimization  
Waterloo, Ontario N2L 3G1, Canada

**Key words:** Trust Regions, Semidefinite programming, Duality, Unconstrained Minimization.

## Abstract

The trust region subproblem (the minimization of a quadratic objective subject to one quadratic constraint and denoted TRS) has many applications in diverse areas, e.g. function minimization, sequential quadratic programming, regularization, ridge regression, and discrete optimization. In particular, it determines the step in trust region algorithms for function minimization. Trust region algorithms are popular for their strong convergence properties. However, a drawback has been the inability to exploit sparsity as well as the difficulty in dealing with the so-called hard case. These concerns have been

---

<sup>\*</sup>McGill University, Department of Mathematics and Statistics, Montréal, Québec, Canada, H3A 2T5. Research supported by a ES-A scholarship of the Natural Sciences and Engineering Research Council of Canada and a doctoral scholarship (B2) from the Fonds de recherche sur la nature et les technologies du Québec. E-mail [fortin@math.mcgill.ca](mailto:fortin@math.mcgill.ca).

<sup>†</sup>University of Waterloo, Department of Combinatorics and Optimization, Waterloo, Ontario, Canada, N2L 3G1. E-mail [hwolkowicz@uwaterloo.ca](mailto:hwolkowicz@uwaterloo.ca).

<sup>0</sup> This report is based on a longer version CORR 2002-22, [13]. Both this report and the longer version are available with URL:  
[orion.math.uwaterloo.ca/~hwolkowi/henry/reports/ABSTRACTS.html](http://orion.math.uwaterloo.ca/~hwolkowi/henry/reports/ABSTRACTS.html)

addressed by recent advances in the theory and algorithmic development.

This paper provides an in depth study of TRS and its properties as well as a survey of recent advances. We emphasize large scale problems and robustness. This is done using semidefinite programming (SDP) and the modern primal-dual approaches as a unifying framework. The SDP framework arises naturally and solves TRS efficiently. In addition, it shows that TRS is always a well-posed problem, i.e. the optimal value and an optimum can be calculated to a given tolerance. We provide both theoretical and empirical evidence to illustrate the strength of the SDP and duality approach. In particular, this includes new insights and techniques for handling the hard case, as well as numerical results on *large* test problems.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Outline . . . . .	5
<b>2</b>	<b>Optimality Conditions</b>	<b>6</b>
2.1	The Hard Case . . . . .	6
2.1.1	Shift, Deflation, Robustness . . . . .	7
<b>3</b>	<b>A Dual Algorithm; The Moré-Sorensen (MS) Algorithm</b>	<b>12</b>
3.1	Handling the Hard Case in MS Algorithm . . . . .	14
<b>4</b>	<b>The Generalized Lanczos Trust Region (GLTR) Algorithm</b>	<b>16</b>
4.1	Handling the Near Hard Case in GLTR Algorithm . . . . .	16
<b>5</b>	<b>Duality and a Semidefinite Framework for the Trust Region Subproblem</b>	<b>17</b>
5.1	Lagrangian Duality and SDP . . . . .	17
5.2	Unconstrained and Linear Duals . . . . .	18
<b>6</b>	<b>Semidefinite Framework</b>	<b>20</b>
6.1	A Semidefinite Framework for the Moré-Sorensen Algorithm . . . . .	20
6.2	A Semidefinite Framework for the GLTR method . . . . .	21

<b>7</b>	<b>The Rendl-Wolkowicz (RW) Algorithm, Modified</b>	<b>22</b>
7.1	Three Useful Functions . . . . .	22
7.1.1	$k(t) = (s^2 + 1)\lambda_1(D(t)) - t$ . . . . .	22
7.1.2	$k'(t) = (s^2 + 1)y_0(t)^2 - 1$ . . . . .	23
7.1.3	$\psi(t) = \sqrt{s^2 + 1} - \frac{1}{y_0(t)}$ . . . . .	25
7.2	Flowchart . . . . .	25
7.3	Techniques Used in the Algorithm . . . . .	26
7.3.1	Newton's Method on $k(t) - M_t = 0$ . . . . .	26
7.3.2	Triangle Interpolation . . . . .	27
7.3.3	Vertical Cut . . . . .	27
7.3.4	Inverse Interpolation . . . . .	27
7.3.5	Recognizing an unconstrained minimum . . . . .	28
7.3.6	Shift and Deflate . . . . .	28
7.3.7	Taking a Primal Step to the Boundary . . . . .	29
<b>8</b>	<b>Numerical Experiments</b>	<b>30</b>
8.1	The Hard Case . . . . .	30
8.2	RW and GLTR Algorithms in TR Framework . . . . .	30
8.3	Accuracy of TRS in a Trust Region Method . . . . .	34
8.4	Large Sparse TRS . . . . .	36
<b>9</b>	<b>Conclusion</b>	<b>37</b>

# 1 Introduction

We are concerned with the following quadratic minimization problem:

$$\begin{aligned}
 \text{(TRS)} \quad q^* = \min \quad & q(x) := x^T A x - 2a^T x \\
 \text{s.t.} \quad & \|x\| \leq s.
 \end{aligned}$$

Here,  $A$  is an  $n \times n$  symmetric (possibly indefinite) matrix,  $a$  is an  $n$ -vector,  $s$  is a positive scalar and  $x$  is the  $n$ -vector of unknowns. All matrix and vector entries are real. This problem is referred as the trust region subproblem. This problem has many applications in e.g.: forming subproblems for constrained optimization [5], regularization of ill-posed problems [36], and regularization for ill-conditioned linear regression problems (called ridge regression, [20]). In particular, it is important in a class of optimization methods called *trust region (TR) methods* for minimization where, at each

iteration of the method, the algorithm determines a step by (approximately) finding the minimum of a quadratic function (a local quadratic model of a given function  $f$ ) restricted to a given ball of radius  $s$ . (This is called the spherical trust region. We do not discuss scaled, ellipsoidal, box, or other trust regions.) The radius  $s$  increases or decreases depending on how well the decrease in the quadratic model predicts the true decrease in  $f$ . The data,  $A$  and  $a$ , respectively, represent the Hessian and the gradient of the modelled function. These methods have advantages over e.g., quasi-Newton methods. Under mild assumptions, they produce a sequence of iterates with an accumulation point that satisfies *both* first and second order necessary optimality conditions (e.g. [10]). Furthermore, if the accumulation point satisfies the second order sufficient optimality conditions, the method reduces to Newton's method locally and convergence is q-quadratic. (For more details see e.g. the recent books [27, 5].)

However, the popularity of TR methods for unconstrained minimization has lagged behind quasi-Newton methods. Numerical difficulties in standard algorithms for TRS can arise when  $a$  is (approximately) perpendicular to the eigenspace of the smallest eigenvalue of  $A$ . This is referred to as the (*near*) *hard case* in the literature. In addition, algorithms for TRS were based on the Cholesky factorization of the Hessian of the Lagrangian, thus sparsity could not always be exploited efficiently. On the other hand, algorithms such as *limited memory quasi-Newton methods* proved to be successful, e.g. [24, 26].

Though TRS appears to be a simple problem, there is a long history of elegant theory and algorithms. (The recent books [3, 5] contain extensive bibliographies. See also the bibliographical database for [5] at URL [www.fundp.ac.be/~phtoint/pht/trbook.bib](http://www.fundp.ac.be/~phtoint/pht/trbook.bib).) In this paper, we emphasize the modern primal-dual approaches. In particular, we study three methods that consider the above mentioned concerns, i.e. the dual based algorithm of Moré-Sorensen 1983 (MS) [25], the semidefinite programming (SDP) based algorithm of Rendl-Wolkowicz 1997 (RW) [31], and the *generalized Lanczos trust region method* 1999 (GLTR) of Gould, Lucidi, Roma and Toint [17].

The classical (MS) algorithm [25] was the first algorithm able to handle the hard case efficiently. (The algorithm of Gay [14] also treats the hard case.) We revisit and modify the RW primal-dual algorithm [31] which is based on a parametric eigenvalue problem, SDP and duality. It is designed specifically to handle large sparse problems; it also handles the hard case efficiently. GLTR (or coincidentally GLRT) is the last algorithm we look at, see [17]. This algorithm uses the Lanczos procedure to obtain a restricted

TRS problem with a tridiagonal matrix. This subproblem can be solved quickly using the MS algorithm.

Several other recent approaches deserve mention. The method by Sorensen [6] is similar to the RW algorithm in that it uses a parametric eigenvalue approach. The DC (difference of convex functions) method of An and Tao [35] and the method of Hager [18] are both designed to exploit sparsity. The method in [18] is similar in spirit to GLTR, i.e. they both solve a sequence of subproblems where TRS is restricted to a special Krylov subspace. The method of Ye [41] exploits a new efficient line search technique.

The main contribution of this paper is the use of SDP and the modern primal-dual approach to motivate, view, and modify existing algorithms for TRS. In addition, we present a novel approach to handling the hard case using a shift of the eigenvalues and deflation. We also include numerical comparisons between the algorithms and specific examples that illustrate the performance on the hard case. In particular, we try to answer questions posed in [17] about the desired accuracy in solving the TRS within a TR minimization algorithm. We include numerical tests, completely in a MATLAB framework, on problems with dimensions of order  $n = 10^5$  for the TR method, and order  $n = 10^6$  for TRS problems.

## 1.1 Outline

We continue in Section 2 with the optimality conditions and definitions of the easy and hard cases for TRS. In particular, Section 2.1 describes the shift process that yields an equivalent well-posed convex program. That TRS is well-posed also follows using SDP, see Theorem 5.2 and Remark 5.1. In Section 3, we use duality to motivate and describe the MS algorithm. The GLTR algorithm is described in Section 4. In Section 5 we present several dual programs to TRS exploiting the strong Lagrangian duality for TRS. These provide the unifying framework for the different algorithms, see Section 6. The RW algorithm, with our modifications, is presented in detail in Section 7. The numerical tests appear in Section 8. Concluding remarks are given in Section 9.

## 2 Optimality Conditions

It is known (see [15] and [32]) that  $x^*$  is a solution to TRS if and only if

$$\left. \begin{aligned} (A - \lambda^* I)x^* &= a && \text{dual feasibility} \\ A - \lambda^* I &\succeq 0, \lambda^* \leq 0 && \\ \|x^*\|^2 &\leq s^2 && \text{primal feasibility} \\ \lambda^*(s^2 - \|x^*\|) &= 0, && \text{complementary slackness} \end{aligned} \right\} \quad (2.1)$$

for some (Lagrange multiplier)  $\lambda^*$ . These conditions are surprising in two respects. First, the conditions characterize optimality of a possibly nonconvex problem, i.e. they are necessary and sufficient. Second, the usual second order positive semidefinite necessary conditions hold on all of  $\mathbb{R}^n$  rather than just the tangent plane at the optimal point.

We have added the descriptive three phrases in (2.1) since this coincides with the framework in [31] and with the modern primal-dual optimization approach, though no dual program appeared in the earlier papers [15, 32].

### 2.1 The Hard Case

If  $A - \lambda^* I \succ 0$  in (2.1), then  $x^*$  is the unique solution to TRS (this is true generically), i.e.  $x^* = (A - \lambda^* I)^{-1}a$ . In general, we denote

$$x(\lambda) = (A - \lambda I)^\dagger a, \quad (2.2)$$

where  $\cdot^\dagger$  denotes the Moore-Penrose generalized inverse. Singularity (or near singularity) of  $A - \lambda I$  results in difficulties in using  $x(\lambda)$ , see Table 2.1.

1. Easy case	2.(a) Hard case (case 1)	2.(b) Hard case (case 2)
$a \notin \mathcal{N}(A - \lambda_1(A)I)$  (implies $\lambda^* < \lambda_1(A)$ )	$a \perp \mathcal{N}(A - \lambda_1(A)I)$ and $\lambda^* < \lambda_1(A)$	$a \perp \mathcal{N}(A - \lambda_1(A)I)$ and $\lambda^* = \lambda_1(A)$ (i) $\ (A - \lambda^* I)^\dagger a\  = s$ or $\lambda^* = 0$ (ii) $\ (A - \lambda^* I)^\dagger a\  < s, \lambda^* < 0$

Table 2.1: The three different cases for the trust region subproblem. We include two subcases (i) and (ii) for the hard case (case 2).

### 2.1.1 Shift, Deflation, Robustness

First, we note that  $\lambda_1(A) > 0$  implies that  $\lambda^* \leq 0 < \lambda_1(A)$ , i.e. the hard case (case 2) cannot hold. Second, if  $\lambda^* = 0$ , then  $A \succeq 0$  and  $\|x^*\| = \|A^\dagger a\| \leq s$ . These two situations can be handled by our algorithm in a standard way. The following deflation technique forms the basis for our approach to handling the hard case. It shows that we can *deflate and/or shift* eigenspaces that are orthogonal to the linear term  $a$  and, thus, avoid the hard case. For example, we could first use Lemma 2.1 Part 3 to ensure that  $A \succeq 0$ . Then, if the possible hard case is detected, we can use the shift in Lemma 2.1 Part 2. Lemma 2.1 Part 1 shows that performing the shift when the hard case did not exist does not cause harm. In many of our numerical tests, our heuristics detected an unconstrained problem after the shifts. This latter problem was solved using PCG, preconditioned conjugate gradients.

**Lemma 2.1.** *Let:  $A = \sum_{i=1}^n \lambda_i(A) v_i v_i^T = P \Lambda P^T$  be the spectral decomposition of  $A$ , with  $v_i$  orthonormal eigenvectors and  $P = [v_1 \ v_2 \ \dots \ v_n]$  an orthogonal matrix. Set the vector  $\bar{a} := P^T a$  and the sets*

$$\begin{aligned} S_1 &= \{i : \bar{a}_i \neq 0, \lambda_i(A) > \lambda_1(A)\} \\ S_2 &= \{i : \bar{a}_i = 0, \lambda_i(A) > \lambda_1(A)\} \\ S_3 &= \{i : \bar{a}_i \neq 0, \lambda_i(A) = \lambda_1(A)\} \\ S_4 &= \{i : \bar{a}_i = 0, \lambda_i(A) = \lambda_1(A)\}. \end{aligned}$$

For  $k = 1, 2, 3, 4$ : the matrices  $A_k := \sum_{i \in S_k} \lambda_i(A) v_i v_i^T$ ; and the ( $A$ -invariant subspace) orthogonal projections  $P_k := \sum_{i \in S_k} v_i v_i^T$ , where  $A_k = P_k = 0$ , if  $S_k = \emptyset$ . Then the following holds.

1. Suppose  $S_3 \neq \emptyset$  (easy case),  $\alpha > 0$ , and  $i \in S_2 \cup S_4$ . Then

$$(x^*, \lambda^*) \text{ solves TRS}$$

**iff**

$$(x^*, \lambda^*) \text{ solves TRS when } A \text{ is replaced by } A + \alpha v_i v_i^T.$$

2. Let  $u^* = (A - \lambda^* I)^\dagger a$  with  $\|u^*\| < s$  and suppose that  $i \in S_2 \cup S_4$  and  $\alpha > 0$ . Then

$$(x^* = u^* + z, \lambda^*), z \in \mathcal{N}(A - \lambda^* I) \text{ solves TRS}$$

**iff**

$$(x^* = u^* + z, \lambda^*), z \in \mathcal{N}(A + \alpha v_i v_i^T - \lambda^* I) \text{ solves TRS when } A \text{ is replaced by } A + \alpha v_i v_i^T.$$

3. Let  $u^* = (A - \lambda^* I)^\dagger a$ . Then

$$(x^*, \lambda^*), \text{ with } x^* = u^* + z, z \in \mathcal{N}(A - \lambda^* I) \text{ solves TRS}$$

**iff**

$$(u^*, \lambda^* - \lambda_1(A)) \text{ solves TRS when } A \text{ is replaced by } A - \lambda_1(A)I.$$

4.

$$x^* \text{ solves TRS and } v_i^T x^* \neq 0, \text{ for some } i \in S_4$$

**iff**

$$\text{the hard case (case 2(ii)) holds.}$$

**PROOF:** The set  $S_3$  can be used to define the hard case, i.e.  $S_3 \neq \emptyset$  if and only if  $a$  is *not* orthogonal to  $\mathcal{N}(A - \lambda_1(A)I)$  if and only if the easy case holds. Note that  $S_3 = \emptyset \Rightarrow S_4 \neq \emptyset$ .

Consider the equivalent problem to TRS obtained after the rotation by  $P^T$  and diagonalization of  $A$ :

$$\begin{aligned} (\text{TRS}_P) \quad q^* = \min \quad & (P^T x)^T \Lambda (P^T x) - 2\bar{a}^T (P^T x) = w^T \Lambda w - 2\bar{a}^T w \\ \text{s.t.} \quad & \|w\| \leq s, \quad w = P^T x. \end{aligned} \quad (2.3)$$

Note that the  $P_k$  form a resolution of the identity,  $I = \sum_{k=1}^4 P_k$ . Moreover,  $x^*$  solves TRS if and only if  $w^* = P^T x^*$  solves (2.3). We set  $w^* = P^T x^*$  and, for  $k = 1, 2, 3, 4$ ,  $E_k = \sum_{i \in S_k} e_i e_i^T$ ,  $\Lambda_k = \sum_{i \in S_k} \lambda_i e_i e_i^T$ ,  $x_k^* := P_k x^*$ ,  $w_k^* := E_k w^*$ , where the  $e_i$  are unit vectors, and  $E_k = \Lambda_k = 0$ ,  $x_k^* = w_k^* = 0$ , if  $S_k = \emptyset$ . In addition,

$$\Lambda = \sum_{i=1}^4 \Lambda_k, w^* = \sum_{i=1}^4 w_k, I = \sum_{i=1}^4 E_k, E_k = P^T P_k P, w_k^* = P^T x_k^*.$$

For simplification, we prove the results for this diagonalized equivalent program.

1. Since  $S_3 \neq \emptyset$ , the definitions imply that the easy case holds and

$$w^* = (\Lambda - \lambda^* I)^{-1} \bar{a} = (\Lambda + \alpha e_i e_i^T - \lambda^* I)^{-1} \bar{a},$$

i.e. the optimality conditions are unchanged after the replacement. This completes the proof of Item 1.



2. As in the above proof, the definitions imply

$$P^T u^* = (\Lambda - \lambda^* I)^\dagger \bar{a} = (\Lambda + \alpha e_i e_i^T - \lambda^* I)^\dagger \bar{a}.$$

In this case, the optimality conditions are unchanged except for the change in the conditions for  $z$ . This completes the proof of Item 2.

3. Note that  $u^* \in \mathcal{R}(A - \lambda^* I) \perp \mathcal{N}(A - \lambda^* I)$ .

Necessity: Assume that  $(x^*, \lambda^*)$  with  $x^* = u^* + z, z \in \mathcal{N}(A - \lambda^* I)$  solves TRS. Then  $w^* = (\Lambda - \lambda^* I)^\dagger \bar{a} + P^T z, P^T z \in \mathcal{N}(\Lambda - \lambda^* I)$ . It follows from the optimality conditions that

$$\begin{aligned} w^* &= (\Lambda - \lambda^* I)^\dagger \bar{a} + P^T z = P^T u^* + P^T z, \\ \lambda^* (\|P^T u^*\|^2 + \|P^T z\|^2 - s^2) &= 0, \quad (\Lambda - \lambda^* I) \succeq 0. \end{aligned} \quad (2.4)$$

By adding and subtracting  $\lambda_1(A)$ , we see that  $(P^T u^*, \lambda^* - \lambda_1(A))$  is optimal for  $\text{TRS}_P$  if we replace  $\Lambda$  by  $\Lambda - \lambda_1(A)I$ .

Conversely, suppose that  $(u^*, \lambda^* - \lambda_1(A))$  solves TRS when  $A$  is replaced by  $A - \lambda_1(A)I$ . Then

$$P^T u^* = (\Lambda - \lambda^* I)^\dagger \bar{a}, \quad \|u^*\| \leq s, \quad (\Lambda - \lambda_1(A)I) \succeq 0. \quad (2.5)$$

We can find an appropriate  $z \in \mathcal{N}(A - \lambda_1(A)I)$  if needed so that  $\|u^*\|^2 + \|z\|^2 = s^2$ , i.e.  $(x^* = u^* + z, \lambda^*)$  solves TRS. This complete the proof of Item 3.

4. Assume that  $x^* = Pw^*$  solves TRS and  $v_i^T x^* \neq 0$ , for some  $i \in S_4$ . Equivalently,  $e_i^T w^* \neq 0$ , for some  $i \in S_4$ . From the definitions,  $\bar{a}_i = 0, \forall i \in S_2 \cup S_4$ . Therefore  $w^* = (\Lambda - \lambda^* I)^\dagger \bar{a} + E_4 v$ , for some  $v \in \mathcal{R}^n$ . The assumption implies that  $E_4 v \neq 0$ .

Conversely, suppose that the hard case (case 2(ii)) holds. Then  $w^* = (\Lambda - \lambda^* I)^\dagger \bar{a} + E_4 v$ , for some  $v \in \mathcal{R}^n$  with  $E_4 v \neq 0$ . This complete the proof of Item 4. ■

Numerically, we cannot distinguish between the hard case and the near hard case. This is handled using the following.

**Lemma 2.2.** *Suppose that  $x^*$  solves TRS and  $\|x^*\| = s$ . Let  $\epsilon > 0$  and  $v \in \mathbb{R}^n$  with  $\|v\| = 1$ . Let  $\mu^*(\epsilon)$  be the optimal value of TRS when  $a$  is perturbed to  $a + \epsilon v$ . Then*

$$-2s\epsilon \leq \mu^* - \mu^*(\epsilon) \leq 2s\epsilon.$$

PROOF:

$$\begin{aligned} \mu^*(\epsilon) &= \min_{\|x\|=s} q(x) - 2\epsilon v^T x \\ &\geq \min_{\|x\|=s} q(x) + \min_{\|x\|=s} -2\epsilon v^T x \\ &= \mu^* - 2\epsilon s. \end{aligned}$$

This proves the right-hand-side inequality.

Since  $x^*$  is optimal for TRS and on the boundary of the ball, we get

$$\begin{aligned} \mu^*(\epsilon) &\leq \mu^* - 2\epsilon v^T x^* \\ &\leq \mu^* + 2\epsilon s. \end{aligned}$$

■

The literature often labels the hard case (case 2) as an ill-posed or degenerate problem, e.g. [17, 18]. Adding a norm constraint to an ill-posed problem is a well-known regularization technique, e.g. [36]. Thus it would appear to be contradictory for TRS to be an ill-posed problem. In fact, we can orthogonally diagonalize the quadratic form; and, we note that the symmetric eigenvalue problem has a condition number of 1, [8]. Then, TRS can be shown to be equivalent to a linearly constraint convex programming problem, see [2], a problem that can be solved efficiently and robustly.

The following lemma and example illustrate that TRS is always a stable convex program. This also follows from the equivalent SDP dual pair in Theorem 5.2 below.

**Lemma 2.3.** *Suppose that the hard case (case 2 (ii)) holds for TRS. Let  $u^* = (A - \lambda^* I)^\dagger a$ ,  $x^* = u^* + z$ ,  $z \in \mathcal{N}(A - \lambda^* I)$ . Then  $z \neq 0$ ,  $\lambda_1(A) \leq 0$  and TRS is equivalent to the following stable convex program*

$$(TRS_s) \quad \begin{aligned} q_s^* &:= \min_{\|x\| \leq s} q_s(x) := x^T (A - \lambda_1(A) I) x - 2a^T x + s^2 \lambda_1(A) \end{aligned} \tag{2.6}$$

*The equivalence is in the sense that the optimal values satisfy  $q^* = q_s^*$ ; and  $(x^*, \lambda^*)$  solves TRS if and only if  $(x^*, 0)$  solves  $TRS_s$ .*

PROOF: Since the hard case (case 2 (ii)) holds, we get  $z \neq 0$  and  $\lambda_1(A) = \lambda^* \leq 0$ . From Lemma 2.1, Item 3 we get  $(u^*, 0)$  solves  $\text{TRS}_s$ . We can then add  $z \in \mathcal{N}(A - \lambda^*I)$  to get  $\|u^* + z\| = s$ .  $\blacksquare$

Convex programs for which Slater's CQ holds are called *stable*, e.g. [16, 9]. They are equivalent to convex programs for which the optimal dual solutions form a convex compact set which further implies that the perturbation function (optimal value function subject to linear perturbations in the data) is convex and Lipschitz continuous. In our case we have the additional strong linear independence CQ which implies that the optimal dual solution is unique and the perturbation function is differentiable. We also have a compact convex feasible set. (See e.g. [9, 16].)

**Example 2.1.** (*Hard Case, Case 2(ii)*) Let

$$A = \begin{bmatrix} 1 + \gamma & 0 \\ 0 & -1 + \delta \end{bmatrix}, \quad a = \begin{bmatrix} 2 + \alpha \\ \beta \end{bmatrix}, \quad s = \sqrt{2},$$

where  $\alpha, \beta, \gamma, \delta$  are perturbations in the data. First suppose that the perturbations are all 0. Then the hard case (case 2(ii)) holds; the optimal Lagrange multiplier is  $\lambda^* = \lambda_1(A) = -1$ ; and the best least squares solution is  $\bar{x} = (A - \lambda^*I)^\dagger a = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  with  $\|\bar{x}\| = 1 < s$ . The optimal solution is obtained from

$$x^* = x^*(0) = \bar{x} + \begin{bmatrix} 0 \\ \pm 1 \end{bmatrix} = \begin{bmatrix} 1 \\ \pm 1 \end{bmatrix}. \quad (2.7)$$

For (small) nonzero perturbations, the optimal Lagrange multiplier  $\lambda^*$  is still unique and  $-1 + \delta$  is the smallest eigenvalue. If  $\beta = 0$ , then the hard case still holds; the optimum is obtained from  $\lambda^* = -1 + \delta$ ; and

$$x^* = x^*(\alpha, \gamma, \delta) = \begin{bmatrix} \frac{2+\alpha}{1+\gamma-\lambda^*} \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \epsilon \pm 1 \end{bmatrix},$$

where  $\epsilon$  is chosen to obtain  $\|x^*\| = s$ , e.g. with  $+1$

$$(1 + \epsilon)^2 + \left( \frac{2 + \alpha}{2 + \gamma - \delta} \right)^2 = s^2 = 2.$$

Depending on the choice of sign, these solutions converge to a solution in (2.7), as the perturbations converge to zero. Moreover, a Taylor series expansion shows that  $\|x^*(0) - x^*(\alpha, \gamma, \delta)\| \leq 2(|\alpha| + |\gamma| + |\delta|)$  for small perturbations, i.e. we have a bounded condition number.

If  $\beta \neq 0$ , then we have the easy case. The unique optimal Lagrange multiplier  $\lambda^* < -1 + \delta$  and the unique optimum is obtained from

$$x^* = \begin{bmatrix} \frac{2+\alpha}{1+\gamma-\lambda^*} \\ \frac{\beta}{-1+\delta-\lambda^*} \end{bmatrix},$$

where  $\lambda^*$  satisfies the positive definiteness condition as well as  $\|x^*\| = s$ . This implies that

$$\left(\frac{2+\alpha}{1+\gamma-\lambda^*}\right)^2 + \left(\frac{\beta}{-1+\delta-\lambda^*}\right)^2 = 2.$$

Since  $\lambda^* \rightarrow -1$  and  $\frac{2+\alpha}{1+\gamma-\lambda^*} \rightarrow 1$ , as the perturbations go to 0, we see that the optimal solutions converge appropriately again.

### 3 A Dual Algorithm; The Moré-Sorensen (MS) Algorithm

We motivate the MS algorithm using duality and illustrate how SDP arises naturally from this setting. The *Lagrangian dual* of TRS is

$$\begin{aligned} q^* = \nu^* &:= \max_{\lambda \leq 0} \min_x x^T (A - \lambda I)x - 2a^T x + \lambda s^2 \\ &= \max_{\lambda \leq 0} h(\lambda) \end{aligned}$$

where the *Lagrangian* is  $L(x, \lambda) := x^T (A - \lambda I)x - 2a^T x + \lambda s^2$  and the *dual functional* is  $h(\lambda) := \min_x L(x, \lambda)$ . Strong duality ( $q^* = \nu^*$  and dual attainment, [34]) holds. Since the inner minimization is unconstrained, we have a *hidden semidefinite constraint* that the Hessian,  $\nabla^2 L(x, \lambda) \succeq 0$ , is positive definite. The dual functional  $h$  is concave with domain within the semidefinite constraint.

Therefore, we can replace TRS with the simpler root finding problem

$$h'(\lambda) = 0 \tag{3.1}$$

(if  $h$  is differentiable), with the restrictions:

$$\lambda \leq 0, \quad \nabla^2 L(x, \lambda) = A - \lambda I \succeq 0. \quad (3.2)$$

We see that semidefiniteness (convexity of the Lagrangian) arises naturally from the duality setting. Note that  $h'(\lambda) = s^2 - x(\lambda)^T x(\lambda)$ , when it exists. If we use the optimality conditions in (2.1) with the descriptive phrases, then we see that the MS algorithm maintains dual feasibility and complementary slackness, while trying to attain primal feasibility. (cf. the dual simplex method for linear programming.)

Though Newton's method has asymptotic q-quadratic convergence, the Newton step does not take into account the semidefinite restrictions, which can result in many backtracking steps. This illustrates the weakness of a dual method compared to a primal-dual method.

The main work in the iterations is a Cholesky factorization used in the evaluation of the derivatives and the Newton step for  $\lambda$ , as well as in the safeguarding and updating scheme that produces either a point  $\lambda$  from which quadratic convergence ensues, or reduces the interval of uncertainty for the optimal  $\lambda$ . If the possible hard case is detected, optimality is reached by taking primal steps to the boundary of the ball. In both cases, given two parameters  $\sigma_1$  and  $\sigma_2$  in  $(0, 1)$ , the algorithm terminates in a finite number of iterations with an approximate solution  $\bar{x}$  which satisfies

$$q(\bar{x}) - q^* \leq \sigma_1(2 - \sigma_1) \max\{|q^*|, \sigma_2\} \quad , \quad \|\bar{x}\| \leq (1 + \sigma_1)\Delta. \quad (3.3)$$

One additional innovation makes the MS algorithm fast. Assume  $A - \lambda I \succ 0$ . There are disadvantages in applying Newton's method to find a root of the function  $h$  or equivalently of  $\psi(\lambda) := \|x(\lambda)\| - s$ . For  $\lambda < \lambda_1(A)$  and close to  $\lambda_1(A)$ , the orthogonal diagonalization of  $A$  shows that

$$\psi(\lambda) = \|Q(\Lambda - \lambda I)^{-1}Q^T a\| - s \approx \frac{c_1}{\lambda_1(A) - \lambda} + d,$$

for some constants  $c_1 > 0, d$ . This function is highly nonlinear for values of  $\lambda$  near  $\lambda_1(A)$ , which equates to slow convergence for Newton's method. Moré-Sorensen solve the equivalent so-called secular equation

$$\phi(\lambda) := \frac{1}{s} - \frac{1}{\|x(\lambda)\|} = 0. \quad (3.4)$$

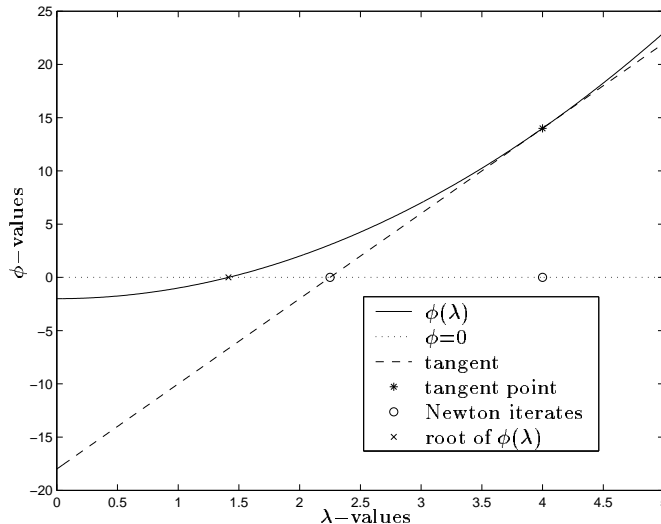


Figure 3.1: Newton's method with the secular function,  $\phi(\lambda)$ .

(See e.g. [28, 29, 30, 11].) The rational structure of  $\|x(\lambda)\|^2$ , shows that this function is less nonlinear, i.e.

$$\phi(\lambda) \approx \frac{1}{s} - \frac{\lambda_1(A) - \lambda}{c_2},$$

for some  $c_2 > 0$ . Therefore, Newton's method applied to this function will be more efficient. One can also show  $\phi(\lambda)$  is a convex function strictly increasing on  $(-\infty, \lambda_1(A))$ . (See [25, Pg 562] for the details.)

### 3.1 Handling the Hard Case in MS Algorithm

From Figure 3.1, we get the following indicator of the easy case for TRS.

**Lemma 3.1.** *Suppose that  $\lambda_0 \leq \min\{0, \lambda_1(A)\}$  and  $\phi(\lambda_0) > 0$  (equivalently  $\|x(\lambda_0)\| > s$ ). Then the hard case (case 2) cannot occur for TRS, i.e.  $\lambda^* < \lambda_1(A)$ . ■*

However, Figure 3.1 also shows that Newton's method can provide a poor prediction for  $\lambda^*$  if  $\phi(\lambda) < 0$  and  $\lambda^*$  is close to 0 and/or  $\lambda_1(A)$ . This would result in many backtracking steps to find the above  $\lambda_0$  (each of which involves

an attempted Cholesky factorization) and make the algorithm inefficient. As discussed above in Section 2.1, in the hard case (case 2), a solution to TRS can be obtained by first finding a solution  $x(\lambda_1(A))$  to the system

$$(A - \lambda_1(A)I)x = a \tag{3.5}$$

with  $\|x\| \leq s$ . If strict inequality holds,  $\|x\| < s$ , then we need an eigenvector  $z \in \mathcal{N}(A - \lambda_1(A)I)$ , and  $\tau \in \mathbb{R}$ , such that  $\|x(\lambda_1(A)) + \tau z\| = s$ , i.e.

$$x^* = x(\lambda_1(A)) + \tau z \tag{3.6}$$

satisfies the optimality conditions. The following lemma by Moré-Sorensen [25] is the key to implementing this idea numerically.

**Lemma 3.2 (Primal step to the boundary).** *Let  $0 < \sigma < 1$  be given and suppose that*

$$A - \lambda I = R^T R, \quad (A - \lambda I)x = a, \quad \lambda \leq 0. \tag{3.7}$$

Let  $z \in \mathbb{R}^n$  satisfy

$$\|x + z\|^2 = s^2 \tag{3.8}$$

and

$$\|Rz\|^2 \leq \sigma(\|Rx\|^2 - \lambda s^2). \tag{3.9}$$

Then

$$|q(x + z) - q(x^*)| \leq \sigma |q(x^*)|. \tag{3.10}$$

where  $x^*$  is optimal for TRS. ■

We will get back to this Lemma in Section 6.1 below, where we show that this lemma is measuring a duality gap in a primal-dual pair of SDPs. Note that (3.7) and (3.8) guarantee the dual and primal feasibility constraints in the optimality conditions (2.1). Comparisons with the RW algorithm are included in Section 7 below.

## 4 The Generalized Lanczos Trust Region (GLTR) Algorithm

As mentioned above, current attempts to solve TRS focus on exploiting sparsity. The Cholesky factorization can be a bottleneck in the MS algorithm for sparse problems without special structure. The GLTR algorithm involves a Lanczos tridiagonalization of the matrix  $A$ , which then allows for fast Cholesky factorizations. The method requires only matrix-vector multiplications and exploits sparsity in  $A$ . The method solves a *sequence* of restricted problems

$$\begin{aligned} \min \quad & q(x) \\ \text{s.t.} \quad & \|x\| \leq s \\ & x \in S \equiv \mathcal{K}_k, \end{aligned} \tag{4.1}$$

where  $\mathcal{K}_k := \text{span}\{a, Aa, A^2a, A^3a, \dots, A^ka\}$  are specially chosen (Krylov) subspaces of  $\mathbb{R}^n$ . (A similar approach based on Krylov subspaces is presented in [18].) The way  $S$  is chosen is inspired by the Steihaug algorithm of [33] (see also [17] and [37]). The authors use heuristics for the MS algorithm that find a starting value of  $\lambda$  on the correct side of  $\lambda^*$  to ensure asymptotic q-quadratic convergence of the Newton iteration  $\lambda^*$ .

### 4.1 Handling the Near Hard Case in GLTR Algorithm

If the near hard case (case 2) occurs for the tridiagonal subproblem, then finding a  $\lambda$  that guarantees the q-quadratic convergence may be difficult and time consuming. (The GLTR algorithm fails if the hard case (case 2) occurs.) In addition, ill-conditioning will slow down both the MS algorithm and therefore the GLTR method (see [17], pp. 515).

Further details are included within the semidefinite framework in Section 6.2 below.

As mentioned above, the method in [18] is similar in spirit to GLTR. The main difference is the choice of vectors to put into the Krylov subspaces  $\mathcal{K}_k$ . The vectors include the direction found from a sequential quadratic programming (SQP) model for TRS.



## 5 Duality and a Semidefinite Framework for the Trust Region Subproblem

Duality plays a central role in designing optimization algorithms, as illustrated in our motivation for the MS algorithm in Section 3. In this section, we focus our attention on different dual programs associated with TRS. We will further see below the role played by duality in both the MS and GLTR algorithms, i.e. they are both *dual based*, and so they exhibit slow convergence and lack of robustness, characteristics of dual algorithms.

For simplicity, we restrict ourselves to the equality TRS, i.e. we minimize over the sphere rather than the ball. (To extend to the standard TRS we would need to add the dual constraint  $\lambda \leq 0$ .) Precisely, consider the slightly different problem

$$\text{(TRS}_=\text{)} \quad q^* = \min_{\|x\|=s} q(x) \quad (5.1)$$

### 5.1 Lagrangian Duality and SDP

From [34], we know that strong Lagrangian duality holds for TRS<sub>=</sub>, i.e.

$$q^* = \nu^* := \max_{\lambda} \min_x L(x, \lambda). \quad (5.2)$$

Since  $L(x, \lambda)$  is a quadratic, the inner min problem is unbounded below unless the hidden constraints,

$$A - \lambda I \succeq 0, \quad a \in \mathcal{R}(A - \lambda I), \quad (5.3)$$

hold. (This can be seen by moving in a direction of an eigenvector corresponding to a negative eigenvalue or, if  $A - \lambda I \succeq 0$ ,  $a \notin \mathcal{R}(A - \lambda I)$ , then moving in a direction  $d \in \mathcal{N}(A - \lambda I)$  such that  $d^T a > 0$ .) This yields the equivalent dual problem

$$q^* = \max_{\substack{A - \lambda I \succeq 0, \\ a \in \mathcal{R}(A - \lambda I)}} \min_x L(x, \lambda).$$

The inner minimum is attained if bounded. We get the following.

**Theorem 5.1.** ([34]) *The Lagrangian dual for  $\text{TRS}_=$  is*

$$(D) \quad q^* = \sup_{A - \lambda I \succ 0} h(\lambda), \quad (5.4)$$

where

$$h(\lambda) := \lambda s^2 - a^T(A - \lambda I)^\dagger a,$$

where  $h$  is a concave function on the feasible set. In the easy case and hard case (case 1), the sup can be replaced by a max. ■

**Corollary 5.1.** *The Lagrangian dual for  $\text{TRS}$  is equivalent to*

$$(D) \quad q^* = \sup_{\substack{A - \lambda I \succ 0 \\ \lambda \leq 0}} h(\lambda). \quad (5.5)$$

In the easy case and hard case (case 1), the sup can be replaced by a max. ■

## 5.2 Unconstrained and Linear Duals

We now present an unconstrained concave maximization problem and a pair of linear SDPs, all of which are equivalent to  $\text{TRS}_=$ , see [31] for the details. Define

$$D(t) = \begin{pmatrix} t & -a^T \\ -a & A \end{pmatrix}, \quad k(t) := (s^2 + 1)\lambda_1(D(t)) - t. \quad (5.6)$$

Then we have the following unconstrained dual problem for  $\text{TRS}_=$ :

$$q^* = \max_t k(t). \quad (5.7)$$

It is well known that  $\lambda_1(D(\cdot))$  is a concave function, and therefore  $k(\cdot)$  is concave as well. Thus, using duality,  $\text{TRS}_=$  is equivalent to an *unconstrained* concave maximization problem in *one* variable. We can also rewrite (5.7) in the following way so that it becomes a linear semidefinite program:

$$\max_{D(t) \succeq \lambda I} (s^2 + 1)\lambda - t. \quad (5.8)$$

Equivalently,

$$\begin{aligned} q^* = \max \quad & (s^2 + 1)\lambda - t \\ \text{s.t.} \quad & \lambda I - tE_{00} \preceq D(0), \end{aligned} \tag{5.9}$$

where  $E_{00}$  is the zero matrix except for 1 in the top left corner. Because Slater's constraint qualification holds for this problem, one can take the Lagrangian dual and get a semidefinite equivalent for  $\text{TRS}_=$

$$\begin{aligned} q^* = \min \quad & \text{trace } D(0)Y \\ \text{s.t.} \quad & \text{trace } Y = s^2 + 1 \\ & -Y_{00} = -1 \\ & Y \succeq 0. \end{aligned} \tag{5.10}$$

**Theorem 5.2.** *The three programs: (5.7), (5.9), (5.10), are equivalent to  $\text{TRS}_=$ . Moreover, the Slater constraint qualification and strict complementarity hold for the primal-dual SDP pair (5.9) and (5.10).*

PROOF: The equivalence with  $\text{TRS}_=$  was already shown above.

That Slater's CQ holds is clear, i.e. choose  $\lambda$  and  $Y$  appropriately.

Now suppose that  $\lambda, t, Y$  are optimal for the SDP pair (5.9) and (5.10). Then  $Z := D(t) - \lambda I \succeq 0$  and singular. Let  $k$  be the multiplicity of  $\lambda_1(D(t))$  and  $y_1, \dots, y_k$  be an orthonormal basis for its eigenspace. Set  $V = [y_1 \ \dots \ y_k]$  and redefine  $Y \leftarrow Y + VV^T$ . We can scale  $DYD$  using a diagonal matrix  $D$  to guarantee that  $Y$  is feasible for (5.10). By construction

$$ZY = 0, Z + Y \succ 0.$$

■

**Corollary 5.2.** *The SDP (5.10) has a rank one optimal solution  $Y^*$ .*

PROOF: Let  $x^*$  be an optimum for  $\text{TRS}$  and

$$y^* = \begin{pmatrix} 1 \\ x^* \end{pmatrix}, \quad Y^* = y^*(y^*)^T.$$

■

**Remark 5.1.** *The primal-dual linear SDP pair can be solved to any desired accuracy in polynomial time, see e.g. [39]. This emphasizes that  $\text{TRS}$  is a well-posed problem.*

## 6 Semidefinite Framework

From above (Theorem 5.2), we saw that  $\text{TRS}_=$  is equivalent to a primal-dual pair of SDPs which could be solved using primal-dual interior-point methods. These methods have revolutionized our view of optimization during the last 15 years. In particular, path-following methods have proven to be an efficient approach for many classes of optimization problems. The main idea for these methods is to apply Newton's method to a perturbation of the primal-dual optimality conditions. Using *both* the primal and dual equations and variables makes for efficient, robust algorithms. (The recent books [38, 40] describe this approach for both linear and semidefinite programming.) Often, compromises have to be made to deal with large sparse problems. In particular, for SDP one often uses a dual based method to exploit sparsity, since the primal variable is often large and dense, see e.g. [19, 1].

We previously motivated the MS algorithm using duality. We now describe the MS and GLTR Algorithms using SDP and the modern primal-dual approach. We see that compromises are made and a full primal-dual path-following method is not used.

### 6.1 A Semidefinite Framework for the Moré-Sorensen Algorithm

For simplicity, we consider the case  $a \neq 0$ . We follow Section 5.2 (see also [31]) and use the following pair of SDP dual programs. (D) is the dual of TRS and (DD) is the dual of (D):

$$(D) \quad q^* = \sup_{A - \lambda I \succ 0} h(\lambda) \quad (6.1)$$

$$(DD) \quad \begin{aligned} q^* = \inf \quad & h(\lambda) + \text{trace}(X(A - \lambda I)) \\ \text{s.t.} \quad & s^2 - a^T((A - \lambda I)^\dagger)^2 a - \text{trace} X = 0 \\ & \lambda < \lambda_l(A) \\ & \text{trace} X \leq s^2 \\ & X \succeq 0, \end{aligned} \quad (6.2)$$

where  $\lambda_l(A)$  is the smallest eigenvalue such that  $a \notin \mathcal{N}(A - \lambda_l(A)I)$ .

In the easy case and the hard case (case 1), we use the dual program (D). The supremum in (6.1) is attained at the stationary point  $\lambda^* \in (-\infty, \lambda_1(A))$ ,

$$h'(\lambda^*) = -a^T((A - \lambda^*I)^{-1})^2 a + s^2 = -\|(A - \lambda^*I)^{-1}a\|^2 + s^2 = 0.$$

Newton's method is applied to the equivalent root finding problem  $\phi(\lambda) = 0$ . Safeguarding guarantees that  $\lambda^*$  stays in the proper interval. Though Newton's method guarantees q-quadratic convergence, this may only happen after many Newton and backtracking steps. The semidefinite constraint is not used explicitly in choosing the Newton direction or the step length.

Things are different in the hard (or near hard) case (case 2), i.e. this is the case when the current estimates satisfy primal feasibility  $\|x(\lambda)\| < s$ . In this case, MS uses a dual-primal approach. Given such a  $\lambda$ , we now use (DD) to reduce the duality gap,  $\text{trace}(X(A - \lambda I))$ , between (D),(DD); and we simultaneously reduce the objective value of TRS. To do this we find  $z$  to move to the boundary (i.e. the primal step is  $\|x + z\|^2 = s^2$ ). The SDP (6.2) suggests how such a  $z$  should be chosen.

$$\begin{aligned} q(x + z) &= (x + z)^T A(x + z) - 2a^T(x + z) + \lambda s^2 - \lambda \|x + z\|^2 \\ &= \lambda s^2 + x^T(A - \lambda I)x + 2x^T(A - \lambda I)z + z^T(A - \lambda I)z - 2a^T x - 2a^T z \\ &= h(\lambda) + z^T(A - \lambda I)z \\ &= h(\lambda) + \text{trace}(zz^T(A - \lambda I)). \end{aligned}$$

To summarize, note that in the MS algorithm,  $\|Rz\|^2 = z^T(A - \lambda I)z$ . Therefore, when a  $z$  is found such that  $\|x + z\|^2 = s^2$  and  $\|Rz\|$  is small, the algorithm is trying to reduce the duality gap between (6.1) and (6.2), while maintaining feasibility for (6.2).

## 6.2 A Semidefinite Framework for the GLTR method

As in the MS algorithm, we now show that the GLTR Algorithm can also be explained using the Lagrangian dual (6.1). Here we outline how their stopping criteria is in fact measuring the duality gap between TRS<sub>=</sub> and this Lagrangian dual. (A detailed discussion is given in [13].)

Each iteration of the algorithm returns a feasible point  $x_k$  for TRS and a corresponding Lagrange multiplier  $\lambda_k$  which are optimal for the subproblem (4.1). The algorithm stops when stationarity is satisfied up to a tolerance, i.e. when

$$\|(A - \lambda_k I)x_k - a\| \tag{6.3}$$

becomes small. When the  $\lambda_k$  are feasible for (6.1) and bounded, then it is possible to show

$$q(x_k) - h(\lambda_k) = O(\|(A - \lambda_k I)x_k - a\|^2),$$

i.e. the duality gap is bounded by a quantity proportional to the square of (6.3). Therefore, convergence of (6.3) to zero implies a zero duality gap.

Though the GLTR Algorithm appears to be a primal algorithm, since simpler primal problems are solved to approximate the solution to TRS, the strength of the approximation is directly linked to the duality gap between TRS and the dual problem (6.1).

## 7 The Rendl-Wolkowicz (RW) Algorithm, Modified

This algorithm both exploits the sparsity of  $A$  and handles the hard case. The algorithm is based on using various primal and dual steps to reduce the interval of uncertainty for the optimum (maximum) of the unconstrained dual program (5.7). We exploit the properties of the eigenvalues and eigenvectors of the parametric matrix  $D(t)$ . Many ideas from the MS algorithm are transformed to the large sparse case, e.g. the primal step to the boundary and the secular function. We also exploit information from the primal-dual pair of linear SDPs (5.9),(5.10). We outline the algorithm with a flowchart in Section 7.2. In Section 7.3 we list new heuristics that take advantage of the structure of  $k(\cdot)$ , accelerate convergence, and facilitate the handling of the hard case.

### 7.1 Three Useful Functions

Graphs, illustrating the properties of these functions, appear in [13].

$$7.1.1 \quad k(t) = (s^2 + 1)\lambda_1(D(t)) - t$$

This is the function that we (implicitly) maximize to solve TRS, see (5.7). Since

$$\lim_{t \rightarrow \infty} \lambda_1(D(t)) = \lambda_1(A) \text{ and } \lim_{t \rightarrow -\infty} (\lambda_1(D(t)) - t) = 0,$$

the asymptotic behavior of  $k(t)$  is

$$\begin{aligned} k(t) &\sim (s^2 + 1)\lambda_1(A) - t, & \text{as } t \rightarrow \infty & \quad (\text{linear with slope } -1), \\ k(t) &\sim s^2 t, & \text{as } t \rightarrow -\infty & \quad (\text{linear with slope } s^2), \end{aligned}$$

i.e.  $k(t)$  is linear as  $|t| \rightarrow \infty$ . Since  $\lambda_1(D(t))$  is concave, so is  $k(t)$ . In the easy case, the function is differentiable and strictly concave. In the hard case, loss of differentiability occurs when the multiplicity of the smallest eigenvalue for  $\lambda_1(D(t))$  changes. The following theorem, based on [31, Proposition 8, Lemma 9, Lemma 15], tells us when this happens.

**Theorem 7.1.** *Let  $A = P\Lambda P^T$  be an orthogonal diagonalization of  $A$ . Let  $\lambda_1(A)$  have multiplicity  $i$  and define*

$$t_0 := \lambda_1(A) + \sum_{j \in \{k | (P^T a)_k \neq 0\}} \frac{(P^T a)_j^2}{\lambda_j(A) - \lambda_1(A)}.$$

Then:

1. In the easy case, for all  $t \in \mathbb{R}$ ,  $\lambda_1(D(t)) < \lambda_1(A)$  and  $\lambda_1(D(t))$  has multiplicity 1.
2. In the hard case:
  - (a) for  $t < t_0$ ,  $\lambda_1(D(t)) < \lambda_1(A)$  and  $\lambda_1(D(t))$  has multiplicity 1;
  - (b) for  $t = t_0$ ,  $\lambda_1(D(t)) = \lambda_1(A)$  and  $\lambda_1(D(t))$  has multiplicity  $1 + i$ ;
  - (c) for  $t > t_0$ ,  $\lambda_1(D(t)) = \lambda_1(A)$  and  $\lambda_1(D(t))$  has multiplicity  $i$ .

■

**Corollary 7.1.** *In the hard case, for  $t \geq t_0$ ,  $k(t) = (s^2 + 1)\lambda_1(A) - t$ .*

■

Note that the maximum  $t^* \leq t_0$ . Difficulties with differentiability arises when  $t^*$  is close to  $t_0$ .

**7.1.2**  $k'(t) = (s^2 + 1)y_0(t)^2 - 1$

Maximizing the concave function  $k(t)$  is equivalent to finding  $0 \in \partial k(t)$  (subgradient). Recall that  $y(t)$  is a normalized eigenvector for  $\lambda_1(D(t))$  and  $y_0(t)$  is its first component. If  $y(t) = \begin{pmatrix} y_0(t) \\ x(t) \end{pmatrix}$ , then, in the differentiable case,

$$\frac{1}{y_0(t)^2} \|x(t)\|^2 = \frac{1 - y_0(t)^2}{y_0(t)^2} = s^2 \text{ if and only if } k'(t) = 0,$$

i.e. this is equivalent to primal feasibility (cf  $h'(\lambda) = 0$  in MS algorithm). To obtain conditions for  $y_0(t) \neq 0$ , we use the following from [31, Lemma 12, Lemma 15].

**Theorem 7.2.** *Let  $y(t)$  be a normalized eigenvector for  $\lambda_1(D(t))$  and let  $y_0(t)$  be its first component. Then:*

1. **In the easy case:** for  $t \in \mathbb{R}$ ,  $y_0(t) \neq 0$ ;
2. **In the hard case:**
  - (a) **for  $t < t_0$ :**  $y_0(t) \neq 0$ ;
  - (b) **for  $t > t_0$ :** there exists a basis of eigenvectors for the eigenspace of  $\lambda_1(D(t))$  such that each eigenvector in the basis has a zero first component ( $y_0(t) = 0$ ) and the vector composed of the last  $n$  components is an eigenvector for  $\lambda_1(A)$ ;
  - (c) **for  $t = t_0$ :** there exists a basis of eigenvectors for the eigenspace of  $\lambda_1(D(t_0))$ , such that one eigenvector of this basis,  $\omega$ , has a non-zero first component ( $\omega_0 \neq 0$ ) and each of the other eigenvectors in the basis has a zero first component ( $y_0(t) = 0$ ) and the vector composed of the last  $n$  components is an eigenvector for  $\lambda_1(A)$ .

■

It is known that the function  $\lambda_1(D(t))$  is differentiable at points where the multiplicity of the eigenvalue is 1. Its derivative is given by  $y_0(t)^2$ , where  $y(t)$  is a normalized eigenvector for  $\lambda_1(D(t))$ , i.e.  $\|y(t)\| = 1$  (see [21]). Therefore, Theorems 7.1 and 7.2 yield the following.

**Corollary 7.2.** 1. *In the easy case:  $k(\cdot)$  is differentiable and  $k'(t) = (s^2 + 1)y_0(t)^2 - 1$ .*

2. *In the hard case:*

- (a) *for  $t < t_0$ ,  $k(\cdot)$  is differentiable and  $k'(t) = (s^2 + 1)y_0(t)^2 - 1$ ;*
- (b) *for  $t = t_0$ ,  $k(\cdot)$  is non-differentiable and the directional derivatives from the left and right are, respectively,  $k'_-(t_0) = (s^2 + 1)\omega_0^2 - 1$  and  $k'_+(t_0) = -1$ ;*
- (c) *for  $t > t_0$ ,  $k(\cdot)$  is differentiable and  $k'(t) = -1$ .*



■

The structure of the eigenvectors along with the shift and deflation Lemma 2.1 can be used to avoid the hard case, i.e. if  $y_0(t)$  is *small*, then we can deflate using the corresponding eigenvector. Lemma 2.2 shows that the deviation from the original problem is *small*.

$$7.1.3 \quad \psi(t) = \sqrt{s^2 + 1} - \frac{1}{y_0(t)}$$

Solving  $\psi(t) = 0$  is equivalent to solving  $k'(t) = 0$ . The advantage is that  $\psi(t)$  is less nonlinear (cf replacing  $h'$  with  $\phi$  in the MS algorithm). It can be shown that  $\psi(t)$  is strictly decreasing and converges to  $\sqrt{s^2 + 1} - 1$  as  $t \rightarrow -\infty$ . In the easy case,  $\psi(t)$  goes to  $-\infty$  as  $t \rightarrow \infty$ . In the hard case,  $\psi(t)$  is undefined for  $t > t_0$ .

## 7.2 Flowchart

In the sequel, the set of  $\{t : k'(t) < 0\}$  is referred to as the *easy side* and its complement as the *hard side*. The details in the flowchart follow in Section 7.3 below.

### • INITIALIZATION:

1. Compute  $\lambda_1 = \lambda_1(A)$  and corresponding eigenvector  $v_1$ . If  $\lambda_1(A) < 0$ , shift  $A \leftarrow A - \lambda_1(A)I$ . ( $\lambda_1(A)\|x^*\|^2$  added back to objective value at end.)

If  $a^T v_1$  is small (near hard case), then deflate, i.e. add the vector  $y_1 = \begin{pmatrix} 0 \\ v_1 \end{pmatrix}$  to the set  $\mathcal{Y}$ .

2. Obtain bounds on  $q^*, \lambda^*$  and  $t^*$ .

If  $\lambda_1 > 0$ , **EXIT** if the optimum is an unconstrained minimizer.

3. Initialize parameters and the stopping criteria; this is based on the optimality conditions, duality gap, and intervals of uncertainty.

• **ITERATION LOOP:** (until convergence to the desired tolerance or until we find the solution is the unconstrained minimizer.)

## 1. FIND a NEW VALUE of $t$ .

- (a) Set  $t$  using Newton's method on  $k(t) - M_t$  if possible; otherwise set it to the the midpoint (default) of the interval of uncertainty for  $t$ .
- (b) If points from the hard and easy side are available:
  - i. Do TRIANGLE INTERPOLATION (Update upper bound on  $q_*$  and set  $t$ , if possible.)
  - ii. Do VERTICAL CUT (Update lower or upper bound for interval of uncertainty for  $t$ .)
- (c) Do INVERSE INTERPOLATION (Set  $t$ , if possible)

## 2. UPDATE

- (a) With new  $t$ , compute (with restarts using a previous eigenvector)  $\lambda = \lambda_1(D(t))$  and corresponding eigenvector  $y$  with  $y_0 \geq 0$ . (Use  $y$  orthogonal to the vectors in  $\mathcal{Y}$ , if possible.)
- (b) If  $\lambda > 0$  and  $y_0^2 > 1/(s^2 + 1)$  then the solution is the unconstrained minimizer. Use Conjugate Gradients and **EXIT**.
- (c) Update bounds on interval of uncertainty of  $q^*$ .
- (d)
  - i. If  $y_0$  is small, then deflate, i.e. add  $y$  to  $\mathcal{Y}$ .
  - ii. elseif  $t$  is on the easy side, update parameters. Take a primal step to the boundary if a hard side point exists.
  - iii. elseif  $t$  is on the hard side, update parameters. Take a primal step to the boundary from this hard side point.
- (e) Save new bounds and update stopping criteria.

## • END LOOP

## 7.3 Techniques Used in the Algorithm

### 7.3.1 Newton's Method on $k(t) - M_t = 0$

We use the upper and lower bounds on  $k(t)$  and the Newton type method presented in [22]. Note that Newton's method applied to  $k(t) - M_t = 0$  at  $t_c$

yields

$$t_+ = t_c - \frac{k(t_c) - M_t}{k'(t_c)} = \frac{(s^2 + 1)(t_c y_0^2(t_c) - \lambda_1(D(t_c))) + M_t}{(s^2 + 1)y_0(t_c)^2 - 1}.$$

One advantage of this method is that the second derivative  $k''$  is not needed. We use this iteration for appropriate choices of  $M_t$  in cases where the inverse iteration on  $\psi$  fails, i.e. if the hard case holds.

### 7.3.2 Triangle Interpolation

Given  $k(t)$ , if we have values of  $t$  from the easy and hard sides,  $t_e$  and  $t_h$ , then we try to find a better approximation  $t_{\text{new}}$  to the maximum of  $k(t)$  using a technique we call *triangle interpolation*, i.e. we find the coordinate of the point where the secant lines intersect. (We use a tangent line on the side where there is only one point.) We also obtain an upper bound  $q_{\text{up}}$  to  $q^* := k(t^*)$  from the point where the secant lines intersect.

### 7.3.3 Vertical Cut

Suppose we have two values of  $t$ ,  $t_e$  and  $t_h$ , with  $k(t_e) < k(t_h)$ . (A similar argument holds for the reverse inequality.) Then we can use the concavity of  $k$  to reduce the interval of uncertainty for  $t$ . We find the intersection of the horizontal line through  $(t_h, k(t_h))$  with the tangent line at the point  $(t_e, k(t_e))$ , i.e.

$$t_{\text{high}} = t_e + (k(t_h) - k(t_e))/k'(t_e),$$

where  $t_{\text{high}}$  is the upper bound on  $t^*$ .

### 7.3.4 Inverse Interpolation

We use (quadratic or linear) inverse interpolation on  $\psi(t) = 0$ , in the case that  $y_0(t) \neq 0$ . Since  $\psi(t)$  is a strictly decreasing function, we can consider its inverse function, say  $t(\psi)$ . We use (concave) quadratic interpolation when possible, i.e. suppose the points  $(\psi_i, t_i)$ ,  $i = 1, 2, 3$ . Then we solve the system

$$\begin{bmatrix} \psi_1^2 & \psi_1 & 1 \\ \psi_2^2 & \psi_2 & 1 \\ \psi_3^2 & \psi_3 & 1 \end{bmatrix} \begin{bmatrix} a \\ b \\ t_{\text{new}} \end{bmatrix} = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix}$$

and get the new estimate  $t_{\text{new}}$ . We use the top right  $2 \times 2$  system in the linear interpolation case.

In the hard case (case 2), inverse interpolation does not provide us with useful information on  $t^*$ . To avoid an erroneous step, we maintain  $\lambda$  on the correct side of  $\lambda_1(A)$ .

### 7.3.5 Recognizing an unconstrained minimum

Problem (TRS) has an inequality constraint, but generally the optimum lies on the boundary and we can solve the same problem with the equality constraint to get the optimum. This does not hold if and only if the matrix  $A$  is positive definite and the unconstrained minimum lies inside the trust region. The following theorem is the key to recognizing this case.

**Theorem 7.3.** *Let  $\bar{x}$  be a solution to  $(A - \lambda I)x = a$  with  $(A - \lambda I)$  positive semidefinite. If  $\lambda \leq 0$ , then  $\bar{x}$  is a solution to  $\min\{x^T Ax - 2a^T x : \|x\| \leq \|\bar{x}\|\}$ . If  $\lambda \geq 0$ , then  $\bar{x}$  is a solution to  $\min\{x^T Ax - 2a^T x : \|x\| \geq \|\bar{x}\|\}$ .*

**PROOF:** The first part follows from the necessary and sufficient optimality conditions and the second part follows easily knowing that the sign of  $\lambda$  plays no role in the positive semidefiniteness of the matrix  $A - \lambda I$  when proving these optimality conditions. ■

In our algorithm, we successively obtain solutions  $x_k$  to  $(A - \lambda_k I)x_k = a$  with  $A - \lambda_k I \succeq 0$ . Therefore each  $x_k$  is a solution to  $\min\{x^T Ax - 2a^T x : \|x\| = \|x_k\|\}$ . Checking the sign of the multiplier  $\lambda_k$  tells us if  $x_k$  is a solution to  $\min\{x^T Ax - 2a^T x : \|x\| \leq \|x_k\|\}$  or  $\min\{x^T Ax - 2a^T x : \|x\| \geq \|x_k\|\}$ . If the later case holds and  $\|x_k\| \leq s$ , then we know the unconstrained minimum lies in the trust region.

### 7.3.6 Shift and Deflate

Let us consider the case when the optimum is on the boundary of the ball. Once the smallest eigenvalue  $\lambda_1(A)$  is found in the initialization step, we can use the shift in Lemma 2.1 Item 3 and Lemma 2.3. Therefore, for simplicity, we can assume that  $\lambda_1(A) = 0$ .

During the algorithm we *deflate* eigenvectors  $y = (y_0 \ v^T)^T$  if  $y_0$  is small (essentially 0). This indicates that  $a^T v$  is small. We perturb  $a \leftarrow a - a^T v v$  and deflate using  $A \leftarrow A + \alpha v v^T$ .

### 7.3.7 Taking a Primal Step to the Boundary

The interpolation and heuristics are used to find a new point  $t$  and then a corresponding  $\lambda$  for the dual problem, i.e. they are used to take a dual step. Once the  $\lambda$  is found, we can find a corresponding primal point  $x(\lambda)$  for the primal problem. This point will be primal feasible if  $t$  (equivalently  $\lambda$ ) is on the hard side and it will be primal infeasible on the other (easy) side. In either case, we now show that we can take an inexpensive primal step, i.e. move to the boundary and improve the objective value. Thus, we get a primal-dual algorithm.

In the easy case, the step is motivated by the following lemma (for a proof see [13] or [12]).

**Lemma 7.1.** *Let  $0 < s_1 < s < s_2$  and*

$$x_h \in \operatorname{argmin}\{q(x) : \|x\| \leq s_1^2\}, \quad x_e \in \operatorname{argmin}\{q(x) : \|x\| \leq s_2^2\}.$$

*Suppose that  $\|x_h\| = s_1$ ,  $\|x_e\| = s_2$ ,  $x_h^T(x_e - x_h) \neq 0$  and the Lagrange multiplier  $\lambda_h$  for  $x_h$  satisfies  $A - \lambda_h I \succ 0$ . Furthermore, let  $m(\alpha) := q(x_h + \alpha(x_e - x_h))$ . Then  $m'(\alpha) \leq 0$ , for  $\alpha \in [0, 1]$ . ■*

Thus, we find  $\bar{\alpha}$  so that  $\|x_h + \bar{\alpha}(x_e - x_h)\| = s$ . We use two values,  $t_h$  and  $t_e$ , respectively on the hard side and the easy side, with

$$x_h := \frac{1}{y_0(t_h)}x(t_h), \quad x_e := \frac{1}{y_0(t_e)}x(t_e).$$

Then if  $x_h^T(x_e - x_h) \neq 0$  and  $\bar{\alpha}$  is defined as in the above lemma, taking a step to the boundary from  $x_h$  to  $x_h + \bar{\alpha}(x_e - x_h)$  will decrease the objective function. We assumed the easy case so that  $y_0(t_e) \neq 0$ . However, in the hard case (case 2)  $y_0(t_e) = 0$ . As in the MS algorithm, we take a step to the boundary. From  $x_h$ , we use an eigenvector  $z$  for the eigenvalue  $\lambda_1(A)$  as the direction to the boundary. This choice is motivated by Lemma 3.2 and the desire to make the quadratic form  $z^T(A - \lambda_1(D(t_h))I)z$  small. We take the step  $x_h + \tau z$  with  $\tau$  chosen to reduce the objective function and satisfy  $\|x_h + \tau z\| = s$ . The explicit expression for  $\tau$  is

$$\tau = \frac{s^2 - \|x_h\|^2}{x_h^T z + \operatorname{sgn}(x_h^T z) \sqrt{(x_h^T z)^2 + (s^2 - \|x_h\|^2)}},$$

where  $\text{sgn}(\cdot)$  equals 1 if its argument is nonnegative and -1 otherwise. Given a direction  $z$ , there are two values of  $\tau$  for which  $x_h + \tau z$  reaches the boundary. [25] proves that to improve the objective, we should pick the one with smallest magnitude.

## 8 Numerical Experiments

### 8.1 The Hard Case

We now provide numerical evidence that our modified RW Algorithm is better suited to handle the hard case (case 2) than the MS Algorithm. It is stated in [25] that the latter algorithm requires few iterations (2-3) in the hard case. However, this appears to hold only when the desired accuracy is low. Many more iterations are required when higher accuracy is desired. Our tests were done using MATLAB 6.1 on a SUNW Ultra-5\_10 with 1 GIG RAM.

Let  $q^*$  be the optimal objective value of TRS and  $\tilde{q}$  be an approximation for  $q^*$ . The Moré-Sorensen Algorithm returns an approximate solution that satisfies  $\tilde{q} \leq (1 - \sigma)^2 q^*$ , where  $0 \leq \sigma < 1$  is an input to the algorithm, and the approximate solution of the RW Algorithm satisfies  $\tilde{q} \leq \frac{1}{1+2\text{dgaptol}} q^*$ , where  $\text{dgaptol}$  is the desired relative duality gap tolerance. Hence, to get equivalent accuracy, we choose  $\sigma = 1 - \sqrt{\frac{1}{1+2\text{dgaptol}}}$ .

We used randomly generated sparse hard case (case 2(ii)) trust region subproblems, where the density is order  $1/(20n \log n)$ . The tolerance parameter  $\text{dgaptol}$  was set to  $10^{-12}$ . Each row in Table 8.1 gives the average number of iterations and cpu time for 10 problems of size  $n$ . We could not go beyond  $n = 640$  for the MS Algorithm due to the large computation times that arise. The results, given in Table 8.1 illustrate the improved performance in both the number of iterations and the computation time.

Note that the GLTR Algorithm does not appear in the above comparison since the algorithm was not designed to handle this case.

### 8.2 RW and GLTR Algorithms in TR Framework

In [17], the the GLTR method is stopped early after a limited extra number of iterations, say  $N$ , once the solution is known to lie on the boundary of the trust region. More precisely, the algorithm stops if the subspace  $S$  in

Dim. n	MS iters.	RW iters.	MS cpu	RW cpu
40	36.4	6.4	0.79	0.55
80	34.4	7.6	1.0	0.57
160	39.2	7.2	6.49	0.61
320	33.8	7.4	23.36	0.77
640	37.8	5.0	149.36	0.78
1280	-	7.6	-	2.06
2560	-	5.0	-	3.18

Table 8.1: Modified-RW and MS Algorithms; hard case (case 2(ii)).

(4.1) is increased in dimension by  $N$  once the solution is known to be on the boundary of the trust region and problems of the type (4.1) are solved. The reason for limiting the size of the subspaces  $S$  once the boundary has been reached is motivated by the fact that the authors in [17] question whether high accuracy is needed for the TRS within a trust region framework. We argue that increased accuracy is needed for TRS just as for the solution of the Newton equation when using inexact Newton methods, as the iterates approach a stationary point, see e.g. [7, 24].

For the upcoming test problems, we used  $N = 2, 6$  and  $n$ , where  $n$  is the problem dimension. Our test problems are of the following form

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{x^T A x}{x^T B x}, \quad (8.1)$$

where  $B$  is a positive definite matrix and  $A$  and  $B$  are generated randomly. The minimum is attained at  $x^*$ , a generalized eigenvector corresponding to the smallest eigenvalue of the generalized eigenvalue problem

$$Ax = \lambda Bx.$$

The optimal value is equal to  $\lambda_1 (B^{-1/2} A B^{-1/2})$ .

To solve each problem we used the trust region method described by Algorithm 8.1 below on the same machine as in Section (8.1). In this algorithm,  $f$  represents the function to be minimized,  $x_j$  is an approximation of a minimizer after  $j$  iterations and  $s_j$  is the radius of the trust region at iteration  $j$ .

**Algorithm 8.1.** (Trust Region Method (TR))

1. Given  $x_j$  and  $s_j$ , calculate  $\nabla f(x_j)$  and  $\nabla^2 f(x_j)$ . Stop if

$$\frac{\|\nabla f(x_j)\|}{1 + |f(x_j)|} < \text{gradtol}. \quad (8.2)$$

2. Find  $\delta_j$  to a given tolerance in the TRS

$$\begin{aligned} \delta_j \in \arg \min & \quad q_j(\delta) := \nabla f(x_j)^T \delta + \frac{1}{2} \delta^T \nabla^2 f(x_j) \delta \\ \text{s.t.} & \quad \|\delta\|^2 \leq s_j^2. \end{aligned} \quad (8.3)$$

3. Evaluate  $r_j = \frac{f(x_j) - f(x_j + \delta_j)}{q_j(0) - q_j(\delta_j)}$ .

4. (a) If  $r_j > 0.95$ , set  $s_{j+1} = 2s_j$  and  $x_{j+1} = x_j + \delta_j$ .

(b) If  $0.01 \leq r_j < 0.95$ , set  $s_{j+1} = s_j$  and  $x_{j+1} = x_j + \delta_j$ .

(c) If  $r_j < 0.01$ , set  $s_{j+1} = 0.5s_j$  and  $x_{j+1} = x_j$ .

Except for the stopping criteria (8.2) which has been scaled here, this algorithm is Algorithm 6.1 in [17]. We chose  $x_0$  randomly and have fixed  $s_0 = 1$ ,  $\text{gradtol} = 10^{-2}$ . We ran five random problems for each problem size  $n = 20, 25$  and  $30$ , where  $n$  is the size of the square matrices  $A$  and  $B$ . If the RW Algorithm solves (8.3) and the solution is on the boundary of the trust region, we stop if the duality gap,  $\text{dgaptol}$ , (between TRS and (5.7)) satisfies

$$\sqrt{\text{dgaptol}} \leq \min\{0.1, \max\{10^{-8}, 10^{-5} \|\nabla f(x_j)\|\}\}. \quad (8.4)$$

Otherwise, the solution is in the interior and we stop with an approximate solution  $\delta_j$  which satisfies

$$\|\nabla f(x_j) + \nabla^2 f(x_j) \delta_j\| \leq \min\{0.1, \max\{10^{-8}, 10^{-5} \|\nabla f(x_j)\|\}\}. \quad (8.5)$$

If the GLTR Algorithm is used, we stop within this algorithm if  $N$  iterations have been done after knowing the solution lies on the boundary of the trust region (see [17]) or if

$$\|(A - \lambda_k I)x_k - a\| < \min\{0.1, \max\{10^{-8}, 10^{-5} \|\nabla f(x_j)\|\}\} \quad (8.6)$$

(see (6.3)) or (8.5) is satisfied, depending if the solution is on the boundary of the trust region or not.



Using (8.4) and (8.6) yields approximately the same accuracy in terms of the duality gap. Recall from Section 6.2 that  $\|(A - \lambda_k I)x_k - a\|$  is an approximation for the square root of the duality gap between TRS and its dual (6.1). The stopping criteria (8.4) and (8.6) are set to reflect this relationship.

For each problem, we give the number of iterations (*iter*) taken by the trust region method 8.1. If the GLTR Algorithm is used to solve the TRS (8.3), we give as well the number of iterations within Algorithm 8.1 where the GLTR Algorithm failed to solve (8.3) because it was unable to solve the restricted problem (4.1). This last output (*hc2*) is an indicator of the (almost) hard case (case 2).

Algorithm used for solving the TRS (8.3)							
	RW	GLTR with N=2		GLTR with N=6		GLTR with N=n	
problem	iter	iter	hc2	iter	hc2	iter	hc2
1	16	36	5	25	0	16	0
2	33	22	0	55	19	48	16
3	50	21	0	15	0	52	16
4	41	39	8	45	16	32	3
5	25	45	10	19	0	25	0

Table 8.2: RW and GLTR; TR framework; size n=20.

Algorithm used for solving the TRS (8.3)							
	RW	GLTR with N=2		GLTR with N=6		GLTR with N=n	
problem	iter	iter	hc2	iter	hc2	iter	hc2
1	31	48	17	34	5	32	2
2	38	26	0	27	1	32	2
3	22	25	0	22	0	22	0
4	20	39	4	31	1	20	0
5	25	26	0	22	0	22	0

Table 8.3: RW and GLTR; TR framework; size n=25.

Algorithm used for solving the TRS (8.3)							
	RW	GLTR with N=2		GLTR with N=6		GLTR with N=n	
problem	iter	iter	hc2	iter	hc2	iter	hc2
1	61	14	0	36	7	57	24
2	38	50	17	35	2	37	6
3	34	22	0	27	1	45	17
4	34	19	14	25	0	36	8
5	36	38	7	26	0	32	3

Table 8.4: RW and GLTR; TR framework; size n=30.

The results given in Tables 8.2, 8.3, 8.4 show that the (almost) hard case (case 2) occurs in many problems and Algorithm 8.1, using the RW Algorithm for (8.3), takes fewer iterations to converge compared to the GLTR Algorithm. This is independent of N. This suggests that handling the hard case (case 2) should be an essential feature for a robust trust region method. However, we reach the same conclusions mentioned in [17] when the hard case (case 2) does not occur. We observe the surprising fact that inexact solutions may indeed lead to less iterations in some cases. This may be due to the fact that the trust region becomes inactive early and Newton’s method takes over. This clearly requires more study.

As we may expect, when the hard case (case 2) does not occur and  $N=n$ , a trust region method, using either the RW Algorithm or the GLTR Algorithm to solve the TRS (8.3), takes more or less the same number of iterations. This should be the case since we are asking for the same accuracy in (8.4) and (8.6).

### 8.3 Accuracy of TRS in a Trust Region Method

Inexact Newton methods can obtain q-superlinear and even q-quadratic convergence rates if the accuracy of the Newton equation increases appropriately as the iterates approach a stationary point, e.g. [27]. Trust region methods such as Algorithm 8.1 are expected to reduce to Newton’s method asymptotically, i.e. the trust region constraint is expected to become inactive for most problems. This happens for example when the second order sufficient

optimality conditions (positive definite Hessian) holds at the limit point. In either case, i.e. whether or not the trust region constraint becomes inactive, the accuracy for solving TRS must increase as we approach the stationary point. We now investigate the number of iterations the trust region Algorithm 8.1 takes to solve an unconstrained minimization problem as the accuracy of the solutions of the TRS (8.3) varies using MATLAB 6.5 on a Sun Fire 280R (UltraSPARC-III) with 2 GIGs RAM.

We solve the problem  $\min_{x \in \mathbb{R}^n} f(x)$  to accuracy given by varying values of  $gradtol$  in the inequality (8.2). The TRS (8.3) is solved using the modified RW Algorithm to accuracy

$$dgaptol = \max\{tol, 10^{-6} \min\{1, \frac{\|\nabla f(x_j)\|}{1 + |f(x_j)|}\}^{1/2}\}, \quad (8.7)$$

for varying values of  $tol$ . We also terminate algorithm 8.1 if more than 1000 iterations are necessary or if the trust region radius  $s_j$  becomes smaller than  $10^{-10}$  (this case is indicated by  $s \cong 0$ ). The results for our two examples are given in Tables 8.5 and 8.6, where the entries are the number of iterations taken by Algorithm 8.1.

	gradtol						
tol	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$	$10^{-8}$	$10^{-9}$
$10^{-1}$	42	117	$\geq 1000$	$\geq 1000$	$\geq 1000$	$\geq 1000$	$\geq 1000$
$10^{-3}$	15	21	50	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$
$10^{-5}$	12	16	18	88	$s \cong 0$	$s \cong 0$	$s \cong 0$
$10^{-7}$	12	22	38	52	67	$s \cong 0$	$s \cong 0$
$10^{-9}$	12	16	18	22	24	50	$s \cong 0$
$10^{-11}$	12	22	38	53	68	83	96

Table 8.5: TR on  $f(x) = \sin(x_1 - 1) + \sum_{i=2}^{1000} 100 \sin(x_i - x_{i-1}^2)$  (see [4]).

From these results we observe two things. First, as the accuracy on the norm of the gradient is decreasing, the trust region method (8.1) using low accuracy solutions for the TRS is eventually outperformed by higher accuracy solutions. Second, the results of Table 8.5 indicate, for a fixed tolerance  $gradtol$  on the norm of the gradient, that more accurate solution of the TRS

	gradtol						
tol	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$	$10^{-8}$	$10^{-9}$
$10^{-1}$	9	10	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$
$10^{-3}$	4	5	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$
$10^{-5}$	4	5	6	$s \cong 0$	$s \cong 0$	$s \cong 0$	$s \cong 0$
$10^{-7}$	4	5	6	7	8	$s \cong 0$	$s \cong 0$
$10^{-9}$	4	5	6	7	8	$s \cong 0$	$s \cong 0$
$10^{-11}$	4	5	6	7	8	$s \cong 0$	$s \cong 0$

Table 8.6: TR on  $f(x) = \sum_{i=1}^{1000} 6x_i^2 + \sum_{i=503}^n x_i + \sum_{i=1}^{998} (x_i x_{i+2} - 4x_i x_{i+1}) - 3x_1 + x_2 + x_{499} - 3x_{500} + 4x_{501}$  (see [23]).

does not necessarily imply fewer iterations. Therefore, for robustness and as for inexact Newton methods, it is beneficial to use increased accuracy for the TRS when approaching the minimum of the objective function  $f$ .

## 8.4 Large Sparse TRS

The results of Table 8.7 show how we used the RW Algorithm to solve problems of size  $n = 100,000$  of different density. Precisely, each row in this table corresponds to the density of the problems (for example, if the density is  $10^{-6}$ , then at most  $n^2 \times 10^{-6}$  entries in the matrix A are non-zero) and each column to the value of the parameter  $dgaptol$ . In each entry of the table we give the average taken over 5 random trust region subproblems of the computation time (cpu) in seconds, the number of matrix-vector multiplications (mvm) and the number of iterations (it) taken by the RW Algorithm to find an approximate solution. We have been using for this section MATLAB 6.1 on a Pentium III with 4 GIGs RAM.

As we may expect, the computation time and the matrix-vector multiplications increase as the density increases and the duality gap tolerance decreases. Furthermore, considering the reasonable length of the computation time taken to solve such TRS, we conclude that it is now within our reach to use trust region methods to minimize functions with hundreds of thousands of variables assuming the Hessian has a sparse structure. In Table 8.8, we considered a TRS of size  $n = 10^6$  with 11 million nonzeros in

density	dgaptol		
	$10^{-12}$	$10^{-10}$	$10^{-8}$
$10^{-8}$	cpu : 18.2 mvm : 185.6 iter : 6.2	cpu : 14.5 mvm : 150.0 iter : 5.4	cpu : 13.5 mvm : 140.0 iter : 4.8
$10^{-6}$	cpu : 21.1 mvm : 210.0 iter : 6.4	cpu : 19.2 mvm : 196.0 iter : 5.4	cpu : 20.0 mvm : 204.0 iter : 5.6
$10^{-4}$	cpu : 91.7 mvm : 341.6 iter : 5.8	cpu : 78.8 mvm : 294.0 iter : 5.0	cpu : 76.0 mvm : 276.0 iter : 5.6

Table 8.7: Modified RW Algorithm on TRS; n=100,000.

the sparse matrix A. The results show that higher accuracy solutions require minimal extra iterations of the RW Algorithm.

	dgaptol					
	$10^{-2}$	$10^{-4}$	$10^{-6}$	$10^{-8}$	$10^{-10}$	$10^{-12}$
cpu	985.1	985.1	1082.8	1082.8	1183.6	1282.9
mvm	240	240	260	260	280	300
iter	3	3	4	4	5	6

Table 8.8: Modified RW Algorithm on TRS; n=1,000,000.

## 9 Conclusion

In this paper we have studied the trust region subproblem, TRS, with emphasis on robustness and solving large sparse problems. We focused on three dual based algorithms: the classical MS algorithm and the recent RW and GLTR Algorithms, designed to solve large and sparse TRS.

We also studied many duals to TRS which can be formulated as semidefinite programs. We have seen how SDP arises naturally for TRS and provides a clear and simple unifying analysis between the different algorithms. In addition, this framework provides insights to the strengths and weaknesses of the algorithms.

In addition, we presented a modified/enhanced RW algorithm with new heuristics and techniques, in particular for taking a primal step to the boundary. However, the main improvement came from a new way of treating the (near) hard case based on Lemma 2.1. Surprisingly, the Lemma shows that for each TRS, it is possible to consider an equivalent TRS where the hard case (case 2) does not occur.

Our final section included numerics which showed the advantage of using the modified RW Algorithm over the MS Algorithm in treating the hard case when high accuracy approximations are needed. We have also shown that handling the hard case in the TRS within a trust region method may have an impact on the total number of iterations if the hard case occurs frequently enough. Thus, the robustness of a TRS Algorithm is indeed an important feature, in particular when the trust region constraint stays active close to the optimal solution. Finally, we showed it is possible to solve large sparse TRS to high accuracy with hundreds of thousands of variables in a small number of iterations.

## References

- [1] S. J. BENSON, Y. YE, and X. ZHANG. Solving large-scale sparse semidefinite programs for combinatorial optimization. *SIAM J. Optim.*, 10(2):443–461 (electronic), 2000.
- [2] A. BEN-TAL and M. TEBoulLE. Hidden convexity in some nonconvex quadratically constrained quadratic programming. *Math. Programming*, 72(1, Ser. A):51–63, 1996.
- [3] Å. BJÖRCK. *Numerical methods for least squares problems*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996.
- [4] A. BOURIACHA. private communication.

- [5] A.R. CONN, N.I.M. GOULD, and P.L. TOINT. *Trust-region methods*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2000.
- [6] D.C. SORENSEN. Minimization of a large-scale quadratic function subject to a spherical constraint. *SIAM Journal on Optimization*, 7(1):141–161, 1997.
- [7] R.S. DEMBO and T. STEihaug. Truncated Newton algorithms for large scale unconstrained optimization. *Mathematical Programming*, 26(72):190–212, 1983.
- [8] J.W. DEMMEL. *Applied numerical linear algebra*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997.
- [9] A.V. FIACCO. *Introduction to Sensitivity and Stability Analysis in Non-linear Programming*, volume 165 of *Mathematics in Science and Engineering*. Academic Press, 1983.
- [10] R. FLETCHER. *Practical Methods of Optimization*. John Wiley and Sons, second edition, 1987.
- [11] G.E. FORSYTHE and G.H. GOLUB. On the stationary values of a second-degree polynomial on the unit sphere. *J. Soc. Indust. Appl. Math.*, 13:1050–1068, 1965.
- [12] C. FORTIN. A survey of the trust region subproblem within a semidefinite framework. Master’s thesis, University of Waterloo, 2000.
- [13] C. FORTIN and H. WOLKOWICZ. A survey of the trust region subproblem within a semidefinite framework. Technical Report CORR 2002-22, University of Waterloo, Waterloo, Canada, 2002. URL:[orion.math.uwaterloo.ca:80/~hwolkowi/henry/reports/ABSTRACTS.html#surveytrs](http://orion.math.uwaterloo.ca:80/~hwolkowi/henry/reports/ABSTRACTS.html#surveytrs).
- [14] D.M. GAY. Computing optimal locally constrained steps. *SIAM J. Sci. Statist. Comput.*, 2:186–197, 1981.
- [15] D.M. GAY. Computing optimal locally constrained steps. *SIAM Journal on Scientific and Statistical Computing*, 2(2):186–197, 1981.

- [16] E.G. GOL'STEIN. *Theory of Convex Programming*. American Mathematical Society, Providence , RI, 1972.
- [17] N.I.M. GOULD, S. LUCIDI, M. ROMA, and P.L. TOINT. Solving the trust-region subproblem using the lanczos method. *SIAM Journal on Optimization*, 9(2):504–525, 1999.
- [18] W.W. HAGER. Minimizing a quadratic over a sphere. Technical report, University of Florida, Gainesville, Fa, 2000.
- [19] C. HELMBERG and F. RENDL. A spectral bundle method for semidefinite programming. *SIAM Journal on Optimization*, 10(3):673 – 696, 2000.
- [20] A.E. HOERL and R.W. KENNARD. Ridge regression: Biased estimation of nonorthogonal problems. *Technometrics*, 12:55–67, 1970.
- [21] R.A. HORN and C.R. JOHNSON. *Matrix Analysis*. Cambridge University Press, 1987.
- [22] Y. LEVIN and A. BEN-ISRAEL. The newton bracketing method for convex minimization. *Comput. Optimiz. Appl.*, 21:213–229, 2002.
- [23] Y. LIN and J. PANG. Iterative methods for large convex quadratic programs: a survey. *SIAM Journal on Control and Optimization*, 25, 1987.
- [24] J.L. MORALES and J. NOCEDAL. Automatic preconditioning by limited memory quasi-Newton updating. *SIAM J. Optim.*, 10(4):1079–1096 (electronic), 2000.
- [25] J.J. MORÉ and D.C. SORESENSEN. Computing a trust region step. *SIAM Journal on Scientific and Statistical Computing*, 4(3):553–572, 1983.
- [26] S.G. NASH and J. NOCEDAL. A numerical study of the limited memory BFGS method and the truncated-Newton method for large scale optimization. *SIAM J. Optim.*, 1(3):358–372, 1991.
- [27] J. NOCEDAL and S.J. WRIGHT. *Numerical optimization*. Springer-Verlag, New York, 1999.



- [28] C. REINSCH. Smoothing by spline functions. *Numerische Mathematik*, 10:177–183, 1967.
- [29] C. REINSCH. Smoothing by spline functions ii. *Numerische Mathematik*, 16:451–454, 1971.
- [30] M.D. REINSCH. An algorithm for minimization using exact second derivatives. Technical Report 515, Harwell Laboratory, Harwell, Oxfordshire, England, 1973.
- [31] F. RENDL and H. WOLKOWICZ. A semidefinite framework for trust region subproblems with applications to large scale minimization. *Mathematical Programming Series B*, 77(2):273–299, 1997.
- [32] D.C. SORENSEN. Newton’s method with a model trust region modification. *SIAM Journal on Numerical Analysis*, 19(2):409–426, 1982.
- [33] T. STEihaug. The conjugate gradient method and trust regions in large scale optimization. *SIAM Journal on Numerical Analysis*, 20(3), 1983.
- [34] R.J. STERN and H. WOLKOWICZ. Indefinite trust region subproblems and nonsymmetric eigenvalue perturbations. *SIAM Journal on Optimization*, 5(2):286–313, 1995.
- [35] P.D. TAO and L.T.H. AN. Difference of convex functions optimization algorithms (dca) for globally minimizing nonconvex quadratic forms on euclidean balls and spheres. *Oper. Res. Lett*, 19(5):207–216, 1996.
- [36] A.N. TIKHONOV and V.Y. ARSENIN. *Solutions of Ill-Posed Problems*. V.H. Winston & Sons, John Wiley & Sons, Washington D.C., 1977. Translation editor Fritz John.
- [37] P.L. TOINT. *Towards an efficient sparsity exploiting Newton method for minimization*, volume Sparse Matrices and Their Uses. I.S.Duff, academic press edition, 1981. pp.57-88.
- [38] R.J. VANDERBEI. *Linear Programming: Foundations and Extensions*. Kluwer Acad. Publ., Dordrecht, 1998.

- [39] H. WOLKOWICZ, R. SAIGAL, and L. VANDENBERGHE, editors. *HANDBOOK OF SEMIDEFINITE PROGRAMMING: Theory, Algorithms, and Applications*. Kluwer Academic Publishers, Boston, MA, 2000. xxvi+654 pages.
- [40] S. WRIGHT. *Primal-Dual Interior-Point Methods*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, Pa, 1996.
- [41] Y. YE. Combining binary search and Newton's method to compute real roots for a class of real functions. *Journal of Complexity*, 10:271–280, 1994.