

Globally Solving the Trust Region Subproblem Using Simple First-Order Methods

Amir Beck* and Yakov Vaisbourd*

October 2, 2017

Abstract

We consider the trust region subproblem which is given by a minimization of a quadratic, not necessarily convex, function over the Euclidean ball. Based on the well-known second-order necessary and sufficient optimality conditions for this problem, we present two sufficient optimality conditions defined solely in terms of the primal variables. Each of these conditions corresponds to one of two possible scenarios that occur in this problem, commonly referred to in the literature as the presence or absence of the “hard case”. We consider a family of first-order methods, which includes the projected and conditional gradient methods. We show that any method belonging to this family produces a sequence which is guaranteed to converge to a stationary point of the trust region subproblem. Based on this result and the established sufficient optimality conditions, we show that convergence to an optimal solution can be also guaranteed as long as the method is properly initialized. In particular, if the method is initialized with the zeros vector and reinitialized with a randomly generated feasible point, then the best of the two obtained vectors is an optimal solution of the problem in probability 1.

1 Introduction

The *trust region subproblem* (TRS) is given by

$$\min \{q(\mathbf{x}) : \mathbf{x} \in B\}, \quad (\text{TRS})$$

with

$$q(\mathbf{x}) := \mathbf{x}^T \mathbf{A} \mathbf{x} - 2\mathbf{b}^T \mathbf{x} \quad \text{and} \quad B := \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|^2 \leq 1\},$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{b} \in \mathbb{R}^n$ and $\|\cdot\|$ stands for the Euclidean ℓ_2 -norm. Though this problem emerges in various applications, it is mostly known due to its central role in the trust region methodology for solving optimization problems, see for example the book

*School of Mathematical Sciences, Tel-Aviv University, Tel-Aviv, Israel
(becka@tauex.tau.ac.il, yakov.vaisbourd@gmail.com)

[7] and references therein.

In spite of the fact that problem (TRS) is not necessarily convex, its complexity is known to be polynomial [26]. Moreover, problem (TRS) possesses necessary and sufficient second-order optimality conditions which have motivated much of the methods designed to solve it. Probably one of the most popular among those methods is the one suggested in [17], where a safeguarded Newton method is used for finding the root of a secular equation emerging from the aforementioned optimality conditions. This method is applicable for small to medium-sized problems due to the Cholesky factorization which is performed at each iteration.

In the last decades, schemes that maintain modest computational and storage requirements were suggested for addressing large-scale instances of problem (TRS). The conjugate gradient was one of the first iterative methods proposed independently in [22] and [24] for this task. This method succeeds in retrieving the optimal solution only if the latter resides in the interior of the feasible set. Otherwise, the scheme breaks down at some stage and the method returns only an approximate solution. In order to continue the process of minimizing over the Krylov subspace from the point in which the conjugate gradient method breaks down, in [12] the authors suggested an adaptation of the Lanczos method that relies on solving at each iteration a much easier TRS, one that involves a tridiagonal matrix in the quadratic term and with a significantly lower dimension. Nevertheless, among other disadvantages of this method, it cannot guarantee to produce the optimal solution either. Trying to overcome these pitfalls, other methods were suggested in the literature such as the sequential subspace method (SSM) [13] which was followed by a variation called phased SSM in [8]. Though these methods can guarantee convergence to the global optimum, to do so they require an estimate of the eigenvector that corresponds to the minimal eigenvalue.

Some schemes for finding the optimal solution of problem (TRS) based on its parameterized eigenvalue reformulation were suggested in the literature. In [21] the authors proposed to examine the close relations between the minimal eigenvector characterization of such a reformulation and the optimality conditions of problem (TRS). Based on this relation, an inverse interpolation scheme for iteratively tuning the parameter is suggested where at each iteration an eigenvalue computation is performed. Various improvements that include a better treatment for the so-called “hard case” (see Subsection 2.3) were suggested in [14, 19]. In [18], a similar parametrized eigenvalue approach was developed using different arguments involving duality theory for SDPs and a homogenization technique.

Under the “difference of convex functions” framework, the projected gradient algorithm was proposed in [23] for solving problem (TRS). In general, there are no guarantees that the sequence generated by the projected gradient method will converge to the global optimal solution, and thus a restarting technique that produces a point with a smaller objective value backed up with a theoretical upper bound on the number of required restarts is suggested in order to avoid convergence to a non-optimal stationary point. Such a restart procedure relies on the availability of an eigenvector that corresponds to the minimal eigenvalue.

In this paper our goal is to show that problem (TRS) can be globally solved by simple first-order methods such as the projected and conditional gradient methods, under proper initialization that can be trivially computed. The dominant computation at each iteration of the class of methods that we consider is multiplication of \mathbf{A} with a given vector, and there

is no need to solve or factorize systems of equations or to compute/approximate eigenvectors of \mathbf{A} .

The paper’s layout is as follows. Section 2 reviews several important mathematical facts that will be pertinent to the analysis throughout the paper: characterizations of stationary and optimality conditions, as well as a discussion on the so-called “easy” and “hard” cases of problem (TRS). Section 3 introduces two sufficient optimality conditions, one for each case (“easy” and “hard”). These conditions, which are expressed solely in terms of the primal variables, will be instrumental for the proof of convergence to a global optimal solution of the class of “first-order conic methods” (FOCM) introduced in Section 4. This class of algorithms includes the projected and conditional gradient methods as special instances. We show that the entire sequence generated by an FOCM converges regardless of the choice of the starting point. Moreover, it is shown that in the “easy case”, an FOCM converges to the optimal solution if initialized with the zeros vector, and in the “hard case”, it converges to the global optimal solution as long as the objective function is non-homogenous and the starting point is randomly generated by a continuous distribution whose support is B . A direct consequence is that an optimal solution can be found with probability 1 (even if it is unclear whether the “easy” or “hard” case hold) by employing an FOCM twice—each with a different initialization.

Notation. Vectors are denoted by boldface lowercase letters, e.g., \mathbf{y} , and matrices by boldface uppercase letters, e.g., \mathbf{B} . Given a matrix $\mathbf{B} \in \mathbb{R}^{m \times n}$ and vector $\mathbf{y} \in \mathbb{R}^n$, $\|\mathbf{B}\|_2$ denotes the spectral norm of \mathbf{B} , and $\|\mathbf{y}\|$ denotes the ℓ_2 norm of \mathbf{y} . The vector $|\mathbf{y}|$ stands for the element-wise absolute value of \mathbf{y} . The matrix \mathbf{I} is the identity matrix whose dimension will be clear from the context. The notation $\mathbf{B} \succeq \mathbf{0}$ means that \mathbf{B} is positive semidefinite. The set $\mathbb{N} = \{0, 1, 2, \dots\}$ is the set of natural numbers.

2 Mathematical Preliminaries

2.1 Stationarity Conditions

We begin by recalling (see for example [2, Section 9.1]) that a point $\bar{\mathbf{x}}$ is a *stationary point* of problem (TRS) if it does not possess any feasible descent directions:

$$\nabla q(\bar{\mathbf{x}})^T(\mathbf{x} - \bar{\mathbf{x}}) \geq 0, \text{ for any } \mathbf{x} \in B.$$

It is well-known that the above stationarity condition can be expressed in a more explicit way in the case of problem (TRS).

Theorem 2.1 ([2, Example 9.6]). *A point $\bar{\mathbf{x}}$ is a stationary point of problem (TRS) if and only if there exists some $\lambda_{(\bar{\mathbf{x}})} \geq 0$ such that the following conditions hold:*

$$(\mathbf{A} + \lambda_{(\bar{\mathbf{x}})}\mathbf{I})\bar{\mathbf{x}} = \mathbf{b}, \tag{2.1}$$

$$(\|\bar{\mathbf{x}}\|^2 - 1)\lambda_{(\bar{\mathbf{x}})} = 0. \tag{2.2}$$

A nonnegative number $\lambda_{(\bar{\mathbf{x}})}$ satisfying (2.1) and (2.2) is called a *Lagrange multiplier associated with $\bar{\mathbf{x}}$* . We also note that condition (2.2) is known as the “complementary slackness” condition. The following simple lemma shows that each stationary point has a *unique* Lagrange multiplier. This uniqueness result on the associated Lagrange multiplier shows that the notation “ $\lambda_{(\bar{\mathbf{x}})}$ ”, that will be used throughout the paper, is well-defined.

Lemma 2.2. *Any stationary point $\bar{\mathbf{x}}$ of problem (TRS) admits a unique Lagrange multiplier $\lambda_{(\bar{\mathbf{x}})} \geq 0$.*

Proof. Assume by contradiction that $\bar{\mathbf{x}}$ admits two Lagrange multipliers $\lambda_{(\bar{\mathbf{x}})}^1$ and $\lambda_{(\bar{\mathbf{x}})}^2$ such that $\lambda_{(\bar{\mathbf{x}})}^1 \neq \lambda_{(\bar{\mathbf{x}})}^2$. Then by (2.1)

$$\lambda_{(\bar{\mathbf{x}})}^1 \bar{\mathbf{x}} = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}} = \lambda_{(\bar{\mathbf{x}})}^2 \bar{\mathbf{x}}.$$

Since $\lambda_{(\bar{\mathbf{x}})}^1 \neq \lambda_{(\bar{\mathbf{x}})}^2$ we obtain that $\bar{\mathbf{x}} = \mathbf{0}$, but in this case, condition (2.2) implies that $\lambda_{(\bar{\mathbf{x}})}^1 = \lambda_{(\bar{\mathbf{x}})}^2 = 0$, leading to a contradiction. \square

It is interesting to note that the converse of Lemma 2.2 does not hold true in general, meaning that it might happen that the same Lagrange multiplier will be associated with different stationary points. For example, if $q(\mathbf{x}) := 0$ ($\mathbf{A} = \mathbf{0}$, $\mathbf{b} = \mathbf{0}$), then all the feasible points are stationary points, and $\lambda = 0$ is the Lagrange multiplier associated with all of them.

Obviously, stationarity is a necessary optimality condition [5]. Combining this fact with Lemma 2.2 yields the following known result.

Theorem 2.3. *Any optimal solution $\bar{\mathbf{x}}$ of problem (TRS) is also a stationary point, meaning that $(\bar{\mathbf{x}}, \lambda_{(\bar{\mathbf{x}})})$ satisfies conditions (2.1) and (2.2).*

Problem (TRS) admits necessary and sufficient second-order optimality conditions which are summarized in the following theorem (see e.g., [10, 20]). The conditions state that a point is optimal if and only if it is stationary and satisfies a certain second-order optimality condition.

Theorem 2.4 ([20, Lemmas 2.4 and 2.8]). *Let $\bar{\mathbf{x}}$ be a feasible point of problem (TRS). Then $\bar{\mathbf{x}}$ is an optimal solution of problem (TRS) if and only if it is a stationary point such that*

$$\mathbf{A} + \lambda_{(\bar{\mathbf{x}})} \mathbf{I} \succeq \mathbf{0}. \tag{2.3}$$

2.2 Spectral Characterization of Optimality and Stationarity

Throughout the analysis in the paper we will consider the stationarity and optimality conditions in terms of the eigenvalues and eigenvectors of \mathbf{A} . Let $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T$ be a spectral decomposition of \mathbf{A} such that $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_n)$ with $d_1 \leq d_2 \leq \dots \leq d_n$ being the increasingly sorted eigenvalues of \mathbf{A} and \mathbf{U} being an orthonormal matrix whose columns are associated eigenvectors of \mathbf{A} . **The notations of the matrix \mathbf{U} and $d_1 \leq d_2 \leq \dots \leq d_n$ will be fixed throughout the paper.**

In terms of the spectral decomposition, Theorems 2.1 and 2.4 as well as Lemma 2.2 can be rewritten in the following compact form.

Theorem 2.5. (a) A point $\bar{\mathbf{x}}$ is a stationary point of problem (TRS) if and only if there exists a unique $\lambda_{(\bar{\mathbf{x}})} \geq 0$ such that the following conditions hold:

$$(d_i + \lambda_{(\bar{\mathbf{x}})})(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \mathbf{b})_i, \quad i = 1, 2, \dots, n. \quad (2.4)$$

$$(\|\bar{\mathbf{x}}\|^2 - 1)\lambda_{(\bar{\mathbf{x}})} = 0. \quad (2.5)$$

(b) A feasible point $\bar{\mathbf{x}}$ of problem (TRS) is optimal if and only if it is a stationary point and satisfies

$$\lambda_{(\bar{\mathbf{x}})} \geq -d_1.$$

2.3 The “Hard Case”

The characterization of an optimal solution given by the three conditions (2.1), (2.2) and (2.3) motivated the development of some efficient methods for solving problem (TRS), see for example [17, 20] and the discussion in the introduction. These methods, as well as many algorithms that appear in the literature, provide a special treatment for the so-called “hard case” which we now describe. Denote the eigenspace of \mathbf{A} corresponding to the smallest eigenvalue, which we denote as d_1 , by

$$\mathcal{E}_1 := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = d_1\mathbf{x}\}.$$

There are two classes of trust region subproblems that are frequently analyzed separately. The classes are determined according to whether or not the vector \mathbf{b} is orthogonal to the subspace \mathcal{E}_1

- “easy case” $\mathbf{b} \notin \mathcal{E}_1$.
- “hard case” $\mathbf{b} \perp \mathcal{E}_1$.

We will often use the following notation for the set of indices of eigenvalues corresponding to the minimal eigenvalue:

$$E_1 := \{i \in \{1, 2, \dots, n\} : d_i = d_1\}. \quad (2.6)$$

With the above notation, the condition describing the “easy case” can also be written as

$$(\mathbf{U}^T \mathbf{b})_i \neq 0 \text{ for some } i \in E_1.$$

It is well known that in the “easy case” there exists a unique optimal solution to problem (TRS).

The “hard case” is, in a sense, not likely to occur since, loosely speaking, for randomly generated data (\mathbf{A} and \mathbf{b}), the probability that the “hard case” will occur is 0. Therefore, the condition $\mathbf{b} \notin \mathcal{E}_1$ describing the “easy case” is considered to be mild. However, theoretically, the “hard case” might arise, and it usually requires a more involved and complicated analysis than the one required for the “easy case” (hence, the name).

3 Two Sufficient Optimality Conditions

In this section we will present two sufficient optimality conditions that will be the basis for showing how a global optimal solution of problem (TRS) can be obtained in both the “hard” and “easy” cases.

The sufficient conditions require that a stationary point will reside in some set. In the “easy case” the relevant set is

$$S_E = \{\mathbf{x} \in B : \text{sign}((\mathbf{U}^T \mathbf{b})_i)(\mathbf{U}^T \mathbf{x})_i \geq 0, i \in I_-\}, \quad (3.1)$$

where $\text{sgn}(\alpha)$ denotes the sign function which returns 1 if $\alpha \geq 0$ and -1 otherwise and

$$I_- = \{i \in \{1, 2, \dots, n\} : d_i \leq 0\}.$$

In the “hard case”, the relevant set is

$$S_H := \{\mathbf{x} \in B : \mathbf{x} \notin \mathcal{E}_1\}. \quad (3.2)$$

We are now ready to prove a sufficient optimality condition for each of the scenarios.

Theorem 3.1 (sufficient optimality condition - “easy case”). *Suppose that the “easy case” holds. Let $\bar{\mathbf{x}} \in S_E$ be a stationary point of problem (TRS). Then $\bar{\mathbf{x}}$ is the optimal solution of problem (TRS).*

Proof. If $d_1 \geq 0$, then problem (TRS) is convex, and obviously any stationary point is optimal. Therefore, we will henceforth assume that $d_1 < 0$. Since $\bar{\mathbf{x}}$ is a stationary point of problem (TRS), then by Theorem 2.5(a) for any $i = 1, 2, \dots, n$,

$$(d_i + \lambda_{(\bar{\mathbf{x}})})(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \mathbf{b})_i. \quad (3.3)$$

Since the “easy case” holds, it follows that there exists $i_1 \in E_1$ such that $(\mathbf{U}^T \mathbf{b})_{i_1} \neq 0$, which by (3.3) implies that $(\mathbf{U}^T \bar{\mathbf{x}})_{i_1} \neq 0$. Combining this with the fact that $\bar{\mathbf{x}} \in S_E$ and $d_1 < 0$, we obtain that $(\mathbf{U}^T \bar{\mathbf{x}})_{i_1}(\mathbf{U}^T \mathbf{b})_{i_1} > 0$. Therefore, by (3.3) with $i = i_1$,

$$\lambda_{(\bar{\mathbf{x}})} = -d_{i_1} + \frac{(\mathbf{U}^T \mathbf{b})_{i_1}}{(\mathbf{U}^T \bar{\mathbf{x}})_{i_1}} > -d_{i_1} = d_1,$$

which in light of Theorem 2.5(b) implies that $\bar{\mathbf{x}}$ is an optimal solution of problem (TRS). \square

Remark 3.2. The condition defining the set S_E , namely “ $\text{sgn}((\mathbf{U}^T \mathbf{b})_i)(\mathbf{U}^T \bar{\mathbf{x}})_i \geq 0$ ” corresponds to a property that is known to hold for optimal solutions of problem (TRS), and is the basis for its so-called “hidden convexity” property, see [4] as well as [2, Lemma 8.7].

Theorem 3.3 (sufficient optimality condition - “hard case”). *Suppose that the “hard case” holds. Let $\bar{\mathbf{x}} \in S_H$ be a stationary point of problem (TRS). Then $\bar{\mathbf{x}}$ is an optimal solution of problem (TRS).*

Proof. Since $\bar{\mathbf{x}} \in S_H$, it holds that $\bar{\mathbf{x}} \notin \mathcal{E}_1$, i.e., there exists $i \in E_1$ such that $(\mathbf{U}^T \bar{\mathbf{x}})_i \neq 0$. For such an i , the “hard case” assumption implies that $(\mathbf{U}^T \mathbf{b})_i = 0$. By Theorem 2.5(a), we obtain that

$$(\lambda_{(\bar{\mathbf{x}})} + d_i)(\mathbf{U}^T \bar{\mathbf{x}})_i = 0,$$

which combined with the fact that $(\mathbf{U}^T \bar{\mathbf{x}})_i \neq 0$ implies that $\lambda_{(\bar{\mathbf{x}})} = -d_i = -d_1$, and hence, by Theorem 2.5(b), $\bar{\mathbf{x}}$ is an optimal solution of problem (TRS). \square

In the next section we will show how the two sufficient conditions described in Theorems 3.1 and 3.3 form the basis for the convergence to a global optimal solution of a class of simple first-order methods that includes, among others, the projected gradient and conditional gradient methods. Specifically, it will be shown that in the “easy case”, an algorithm that belongs to this class of methods, initialized at a point in S_E , converges to a stationary point in S_E , which by Theorem 3.1, is also optimal. In the “hard case”, based on Theorem 3.3, it will be shown that convergence to an optimal solution is guaranteed as long as the initial point belongs to S_H . If it is not clear (as is the common situation) whether the hard or easy cases hold, then a global optimal solution can be obtained by running the method twice—once with $\mathbf{x}^0 \in S_E$ and the second time with $\mathbf{x}^0 \in S_H$. Note that although both S_E and S_H are described in terms of the spectral decomposition of \mathbf{A} , an operation that should be avoided, we will show that there is a simple way to choose the initial vector \mathbf{x}^0 without the need to access any spectral information.

4 First-Order Conic Methods

4.1 Definition and Examples

We begin with the definition of a “first-order conic method”.

Definition 4.1 (first-order conic method). A method which generates a sequence of points $\{\mathbf{x}^k\}_{k=0}^{\infty}$ according to the general step

$$\mathbf{x}^{k+1} = \theta_1^k \mathbf{x}^k + \theta_2^k (-\nabla q(\mathbf{x}^k)). \quad (4.1)$$

with some $\mathbf{x}^0 \in B$ is called a *first-order conic method* (FOCM) if it satisfies the following conditions:

- (A) The method generates a feasible sequence, i.e., $\mathbf{x}^k \in B$ for any $k \in \mathbb{N}$.
- (B) $q(\mathbf{x}^{k+1}) \leq q(\mathbf{x}^k)$ for all $k \in \mathbb{N}$.
- (C) All limit points of the sequence $\{\mathbf{x}^k\}_{k=0}^{\infty}$ are stationary points of problem (TRS).
- (D) $\theta^k \in \mathbb{R}_+^2 \setminus \{\mathbf{0}\}$ and $\theta_1^k \in [0, 1]$ for any $k \in \mathbb{N}$.

Property (D) in the definition of an FOCM above implies that

$$\mathbf{x}^{k+1} \in B \cap \text{cone}(\{\mathbf{x}^k, -\nabla q(\mathbf{x}^k)\}),$$

which explains the name of the class of methods. Note that Definition 4.1 is quite general. Following are two examples of algorithms which are FOCMs.

4.1.1 Projected Gradient

The projected gradient (PG) method is one of the most fundamental methods for solving differentiable constrained optimization problems [5, 15]. In the case of problem (TRS), the general step of the PG method is given by

Projected Gradient (PG) Method

For any $k = 0, 1, 2, \dots$, execute the following:

$$\mathbf{x}^{k+1} := P_B(\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)) = \begin{cases} \mathbf{x}^k - t_k \nabla q(\mathbf{x}^k), & \|\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)\| \leq 1, \\ \frac{\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)}{\|\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)\|}, & \text{otherwise,} \end{cases}$$

where $t_k \geq 0$ is a corresponding step-size and

$$P_B(\mathbf{y}) := \operatorname{argmin}\{\|\mathbf{x} - \mathbf{y}\|^2 : \mathbf{x} \in B\} = \begin{cases} \mathbf{y}, & \|\mathbf{y}\| \leq 1, \\ \frac{\mathbf{y}}{\|\mathbf{y}\|}, & \|\mathbf{y}\| > 1. \end{cases}$$

is the orthogonal projection operator. This method, with a constant step size, was previously suggested for solving problem (TRS) in [23] where a restarting procedure, that involves an eigenvector computation, was used in order to guarantee convergence to the optimal solution.

Property (A) in the definition of an FOCM (feasibility) is surely satisfied because of the projection step. Property (D) also holds since $t_k \geq 0$ and $\theta_1^k \in [0, 1]$ since it is either equal to 1 if $\|\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)\| \leq 1$ and if $\|\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)\| > 1$, then $\theta_1^k = \frac{1}{\|\mathbf{x}^k - t_k \nabla q(\mathbf{x}^k)\|} < 1$. The validity of properties (B) and (C) can be shown under some step-size regimes. Here we consider the following two popular step-size strategies:

- **Constant.** $t_k = \frac{1}{\bar{L}}$, where $\bar{L} \in (\|\mathbf{A}\|_2, \infty)$.
- **Backtracking.** The procedure requires three parameters (s, γ, η) where $s > 0$, $\gamma \in (0, 1)$ and $\eta > 1$. The step-size is chosen as $t_k = \frac{1}{L_k}$ where L_k is picked as follows: first, L_k is set to be equal to the initial guess s . Then, while

$$q(\mathbf{x}^k) - q\left(P_B\left(\mathbf{x}^k - \frac{1}{L_k} \nabla q(\mathbf{x}^k)\right)\right) < \gamma L_k \left\| \mathbf{x}^k - P_B\left(\mathbf{x}^k - \frac{1}{L_k} \nabla q(\mathbf{x}^k)\right) \right\|^2,$$

we set $L_k := \eta L_k$.

It is well known (see for example [3, Theorem 10.15]) that properties (B) and (C) hold for the PG method under the constant and backtracking strategies as described above. We can thus conclude the following.

Theorem 4.2. *The PG method with either constant or backtracking step-size strategies is an FOCM.*

4.1.2 Conditional Gradient

Another fundamental method for solving differentiable and constrained optimization problems is the conditional gradient (CG) method [5, 9, 15]. In the case of problem (TRS), the CG method takes the following form.

Conditional Gradient (CG) Method

For any $k = 0, 1, 2, \dots$, execute the following:

- (a) compute $\mathbf{p}^k \in \operatorname{argmin}_{\mathbf{p} \in B} \langle \mathbf{p}, \nabla q(\mathbf{x}^k) \rangle$;
- (b) choose $t_k \in [0, 1]$ and set $\mathbf{x}^{k+1} = \mathbf{x}^k + t_k(\mathbf{p}^k - \mathbf{x}^k)$.

Step (a) of the method is not well-defined when $\nabla q(\mathbf{x}^k) = \mathbf{0}$ since in this case we can choose \mathbf{p}^k to be any point in B . To make the method well-defined, and in order to avoid unnecessary pathological situations, we will set $\mathbf{p}^k = \mathbf{0}$ and $t_k = 0$ in this case. With this convention, we have

$$\mathbf{p}^k = \begin{cases} -\frac{\nabla q(\mathbf{x}^k)}{\|\nabla q(\mathbf{x}^k)\|}, & \nabla q(\mathbf{x}^k) \neq \mathbf{0}, \\ \mathbf{0}, & \nabla q(\mathbf{x}^k) = \mathbf{0}. \end{cases}$$

and thus, the general update step is given by

$$\mathbf{x}^{k+1} = \begin{cases} (1 - t_k)\mathbf{x}^k - \frac{t_k}{\|\nabla q(\mathbf{x}^k)\|} \nabla q(\mathbf{x}^k), & \nabla q(\mathbf{x}^k) \neq \mathbf{0}, \\ \mathbf{x}^k, & \nabla q(\mathbf{x}^k) = \mathbf{0}. \end{cases}$$

The above representation shows that the CG method satisfies properties (A) and (D) required from an FOCM. We will consider the following two step-size strategies:

- **Adaptive.** $t_k = \min \left\{ 1, \frac{S(\mathbf{x}^k)}{2\|\mathbf{A}\|_2\|\mathbf{p} - \mathbf{x}^k\|^2} \right\}$ where $S(\mathbf{x}) := \langle \nabla q(\mathbf{x}), \mathbf{x} \rangle + \|\nabla q(\mathbf{x})\|$.
- **Exact line search.** $t_k \in \operatorname{argmin}_{t \in [0, 1]} q(\mathbf{x}^k + t(\mathbf{p}^k - \mathbf{x}^k))$.

The CG method with either of the above step-size strategies is known to satisfy properties (B) and (C) in the definition of an FOCM (see for example [3, Chapter 13]), and we can thus conclude that in these settings, the CG method is an FOCM.

Theorem 4.3. *The CG method with either an adaptive or an exact line-search strategies is an FOCM.*

4.2 Convergence of the Sequence Generated by an FOCM

Our main objective will be to show how to achieve the global optimal point of problem (TRS) using any FOCM method. Before that, we will show that the sequence generated by an FOCM converges to a stationary point. This is a stronger property than condition (C) in the definition of an FOCM (Definition 4.1), which only warrants that limit points of the generated sequence are stationary point. This result will also be useful in showing later on how convergence to an optimal solution can be achieved under proper initialization.

Theorem 4.4. *Let $\{\mathbf{x}^k\}_{k=0}^{\infty}$ be a sequence generated by an FOCM. Then $\{\mathbf{x}^k\}_{k=0}^{\infty}$ converges to a stationary point of problem (TRS).*

Proof. Due to the monotonicity of $\{q(\mathbf{x}^k)\}_{k=0}^{\infty}$ (property (B) in Definition 4.1) and the fact that problem (TRS) admits an optimal solution, we can conclude that there is some \bar{q} such that

$$\lim_{k \rightarrow \infty} q(\mathbf{x}^k) = \bar{q}. \quad (4.2)$$

Since by the definition of an FOCM (specifically, property (C) in Definition 4.1), any limit point of the sequence is a stationary point, all we need to prove is that the sequence converges. Assume by contradiction that the sequence does not converge. Let then $\bar{\mathbf{x}}, \bar{\mathbf{y}} \in B$, $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$ be two different accumulation points of the sequence $\{\mathbf{x}^k\}_{k=0}^{\infty}$. By (4.2) and the continuity of q , we obtain that $q(\bar{\mathbf{x}}) = q(\bar{\mathbf{y}}) = \bar{q}$. In addition, by property (C) in the definition of an FOCM (Definition 4.1), both points are stationary points of problem (TRS). We will show that $\lambda_{(\bar{\mathbf{x}})} = \lambda_{(\bar{\mathbf{y}})}$. Indeed,

$$\begin{aligned} q(\bar{\mathbf{x}}) - q(\bar{\mathbf{y}}) &= \bar{\mathbf{x}}^T \mathbf{A} \bar{\mathbf{x}} - 2\mathbf{b}^T \bar{\mathbf{x}} - \bar{\mathbf{y}}^T \mathbf{A} \bar{\mathbf{y}} + 2\mathbf{b}^T \bar{\mathbf{y}} \\ &= (\bar{\mathbf{x}} - \bar{\mathbf{y}})^T \mathbf{A} (\bar{\mathbf{x}} - \bar{\mathbf{y}}) + 2(\mathbf{A} \bar{\mathbf{y}} - \mathbf{b})^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) \\ &= (\bar{\mathbf{x}} - \bar{\mathbf{y}})^T \mathbf{A} (\bar{\mathbf{x}} - \bar{\mathbf{y}}) - 2\lambda_{(\bar{\mathbf{y}})} \bar{\mathbf{y}}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}), \end{aligned}$$

where in the last equality we have used (2.1) due to the stationarity of $\bar{\mathbf{y}}$. Similarly, we obtain that

$$q(\bar{\mathbf{y}}) - q(\bar{\mathbf{x}}) = (\bar{\mathbf{y}} - \bar{\mathbf{x}})^T \mathbf{A} (\bar{\mathbf{y}} - \bar{\mathbf{x}}) - 2\lambda_{(\bar{\mathbf{x}})} \bar{\mathbf{x}}^T (\bar{\mathbf{y}} - \bar{\mathbf{x}}).$$

Since we already established that $q(\bar{\mathbf{x}}) = q(\bar{\mathbf{y}})$, by combining the last two results we obtain that

$$\lambda_{(\bar{\mathbf{x}})} (\bar{\mathbf{y}}^T \bar{\mathbf{x}} - \|\bar{\mathbf{x}}\|^2) = \lambda_{(\bar{\mathbf{y}})} (\bar{\mathbf{y}}^T \bar{\mathbf{x}} - \|\bar{\mathbf{y}}\|^2). \quad (4.3)$$

If $\|\bar{\mathbf{x}}\|^2 = \|\bar{\mathbf{y}}\|^2$ then due to our assumption that $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$, we obtain that¹ $\bar{\mathbf{y}}^T \bar{\mathbf{x}} < \|\bar{\mathbf{x}}\|^2 = \|\bar{\mathbf{y}}\|^2$ and thus $\lambda_{(\bar{\mathbf{x}})} = \lambda_{(\bar{\mathbf{y}})}$. Otherwise, assume without loss of generality that $\|\bar{\mathbf{x}}\|^2 < \|\bar{\mathbf{y}}\|^2$ which

¹Since $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$ and $\|\bar{\mathbf{x}}\| = \|\bar{\mathbf{y}}\|$, it follows that $\bar{\mathbf{x}}, \bar{\mathbf{y}} \neq \mathbf{0}$. To prove the strict inequality, note that by the Cauchy-Schwarz inequality $\bar{\mathbf{y}}^T \bar{\mathbf{x}} \leq \|\bar{\mathbf{y}}\| \cdot \|\bar{\mathbf{x}}\| = \|\bar{\mathbf{x}}\|^2$. Equality will hold only if $\bar{\mathbf{y}} = \alpha \bar{\mathbf{x}}$ for some nonnegative α , but since $\|\bar{\mathbf{x}}\| = \|\bar{\mathbf{y}}\|$, α has to be equal to 1, leading to a contradiction to the assumption that $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$.

implies by the complementary slackness condition (2.5) that $\lambda_{(\bar{\mathbf{x}})} = 0$. Then (4.3) boils down to $\lambda_{(\bar{\mathbf{y}})}(\bar{\mathbf{y}}^T \bar{\mathbf{x}} - \|\bar{\mathbf{y}}\|^2) = 0$ and since $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$ we obtain that² $\lambda_{(\bar{\mathbf{y}})} = 0$.

To summarize, up to this point we assumed by contradiction that the sequence $\{\mathbf{x}^k\}_{k=0}^\infty$ has two different accumulation points $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$ and showed that $\lambda_{(\bar{\mathbf{x}})} = \lambda_{(\bar{\mathbf{y}})}$. We will denote the common Lagrange multiplier value by $\lambda := \lambda_{(\bar{\mathbf{x}})} = \lambda_{(\bar{\mathbf{y}})}$. Now, since both $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are stationary points of problem (TRS), then by Theorem 2.5(a)

$$(\lambda + d_i)(\mathbf{U}^T \mathbf{z})_i = (\mathbf{U}^T \mathbf{b})_i, \quad \mathbf{z} \in \{\bar{\mathbf{x}}, \bar{\mathbf{y}}\}, \quad i = 1, 2, \dots, n. \quad (4.4)$$

If $\lambda \neq -d_i$ for all $i = 1, 2, \dots, n$, then by (4.4)

$$(\mathbf{U}^T \bar{\mathbf{x}})_i = \frac{(\mathbf{U}^T \mathbf{b})_i}{\lambda + d_i} = (\mathbf{U}^T \bar{\mathbf{y}})_i, \quad i = 1, 2, \dots, n,$$

implying that $\bar{\mathbf{x}} = \bar{\mathbf{y}}$, which is a contradiction to the assumption that $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$. We will therefore assume from now on that there is some $m \in \{1, 2, \dots, n\}$ for which $\lambda_{(\bar{\mathbf{x}})} = -d_m$. We denote (expanding the notation “ E_1 ”, see (2.6)),

$$E_m = \{\ell \in \{1, 2, \dots, n\} : d_\ell = d_m\}.$$

For any $i \notin E_m$, by (4.4),

$$(\mathbf{U}^T \bar{\mathbf{x}})_i = \frac{(\mathbf{U}^T \mathbf{b})_i}{\lambda + d_i} = (\mathbf{U}^T \bar{\mathbf{y}})_i.$$

We will show that $(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \bar{\mathbf{y}})_i$ also for any $i \in E_m$. Let then $i \in E_m$. Then due to (4.4), since $\lambda = -d_m$, it must hold that $(\mathbf{U}^T \mathbf{b})_i = 0$, and thus, the update step (4.1) of the FOCM method is (recalling that $d_i = d_m$ for $i \in E_m$),

$$\begin{aligned} (\mathbf{U}^T \mathbf{x}^{k+1})_i &= \theta_1^k (\mathbf{U}^T \mathbf{x}^k)_i + 2\theta_2^k (-d_i (\mathbf{U}^T \mathbf{x}^k)_i) \\ &= \left[\prod_{j=0}^k (\theta_1^j - 2\theta_2^j d_m) \right] (\mathbf{U}^T \mathbf{x}^0)_i, \end{aligned}$$

meaning that

$$(\mathbf{U}^T \mathbf{x}^{k+1})_i = \pi_k (\mathbf{U}^T \mathbf{x}^0)_i, \quad (4.5)$$

where

$$\pi_k := \prod_{j=0}^k (\theta_1^j - 2\theta_2^j d_m).$$

Note that $\pi_k \geq 0$ for any k since $-d_m = \lambda \geq 0$ and $\theta_1^j, \theta_2^j \geq 0$.

We will now consider two cases.

(a) If $\lambda = 0$, then for any $i \in E_m$ (4.5) is simply

$$(\mathbf{U}^T \mathbf{x}^{k+1})_i = \pi_k (\mathbf{U}^T \mathbf{x}^0)_i,$$

where here $\pi_k = \prod_{j=0}^k \theta_1^j$. Since $\theta_1^k \in [0, 1]$ for any $k \in \mathbb{N}$, the sequence $\{\pi_k\}_{k=0}^\infty$ is non-increasing and bounded from below by zero, and thus converges to some β . Hence, $\{(\mathbf{U}^T \mathbf{x}^k)_i\}_{k=0}^\infty$ converges to $\beta (\mathbf{U}^T \mathbf{x}^0)_i$. Since $(\mathbf{U}^T \bar{\mathbf{x}})_i$ and $(\mathbf{U}^T \bar{\mathbf{y}})_i$ are limit points of the convergent sequence $\{(\mathbf{U}^T \mathbf{x}^k)_i\}_{k=0}^\infty$, it follows that $(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \bar{\mathbf{y}})_i$.

²Indeed, by the Cauchy-Schwarz inequality $\bar{\mathbf{y}}^T \bar{\mathbf{x}} \leq \|\bar{\mathbf{y}}\| \cdot \|\bar{\mathbf{x}}\| < \|\bar{\mathbf{y}}\|^2$, where the strict inequality is due to the facts that $\bar{\mathbf{y}} \neq \mathbf{0}$ and $\|\bar{\mathbf{x}}\| < \|\bar{\mathbf{y}}\|$.

(b) If $\lambda \neq 0$, then due to the complementary slackness condition and the fact, previously established, that $(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \bar{\mathbf{y}})_i$ for any $i \notin E_m$, we obtain that

$$\sum_{i \in E_m} (\mathbf{U}^T \bar{\mathbf{x}})_i^2 = \sum_{i \in E_m} (\mathbf{U}^T \bar{\mathbf{y}})_i^2. \quad (4.6)$$

By (4.5) it follows that $(\mathbf{U}^T \bar{\mathbf{x}})_i (\mathbf{U}^T \bar{\mathbf{y}})_i \geq 0$ for any $i \in E_m$ and consequently, the fact that $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$ implies that $|\mathbf{U}^T \bar{\mathbf{x}}| \neq |\mathbf{U}^T \bar{\mathbf{y}}|$. Combining this with (4.6), it follows that there must be some $i_1, i_2 \in E_m$ and some $\epsilon > 0$ such that

$$|(\mathbf{U}^T \bar{\mathbf{x}})_{i_1}| - |(\mathbf{U}^T \bar{\mathbf{y}})_{i_1}| > 2\epsilon \quad \text{and} \quad |(\mathbf{U}^T \bar{\mathbf{y}})_{i_2}| - |(\mathbf{U}^T \bar{\mathbf{x}})_{i_2}| > 2\epsilon. \quad (4.7)$$

Since $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are accumulation points of the sequence $\{\mathbf{x}^k\}_{k=0}^\infty$, there are two subsequences $\{\mathbf{x}^{k_t^1}\}_{t=0}^\infty$ and $\{\mathbf{x}^{k_t^2}\}_{t=0}^\infty$ such that $\mathbf{x}^{k_t^1} \rightarrow \bar{\mathbf{x}}$ and $\mathbf{x}^{k_t^2} \rightarrow \bar{\mathbf{y}}$ as $t \rightarrow \infty$. Hence there exists $t_1 \in \mathbb{N}$ such that

$$|(\mathbf{U}^T \mathbf{x}^{k_t^1})_{i_1} - (\mathbf{U}^T \bar{\mathbf{x}})_{i_1}| < \epsilon \quad \text{and} \quad |(\mathbf{U}^T \mathbf{x}^{k_t^1})_{i_2} - (\mathbf{U}^T \bar{\mathbf{x}})_{i_2}| < \epsilon \quad \text{for any } t \geq t_1. \quad (4.8)$$

Similarly, there exists $t_2 \in \mathbb{N}$ such that

$$|(\mathbf{U}^T \mathbf{x}^{k_t^2})_{i_1} - (\mathbf{U}^T \bar{\mathbf{y}})_{i_1}| < \epsilon \quad \text{and} \quad |(\mathbf{U}^T \mathbf{x}^{k_t^2})_{i_2} - (\mathbf{U}^T \bar{\mathbf{y}})_{i_2}| < \epsilon \quad \text{for any } t \geq t_2. \quad (4.9)$$

Assume, without the loss of generality, that $k_{t_2}^2 > k_{t_1}^1$. Combining (4.7) and (4.8), we obtain

$$|(\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_1}| - |(\mathbf{U}^T \bar{\mathbf{y}})_{i_1}| > \epsilon, \quad (4.10)$$

$$|(\mathbf{U}^T \bar{\mathbf{y}})_{i_2}| - |(\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_2}| > \epsilon. \quad (4.11)$$

Due to (4.5) we also have that

$$(\mathbf{U}^T \mathbf{x}^{k_{t_2}^2})_{i_1} = \eta (\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_1} \quad \text{and} \quad (\mathbf{U}^T \mathbf{x}^{k_{t_2}^2})_{i_2} = \eta (\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_2}, \quad (4.12)$$

where

$$\eta = \left[\prod_{j=k_{t_1}^1}^{k_{t_2}^2-1} (\theta_1^j - 2\theta_2^j d_m) \right].$$

Thus, (4.9) and (4.12) yield

$$|\eta (\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_1} - (\mathbf{U}^T \bar{\mathbf{y}})_{i_1}| < \epsilon, \quad (4.13)$$

$$|\eta (\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_2} - (\mathbf{U}^T \bar{\mathbf{y}})_{i_2}| < \epsilon. \quad (4.14)$$

Note that in addition, (4.5) implies that

$$(\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_1} (\mathbf{U}^T \bar{\mathbf{y}})_{i_1} \geq 0, \quad (4.15)$$

$$(\mathbf{U}^T \mathbf{x}^{k_{t_1}^1})_{i_2} (\mathbf{U}^T \bar{\mathbf{y}})_{i_2} \geq 0. \quad (4.16)$$

Combining (4.10), (4.13) and (4.15) implies that $\eta < 1$. This is a direct consequence of the following simple result on real numbers: if $a, b \in \mathbb{R}$ satisfy (for some $\epsilon > 0$) $ab \geq 0, |a| - |b| > \epsilon, |\eta a - b| < \epsilon$, then $\eta < 1$. Similarly, (4.11), (4.14) and (4.16) imply that $\eta > 1$. We thus obtained a contradiction to our assumption that $\bar{\mathbf{x}} \neq \bar{\mathbf{y}}$. \square

Remark 4.5. Convergence of the sequence generated by the projected gradient method can be established using the Kurdyka-Łojasiewicz property under the condition that the step-sizes t_k reside in the interval $(2\|\mathbf{A}\|_2, M)$ for some $M > 2\|\mathbf{A}\|_2$, see [1] and the precise statement in [6, Proposition 3].

4.3 Convergence to the Optimal Solution – the “Easy Case”

So far, we have established the convergence of the sequence generated by an FOCM, but we did not prove that it converges to the global optimal solution of problem (TRS). This is actually not correct in general. What we can prove is that both in the easy and hard cases, certain choices of the initial vector \mathbf{x}^0 will guarantee convergence to the optimal solution.

Theorem 4.6 below shows that in the setting of the “easy case”, if $\mathbf{x}^0 \in S_E$, then an FOCM will converge to the optimal solution of problem (TRS). This is done by showing that the limit of the sequence is a stationary point in S_E , which by the sufficient condition described in Theorem 3.1, implies that the limit is an optimal solution.

Theorem 4.6 (convergence to the optimal solution in the “easy case”). *Let $\{\mathbf{x}^k\}_{k=0}^\infty$ be a sequence generated by an FOCM. If $\mathbf{x}^0 \in S_E$, then the sequence $\{\mathbf{x}^k\}_{k=0}^\infty$ converges to the optimal solution of problem (TRS).*

Proof. By the definition of an FOCM, for any $k \in \mathbb{N}$ there exists $\theta^k \in \mathbb{R}_+^2 \setminus \{0\}$ such that

$$\mathbf{x}^{k+1} = \theta_1^k \mathbf{x}^k + \theta_2^k (-\nabla q(\mathbf{x}^k)),$$

and $\mathbf{x}^{k+1} \in B$. We will show that the whole sequence $\{\mathbf{x}^k\}_{k=0}^\infty$ is in S_E by induction.

By the statement of the theorem, $\mathbf{x}^0 \in S_E$. Assume that $\mathbf{x}^k \in S_E$. Then for any $i \in I_-$,

$$\begin{aligned} \operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (-\mathbf{U}^T \nabla q(\mathbf{x}^k))_i &= \operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (-2d_i(\mathbf{U}^T \mathbf{x}^k)_i + 2(\mathbf{U}^T \mathbf{b})_i) \\ &= -2d_i \operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (\mathbf{U}^T \mathbf{x}^k)_i + 2|(\mathbf{U}^T \mathbf{b})_i| \\ &\geq 0. \end{aligned}$$

where the inequality follows by the inclusion $\mathbf{x}^k \in S_E$ and the fact that $d_i \leq 0$ for any $i \in I_-$.

We have thus shown that if $\mathbf{x}^k \in S_E$, then $-\nabla q(\mathbf{x}^k) \in S_E$. This result and the fact that $\theta^k \in \mathbb{R}_+^2 \setminus \{0\}$ imply that for any $i \in I_-$,

$$\operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (\mathbf{U}^T \mathbf{x}^{k+1})_i = \theta_1^k \operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (\mathbf{U}^T \mathbf{x}^k)_i + \theta_2^k \operatorname{sgn}((\mathbf{U}^T \mathbf{b})_i) (-\mathbf{U}^T \nabla q(\mathbf{x}^k))_i \geq 0.$$

Thus, $\mathbf{x}^{k+1} \in S_E$, showing that $\mathbf{x}^k \in S_E$ for any $k \in \mathbb{N}$. Therefore, since $\{\mathbf{x}^k\}_{k=0}^\infty$ converges to some stationary point $\bar{\mathbf{x}}$ (Theorem 4.4) and since S_E is a closed set, it follows that $\bar{\mathbf{x}} \in S_E$. We have thus shown that the limit point $\bar{\mathbf{x}}$ of the the sequence is a stationary point in S_E , and thus, by Theorem 3.1, it is an optimal solution of problem (TRS). \square

The set S_E as defined in (3.1) is given in terms of the spectral decomposition of \mathbf{A} . Thus, the task of finding a point $\mathbf{x}^0 \in S_E$ seems to require some spectral knowledge; however, such knowledge is not really needed by the trivial observation that $\mathbf{0} \in S_E$. Consequently, Theorem 4.6 ensures that in the “easy case”, if an FOCM is initialized with $\mathbf{x}^0 = \mathbf{0}$, then it is guaranteed to converge to the optimal solution.

Corollary 4.7. *Suppose that the “easy case” holds and let $\{\mathbf{x}^k\}_{k=0}^\infty$ be a sequence generated by an FOCM initialized with $\mathbf{x}^0 = \mathbf{0}$. Then the sequence converges to an optimal solution of problem (TRS).*

4.4 Convergence to the Optimal Solution – the “Hard Case”

In the “hard case”, we will obtain a result that resembles a well-known result (see for example [11, Section 8.2.1]) for problem (TRS) with $\mathbf{b} = \mathbf{0}$, which shows that the so-called “power iteration method” converges to the dominant eigenvector of the matrix \mathbf{A} under the condition that the dominant eigenvalue is unique³ and the mild condition that the initial point \mathbf{x}^0 satisfies

$$\mathbf{x}^0 \notin \mathcal{E}_1. \quad (4.17)$$

The above condition is the same as the relation $\mathbf{x}^0 \in S_H$. We note that the power iteration method is known to be equivalent to the conditional gradient method (see [16]), and as such, it is an FOCM. Theorem 4.8 below shows that if the “hard case” holds and $\mathbf{b} \neq \mathbf{0}$, then under the same mild condition on the initial vector (4.17), convergence of any FOCM to the optimal solution of problem (TRS) is warranted. The proof is based on the sufficient optimality condition given in Theorem 3.3.

Theorem 4.8 (convergence to an optimal solution in the “hard case”). *Suppose that the “hard case” holds and that $\mathbf{b} \neq \mathbf{0}$. Let $\{\mathbf{x}^k\}_{k=0}^\infty$ be a sequence generated by an FOCM with $\mathbf{x}^0 \in S_H$. Then the sequence $\{\mathbf{x}^k\}_{k=0}^\infty$ converges to an optimal solution of problem (TRS).*

Proof. By Theorem 4.4 $\{\mathbf{x}^k\}_{k=0}^\infty$ converges, and we will denote its limit by $\bar{\mathbf{x}}$. By the definition of an FOCM, $\bar{\mathbf{x}}$ is a stationary point of problem (TRS). As usual, the associated Lagrange multiplier of $\bar{\mathbf{x}}$ will be denoted by $\lambda_{(\bar{\mathbf{x}})}$. By Theorem 2.5(a), $(\bar{\mathbf{x}}, \lambda_{(\bar{\mathbf{x}})})$ satisfies the equations

$$(\lambda_{(\bar{\mathbf{x}})} + d_i)(\mathbf{U}^T \bar{\mathbf{x}})_i = (\mathbf{U}^T \mathbf{b})_i, \quad i = 1, 2, \dots, n. \quad (4.18)$$

We will prove that $\bar{\mathbf{x}}$ is an optimal solution of (TRS). If $d_1 \geq 0$, then problem (TRS) is convex, and hence $\bar{\mathbf{x}}$ must be an optimal solution, and we are done. We will henceforth consider the case where $d_1 < 0$.

Assume by contradiction that $\bar{\mathbf{x}}$ is not an optimal solution of problem (TRS). By Theorem 3.3, this entails that $\bar{\mathbf{x}} \notin S_H$, which by the definition of S_H (see (3.2)) implies that $(\mathbf{U}^T \bar{\mathbf{x}})_i = 0$ for any $i \in E_1$, and in particular that

$$(\mathbf{U}^T \bar{\mathbf{x}})_1 = 0. \quad (4.19)$$

Another consequence of the fact that $\bar{\mathbf{x}}$ is a non-optimal stationary point of problem (TRS) and Theorem 2.5(b) is that

$$\lambda_{(\bar{\mathbf{x}})} < -d_1. \quad (4.20)$$

Since $\mathbf{x}^0 \in S_H$, it follows that there exists $i \in E_1$ for which $(\mathbf{U}^T \mathbf{x}^0)_i \neq 0$ and we assume w.l.o.g. that $(\mathbf{U}^T \mathbf{x}^0)_1 \neq 0$.

Since $(\mathbf{U}^T \mathbf{b})_1 = 0$, the update formula of the method (4.1) can be written as

$$(\mathbf{U}^T \mathbf{x}^{k+1})_1 = \theta_1^k (\mathbf{U}^T \mathbf{x}^k)_1 + 2\theta_2^k (-d_1 (\mathbf{U}^T \mathbf{x}^k)_1) = [\theta_1^k - 2\theta_2^k d_1] (\mathbf{U}^T \mathbf{x}^k)_1,$$

³This condition can be relaxed to some extent [25, Chapter 9].

meaning that

$$(\mathbf{U}^T \mathbf{x}^{k+1})_1 = \alpha_k (\mathbf{U}^T \mathbf{x}^k)_1, \quad (4.21)$$

where $\alpha_k = \theta_1^k - 2\theta_2^k d_1$. Note that $\alpha_k > 0$ since $d_1 < 0$ and $\boldsymbol{\theta}^k \in \mathbb{R}_+^2 \setminus \{\mathbf{0}\}$ for all k . Combining the fact that $(\mathbf{U}^T \mathbf{x}^0)_1 \neq 0$ and the validity of (4.21) with $\alpha_k > 0$, we can conclude that

$$(\mathbf{U}^T \mathbf{x}^k)_1 \neq 0 \text{ for any } k \in \mathbb{N}. \quad (4.22)$$

Since $\mathbf{b} \neq \mathbf{0}$, there exists $i \in \{1, 2, \dots, n\}$ such that $(\mathbf{U}^T \mathbf{b})_i \neq 0$, which by (4.18) implies that $\lambda_{(\bar{\mathbf{x}})} \neq -d_i$ and

$$(\mathbf{U}^T \bar{\mathbf{x}})_i \neq 0. \quad (4.23)$$

Since $(\mathbf{U}^T \mathbf{x}^k)_i \rightarrow (\mathbf{U}^T \bar{\mathbf{x}})_i$ as $k \rightarrow \infty$, it follows by (4.18) and (4.23) that there exists $K_1 \in \mathbb{N}$ such that

$$(\mathbf{U}^T \mathbf{x}^k)_i [(\lambda_{(\bar{\mathbf{x}})} + d_i)(\mathbf{U}^T \mathbf{b})_i] > 0 \text{ for all } k \geq K_1. \quad (4.24)$$

The last equation implies further that

$$(\mathbf{U}^T \mathbf{x}^k)_i (\mathbf{U}^T \mathbf{x}^{k+1})_i > 0 \text{ for all } k \geq K_1. \quad (4.25)$$

For any $k \geq K_1$, we can write

$$(\mathbf{U}^T \mathbf{x}^{k+1})_i = \theta_1^k (\mathbf{U}^T \mathbf{x}^k)_i + 2\theta_2^k (-d_i (\mathbf{U}^T \mathbf{x}^k)_i + (\mathbf{U}^T \mathbf{b})_i),$$

meaning that

$$(\mathbf{U}^T \mathbf{x}^{k+1})_i = \eta_k (\mathbf{U}^T \mathbf{x}^k)_i, \quad (4.26)$$

where

$$\eta_k = \theta_1^k + 2\theta_2^k \left(-d_i + \frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \mathbf{x}^k)_i} \right).$$

Combining (4.25) and (4.26) we can also conclude that $\eta_k > 0$ for all $k \geq K_1$.

We now consider two cases:

- (i) $\lambda_{(\bar{\mathbf{x}})} < -d_i$. In this case, by (4.24), it follows that $\frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \mathbf{x}^k)_i} < 0$ for all $k \geq K_1$. We can therefore conclude that for any $k \geq K_1$

$$0 \leq \eta_k = \theta_1^k + 2\theta_2^k \left(-d_i + \frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \mathbf{x}^k)_i} \right) \leq \theta_1^k + 2\theta_2^k (-d_i) \leq \theta_1^k + 2\theta_2^k (-d_1) = \alpha_k.$$

- (ii) $\lambda_{(\bar{\mathbf{x}})} > -d_i$. In this case, by (4.18), it follows that

$$|(\mathbf{U}^T \bar{\mathbf{x}})_i| = \frac{|(\mathbf{U}^T \mathbf{b})_i|}{\lambda_{(\bar{\mathbf{x}})} + d_i} > \frac{|(\mathbf{U}^T \mathbf{b})_i|}{-d_1 + d_i},$$

where the inequality is due to (4.20). The above inequality can be also rewritten as

$$-d_1 > -d_i + \frac{|(\mathbf{U}^T \mathbf{b})_i|}{|(\mathbf{U}^T \bar{\mathbf{x}})_i|} = -d_i + \frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \bar{\mathbf{x}})_i},$$

where the equality is due to (4.18). Since $(\mathbf{U}^T \mathbf{x}^k)_i \rightarrow (\mathbf{U}^T \bar{\mathbf{x}})_i$ as $k \rightarrow \infty$ we obtain that there is some K_2 such that

$$-d_1 > -d_i + \frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \mathbf{x}^k)_i} \text{ for all } k \geq K_2.$$

We can thus conclude that in this case, for any $k \geq \max\{K_1, K_2\}$, it holds that

$$0 \leq \eta_k = \theta_1^k + 2\theta_2^k \left(-d_i + \frac{(\mathbf{U}^T \mathbf{b})_i}{(\mathbf{U}^T \mathbf{x}^k)_i} \right) \leq \theta_1^k + 2\theta_2^k(-d_1) = \alpha_k.$$

In both cases, we showed that there exists $K \in \mathbb{N}$ ($K = K_1$ in the first case and $K = \max\{K_1, K_2\}$ in the second case) such that for all $k \geq K$ it holds that $0 \leq \eta_k \leq \alpha_k$. We can now conclude that for any $k \geq K$,

$$\begin{aligned} |(\mathbf{U}^T \mathbf{x}^{k+1})_i| &= \prod_{j=K}^k \eta_j |(\mathbf{U}^T \mathbf{x}^K)_i| && [(4.26)] \\ &\leq \prod_{j=K}^k \alpha_j |(\mathbf{U}^T \mathbf{x}^K)_i| && [0 \leq \eta_k \leq \alpha_k] \\ &= \prod_{j=K}^k \alpha_j |(\mathbf{U}^T \mathbf{x}^K)_1| \frac{|(\mathbf{U}^T \mathbf{x}^K)_i|}{|(\mathbf{U}^T \mathbf{x}^K)_1|} && [(4.22)] \\ &= |(\mathbf{U}^T \mathbf{x}^k)_1| \frac{|(\mathbf{U}^T \mathbf{x}^K)_i|}{|(\mathbf{U}^T \mathbf{x}^K)_1|} && [(4.21)]. \end{aligned}$$

Since $\lim_{k \rightarrow \infty} (\mathbf{U}^T \mathbf{x}^k)_1 = (\mathbf{U}^T \bar{\mathbf{x}})_1 = 0$ (equation (4.19)), it follows by taking $k \rightarrow \infty$ in the above inequality that $(\mathbf{U}^T \bar{\mathbf{x}})_i = 0$, which is a contradiction to (4.23), thus establishing the desired result that $\bar{\mathbf{x}}$ is an optimal solution of problem (TRS). \square

Remark 4.9 (finding a point satisfying $\mathbf{x} \in S_H$). The condition $\mathbf{x}^0 \in S_H$, which is also the same as condition (4.17), can be more explicitly written as

$$(\mathbf{U}^T \mathbf{x}^0)_i \neq 0 \text{ for some } i \in E_1.$$

This is an extremely mild condition in the sense that if \mathbf{x}^0 is randomly generated via a continuous distribution with support B (for example, uniform distribution over B), then the probability that it will be satisfied is 1, and in this case, an FOCM will surely converge to the optimal solution of problem (TRS), as guaranteed by Theorem 4.8.

4.5 A Double-Start FOCM Method

Verifying whether the “easy” or “hard” cases hold requires spectral information which we do not assume to possess, and thus in this common case of uncertainty, we can employ an FOCM twice with two starting points: $\mathbf{x}^0 = \mathbf{0}$ and \mathbf{x}^0 which is randomly generated via a continuous distribution whose support is B . The best of the two resulting points (in terms of the objective function q) will be an optimal solution of problem (TRS) in probability 1.

Double-Start FOCM

Given an FOCM, execute the following:

- (a) Employ the FOCM with initial point $\mathbf{x}^0 = \mathbf{0}$ and obtain an output $\bar{\mathbf{x}}$.
- (b) Employ the FOCM with initial point \mathbf{x}^0 chosen via a continuous distribution function over B and obtain an output $\tilde{\mathbf{x}}$.
- (c) The output of the method is $\mathbf{x}^* \in \underset{\mathbf{z}}{\operatorname{argmin}}\{q(\mathbf{z}) : \mathbf{z} \in \{\bar{\mathbf{x}}, \tilde{\mathbf{x}}\}\}$.

References

- [1] Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized gauss–seidel methods. *Math. Program.*, 137(1):91–129, Feb 2013.
- [2] Amir Beck. *Introduction to Nonlinear Optimization: Theory, Algorithms, and Applications with MATLAB*, volume 19. SIAM, 2014.
- [3] Amir Beck. *First-Order Methods in Optimization*, volume 25. SIAM, 2017.
- [4] Aharon Ben-Tal and Marc Teboulle. Hidden convexity in some nonconvex quadratically constrained quadratic programming. *Math. Program.*, 72(1):51–63, 1996.
- [5] Dimitri P. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2nd edition, September 1999.
- [6] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Math. Program.*, 146(1):459–494, Aug 2014.
- [7] Andrew R. Conn, Nicholas I.M. Gould, and Philippe L. Toint. *Trust region methods*, volume 1. SIAM, 2000.
- [8] Jennifer B. Erway, Philip E. Gill, and Joshua D. Griffin. Iterative methods for finding a trust-region step. *SIAM J. Optim.*, 20(2):1110–1131, 2009.
- [9] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval Res. Logist. Quart.*, 3(1-2):95–110, 1956.
- [10] David M. Gay. Computing optimal locally constrained steps. *SIAM J. Sci. Comput.*, 2(2):186–197, 1981.
- [11] Gene H. Golub and Charles F. Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.

- [12] Nicholas I.M. Gould, Stefano Lucidi, Massimo Roma, and Philippe L. Toint. Solving the trust-region subproblem using the lanczos method. *SIAM J. Optim.*, 9(2):504–525, 1999.
- [13] William W. Hager. Minimizing a quadratic over a sphere. *SIAM J. Optim.*, 12(1):188–208, 2001.
- [14] Jörg Lampe, Marielba Rojas, Danny C. Sorensen, and Heinrich Voss. Accelerating the lstrs algorithm. *SIAM J. Sci. Comput.*, 33(1):175–194, 2011.
- [15] Evgeny S. Levitin and Boris T. Polyak. Constrained minimization methods. *U.S.S.R. Comput. Math. Math. Phys.*, 6(5):787–823, 1966.
- [16] Ronny Luss and Marc Teboulle. Conditional gradient algorithms for rank-one matrix approximations with a sparsity constraint. *SIAM Rev.*, 55(1):65–98, 2013.
- [17] Jorge J. Moré and Danny C. Sorensen. Computing a trust region step. *SIAM J. Sci. Comput.*, 4(3):553–572, 1983.
- [18] Franz Rendl and Henry Wolkowicz. A semidefinite framework for trust region subproblems with applications to large scale minimization. *Math. Program.*, 77(1):273–299, 1997.
- [19] Marielba Rojas, Sandra A. Santos, and Danny C. Sorensen. A new matrix-free algorithm for the large-scale trust-region subproblem. *SIAM J. Optim.*, 11(3):611–646, 2001.
- [20] Danny C. Sorensen. Newton’s method with a model trust region modification. *SIAM J. Numer. Anal.*, 19(2):409–426, 1982.
- [21] Danny C. Sorensen. Minimization of a large-scale quadratic function subject to a spherical constraint. *SIAM J. Optim.*, 7(1):141–161, 1997.
- [22] Trond Steihaug. The conjugate gradient method and trust regions in large scale optimization. *SIAM J. Numer. Anal.*, 20(3):626–637, 1983.
- [23] Pham Dinh Tao and Le Thi Hoai An. A dc optimization algorithm for solving the trust-region subproblem. *SIAM J. Optim.*, 8(2):476–505, 1998.
- [24] Philippe L. Toint. Towards an efficient sparsity exploiting newton method for minimization. In *Sparse Matrices and Their Uses*, pages 57–88. Academic Press, London, New York, 1981.
- [25] James Hardy Wilkinson. *The algebraic eigenvalue problem*, volume 87. Clarendon Press Oxford, 1965.
- [26] Yinyu Ye. A new complexity result on minimization of a quadratic function with a sphere constraint. In *Recent Advances in Global Optimization*, pages 19–31. Princeton University Press, 1992.