

# Supplementary Material

## On Asymptotic Linear Convergence of Projected Gradient Descent for Constrained Least Squares

Trung Vu, *Graduate Student Member, IEEE*, Raviv Raich, *Senior Member, IEEE*

### I. RELATED WORKS

In this section, we review existing approaches to convergence analysis of iterative first-order methods in optimization including projected gradient descent. We present several aspects of convergence, namely, convergence to a global versus a local optimum and speed of convergence. Finally, we clarify our contribution in this work with regard to previous works in the literature.

#### A. Convergence of Iterative First-Order Methods

Convergence properties of iterative algorithms such as PGD often involve two key aspects: the quality of convergent points and the speed of convergence. On the one hand, the quality of convergent points provides useful insights into when the algorithm converges, whether it converges to a stationary point or a set of stationary points of the problem, and how big is the gap between the objective function at the convergent point and the optimal objective value. On the other hand, the speed of convergence concerns the order of convergence, the rate of convergence, and the number of iterations required to obtain sufficiently small errors. Let  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  be the sequence of updates generated by a certain iterative first-order method (e.g., PGD). In order to prove the convergence of the algorithm, it is common [1]–[5] to consider the convergence of the following quantities to  $\mathbf{0}$  as  $k \rightarrow \infty$ : (i) the norm of the generalized gradient ( $\|\frac{1}{\eta}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})\|$ ), (ii) the gap between current objective function and the optimal value ( $|f(\mathbf{x}^{(k)}) - f^*|$ ), and (iii) the distance to a convergent point ( $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ ). Here, we note that  $f^*$  and  $\mathbf{x}^*$  are the limiting points of the objective function  $f(\mathbf{x}^{(k)})$  and the parameter  $\mathbf{x}^{(k)}$  as the number of iterations  $k$  goes to infinity, respectively. In (i), the convergence of the generalized gradient norm to 0 implies the stationarity condition of the constrained problem is satisfied. It follows that the algorithm converges to a set of stationary points of the problem. In (ii), the convergence on the function side is often obtained via the monotonicity of the objective-value sequence  $\{f(\mathbf{x}^{(k)})\}_{k=0}^{\infty}$  (e.g., decreasing to a limiting value  $f^*$ ). This in turn implies the sequence  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  converges to a set of local optima that yields the

same objective function value  $f^*$ .<sup>1</sup> In (iii), the convergence of  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$  implies convergence to a unique point that is often an *isolated local optimum point* of the problem. Typically, convergence on the domain side is used in linear convergence proofs for strongly convex settings.

#### B. Convergence to a Global Optimum

In general, a stationary point can be a saddle point, a local/global minimum, or local/global maximum of the problem. When both the objective function and the constraint set are convex, it is well-known that all stationary points are also global optima of the problem. Convergence analysis of iterative algorithm (e.g., PGD) in convex optimization therefore focus on providing a universal upper bound on the distance to the global solutions. Analysis on the domain side (iii) is usually used in the presence of *strong convexity* that guarantees the *uniqueness* of the global optimum [1]-Section 8.6. Without the strong convexity, one may resort to analysis on the function side (ii) in order to prove convergence to a set of global optima [6]-Section 10.4.3. When convexity is not guaranteed, due to a non-convex objective and/or a non-convex constraint set, convergence analysis has recourse to a set of stationary points by bounding the generalized gradient norm through iterations (i) [3]-Section 2.3.2. Notwithstanding, recent advances in structured non-convex optimization have shed light on convergence guarantees to global solutions of the problem. By exploiting the special structure of some classes of non-convex problems and using appropriate initialization, PGD can be shown to converge to a unique global optimum despite the non-convexity of these problems. Examples of such powerful results include sparse recovery with restricted isometry properties [7], matrix completion with incoherence properties [8], empirical risk minimization with restricted strong convexity and smoothness properties [9], and spherically constrained quadratic minimization with hidden convexity [10].

#### C. Convergence to a Local Optimum

In general non-convex settings, domain-side convergence analysis is restricted to the local region around the convergence point  $\mathbf{x}^*$ . Such points can be a saddle point, a

<sup>1</sup>An example for such scenario is minimizing a convex but not strongly convex function  $f(\mathbf{x}) = \|\mathbf{x}\|_1$  subject to  $\mathbf{x} \in \mathbb{R}^n$  and  $\|\mathbf{x}\|_2^2 = 1$ . The  $2n$  vectors  $\{\mathbf{e}_i\}_{i=1}^n$  and  $\{-\mathbf{e}_i\}_{i=1}^n$  are local minimizers that obtain the same objective function value. It is worthwhile mentioning that they are also the global solutions of the foregoing problem.

local minimum, or a local maximum of the problem. The ROC associated with  $\mathbf{x}^*$  is the neighborhood in which the algorithm (e.g., PGD) is guaranteed to converge to  $\mathbf{x}^*$  when initialized inside this region. To a certain extent, the ROC in the aforementioned global convergence analysis is the entire feasible space. However, while global convergence analysis does not require the initialization to be close to the global solution, it often ignores the local structure near the solution needed for establishing sharp bounds on the speed of convergence. In particular, bounding techniques employed in global convergence analysis hold universally, including worst-case scenarios. Thus, in many problem-specific settings where the solution lies in a benign neighborhood, the global analysis could lead to conservative convergence rate bounds. As an illustration, in minimizing a smooth and strongly convex function  $f$ , gradient descent with a fixed step size achieves the rate of convergence at most  $(\kappa - 1)/(\kappa + 1)$  [11], where  $\kappa$  is the (global) condition number of  $f$ . Recall that the condition number of a differentiable convex function is the ratio of its smoothness  $L$  to strong convexity  $\mu$  [4]. For any quadratic function, this global bound is also an exact and attainable estimate thanks to the fact that the objective curvature is unchanged everywhere. For non-quadratic objectives, on the other hand, this global bound may be loose as  $\kappa$  takes into account the worst-case scenario, in which the objective function is most ill-conditioned. The asymptotic behavior of gradient descent near the solution indeed relies on the condition number of the local Hessian  $\kappa(\mathbf{x}^*)$  of the objective function, defining as  $\lambda_{\max}(\nabla^2 f(\mathbf{x}^*))/\lambda_{\min}(\nabla^2 f(\mathbf{x}^*))$ . Generally, we have  $\mu \leq \lambda_{\min}(\nabla^2 f(\mathbf{x})) \leq \lambda_{\max}(\nabla^2 f(\mathbf{x})) \leq L$ , for any  $\mathbf{x}$  in the domain of  $f$ , which implies  $\kappa(\mathbf{x}^*) \leq \kappa$ . This local condition number  $\kappa(\mathbf{x}^*)$  can be significantly smaller than the global condition number  $\kappa$  and hence, *a local convergence analysis can yield a tighter bound that reflects the actual convergence speed of the algorithm near the solution*. Similar situation also occurs for constrained least squares in which the Hessian restricted to the constrained set can depend on the local structure of the set.

#### D. Speed of Convergence

To illustrate the concept of convergence speed, let us consider the convergence on the domain side, i.e., the distance  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ . Let  $\mu$  be a number between 0 and 1. The convergence of  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  to  $\mathbf{x}^*$  is said to be at rate  $\mu \triangleq \mu(\{\mathbf{x}^{(k)}\}_{k=0}^{\infty})$  if  $\mu = \inf_{\{\epsilon_k\}_{k=0}^{\infty}} \lim_{k \rightarrow \infty} \epsilon_{k+1}/\epsilon_k$ , for any monotonically decreasing sequence  $\{\epsilon_k\}_{k=0}^{\infty}$  satisfying  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \epsilon_k$  for all index  $k$ . The asymptotic rate of convergence of gradient descent to  $\mathbf{x}^*$ , denoted by  $\rho$ , is defined by the worst-case rate of convergence among all possible sequences  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  that are generated by the algorithm and converge to  $\mathbf{x}^*$ , i.e.,  $\rho = \sup_{\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}} \mu(\{\mathbf{x}^{(k)}\}_{k=0}^{\infty})$ . Depending on the value of  $\rho$  in the interval  $[0, 1]$ , the convergence is said to be *sublinear* when  $\rho = 1$ , *linear* when  $0 < \rho < 1$ , or *superlinear* when  $\rho = 0$ . The lower the value of  $\rho$  is, the faster the speed of convergence is and the fewer the number of iterations needed is to obtain a close approximation of the solution. Thus, analytical estimation of the convergence

rate plays a pivotal role in convergence analysis. We would like to note two distinct methods for linear convergence rate analysis dating back to the 1960s. The first approach was proposed by Polyak [2], based on his earlier study into nonlinear difference equations [12]. The author analyzed the asymptotic convergence of gradient descent for minimizing some objective function  $f$ . Assuming  $\mathbf{x}^*$  is a non-singular local minimum of  $f$ , Polyak showed that for any  $\delta > 0$ , there exists  $\epsilon > 0$  such that if  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| < \epsilon$  then the sequence  $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$  generated by gradient descent satisfies

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \|\mathbf{x}^{(0)} - \mathbf{x}^*\| (\rho + \delta)^k, \quad (1)$$

where  $\rho = \max\{|1 - \eta\lambda_{\max}|, |1 - \eta\lambda_{\min}|\}$  and  $\lambda_{\max}$  and  $\lambda_{\min}$  are the largest and smallest eigenvalues of  $\nabla^2 f(\mathbf{x}^*)$ , respectively. Here we emphasize that  $f$  does not need to be smooth and strongly convex everywhere but only so around  $\mathbf{x}^*$ . By setting  $\eta_{opt} = 2/(\lambda_{\max} + \lambda_{\min})$ , the optimal rate of convergence is given by  $\rho_{opt} = (\kappa^* - 1)/(\kappa^* + 1)$ , where  $\kappa^* = \lambda_{\max}/\lambda_{\min}$  is the condition number of the local Hessian  $\nabla^2 f(\mathbf{x}^*)$ . When  $f$  is a strongly convex quadratic, the local result coincides with the aforementioned global result in [11] ( $\kappa^* = \kappa$ ). The expression of  $\rho$  in (1) is called the **asymptotic convergence rate** of gradient descent with fixed step size  $\eta$ .<sup>2</sup> The second approach was developed by Daniel [13] in 1967, while studying gradient descent with *exact line search*, i.e., choosing  $\eta$  that minimizes the objective at each iteration. Utilizing the Kantorovich inequality [14], the author proved that if  $\mathbf{x}^{(0)}$  is sufficiently close to  $\mathbf{x}^*$ , there exist a constant  $\epsilon$  and a sequence  $\{q_k\}_{k=0}^{\infty}$  such that

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq \epsilon \prod_{i=0}^k q_i, \quad \lim_{k \rightarrow \infty} q_k = (\kappa^* - 1)/(\kappa^* + 1).$$

Note that here the characteristics of convergence are also exploited through the Hessian  $\nabla^2 f(\mathbf{x}^*)$ . This result was then extended to study the asymptotic convergence of projected gradient descent for constrained optimization [15]–[17].

## II. PROOF OF EXAMPLE 1

Our goal in this proof is to establish the Lipschitz differentiability of the projection operator onto the unit sphere  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| = 1\}$ . We start by establishing the Lipschitz differentiability at a point on  $\mathcal{C}$  and then extend it to any nonzero point in  $\mathbb{R}^n$ . For the Lipschitz differentiability on  $\mathcal{C}$ , we introduce the following lemma:

**Lemma 1.** *For any  $\mathbf{x}^* \in \mathcal{C}$ , we have*

$$\sup_{\mathbf{y} \in \Pi_{\mathcal{C}}(\mathbf{x}^* + \delta)} \|\mathbf{y} - \mathbf{x}^* - (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^T)\delta\| \leq 2\|\delta\|^2. \quad (2)$$

*Proof.* We consider two cases:

**Case 1:** If  $\mathbf{x}^* + \delta = \mathbf{0}$ , then  $\Pi_{\mathcal{C}}(\mathbf{0}) = \mathcal{C}$  and  $\|\delta\| = \|\mathbf{x}^*\| = 1$ . For any  $\mathbf{y} \in \mathcal{C}$ , substituting  $\delta = -\mathbf{x}^*$  and then using the fact

<sup>2</sup>It is worthwhile to mention that using a similar technique, Nesterov [4] proved that the asymptotic rate is at most  $\hat{\rho} = (\kappa^* + 1)/(\kappa^* + 3)$ . While this bound also exploits the local information of the optimization problem, we note that it is not as tight as the bound in (1).

that  $\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top$  is the projection onto the null space of  $\mathbf{x}^*$ , we have

$$\begin{aligned} \mathbf{y} - \mathbf{x}^* - (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top)\boldsymbol{\delta} &= \mathbf{y} - \mathbf{x}^* + (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top)\mathbf{x}^* \\ &= \mathbf{y} - \mathbf{x}^*. \end{aligned}$$

Next, taking the norm and using the triangle inequality yield

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}^* - (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top)\boldsymbol{\delta}\| &= \|\mathbf{y} - \mathbf{x}^*\| \\ &\leq \|\mathbf{y}\| + \|\mathbf{x}^*\| = 2\|\boldsymbol{\delta}\|^2, \end{aligned}$$

where the last step stems from  $\|\mathbf{y}\| = \|\mathbf{x}^*\| = \|\boldsymbol{\delta}\| = 1$ . Thus, (2) holds in this case.

**Case 2:** If  $\mathbf{x}^* + \boldsymbol{\delta} \neq \mathbf{0}$ , then  $\Pi_C(\mathbf{x}^* + \boldsymbol{\delta})$  is singleton containing the unique projection

$$\mathcal{P}_C(\mathbf{x}^* + \boldsymbol{\delta}) = \frac{\mathbf{x}^* + \boldsymbol{\delta}}{\|\mathbf{x}^* + \boldsymbol{\delta}\|}.$$

Hence, (2) is equivalent to

$$\left\| \frac{\mathbf{x}^* + \boldsymbol{\delta}}{\|\mathbf{x}^* + \boldsymbol{\delta}\|} - \mathbf{x}^* - (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top)\boldsymbol{\delta} \right\| \leq 2\|\boldsymbol{\delta}\|^2. \quad (3)$$

We prove (3) by (i) showing that for any scalars  $u > 0$  and  $(1-u)^2 \leq v \leq (1+u)^2$ :

$$(17u-2)v^2 - 2u(1-u)^2v + (1-u)^4(u+2) \geq 0, \quad (4)$$

and (ii) showing that (4) is equivalent to (3) with  $u = \|\mathbf{x}^* + \boldsymbol{\delta}\| > 0$  and  $v = \|\boldsymbol{\delta}\|^2 \geq 0$ .

(i) To prove (4), let us consider the following cases:

1) If  $0 < u \leq 2/17$ , then for  $v \leq (1+u)^2$ , we have

$$\begin{aligned} (17u-2)v^2 - 2u(1-u)^2v + (1-u)^4(u+2) &\geq (17u-2)(1+u)^4 - 2u(1-u)^2(1+u)^2 \\ &\quad + (1-u)^4(u+2) \\ &= 16u^2(u+2)(u^2+2u+2) \geq 0. \end{aligned}$$

2) If  $2/17 < u \leq 1/2$ , then for  $(1-u)^2 \leq v \leq (1+u)^2$ , the following holds

$$\begin{aligned} (17u-2)v^2 - 2u(1-u)^2v + (1-u)^4(u+2) &\geq (17u-2)(1-u)^4 - 2u(1-u)^2(1+u)^2 \\ &\quad + (1-u)^4(u+2) \\ &= 8u(1-u)^2(2-u)(1-2u) \geq 0. \end{aligned}$$

3) If  $u > 1/2$ , using the quadratic vertex at  $v = u(1-u)^2/(17u-2)$  as the minimum point, we obtain

$$\begin{aligned} (17u-2)v^2 - 2u(1-u)^2v + (1-u)^4(u+2) &\geq \frac{4(1-u)^4(4u^2+8u-1)}{17u-2} \geq 0. \end{aligned}$$

(ii) Now for  $u = \|\mathbf{x}^* + \boldsymbol{\delta}\| > 0$  and  $v = \|\boldsymbol{\delta}\|^2 \geq 0$ , we have  $(\mathbf{x}^*)^\top\boldsymbol{\delta} = (u^2 - v - 1)/2$  and

$$\begin{aligned} (3) &\Leftrightarrow \left\| \frac{\mathbf{x}^* + \boldsymbol{\delta}}{\|\mathbf{x}^* + \boldsymbol{\delta}\|} - \mathbf{x}^* - (\mathbf{I} - \mathbf{x}^*(\mathbf{x}^*)^\top)\boldsymbol{\delta} \right\| \leq 2\|\boldsymbol{\delta}\|^2 \\ &\Leftrightarrow \|\mathbf{x}^* + \boldsymbol{\delta} - \|\mathbf{x}^* + \boldsymbol{\delta}\|(\mathbf{x}^* + \boldsymbol{\delta} - \mathbf{x}^*(\mathbf{x}^*)^\top\boldsymbol{\delta})\|^2 \\ &\quad \leq 4\|\mathbf{x}^* + \boldsymbol{\delta}\|^2\|\boldsymbol{\delta}\|^4 \\ &\Leftrightarrow \|(1-u)(\mathbf{x}^* + \boldsymbol{\delta}) + u((\mathbf{x}^*)^\top\boldsymbol{\delta})\mathbf{x}^*\|_2^2 \leq 4u^2v^2 \\ &\Leftrightarrow (1-u)^2u^2 + u^2\left(\frac{u^2-v-1}{2}\right)^2 \\ &\quad + 2u(1-u)\frac{u^2-v-1}{2}\frac{u^2-v+1}{2} \leq 4u^2v^2 \\ &\Leftrightarrow (4). \end{aligned}$$

Finally, by the triangle inequality, we have

$$\|\|\mathbf{x}^* + \boldsymbol{\delta}\| - \|\mathbf{x}^*\|\| \leq \|\boldsymbol{\delta}\| \leq \|\mathbf{x}^*\| + \|\mathbf{x}^* + \boldsymbol{\delta}\|,$$

which in turn verifies  $(1-u)^2 \leq v \leq (1+u)^2$ . This completes our proof of the lemma.  $\square$

Next, to extend the result in Lemma 1 to any  $\mathbf{x} \in \mathbb{R} \setminus \{0\}$ , we substitute  $\mathbf{x}^* = \mathbf{x}/\|\mathbf{x}\|$  and  $\boldsymbol{\delta} = \boldsymbol{\delta}/\|\mathbf{x}\|$  into (2) and obtain

$$\sup_{\mathbf{y} \in \Pi_C\left(\frac{\mathbf{x}}{\|\mathbf{x}\|} + \frac{\boldsymbol{\delta}}{\|\mathbf{x}\|}\right)} \left\| \mathbf{y} - \frac{\mathbf{x}}{\|\mathbf{x}\|} - \left(\mathbf{I}_n - \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|^2}\right) \frac{\boldsymbol{\delta}}{\|\mathbf{x}\|} \right\| \leq 2\frac{\|\boldsymbol{\delta}\|^2}{\|\mathbf{x}\|^2}. \quad (5)$$

Since the projection onto the unit sphere is scale-invariant,

$$\Pi_C\left(\frac{\mathbf{x}}{\|\mathbf{x}\|} + \frac{\boldsymbol{\delta}}{\|\mathbf{x}\|}\right) = \Pi_C(\mathbf{x} + \boldsymbol{\delta}). \quad (6)$$

Substituting (6) into (5) yields (6). Thus, by Definition 2, for any  $\mathbf{x} \neq \mathbf{0}$  we obtain

$$\begin{aligned} \nabla\mathcal{P}_C(\mathbf{x}) &= \frac{1}{\|\mathbf{x}\|} \left( \mathbf{I}_n - \frac{\mathbf{x}\mathbf{x}^\top}{\|\mathbf{x}\|^2} \right), \\ c_1(\mathbf{x}) &= \infty, \quad c_2(\mathbf{x}) = \frac{2}{\|\mathbf{x}\|^2}. \end{aligned}$$

### III. DETAILS OF APPLICATION B - ITERATIVE HARD THRESHOLDING FOR SPARSE RECOVERY

#### A. Proof of (42)

In this subsection, we first show that any  $\mathbf{x}^* \in \Phi_{=s}$  and  $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, |x_{[s]}^*|/\sqrt{2})$  share the same index set of  $s$ -largest elements (in magnitude), i.e.,  $\Omega_s(\mathbf{x}^*)$ . Then, we construct a counter-example to demonstrate that  $|x_{[s]}^*|/\sqrt{2}$  is the largest possible radius so that (42) holds.

First, we show that for any  $i \in \Omega_s(\mathbf{x}^*)$  and  $j \in \{1, \dots, n\} \setminus \Omega_s(\mathbf{x}^*)$ ,  $|x_j| < |x_i|$  as follows. In particular, we have

$$\begin{aligned} |x_j - x_j^*| + |x_i - x_i^*| &\leq \sqrt{2((x_j - x_j^*)^2 + (x_i - x_i^*)^2)} \\ &\leq \sqrt{2\|\mathbf{x} - \mathbf{x}^*\|^2} < |x_{[s]}^*|, \end{aligned}$$

where the last inequality stems from the fact that  $\|\mathbf{x} - \mathbf{x}^*\| < |x_{[s]}^*|/\sqrt{2}$ . Now, since  $x_j^* = 0$  for all  $j \in \{1, \dots, n\} \setminus \Omega_s(\mathbf{x}^*)$ , we have

$$\begin{aligned} |x_j| &= |x_j - x_j^*| \\ &< |x_{[s]}^*| - |x_i - x_i^*| \\ &\leq |x_i^*| - |x_i - x_i^*| \\ &\leq |x_i^* + (x_i - x_i^*)| = |x_i|, \end{aligned} \quad (7)$$

Therefore, every  $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, |x_{[s]}^*|/\sqrt{2})$  shares the same index set of  $s$ -largest (in magnitude) elements with  $\mathbf{x}^*$ , i.e.,  $\Omega_s(\mathbf{x}) = \Omega_s(\mathbf{x}^*)$ , which implies (42).

We now construct the counter-example as a point  $\mathbf{x}$  such that  $\Omega_s(\mathbf{x}) \neq \Omega_s(\mathbf{x}^*)$  and  $\mathbf{x}$  is not in  $\mathcal{B}(\mathbf{x}^*, |x_{[s]}^*|/\sqrt{2})$  but arbitrarily close to its boundary. Without loss of generality, assume that  $|x_1^*| \geq \dots \geq |x_s^*| > |x_{s+1}^*| = \dots = |x_n^*| = 0$ . For arbitrarily small  $\epsilon > 0$ , define  $\mathbf{x}$  as

$$x_i = \begin{cases} x_s^*/2 & \text{if } i = s, \\ x_s^*/2 + \epsilon & \text{if } i = s + 1, \\ x_i & \text{otherwise.} \end{cases}$$

Then, since  $x_{s+1} < x_s$ ,  $\mathbf{x}$  does not shares the same index set of  $s$ -largest (in magnitude) elements with  $\mathbf{x}^*$ . On the other hand, as  $\epsilon \rightarrow 0$ , we have

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^*\| &= \sqrt{\sum_{i=1}^n (x_i - x_i^*)^2} \\ &= \sqrt{\left(-\frac{x_s^*}{2}\right)^2 + \left(\frac{x_s^*}{2} + \epsilon\right)^2} \rightarrow \frac{1}{\sqrt{2}}|x_{[s]}^*|. \end{aligned}$$

This means  $\mathbf{x} \notin \mathcal{B}(\mathbf{x}^*, |x_{[s]}^*|/\sqrt{2})$  but it can approach the boundary of the ball as  $\epsilon$  decreases to 0.

### B. Proof of Remark 6

In the following, we show any stationary point  $\mathbf{x}^*$  of (40) is also a local minimum by proving that the objective function does not decrease if we add any perturbation to  $\mathbf{x}^*$  on  $\mathcal{C}$ . Let us consider any perturbation  $\delta$  such that  $\delta \in \mathcal{B}(\mathbf{0}, c_1(\mathbf{x}^*))$  and  $\mathbf{x} = \mathbf{x}^* + \delta \in \mathcal{C}$ . Since  $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, c_1(\mathbf{x}^*))$ , using (7), we have  $|x_{[1]}| \geq \dots \geq |x_{[s]}| > 0$ . On the other hand, since  $\mathbf{x}$  has no more than  $s$  non-zero entries, it must hold that  $|x_{[s+1]}| = \dots = |x_{[n]}| = 0$ . Therefore,  $\mathbf{x} = \mathbf{S}_{\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*}^\top \mathbf{x}$ , which implies  $\delta = \mathbf{S}_{\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*}^\top \delta$ . Now we represent the change in the objective function as

$$\begin{aligned} &\frac{1}{2} \|\mathbf{A}(\mathbf{x}^* + \delta) - \mathbf{b}\|^2 - \frac{1}{2} \|\mathbf{A}\mathbf{x}^* - \mathbf{b}\|^2 \\ &= \frac{1}{2} \delta^\top \mathbf{A}^\top \mathbf{A} \delta + \delta^\top \mathbf{A}^\top (\mathbf{A}\mathbf{x}^* - \mathbf{b}) \\ &= \frac{1}{2} \delta^\top \mathbf{S}_{\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*}^\top \mathbf{A}^\top \mathbf{A} \mathbf{S}_{\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*}^\top \delta + \delta^\top \mathbf{S}_{\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*}^\top \mathbf{A}^\top (\mathbf{A}\mathbf{x}^* - \mathbf{b}) \\ &= \frac{1}{2} \delta^\top \mathbf{S}_{\mathbf{x}^*} (\mathbf{S}_{\mathbf{x}^*}^\top \mathbf{A}^\top \mathbf{A} \mathbf{S}_{\mathbf{x}^*}) \mathbf{S}_{\mathbf{x}^*}^\top \delta \geq 0, \end{aligned} \quad (8)$$

where the last equality uses the stationarity condition in (43). From (8), we conclude  $\mathbf{x}^*$  is a local minimum of (40).

### REFERENCES

- [1] D. G. Luenberger, Y. Ye *et al.*, *Linear and nonlinear programming*. Springer, 1984, vol. 2.
- [2] B. T. Polyak, *Introduction to optimization*. Optimization software. Inc., Publications Division, New York, 1987, vol. 1.
- [3] D. P. Bertsekas, "Nonlinear programming," *J. Oper. Res. Soc.*, vol. 48, no. 3, pp. 334–334, 1997.
- [4] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2003, vol. 87.
- [5] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [6] A. Beck, *First-order methods in optimization*. SIAM, 2017.
- [7] T. Blumensath and M. E. Davies, "Iterative hard thresholding for compressed sensing," *Appl. Comput. Harmon. Anal.*, vol. 27, no. 3, pp. 265–274, 2009.
- [8] C. Ma, K. Wang, Y. Chi, and Y. Chen, "Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution," in *Found. Comput. Math.*, 2020, pp. 451–632.
- [9] R. Khanna and A. Kyrillidis, "IHT dies hard: Provable accelerated iterative hard thresholding," in *Proc. Int. Conf. Artif. Intell. Stat.*, 2018, pp. 188–198.
- [10] A. Beck and Y. Vaisbourd, "Globally solving the trust region subproblem using simple first-order methods," *SIAM J. Optim.*, vol. 28, no. 3, pp. 1951–1967, 2018.
- [11] B. T. Polyak, "Gradient methods for minimizing functionals," *Zhurnal vychislitel'noi matematiki i matematicheskoi fiziki*, vol. 3, no. 4, pp. 643–653, 1963.
- [12] —, "Some methods of speeding up the convergence of iteration methods," *USSR Comput. Math. Math. Phys.*, vol. 4, no. 5, pp. 1–17, 1964.
- [13] J. W. Daniel, "The conjugate gradient method for linear and nonlinear operator equations," *SIAM J. Numer. Anal.*, vol. 4, no. 1, pp. 10–26, 1967.
- [14] L. V. Kantorovich, "Functional analysis and applied mathematics," *Uspekhi Matematicheskikh Nauk*, vol. 3, no. 6, pp. 89–185, 1948.
- [15] D. G. Luenberger, "The gradient projection method along geodesics," *Manag. Sci.*, vol. 18, no. 11, pp. 620–631, 1972.
- [16] A. Lichnewsky, "Minimisation des fonctionnelles définies sur une variété par la méthode du gradient conjugué," Ph.D. dissertation, These de Doctorat d'Etat. Paris: Université de Paris-Sud, 1979.
- [17] D. Gabay, "Minimizing a differentiable function over a differential manifold," *J. Optim. Theory Appl.*, vol. 37, no. 2, pp. 177–219, 1982.