

A NEWTON-CG BASED AUGMENTED LAGRANGIAN METHOD  
FOR FINDING A SECOND-ORDER STATIONARY POINT OF  
NONCONVEX EQUALITY CONSTRAINED OPTIMIZATION WITH  
COMPLEXITY GUARANTEES\*

CHUAN HE<sup>†</sup>, ZHAOSONG LU<sup>†</sup>, AND TING KEI PONG<sup>‡</sup>

**Abstract.** In this paper we consider finding a second-order stationary point (SOSP) of nonconvex equality constrained optimization when a nearly feasible point is known. In particular, we first propose a new Newton-CG method for finding an approximate SOSP of unconstrained optimization and show that it enjoys a substantially better complexity than the Newton-CG method in [C. W. Royer, M. O’Neill, and S. J. Wright, Math. Program., 180 (2020), pp. 451–488]. We then propose a Newton-CG based augmented Lagrangian (AL) method for finding an approximate SOSP of nonconvex equality constrained optimization, in which the proposed Newton-CG method is used as a subproblem solver. We show that under a generalized linear independence constraint qualification (GLICQ), our AL method enjoys a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2} \min\{n, \epsilon^{-3/4}\})$  for finding an  $(\epsilon, \sqrt{\epsilon})$ -SOSP of nonconvex equality constrained optimization with high probability, which are significantly better than the ones achieved by the proximal AL method in [Y. Xie and S. J. Wright, J. Sci. Comput., 86 (2021), pp. 1–30]. Besides, we show that it has a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-5/4}\})$  when the GLICQ does not hold. To the best of our knowledge, all the complexity results obtained in this paper are new for finding an approximate SOSP of nonconvex equality constrained optimization with high probability. Preliminary numerical results also demonstrate the superiority of our proposed methods over the other competing algorithms.

23 **Key words.** Nonconvex equality constrained optimization, second-order stationary point,  
24 augmented Lagrangian method, Newton-conjugate gradient method, iteration complexity, operation  
25 complexity

26 MSC codes. 49M15, 68Q25, 90C06, 90C26, 90C30, 90C60

**1. Introduction.** In this paper we consider nonconvex equality constrained optimization problem

$$29 \quad (1.1) \quad \min_{x \in \mathbb{R}^n} f(x) \quad \text{s. t. } c(x) = 0,$$

30 where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are twice continuously differentiable, and we  
 31 assume that problem (1.1) has at least one optimal solution. Since (1.1) is a nonconvex  
 32 optimization problem, it may have many local but non-global minimizers and finding  
 33 its global minimizer is generally NP-hard. A first-order stationary point (FOSP) of it  
 34 is usually found in practice instead. Nevertheless, a mere FOSP may sometimes not  
 35 suit our needs and a *second-order stationary point* (SOSP) needs to be sought. For  
 36 example, in the context of linear semidefinite programming (SDP), a powerful approach  
 37 to solving it is by solving an equivalent nonconvex equality constrained optimization  
 38 problem [17, 18]. It was shown in [18, 15] that under some mild conditions an SOSP  
 39 of the latter problem can yield an optimal solution of the linear SDP, while a mere  
 40 FOSP generally cannot. It is therefore important to find an SOSP of problem (1.1).

\*Submitted to the editors January 11, 2023.

**Funding:** The work of the second author was partially supported by NSF Award IIS-2211491. The work of the third author was partially supported by a Research Scheme of the Research Grants Council of Hong Kong SAR, China (Project No. T22-504/21R).

<sup>†</sup>Department of Industrial and Systems Engineering, University of Minnesota, USA  
(he000233@umn.edu, zhaosong@umn.edu).

<sup>‡</sup>Department of Applied Mathematics, the Hong Kong Polytechnic University, Hong Kong, People's Republic of China (tk.pong@polyu.edu.hk)

41 In recent years, numerous methods with complexity guarantees have been developed  
 42 for finding an approximate SOSP of several types of nonconvex optimization. For  
 43 example, cubic regularized Newton methods [52, 25, 1, 22], accelerated gradient  
 44 methods [23, 24], trust-region methods [34, 35, 50], quadratic regularization method  
 45 [12], second-order line-search method [57], and Newton-conjugate gradient (Newton-  
 46 CG) method [56] were developed for nonconvex unconstrained optimization. In  
 47 addition, interior-point method [8] and log-barrier method [54] were proposed for  
 48 nonconvex optimization with sign constraints. The interior-point method [8] was also  
 49 generalized in [38] to solve nonconvex optimization with sign constraints and additional  
 50 linear equality constraints. Furthermore, a projected gradient descent method with  
 51 random perturbations was proposed in [47] for nonconvex optimization with linear  
 52 inequality constraints. Iteration complexity was established for these methods for  
 53 finding an approximate SOSP. Besides, operation complexity measured by the amount  
 54 of fundamental operations such as gradient evaluations and matrix-vector products  
 55 was also studied in [1, 23, 34, 41, 24, 57, 22, 56].

56 Several methods including trust-region methods [21, 33], sequential quadratic  
 57 programming method [14], two-phase method [9, 30, 32] and augmented Lagrangian  
 58 (AL) type methods [4, 10, 58, 60] were proposed for finding an SOSP of problem (1.1).  
 59 However, only a few of them have *complexity guarantees* for finding an approximate  
 60 SOSP of (1.1). In particular, the inexact AL method [58] has a worst-case complexity  
 61 in terms of the number of calls to a second-order oracle. Yet its operation complexity,  
 62 measured by the amount of fundamental operations such as gradient evaluations and  
 63 Hessian-vector products, is unknown. To the best of our knowledge, the proximal  
 64 AL method in [60] appears to be the only existing method that enjoys a worst-  
 65 case complexity for finding an approximate SOSP of (1.1) in terms of fundamental  
 66 operations. In this method, given an iterate  $x^k$  and a multiplier estimate  $\lambda^k$  at the  
 67  $k$ th iteration, the next iterate  $x^{k+1}$  is obtained by finding an approximate stochastic  
 68 SOSP of the proximal AL subproblem:

$$69 \quad \min_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda^k; \rho) + \beta \|x - x^k\|^2/2$$

70 for some suitable positive  $\rho$  and  $\beta$  using a Newton-CG method proposed in [56], where  
 71  $\mathcal{L}$  is the AL function of (1.1) defined as

$$72 \quad \mathcal{L}(x, \lambda; \rho) := f(x) + \lambda^T c(x) + \rho \|c(x)\|^2/2.$$

73 Then the multiplier estimate is updated using the classical scheme, i.e.,  $\lambda^{k+1} =$   
 74  $\lambda^k + \rho c(x^{k+1})$  (e.g., see [39, 55]). The authors of [60] studied the worst-case complexity  
 75 of their proximal AL method including: (i) *total inner iteration complexity*, which  
 76 measures the total number of iterations of the Newton-CG method [56] performed in  
 77 their method; (ii) *operation complexity*, which measures the total number of gradient  
 78 evaluations and matrix-vector products involving the Hessian of the AL function that  
 79 are evaluated in their method. Under some suitable assumptions, including that a  
 80 generalized linear independence constraint qualification (GLICQ) holds at all iterates,  
 81 it was established in [60] that their proximal AL method enjoys a total inner iteration  
 82 complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-3/4}\})$  for  
 83 finding an  $(\epsilon, \sqrt{\epsilon})$ -SOSP of problem (1.1) with high probability.<sup>1</sup> Yet, there is a big

<sup>1</sup>In fact, a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7} \min\{n, \epsilon^{-1}\})$  were established in [60] for finding an  $(\epsilon, \epsilon)$ -SOSP of problem (1.1) with high probability; see [60, Theorem 4(ii), Corollary 3(ii), Theorem 5]. Nonetheless, they can be modified to obtain the aforementioned complexity for finding an  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1) with high probability.

84 gap between these complexities and the iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-3/2})$  and the  
 85 operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-3/2} \min\{n, \epsilon^{-1/4}\})$  that are achieved by the methods in  
 86 [1, 24, 57, 56] for finding an  $(\epsilon, \sqrt{\epsilon})$ -SOSP of nonconvex unconstrained optimization  
 87 with high probability, which is a special case of (1.1) with  $c \equiv 0$ . Also, there is a lack  
 88 of complexity guarantees for this proximal AL method when the GLICQ does not  
 89 hold. It shall be mentioned that Newton-CG based AL methods were also developed  
 90 for efficiently solving various convex optimization problems (e.g., see [61, 62]), though  
 91 their complexities remain unknown.

92 In this paper we propose a Newton-CG based AL method for finding an approxi-  
 93 mate SOSP of problem (1.1) with high probability, and study its worst-case complexity  
 94 with and without the assumption of a GLICQ. In particular, we show that this method  
 95 enjoys a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2})$  and an operation complexity  
 96 of  $\tilde{\mathcal{O}}(\epsilon^{-7/2} \min\{n, \epsilon^{-3/4}\})$  for finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1) under the  
 97 GLICQ, which are significantly better than the aforementioned ones achieved by the  
 98 proximal AL method in [60]. Besides, when the GLICQ does not hold, we show that  
 99 it has a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and an operation complexity of  
 100  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-5/4}\})$  for finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1), which fills the  
 101 research gap in this topic. Specifically, our AL method (Algorithm 4.1) proceeds in  
 102 the following manner. Instead of directly solving problem (1.1), it solves a perturbed  
 103 problem of (1.1) with  $c$  replaced by its perturbed counterpart  $\tilde{c}$  constructed by using  
 104 a nearly feasible point of (1.1) (see (4.4) for details). At the  $k$ th iteration, an approxi-  
 105 mate stochastic SOSP  $x^{k+1}$  of the AL subproblem of this perturbed problem is found  
 106 by our newly proposed Newton-CG method (Algorithm 3.1) for a penalty parameter  
 107  $\rho_k$  and a truncated Lagrangian multiplier  $\lambda^k$ , which results from projecting onto a  
 108 Euclidean ball the standard multiplier estimate  $\tilde{\lambda}^k$  obtained by the classical scheme  
 109  $\tilde{\lambda}^k = \lambda^{k-1} + \rho_k \tilde{c}(x^k)$ .<sup>2</sup> The penalty parameter  $\rho_{k+1}$  is then updated by the following  
 110 practical scheme (e.g., see [7, Section 4.2]):

$$\rho_{k+1} = \begin{cases} r\rho_k & \text{if } \|\tilde{c}(x^{k+1})\| > \alpha \|\tilde{c}(x^k)\|, \\ \rho_k & \text{otherwise} \end{cases}$$

111 for some  $r > 1$  and  $\alpha \in (0, 1)$ . It shall be mentioned that in contrast with the classical  
 112 AL method, our method has two distinct features: (i) the values of the AL function  
 113 along the iterates are bounded from above; (ii) the multiplier estimates associated  
 114 with the AL subproblems are bounded. In addition, to solve the AL subproblems with  
 115 better complexity guarantees, we propose a variant of the Newton-CG method in [56]  
 116 for finding an approximate stochastic SOSP of unconstrained optimization, whose  
 117 complexity has significantly less dependence on the Lipschitz constant of the Hessian of  
 118 the objective than that of the Newton-CG method in [56], while improving or retaining  
 119 the same order of dependence on tolerance parameter. Given that such a Lipschitz  
 120 constant is typically large for the AL subproblems, our Newton-CG method (Algorithm  
 121 3.1) is a much more favorable subproblem solver than the Newton-CG method in [56]  
 122 that is used in the proximal AL method in [60] from theoretical complexity perspective.  
 123

124 The main contributions of this paper are summarized below.

125 • We propose a new Newton-CG method for finding an approximate SOSP of  
 126 unconstrained optimization and show that it enjoys an iteration and operation

---

<sup>2</sup>The  $\lambda^k$  obtained by projecting  $\tilde{\lambda}^k$  onto a compact set is also called a safeguarded Lagrangian multiplier in the relevant literature [11, 42, 13], which has been shown to enjoy many practical and theoretical advantages (see [11] for discussions).

complexity with a *quadratic* dependence on the Lipschitz constant of the Hessian of the objective that improves the *cubic* dependence achieved by the Newton-CG method in [56], while improving or retaining the same order of dependence on tolerance parameter. In addition, our complexity results are established under the assumption that the Hessian of the objective is Lipschitz continuous in a convex neighborhood of a level set of the objective. This assumption is weaker than the one commonly imposed for the Newton-CG method in [56] and some other methods (e.g., [12, 35]) that the Hessian of the objective is Lipschitz continuous in a convex set containing this neighborhood and also *all the trial points* arising in the line search or trust region steps of the methods (see Section 3 for more detailed discussion).

- We propose a Newton-CG based AL method for finding an approximate SOSP of nonconvex equality constrained optimization (1.1) with high probability, and study its worst-case complexity with and without the assumption of a GLICQ. Prior to our work, there was no complexity study on finding an approximate SOSP of problem (1.1) without imposing a GLICQ. Besides, under the GLICQ and some other suitable assumptions, we show that our method enjoys a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2} \min\{n, \epsilon^{-3/4}\})$  for finding an  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1) with high probability, which are significantly better than the respective complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-3/4}\})$  achieved by the proximal AL method in [60]. To the best of our knowledge, all the complexity results obtained in this paper are new for finding an approximate SOSP of nonconvex equality constrained optimization with high probability.

For ease of comparison, we summarize in Table 1 the total inner iteration and operation complexity of our AL method and the proximal AL method in [60] for finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of problem (1.1) with or without assuming GLICQ.

TABLE 1  
*Total inner iteration and operation complexity of finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1).*

Method	GLICQ	Total inner iteration complexity	Operation complexity
Proximal AL method [60]	✓	$\mathcal{O}(\epsilon^{-11/2})$	$\mathcal{O}(\epsilon^{-11/2} \min\{n, \epsilon^{-3/4}\})$
Proximal AL method [60]	✗	unknown	unknown
Our AL method	✓	$\tilde{\mathcal{O}}(\epsilon^{-7/2})$	$\tilde{\mathcal{O}}(\epsilon^{-7/2} \min\{n, \epsilon^{-3/4}\})$
Our AL method	✗	$\tilde{\mathcal{O}}(\epsilon^{-11/2})$	$\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-5/4}\})$

It shall be mentioned that there are many works other than [60] studying complexity of AL methods for nonconvex constrained optimization. However, they aim to find an approximate FOSP rather than SOSP of the problem (e.g., see [40, 37, 13, 51, 45]). Since our main focus is on the complexity of finding an approximate SOSP by AL methods, we do not include them in the above table for comparison.

The rest of this paper is organized as follows. In Section 2, we introduce some notation and optimality conditions. In Section 3, we propose a Newton-CG method for unconstrained optimization and study its worst-case complexity. In Section 4, we propose a Newton-CG based AL method for (1.1) and study its worst-case complexity. We present numerical results and the proof of the main results in Sections 5 and 6, respectively. In Section 7, we discuss some future research directions.

**2. Notation and preliminaries.** Throughout this paper, we let  $\mathbb{R}^n$  denote the  $n$ -dimensional Euclidean space. We use  $\|\cdot\|$  to denote the Euclidean norm of a vector or the spectral norm of a matrix. For a real symmetric matrix  $H$ , we use  $\lambda_{\min}(H)$

168 to denote its minimum eigenvalue. The Euclidean ball centered at the origin with  
 169 radius  $R \geq 0$  is denoted by  $\mathcal{B}_R := \{x : \|x\| \leq R\}$ , and we use  $\Pi_{\mathcal{B}_R}(v)$  to denote the  
 170 Euclidean projection of a vector  $v$  onto  $\mathcal{B}_R$ . For a given finite set  $\mathcal{A}$ , we let  $|\mathcal{A}|$  denote  
 171 its cardinality. For any  $s \in \mathbb{R}$ , we let  $\text{sgn}(s)$  be 1 if  $s \geq 0$  and let it be  $-1$  otherwise.  
 172 In addition,  $\tilde{\mathcal{O}}(\cdot)$  represents  $\mathcal{O}(\cdot)$  with logarithmic terms omitted.

173 Suppose that  $x^*$  is a local minimizer of problem (1.1) and the linear independence  
 174 constraint qualification holds at  $x^*$ , i.e.,  $\nabla c(x^*) := [\nabla c_1(x^*) \ \nabla c_2(x^*) \ \cdots \ \nabla c_m(x^*)]$   
 175 has full column rank. Then there exists a Lagrangian multiplier  $\lambda^* \in \mathbb{R}^m$  such that

$$176 \quad (2.1) \quad \nabla f(x^*) + \nabla c(x^*)\lambda^* = 0,$$

$$177 \quad (2.2) \quad d^T (\nabla^2 f(x^*) + \sum_{i=1}^m \lambda_i^* \nabla^2 c_i(x^*)) d \geq 0, \quad \forall d \in \mathcal{C}(x^*),$$

178 where  $\mathcal{C}(\cdot)$  is defined as

$$179 \quad (2.3) \quad \mathcal{C}(x) := \{d \in \mathbb{R}^n : \nabla c(x)^T d = 0\}.$$

180 The relations (2.1) and (2.2) are respectively known as the first- and second-order  
 181 optimality conditions for (1.1) in the literature (e.g., see [53]). Note that it is in  
 182 general impossible to find a point that exactly satisfies (2.1) and (2.2). Thus, we  
 183 are instead interested in finding a point that satisfies their approximate counterparts.  
 184 In particular, we introduce the following definitions of an approximate first-order  
 185 stationary point (FOSP) and second-order stationary point (SOSP), which are similar  
 186 to those considered in [4, 10, 60]. The rationality of them can be justified by the study  
 187 of the sequential optimality conditions for constrained optimization [3, 4].

188 **DEFINITION 2.1 ( $\epsilon_1$ -first-order stationary point).** *Let  $\epsilon_1 > 0$ . We say that  
 189  $x \in \mathbb{R}^n$  is an  $\epsilon_1$ -first-order stationary point ( $\epsilon_1$ -FOSP) of problem (1.1) if it, together  
 190 with some  $\lambda \in \mathbb{R}^m$ , satisfies*

$$191 \quad (2.4) \quad \|\nabla f(x) + \nabla c(x)\lambda\| \leq \epsilon_1, \quad \|c(x)\| \leq \epsilon_1.$$

192 **DEFINITION 2.2 (( $\epsilon_1, \epsilon_2$ )-second-order stationary point).** *Let  $\epsilon_1, \epsilon_2 > 0$ . We  
 193 say that  $x \in \mathbb{R}^n$  is an  $(\epsilon_1, \epsilon_2)$ -second-order stationary point (( $\epsilon_1, \epsilon_2$ )-SOSP) of problem  
 194 (1.1) if it, together with some  $\lambda \in \mathbb{R}^m$ , satisfies (2.4) and additionally*

$$195 \quad (2.5) \quad d^T (\nabla^2 f(x) + \sum_{i=1}^m \lambda_i \nabla^2 c_i(x)) d \geq -\epsilon_2 \|d\|^2, \quad \forall d \in \mathcal{C}(x),$$

196 where  $\mathcal{C}(\cdot)$  is defined as in (2.3).

197 **3. A Newton-CG method for unconstrained optimization.** In this section  
 198 we propose a variant of Newton-CG method [56, Algorithm 3] for finding an approxi-  
 199 mate SOSP of a class of unconstrained optimization problems, which will be used as a  
 200 subproblem solver for the AL method proposed in the next section. In particular, we  
 201 consider an unconstrained optimization problem

$$202 \quad (3.1) \quad \min_{x \in \mathbb{R}^n} F(x),$$

203 where the function  $F$  satisfies the following assumptions.

204 *Assumption 3.1.* (a) The level set  $\mathcal{L}_F(u^0) := \{x : F(x) \leq F(u^0)\}$  is compact for  
 205 some  $u^0 \in \mathbb{R}^n$ .

206 (b) The function  $F$  is twice Lipschitz continuously differentiable in a convex open  
 207 neighborhood, denoted by  $\Omega$ , of  $\mathcal{L}_F(u^0)$ , that is, there exists  $L_H^F > 0$  such that

$$208 \quad (3.2) \quad \|\nabla^2 F(x) - \nabla^2 F(y)\| \leq L_H^F \|x - y\|, \quad \forall x, y \in \Omega.$$

209 By Assumption 3.1, there exist  $F_{\text{low}} \in \mathbb{R}$ ,  $U_g^F > 0$  and  $U_H^F > 0$  such that

210 (3.3) 
$$F(x) \geq F_{\text{low}}, \quad \|\nabla F(x)\| \leq U_g^F, \quad \|\nabla^2 F(x)\| \leq U_H^F, \quad \forall x \in \mathcal{L}_F(u^0).$$

211 Recently, a Newton-CG method [56, Algorithm 3] was developed to find an  
 212 approximate stochastic SOSP of problem (3.1), which is not only easy to implement  
 213 but also enjoys a nice feature that the main computation consists only of gradient  
 214 evaluations and Hessian-vector products associated with the function  $F$ . Under the  
 215 assumption that  $\nabla^2 F$  is Lipschitz continuous in a convex open set containing  $\mathcal{L}_F(u^0)$   
 216 and also *all the trial points* arising in the line search steps of this method (see [56,  
 217 Assumption 2]), it was established in [56, Theorem 4, Corollary 2] that the iteration  
 218 and operation complexity of this method for finding a stochastic  $(\epsilon_g, \epsilon_H)$ -SOSP of (3.1)  
 219 (namely, a point  $x$  satisfying  $\|\nabla F(x)\| \leq \epsilon_g$  deterministically and  $\lambda_{\min}(\nabla^2 F(x)) \geq -\epsilon_H$   
 220 with high probability) are

221 (3.4) 
$$\mathcal{O}((L_H^F)^3 \max\{\epsilon_g^{-3} \epsilon_H^3, \epsilon_H^{-3}\}) \text{ and } \tilde{\mathcal{O}}((L_H^F)^3 \max\{\epsilon_g^{-3} \epsilon_H^3, \epsilon_H^{-3}\} \min\{n, (U_H^F/\epsilon_H)^{1/2}\}),$$

222 respectively, where  $\epsilon_g, \epsilon_H \in (0, 1)$  are prescribed tolerances. Yet, this assumption can  
 223 be hard to check because these trial points are *unknown* before the method terminates  
 224 and moreover the distance between the origin and them depends on the tolerance  $\epsilon_H$   
 225 in  $\mathcal{O}(\epsilon_H^{-1})$  (see [56, Lemma 3]). In addition, as seen from (3.4), iteration and operation  
 226 complexity of the Newton-CG method in [56] depend *cubically* on  $L_H^F$ . Notice that  $L_H^F$   
 227 can sometimes be very large. For example, the AL subproblems arising in Algorithm 4.1  
 228 have  $L_H^F = \mathcal{O}(\epsilon_1^{-2})$  or  $\mathcal{O}(\epsilon_1^{-1})$ , where  $\epsilon_1 \in (0, 1)$  is a prescribed tolerance for problem  
 229 (1.1) (see Section 4). The cubic dependence on  $L_H^F$  makes such a Newton-CG method  
 230 not appealing as an AL subproblem solver from theoretical complexity perspective.

231 In the rest of this section, we propose a variant of the Newton-CG method [56,  
 232 Algorithm 3] and show that under Assumption 3.1, it enjoys an iteration and operation  
 233 complexity of

234 (3.5) 
$$\mathcal{O}((L_H^F)^2 \max\{\epsilon_g^{-2} \epsilon_H, \epsilon_H^{-3}\}) \text{ and } \tilde{\mathcal{O}}((L_H^F)^2 \max\{\epsilon_g^{-2} \epsilon_H, \epsilon_H^{-3}\} \min\{n, (U_H^F/\epsilon_H)^{1/2}\}),$$

235 for finding a stochastic  $(\epsilon_g, \epsilon_H)$ -SOSP of problem (3.1), respectively. These complexities  
 236 are substantially superior to those in (3.4) achieved by the Newton-CG method in  
 237 [56]. Indeed, the complexities in (3.5) depend quadratically on  $L_H^F$ , while those in  
 238 (3.4) depend cubically on  $L_H^F$ . In addition, it can be verified that they improve or  
 239 retain the order of dependence on  $\epsilon_g$  and  $\epsilon_H$  given in (3.4).

240 **3.1. Main components of a Newton-CG method.** In this subsection we  
 241 briefly discuss two main components of the Newton-CG method in [56], which will be  
 242 used to propose a variant of this method for finding an approximate stochastic SOSP  
 243 of problem (3.1) in the next subsection.

244 The first main component of the Newton-CG method in [56] is a *capped CG method*  
 245 [56, Algorithm 1], which is a modified CG method, for solving a possibly indefinite  
 246 linear system

247 (3.6) 
$$(H + 2\epsilon I)d = -g,$$

248 where  $0 \neq g \in \mathbb{R}^n$ ,  $\epsilon > 0$ , and  $H \in \mathbb{R}^{n \times n}$  is a symmetric matrix. This capped  
 249 CG method terminates within a finite number of iterations. It outputs either an  
 250 approximate solution  $d$  to (3.6) such that  $\|(H + 2\epsilon I)d + g\| \leq \hat{\zeta}\|g\|$  and  $d^T H d \geq$   
 251  $-\epsilon\|d\|^2$  for some  $\hat{\zeta} \in (0, 1)$  or a sufficiently negative curvature direction  $d$  of  $H$  with  
 252  $d^T H d < -\epsilon\|d\|^2$ . The second main component of the Newton-CG method in [56] is

253 a minimum eigenvalue oracle that either produces a sufficiently negative curvature  
 254 direction  $v$  of  $H$  with  $\|v\| = 1$  and  $v^T Hv \leq -\varepsilon/2$  or certifies that  $\lambda_{\min}(H) \geq -\varepsilon$   
 255 holds with high probability. For ease of reference, we present these two components in  
 256 Algorithms A.1 and B.1 in Appendices A and B, respectively.

---

**Algorithm 3.1** A Newton-CG method for problem (3.1)

---

*Input:* Tolerances  $\epsilon_g, \epsilon_H \in (0, 1)$ , backtracking ratio  $\theta \in (0, 1)$ , starting point  $u^0$ , CG-accuracy parameter  $\zeta \in (0, 1)$ , line-search parameter  $\eta \in (0, 1)$ , probability parameter  $\delta \in (0, 1)$ .

Set  $x^0 = u^0$ ;

**for**  $t = 0, 1, 2, \dots$  **do**

**if**  $\|\nabla F(x^t)\| > \epsilon_g$  **then**

    Call Algorithm A.1 with  $H = \nabla^2 F(x^t)$ ,  $\varepsilon = \epsilon_H$ ,  $g = \nabla F(x^t)$ , accuracy parameter  $\zeta$ , and  $U = 0$  to obtain outputs  $d$ , d\_type;

**if** d\_type=NC **then**

$$(3.7) \quad d^t \leftarrow -\text{sgn}(d^T \nabla F(x^t)) \frac{|d^T \nabla^2 F(x^t) d|}{\|d\|^3} d;$$

**else** {d\_type=SOL}

$$(3.8) \quad d^t \leftarrow d;$$

**end if**

    Go to **Line Search**;

**else**

    Call Algorithm B.1 with  $H = \nabla^2 F(x^t)$ ,  $\varepsilon = \epsilon_H$ , and probability parameter  $\delta$ ;

**if** Algorithm B.1 certifies that  $\lambda_{\min}(\nabla^2 F(x^t)) \geq -\epsilon_H$  **then**

      Output  $x^t$  and terminate;

**else** {Sufficiently negative curvature direction  $v$  returned by Algorithm B.1}

      Set d\_type=NC and

$$(3.9) \quad d^t \leftarrow -\text{sgn}(v^T \nabla F(x^t)) |v^T \nabla^2 F(x^t) v| v;$$

      Go to **Line Search**;

**end if**

**end if**

**Line Search:**

**if** d\_type=SOL **then**

    Find  $\alpha_t = \theta^{j_t}$ , where  $j_t$  is the smallest nonnegative integer  $j$  such that

$$(3.10) \quad F(x^t + \theta^j d^t) < F(x^t) - \eta \epsilon_H \theta^{2j} \|d^t\|^2;$$

**else** {d\_type=NC}

    Find  $\alpha_t = \theta^{j_t}$ , where  $j_t$  is the smallest nonnegative integer  $j$  such that

$$(3.11) \quad F(x^t + \theta^j d^t) < F(x^t) - \eta \theta^{2j} \|d^t\|^3 / 2;$$

**end if**

$x^{t+1} = x^t + \alpha_t d^t$ ;

**end for**

---

257 **3.2. A Newton-CG method for problem (3.1).** In this subsection we propose  
 258 a Newton-CG method in Algorithm 3.1, which is a variant of the Newton-CG method  
 259 [56, Algorithm 3], for finding an approximate stochastic SOSCP of problem (3.1).

260 Our Newton-CG method (Algorithm 3.1) follows the same framework as [56,  
 261 Algorithm 3]. In particular, at each iteration, if the gradient of  $F$  at the current  
 262 iterate is not desirably small, then the capped CG method (Algorithm A.1) is called  
 263 to solve a damped Newton system for obtaining a descent direction and a subsequent  
 264 line search along this direction results in a sufficient reduction on  $F$ . Otherwise, the  
 265 current iterate is already an approximate first-order stationary point of (3.1), and the  
 266 minimum eigenvalue oracle (Algorithm B.1) is then called, which either produces a

267 sufficiently negative curvature direction for  $F$  and a subsequent line search along this  
 268 direction results in a sufficient reduction on  $F$ , or certifies that the current iterate is  
 269 an approximate SOSP of (3.1) with high probability and terminates the algorithm.  
 270 More details about this framework can be found in [56].

271 Despite sharing the same framework, our Newton-CG method and [56, Algorithm 3]  
 272 use different line search criteria. Indeed, our Newton-CG method uses a hybrid line  
 273 search criterion adopted from [59], which is a combination of the quadratic descent  
 274 criterion (3.10) and the cubic descent criterion (3.11). Specifically, it uses the quadratic  
 275 descent criterion (3.10) when the search direction is of type ‘SOL’. On the other hand,  
 276 it uses the cubic descent criterion (3.11) when the search direction is of type ‘NC’.<sup>3</sup>  
 277 In contrast, the Newton-CG method in [56] always uses a cubic descent criterion  
 278 regardless of the type of search directions. As observed from Theorem 3.2 below, our  
 279 Newton-CG method achieves an iteration and operation complexity given in (3.5),  
 280 which are superior to those in (3.4) achieved by [56, Algorithm 3] in terms of the order  
 281 dependence on  $L_H^F$ , while improving or retaining the order of dependence on  $\epsilon_g$  and  
 282  $\epsilon_H$  as given in (3.4). Consequently, our Newton-CG method is more appealing than  
 283 [56, Algorithm 3] as an AL subproblem solver for the AL method proposed in Section  
 284 4 from theoretical complexity perspective.

285 The following theorem states the iteration and operation complexity of Algo-  
 286 rithm 3.1, whose proof is deferred to Section 6.1.

287 **THEOREM 3.2.** *Suppose that Assumption 3.1 holds. Let*

$$288 \quad (3.12) \quad T_1 := \left\lceil \frac{F_{\text{hi}} - F_{\text{low}}}{\min\{c_{\text{sol}}, c_{\text{nc}}\}} \max\{\epsilon_g^{-2}\epsilon_H, \epsilon_H^{-3}\} \right\rceil + \left\lceil \frac{F_{\text{hi}} - F_{\text{low}}}{c_{\text{nc}}} \epsilon_H^{-3} \right\rceil + 1, \quad T_2 := \left\lceil \frac{F_{\text{hi}} - F_{\text{low}}}{c_{\text{nc}}} \epsilon_H^{-3} \right\rceil + 1,$$

289 where  $F_{\text{hi}} = F(u^0)$ ,  $F_{\text{low}}$  is given in (3.3), and

$$290 \quad (3.13) \quad c_{\text{sol}} := \eta \min \left\{ \left[ \frac{4}{4+\zeta+\sqrt{(4+\zeta)^2+8L_H^F}} \right]^2, \left[ \frac{\min\{6(1-\eta), 2\}\theta}{L_H^F} \right]^2 \right\},$$

$$291 \quad (3.14) \quad c_{\text{nc}} := \frac{\eta}{16} \min \left\{ 1, \left[ \frac{\min\{3(1-\eta), 1\}\theta}{L_H^F} \right]^2 \right\}.$$

292 Then the following statements hold.

293 (i) The total number of calls of Algorithm B.1 in Algorithm 3.1 is at most  $T_2$ .  
 294 (ii) The total number of calls of Algorithm A.1 in Algorithm 3.1 is at most  $T_1$ .  
 295 (iii) **(iteration complexity)** Algorithm 3.1 terminates in at most  $T_1 + T_2$  iterations  
 296 with

$$297 \quad (3.15) \quad T_1 + T_2 = \mathcal{O}((F_{\text{hi}} - F_{\text{low}})(L_H^F)^2 \max\{\epsilon_g^{-2}\epsilon_H, \epsilon_H^{-3}\}).$$

298 Also, its output  $x^t$  satisfies  $\|\nabla F(x^t)\| \leq \epsilon_g$  deterministically and  $\lambda_{\min}(\nabla^2 F(x^t))$   
 299  $\geq -\epsilon_H$  with probability at least  $1 - \delta$  for some  $0 \leq t \leq T_1 + T_2$ .

300 (iv) **(operation complexity)** Algorithm 3.1 requires at most

$$301 \quad \tilde{\mathcal{O}}((F_{\text{hi}} - F_{\text{low}})(L_H^F)^2 \max\{\epsilon_g^{-2}\epsilon_H, \epsilon_H^{-3}\} \min\{n, (U_H^F/\epsilon_H)^{1/2}\})$$

302 matrix-vector products, where  $U_H^F$  is given in (3.3).

303 **4. A Newton-CG based AL method for problem (1.1).** In this section we  
 304 propose a Newton-CG based AL method for finding a stochastic  $(\epsilon_1, \epsilon_2)$ -SOSP of  
 305 problem (1.1) for any prescribed tolerances  $\epsilon_1, \epsilon_2 \in (0, 1)$ . Before proceeding, we make  
 306 some additional assumptions on problem (1.1).

---

<sup>3</sup>SOL and NC stand for “approximate solution” and “negative curvature”, respectively.

307     *Assumption 4.1.* (a) An  $\epsilon_1/2$ -approximately feasible point  $z_{\epsilon_1}$  of problem (1.1),  
 308     namely satisfying  $\|c(z_{\epsilon_1})\| \leq \epsilon_1/2$ , is known.  
 309     (b) There exist constants  $f_{\text{hi}}$ ,  $f_{\text{low}}$  and  $\gamma > 0$ , independent of  $\epsilon_1$  and  $\epsilon_2$ , such that

310     (4.1)                     $f(z_{\epsilon_1}) \leq f_{\text{hi}},$

311     (4.2)                     $f(x) + \gamma\|c(x)\|^2/2 \geq f_{\text{low}}, \quad \forall x \in \mathbb{R}^n,$

312     where  $z_{\epsilon_1}$  is given in (a).

313     (c) There exist some  $\delta_f, \delta_c > 0$  such that the set

314     (4.3)                     $\mathcal{S}(\delta_f, \delta_c) := \{x : f(x) \leq f_{\text{hi}} + \delta_f, \|c(x)\| \leq 1 + \delta_c\}$

315     is compact with  $f_{\text{hi}}$  given above. Also,  $\nabla^2 f$  and  $\nabla^2 c_i$ ,  $i = 1, 2, \dots, m$ , are Lipschitz  
 316     continuous in a convex open neighborhood, denoted by  $\Omega(\delta_f, \delta_c)$ , of  $\mathcal{S}(\delta_f, \delta_c)$ .

317     We now make some remarks on Assumption 4.1.

318     *Remark 4.2.* (i) A very similar assumption as Assumption 4.1(a) was con-  
 319     sidered in [31, 37, 49, 60]. By imposing Assumption 4.1(a), we restrict our  
 320     study on problem (1.1) for which an  $\epsilon_1/2$ -approximately feasible point  $z_{\epsilon_1}$   
 321     can be found by an inexpensive procedure. One example of such problem  
 322     instances arises when there exists  $v^0$  such that  $\{x : \|c(x)\| \leq \|c(v^0)\|\}$  is  
 323     compact,  $\nabla^2 c_i$ ,  $1 \leq i \leq m$ , is Lipschitz continuous on a convex neighborhood  
 324     of this set, and the LICQ holds on this set. Indeed, for this instance, a point  
 325      $z_{\epsilon_1}$  satisfying  $\|c(z_{\epsilon_1})\| \leq \epsilon_1/2$  can be computed by applying our Newton-CG  
 326     method (Algorithm 3.1) to the problem  $\min_{x \in \mathbb{R}^n} \|c(x)\|^2$ . As seen from The-  
 327     orem 3.2, the resulting iteration and operation complexity of Algorithm 3.1 for  
 328     finding such  $z_{\epsilon_1}$  are respectively  $\mathcal{O}(\epsilon_1^{-3/2})$  and  $\tilde{\mathcal{O}}(\epsilon_1^{-3/2} \min\{n, \epsilon_1^{-1/4}\})$ , which  
 329     are negligible compared with those of our AL method (see Theorems 4.10 and  
 330     4.14 below). As another example, when the standard error bound condition  
 331      $\|c(x)\|^2 = \mathcal{O}(\|\nabla(\|c(x)\|^2)\|^\nu)$  holds on a level set of  $\|c(x)\|$  for some  $\nu > 0$ ,  
 332     one can find the above  $z_{\epsilon_1}$  by applying a gradient method to the problem  
 333      $\min_{x \in \mathbb{R}^n} \|c(x)\|^2$  (e.g., see [46, 58]). In addition, the Newton-CG based AL  
 334     method (Algorithm 4.1) proposed below is a second-order method with the  
 335     aim to find a second-order stationary point. It is more expensive than a  
 336     first-order method in general. To make best use of such an AL method in  
 337     practice, it is natural to run a first-order method in advance to obtain an  
 338      $\epsilon_1/2$ -first-order stationary point  $z_{\epsilon_1}$  and then run the AL method using  $z_{\epsilon_1}$  as  
 339     an  $\epsilon_1/2$ -approximately feasible point. Therefore, Assumption 4.1(a) is met  
 340     in practice, provided that an  $\epsilon_1/2$ -first-order stationary point of (1.1) can be  
 341     found by a first-order method.

342     (ii) Assumption 4.1(b) is mild. In particular, the assumption in (4.1) holds  
 343     if  $f(x) \leq f_{\text{hi}}$  holds for all  $x$  with  $\|c(x)\| \leq 1$ , which is imposed in [60,  
 344     Assumption 3]. It also holds if problem (1.1) has a known feasible point,  
 345     which is often imposed for designing AL methods for nonconvex constrained  
 346     optimization (e.g., see [49, 31, 48, 37]). Besides, the assumption in (4.2) implies  
 347     that the quadratic penalty function is bounded below when the associated  
 348     penalty parameter is sufficiently large, which is typically used in the study  
 349     of quadratic penalty and AL methods for solving problem (1.1) (e.g., see  
 350     [40, 37, 60, 43]). Clearly, when  $\inf_{x \in \mathbb{R}^n} f(x) > -\infty$ , one can see that (4.2)  
 351     holds for any  $\gamma > 0$ . In general, one possible approach to identifying  $\gamma$  is to  
 352     apply the techniques on infeasibility detection developed in the literature (e.g.,

353 [20, 19, 6]) to check the infeasibility of the level set  $\{x : f(x) + \gamma\|c(x)\|^2/2 \leq$   
 354  $\tilde{f}_{\text{low}}\}$  for some sufficiently small  $\tilde{f}_{\text{low}}$ . Note that this level set being infeasible  
 355 for some  $\tilde{f}_{\text{low}}$  implies that (4.2) holds for the given  $\gamma$  and  $f_{\text{low}} = \tilde{f}_{\text{low}}$ .  
 356 (iii) Assumption 4.1(c) is not too restrictive. Indeed, the set  $\mathcal{S}(\delta_f, \delta_c)$  is compact  
 357 if  $f$  or  $f(\cdot) + \gamma\|c(\cdot)\|^2/2$  is level-bounded. The latter level-boundedness  
 358 assumption is commonly imposed for studying AL methods (e.g., see [37, 60]),  
 359 which is stronger than our assumption.

360 We next propose a Newton-CG based AL method in Algorithm 4.1 for finding a  
 361 stochastic  $(\epsilon_1, \epsilon_2)$ -SOSP of problem (1.1) under Assumption 4.1. Instead of solving  
 362 (1.1) directly, this method solves the perturbed problem:

363 (4.4) 
$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s. t. } \tilde{c}(x) := c(x) - c(z_{\epsilon_1}) = 0,$$

364 where  $z_{\epsilon_1}$  is given in Assumption 4.1(a). Specifically, at the  $k$ th iteration, this method  
 365 applies the Newton-CG method (Algorithm 3.1) to find an approximate stochastic  
 366 SOSP  $x^{k+1}$  of the AL subproblem associated with (4.4):

367 (4.5) 
$$\min_{x \in \mathbb{R}^n} \{\tilde{\mathcal{L}}(x, \lambda^k, \rho_k) := f(x) + (\lambda^k)^T \tilde{c}(x) + \rho_k \|\tilde{c}(x)\|^2/2\}$$

368 such that  $\tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)$  is below a threshold (see (4.6) and (4.7)), where  $\lambda^k$  is a  
 369 truncated Lagrangian multiplier, i.e., the one that results from projecting the standard  
 370 multiplier estimate  $\tilde{\lambda}^k$  onto an Euclidean ball (see step 6 of Algorithm 4.1). The  
 371 standard multiplier estimate  $\tilde{\lambda}^{k+1}$  is then updated by the classical scheme described  
 372 in step 4 of Algorithm 4.1. Finally, the penalty parameter  $\rho_{k+1}$  is adaptively updated  
 373 based on the improvement on constraint violation (see step 7 of Algorithm 4.1). Such  
 374 a practical update scheme is often adopted in the literature (e.g., see [7, 2, 31]).

375 We would like to point out that the truncated Lagrangian multiplier sequence  $\{\lambda^k\}$   
 376 is used in the AL subproblems of Algorithm 4.1 and is bounded, while the standard  
 377 Lagrangian multiplier sequence  $\{\tilde{\lambda}^k\}$  is used in those of the classical AL methods  
 378 and can be unbounded. Therefore, Algorithm 4.1 can be viewed as a safeguarded  
 379 AL method. Truncated Lagrangian multipliers have been used in the literature for  
 380 designing some AL methods [2, 11, 42, 13], and will play a crucial role in the subsequent  
 381 complexity analysis of Algorithm 4.1.

382 *Remark 4.3.* (i) Notice that the starting point  $x_{\text{init}}^0$  of Algorithm 4.1 can  
 383 be different from  $z_{\epsilon_1}$  and it may be rather infeasible, though  $z_{\epsilon_1}$  is a nearly  
 384 feasible point of (1.1). Besides,  $z_{\epsilon_1}$  is used to ensure convergence of Algorithm  
 385 4.1. Specifically, if the algorithm runs into a “poorly infeasible point”  $x^k$ ,  
 386 namely satisfying  $\tilde{\mathcal{L}}(x^k, \lambda^k; \rho_k) > f(z_{\epsilon_1})$ , it will be superseded by  $z_{\epsilon_1}$  (see  
 387 (4.8)), which prevents the iterates  $\{x^k\}$  from converging to an infeasible point.  
 388 Yet,  $x^k$  may be rather infeasible when  $k$  is not large. Thus, Algorithm 4.1  
 389 substantially differs from a funneling or two-phase type algorithm, in which a  
 390 nearly feasible point is found in Phase 1, and then approximate stationarity  
 391 is sought while near feasibility is maintained throughout Phase 2 (e.g., see  
 392 [9, 16, 26, 27, 28, 29, 30, 36]).  
 393 (ii) The choice of  $\rho_0$  in Algorithm 4.1 is mainly for the simplicity of complexity  
 394 analysis. Yet, it may be overly large and lead to highly ill-conditioned AL  
 395 subproblems in practice. To make Algorithm 4.1 practically more efficient, one  
 396 can possibly modify it by choosing a relatively small initial penalty parameter,  
 397 then solving the subsequent AL subproblems by a first-order method until an  
 398  $\epsilon_1$ -first-order stationary point  $\hat{x}$  of (1.1) along with a Lagrangian multiplier  $\hat{\lambda}$

**Algorithm 4.1** A Newton-CG based AL method for problem (1.1)

Let  $\gamma$  be given in Assumption 4.1.

**Input:**  $\epsilon_1, \epsilon_2 \in (0, 1)$ ,  $\Lambda > 0$ ,  $x^0 \in \mathbb{R}^n$ ,  $\lambda^0 \in \mathcal{B}_\Lambda$ ,  $\rho_0 > 2\gamma$ ,  $\alpha \in (0, 1)$ ,  $r > 1$ ,  $\delta \in (0, 1)$ , and  $z_{\epsilon_1}$  given in Assumption 4.1.

- 1: Set  $k = 0$ .
- 2: Set  $\tau_k^g = \max\{\epsilon_1, r^{k \log \epsilon_1 / \log 2}\}$  and  $\tau_k^H = \max\{\epsilon_2, r^{k \log \epsilon_2 / \log 2}\}$ .
- 3: Call Algorithm 3.1 with  $\epsilon_g = \tau_k^g$ ,  $\epsilon_H = \tau_k^H$  and  $u^0 = x_{\text{init}}^k$  to find an approximate solution  $x^{k+1}$  to  $\min_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  such that

$$(4.6) \quad \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k) \leq f(z_{\epsilon_1}), \quad \|\nabla_x \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)\| \leq \tau_k^g,$$

$$(4.7) \quad \lambda_{\min}(\nabla_{xx}^2 \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)) \geq -\tau_k^H \text{ with probability at least } 1 - \delta,$$

where

$$(4.8) \quad x_{\text{init}}^k = \begin{cases} z_{\epsilon_1} & \text{if } \tilde{\mathcal{L}}(x^k, \lambda^k; \rho_k) > f(z_{\epsilon_1}), \\ x^k & \text{otherwise,} \end{cases} \quad \text{for } k \geq 0.$$

- 4: Set  $\tilde{\lambda}^{k+1} = \lambda^k + \rho_k \tilde{c}(x^{k+1})$ .
- 5: If  $\tau_k^g \leq \epsilon_1$ ,  $\tau_k^H \leq \epsilon_2$  and  $\|c(x^{k+1})\| \leq \epsilon_1$ , then output  $(x^{k+1}, \tilde{\lambda}^{k+1})$  and terminate.
- 6: Set  $\lambda^{k+1} = \Pi_{\mathcal{B}_\Lambda}(\tilde{\lambda}^{k+1})$ .
- 7: If  $k = 0$  or  $\|\tilde{c}(x^{k+1})\| > \alpha \|\tilde{c}(x^k)\|$ , set  $\rho_{k+1} = r \rho_k$ . Otherwise, set  $\rho_{k+1} = \rho_k$ .
- 8: Set  $k \leftarrow k + 1$ , and go to step 2.

399 is found, and finally performing the steps described in Algorithm 4.1 but with  
400  $x^0 = \hat{x}$  and  $\lambda^0 = \Pi_{\mathcal{B}_\Lambda}(\hat{\lambda})$ .

401 Before analyzing the complexity of Algorithm 4.1, we first argue that it is well-  
402 defined if  $\rho_0$  is suitably chosen. Specifically, we will show that when  $\rho_0$  is sufficiently  
403 large, one can apply the Newton-CG method (Algorithm 3.1) to the AL subproblem  
404  $\min_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  with  $x_{\text{init}}^k$  as the initial point to find an  $x^{k+1}$  satisfying (4.6) and  
405 (4.7). To this end, we start by noting from (4.1), (4.4), (4.5) and (4.8) that

$$406 \quad (4.9) \quad \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k) \leq \max\{\tilde{\mathcal{L}}(z_{\epsilon_1}, \lambda^k; \rho_k), f(z_{\epsilon_1})\} = f(z_{\epsilon_1}) \leq f_{\text{hi}}.$$

407 Based on the above observation, we show in the next lemma that when  $\rho_0$  is sufficiently  
408 large,  $\tilde{\mathcal{L}}(\cdot, \lambda^k; \rho_k)$  is bounded below and its certain level set is bounded, whose proof is  
409 deferred to Section 6.2.

410 **LEMMA 4.4.** *Suppose that Assumption 4.1 holds. Let  $(\lambda^k, \rho_k)$  be generated at the  
411  $k$ th iteration of Algorithm 4.1 for some  $k \geq 0$ , and  $\mathcal{S}(\delta_f, \delta_c)$  and  $x_{\text{init}}^k$  be defined in  
412 (4.3) and (4.8), respectively, and let  $f_{\text{hi}}$ ,  $f_{\text{low}}$ ,  $\delta_f$  and  $\delta_c$  be given in Assumption 4.1.  
413 Suppose that  $\rho_0$  is sufficiently large such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where*

$$414 \quad (4.10) \quad \delta_{f,1} := \Lambda^2 / (2\rho_0) \quad \text{and} \quad \delta_{c,1} := \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho_0 - 2\gamma} + \frac{\Lambda^2}{(\rho_0 - 2\gamma)^2}} + \frac{\Lambda}{\rho_0 - 2\gamma}.$$

415 Then the following statements hold.

- (i)  $\{x : \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)\} \subseteq \mathcal{S}(\delta_f, \delta_c)$ .
- (ii)  $\inf_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \geq f_{\text{low}} - \gamma - \Lambda \delta_c$ .

416 Using Lemma 4.4, we can verify that the Newton-CG method (Algorithm 3.1),  
417 starting with  $u^0 = x_{\text{init}}^k$ , is capable of finding an approximate solution  $x^{k+1}$  of the  
418 AL subproblem  $\min_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  satisfying (4.6) and (4.7). Indeed, let  $F(\cdot) =$   
419  $\tilde{\mathcal{L}}(\cdot, \lambda^k; \rho_k)$  and  $u^0 = x_{\text{init}}^k$ . By these and Lemma 4.4, one can see that  $\{x : F(x) \leq$   
420  $F(u^0)\} \subseteq \mathcal{S}(\delta_f, \delta_c)$ . It then follows from this and Assumption 4.1(c) that the level set

423  $\{x : F(x) \leq F(u^0)\}$  is compact and  $\nabla^2 F$  is Lipschitz continuous on a convex open  
 424 neighborhood of  $\{x : F(x) \leq F(u^0)\}$ . Thus, such  $F$  and  $u^0$  satisfy Assumption 3.1.  
 425 Based on this and the discussion in Section 3, one can conclude that Algorithm 3.1,  
 426 starting with  $u^0 = x_{\text{init}}^k$ , is applicable to the AL subproblem  $\min_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$ .  
 427 Moreover, it follows from Theorem 3.2 that this algorithm with  $(\epsilon_g, \epsilon_H) = (\tau_k^g, \tau_k^H)$   
 428 can produce a point  $x^{k+1}$  satisfying (4.7) and also the second relation in (4.6). In  
 429 addition, since this algorithm is descent and its starting point is  $x_{\text{init}}^k$ , its output  $x^{k+1}$   
 430 must satisfy  $\tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)$ , which along with (4.9) implies that  
 431  $\tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k) \leq f(z_{\epsilon_1})$  and thus  $x^{k+1}$  also satisfies the first relation in (4.6).

432 The above discussion leads to the following conclusion concerning the *well-definedness of Algorithm 4.1*.

434 **THEOREM 4.5.** *Under the same settings as in Lemma 4.4, the Newton-CG method*  
 435 *(Algorithm 3.1) applied to the AL subproblem  $\min_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  with  $u^0 = x_{\text{init}}^k$*   
 436 *finds a point  $x^{k+1}$  satisfying (4.6) and (4.7).*

437 The following theorem characterizes the *output of Algorithm 4.1*. Its proof is  
 438 deferred to Section 6.2.

439 **THEOREM 4.6.** *Suppose that Assumption 4.1 holds and that  $\rho_0$  is sufficiently large*  
 440 *such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). If Algorithm*  
 441 *4.1 terminates at some iteration  $k$ , then  $x^{k+1}$  is a deterministic  $\epsilon_1$ -FOSP of problem*  
 442 *(1.1), and moreover, it is an  $(\epsilon_1, \epsilon_2)$ -SOSP of (1.1) with probability at least  $1 - \delta$ .*

443 **Remark 4.7.** As seen from this theorem, the output of Algorithm 4.1 is a stochastic  
 444  $(\epsilon_1, \epsilon_2)$ -SOSP of problem (1.1). Nevertheless, one can easily modify Algorithm 4.1  
 445 to seek some other approximate solutions. For example, if one is only interested in  
 446 finding an  $\epsilon_1$ -FOSP of (1.1), one can remove the condition (4.7) from Algorithm 4.1.  
 447 In addition, if one aims to find a deterministic  $(\epsilon_1, \epsilon_2)$ -SOSP of (1.1), one can replace  
 448 the condition (4.7) and Algorithm 3.1 by  $\lambda_{\min}(\nabla_{xx}^2 \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)) \geq -\tau_k^H$  and a  
 449 deterministic counterpart, respectively. The purpose of imposing high probability in  
 450 the condition (4.7) is to enable us to derive operation complexity of Algorithm 4.1  
 451 measured by the number of matrix-vector products.

452 In the rest of this section, we study the worst-case complexity of Algorithm 4.1.  
 453 Since our method has two nested loops, particularly, outer loops executed by the  
 454 AL method and inner loops executed by the Newton-CG method for solving the AL  
 455 subproblems, we consider the following measures of complexity for Algorithm 4.1.

- 456 • *Outer iteration complexity*, which measures the number of outer iterations of  
 457 Algorithm 4.1;
- 458 • *Total inner iteration complexity*, which measures the total number of iterations  
 459 of the Newton-CG method that are performed in Algorithm 4.1;
- 460 • *Operation complexity*, which measures the total number of matrix-vector  
 461 products involving the Hessian of the augmented Lagrangian function that  
 462 are evaluated in Algorithm 4.1.

463 **4.1. Outer iteration complexity of Algorithm 4.1.** In this subsection we  
 464 establish outer iteration complexity of Algorithm 4.1. For notational convenience, we  
 465 rewrite  $(\tau_k^g, \tau_k^H)$  arising in Algorithm 4.1 as

466 (4.11)  $(\tau_k^g, \tau_k^H) = (\max\{\epsilon_1, \omega_1^k\}, \max\{\epsilon_2, \omega_2^k\})$  with  $(\omega_1, \omega_2) := (r^{\log \epsilon_1 / \log 2}, r^{\log \epsilon_2 / \log 2})$ ,

467 where  $\epsilon_1, \epsilon_2$  and  $r$  are the input parameters of Algorithm 4.1. Since  $r > 1$  and  
 468  $\epsilon_1, \epsilon_2 \in (0, 1)$ , it is not hard to verify that  $\omega_1, \omega_2 \in (0, 1)$ . Also, we introduce the

469 following quantity that will be used frequently later:

470 (4.12) 
$$K_{\epsilon_1} := \lceil \min\{k \geq 0 : \omega_1^k \leq \epsilon_1\} \rceil = \lceil \log \epsilon_1 / \log \omega_1 \rceil.$$

471 In view of (4.11), (4.12) and the fact that

472 (4.13) 
$$\log \epsilon_1 / \log \omega_1 = \log \epsilon_2 / \log \omega_2 = \log 2 / \log r,$$

473 we see that  $(\tau_k^g, \tau_k^H) = (\epsilon_1, \epsilon_2)$  for all  $k \geq K_{\epsilon_1}$ . This along with the termination  
474 criterion of Algorithm 4.1 implies that it runs for at least  $K_{\epsilon_1}$  iterations and terminates  
475 once  $\|c(x^{k+1})\| \leq \epsilon_1$  for some  $k \geq K_{\epsilon_1}$ . As a result, to establish outer iteration  
476 complexity of Algorithm 4.1, it suffices to bound such  $k$ . The resulting outer iteration  
477 complexity of Algorithm 4.1 is presented below, whose proof is deferred to Section 6.2.

478 **THEOREM 4.8.** *Suppose that Assumption 4.1 holds and that  $\rho_0$  is sufficiently large  
479 such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). Let*

480 (4.14) 
$$\rho_{\epsilon_1} := \max \{8(f_{\text{hi}} - f_{\text{low}} + \gamma)\epsilon_1^{-2} + 4\Lambda\epsilon_1^{-1} + 2\gamma, 2\rho_0\},$$

481 (4.15) 
$$\bar{K}_{\epsilon_1} := \inf\{k \geq K_{\epsilon_1} : \|c(x^{k+1})\| \leq \epsilon_1\},$$

482 where  $K_{\epsilon_1}$  is defined in (4.12), and  $\gamma$ ,  $f_{\text{hi}}$  and  $f_{\text{low}}$  are given in Assumption 4.1. Then  
483  $\bar{K}_{\epsilon_1}$  is finite, and Algorithm 4.1 terminates at iteration  $\bar{K}_{\epsilon_1}$  with

484 (4.16) 
$$\bar{K}_{\epsilon_1} \leq \left( \frac{\log(\rho_{\epsilon_1}\rho_0^{-1})}{\log r} + 1 \right) \left( \left\lceil \frac{\log(\epsilon_1(2\delta_{c,1})^{-1})}{\log \alpha} \right\rceil + 2 \right) + 1.$$

485 Moreover,  $\rho_k \leq r\rho_{\epsilon_1}$  holds for  $0 \leq k \leq \bar{K}_{\epsilon_1}$

486 **Remark 4.9 (Upper bounds for  $\bar{K}_{\epsilon_1}$  and  $\{\rho_k\}$ ).** As observed from Theorem  
487 4.8, the number of outer iterations of Algorithm 4.1 for finding a stochastic  $(\epsilon_1, \epsilon_2)$ -  
488 SOSP of problem (1.1) is  $\bar{K}_{\epsilon_1} + 1$ , which is at most of  $\mathcal{O}(|\log \epsilon_1|^2)$ . In addition, the  
489 penalty parameters  $\{\rho_k\}$  generated in this algorithm are at most of  $\mathcal{O}(\epsilon_1^{-2})$ .

#### 490 4.2. Total inner iteration and operation complexity of Algorithm 4.1.

491 We present the total inner iteration and operation complexity of Algorithm 4.1 for  
492 finding a stochastic  $(\epsilon_1, \epsilon_2)$ -SOSP of (1.1), whose proof is deferred to Section 6.2.

493 **THEOREM 4.10.** *Suppose that Assumption 4.1 holds and that  $\rho_0$  is sufficiently  
494 large such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). Then  
495 the following statements hold.*

- 496 (i) *The total number of iterations of Algorithm 3.1 performed in Algorithm 4.1 is at  
497 most  $\tilde{\mathcal{O}}(\epsilon_1^{-4} \max\{\epsilon_1^{-2}\epsilon_2, \epsilon_2^{-3}\})$ . If  $c$  is further assumed to be affine, then it is at  
498 most  $\tilde{\mathcal{O}}(\max\{\epsilon_1^{-2}\epsilon_2, \epsilon_2^{-3}\})$ .*
- 499 (ii) *The total number of matrix-vector products performed by Algorithm 3.1 in Al-  
500 gorithm 4.1 is at most  $\tilde{\mathcal{O}}(\epsilon_1^{-4} \max\{\epsilon_1^{-2}\epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1}\epsilon_2^{-1/2}\})$ . If  $c$  is further  
501 assumed to be affine, then it is at most  $\tilde{\mathcal{O}}(\max\{\epsilon_1^{-2}\epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1}\epsilon_2^{-1/2}\})$ .*

502 **Remark 4.11.** (i) Note that the above complexity results of Algorithm 4.1 are  
503 established without assuming any constraint qualification (CQ). In contrast,  
504 similar complexity results are obtained in [60] for a proximal AL method under  
505 a generalized LICQ condition. To the best of our knowledge, our work provides  
506 the first study on complexity for finding a stochastic SOSP of (1.1) without CQ.  
507 (ii) Letting  $(\epsilon_1, \epsilon_2) = (\epsilon, \sqrt{\epsilon})$  for some  $\epsilon \in (0, 1)$ , we see that Algorithm 4.1 achieves  
508 a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and an operation complexity of  
509  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-5/4}\})$  for finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of problem (1.1)  
510 without constraint qualification.

**4.3. Enhanced complexity of Algorithm 4.1 under constraint qualification.** In this subsection we study complexity of Algorithm 4.1 under one additional assumption that a generalized linear independence constraint qualification (GLICQ) holds for problem (1.1), which is introduced below. In particular, under GLICQ we will obtain an enhanced total inner iteration and operation complexity for Algorithm 4.1, which are significantly better than the ones in Theorem 4.10 when problem (1.1) has nonlinear constraints. Moreover, when  $(\epsilon_1, \epsilon_2) = (\epsilon, \sqrt{\epsilon})$  for some  $\epsilon \in (0, 1)$ , our enhanced complexity bounds are also better than those obtained in [60] for a proximal AL method. We now introduce the GLICQ assumption for problem (1.1).

*Assumption 4.12 (GLICQ).*  $\nabla c(x)$  has full column rank for all  $x \in \mathcal{S}(\delta_f, \delta_c)$ , where  $\mathcal{S}(\delta_f, \delta_c)$  is as in (4.3).

*Remark 4.13.* A related yet different GLICQ is imposed in [60, Assumption 2(ii)] for problem (1.1), which assumes that  $\nabla c(x)$  has full column rank for all  $x$  in a level set of  $f(\cdot) + \gamma\|c(\cdot)\|^2/2$ . It is not hard to verify that this assumption is generally stronger than the above GLICQ assumption.

The following theorem shows that under Assumption 4.12, the total inner iteration and operation complexity results presented in Theorem 4.10 can be significantly improved, whose proof is deferred to Section 6.2.

**THEOREM 4.14.** *Suppose that Assumptions 4.1 and 4.12 hold and that  $\rho_0$  is sufficiently large such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). Then the following statements hold.*

- (i) *The total number of iterations of Algorithm 3.1 performed in Algorithm 4.1 is at most  $\tilde{\mathcal{O}}(\epsilon_1^{-2} \max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\})$ . If  $c$  is further assumed to be affine, then it is at most  $\tilde{\mathcal{O}}(\max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\})$ .*
- (ii) *The total number of matrix-vector products performed by Algorithm 3.1 in Algorithm 4.1 is at most  $\tilde{\mathcal{O}}(\epsilon_1^{-2} \max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1/2} \epsilon_2^{-1/2}\})$ . If  $c$  is further assumed to be affine, then it is at most  $\tilde{\mathcal{O}}(\max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1/2} \epsilon_2^{-1/2}\})$ .*

*Remark 4.15.* (i) As seen from Theorem 4.14, when problem (1.1) has nonlinear constraints, under GLICQ and some other suitable assumptions, Algorithm 4.1 achieves significantly better complexity bounds than the ones in Theorem 4.10 without constraint qualification.

- (ii) Letting  $(\epsilon_1, \epsilon_2) = (\epsilon, \sqrt{\epsilon})$  for some  $\epsilon \in (0, 1)$ , we see that when problem (1.1) has nonlinear constraints, under GLICQ and some other suitable assumptions, Algorithm 4.1 achieves a total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2})$  and an operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-7/2} \min\{n, \epsilon^{-3/4}\})$ . They are vastly better than the total inner iteration complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2})$  and the operation complexity of  $\tilde{\mathcal{O}}(\epsilon^{-11/2} \min\{n, \epsilon^{-3/4}\})$  that are achieved by a proximal AL method in [60] for finding a stochastic  $(\epsilon, \sqrt{\epsilon})$ -SOSP of (1.1) yet under a generally stronger GLICQ.

**5. Numerical results.** We conduct some preliminary experiments to test the performance of our proposed methods (Algorithms 3.1 and 4.1), and compare them with the Newton-CG method in [56] and the proximal AL method in [60], respectively. All the algorithms are coded in Matlab and all the computations are performed on a desktop with a 3.79 GHz AMD 3900XT 12-Core processor and 32 GB of RAM.

**5.1. Regularized robust regression.** In this subsection we consider the regularized robust regression problem

$$556 \quad (5.1) \quad \min_{x \in \mathbb{R}^n} \sum_{i=1}^m \phi(a_i^T x - b_i) + \mu \|x\|_4^4,$$

n	m	$\mu$	Objective value		Iterations		CPU time (seconds)	
			Algorithm 1	Newton-CG	Algorithm 1	Newton-CG	Algorithm 1	Newton-CG
100	10	1	5.9	5.9	85.7	116.3	1.4	1.6
100	50	1	45.9	45.9	82.6	158.2	1.0	2.7
100	90	1	84.8	84.8	102.2	224.7	2.0	4.2
500	50	5	42.2	42.5	173.1	344.7	44.2	72.2
500	250	5	243.0	242.9	145.5	362.4	41.9	95.0
500	450	5	442.2	442.2	163.7	425.2	47.6	138.3
1000	100	10	90.1	90.4	162.5	361.0	110.8	259.0
1000	500	10	491.1	491.2	158.3	475.4	129.1	558.4
1000	900	10	891.1	891.1	193.5	300.7	187.0	298.5

TABLE 2  
Numerical results for problem (5.1)

557 where  $\phi(t) = t^2/(1+t^2)$ ,  $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$  for any  $p \geq 1$ , and  $\mu > 0$ .

558 For each triple  $(n, m, \mu)$ , we randomly generate 10 instances of problem (5.1). In  
559 particular, we first randomly generate  $a_i$ ,  $1 \leq i \leq m$ , with all the entries independently  
560 chosen from the standard normal distribution. We then randomly generate  $\bar{b}_i$  according  
561 to the standard normal distribution and set  $b_i = 2m\bar{b}_i$  for  $i = 1, \dots, m$ .

562 Our aim is to find a  $(10^{-5}, 10^{-5/2})$ -SOSP of (5.1) for the above instances by  
563 Algorithm 3.1 and the Newton-CG method in [56] and compare their performance. For a  
564 fair comparison, we use a minimum eigenvalue oracle that returns a deterministic output  
565 for them so that they both certainly output an approximate second-order stationary  
566 point. Specifically, we use the Matlab subroutine  $[v, \lambda] = \text{eigs}(H, 1, 'smallestreal')$  as the  
567 minimum eigenvalue oracle to find the minimum eigenvalue  $\lambda$  and its associated unit  
568 eigenvector  $v$  of a real symmetric matrix  $H$ . Also, for both methods, we choose the  
569 all-ones vector as the initial point, and set  $\theta = 0.8$ ,  $\zeta = 0.5$ , and  $\eta = 0.2$ .

570 The computational results of Algorithm 3.1 and the Newton-CG method in [56]  
571 for the instances randomly generated above are presented in Table 2. In detail, the  
572 value of  $n$ ,  $m$ , and  $\mu$  is listed in the first three columns, respectively. For each triple  
573  $(n, m, \mu)$ , the average CPU time (in seconds), the average number of iterations, and  
574 the average final objective value over 10 random instances are given in the rest of  
575 the columns. One can observe that both methods output an approximate solution  
576 with a similar objective value, while our Algorithm 3.1 substantially outperforms  
577 the Newton-CG method in [56] in terms of CPU time. This is consistent with our  
578 theoretical finding that Algorithm 3.1 achieves a better iteration complexity than the  
579 Newton-CG method in [56] in terms of dependence on the Lipschitz constant of the  
580 Hessian for finding an approximate SOSP.

581 **5.2. Spherically constrained regularized robust regression.** In subsection  
582 we consider the spherically constrained regularized robust regression problem

583 (5.2) 
$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^m \phi(a_i^T x - b_i) + \mu \|x\|_4^4 \quad \text{s. t.} \quad \|x\|_2^2 = 1,$$

584 where  $\phi(t) = t^2/(1+t^2)$ ,  $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$  for any  $p \geq 1$ , and  $\mu > 0$  is a  
585 tuning parameter. For each triple  $(n, m, \mu)$ , we randomly generate 10 instances of  
586 problem (5.2) in the same manner as described in Subsection 5.1.

587 Our aim is to find a  $(10^{-4}, 10^{-2})$ -SOSP of (5.2) for the above instances by  
588 Algorithm 4.1 and the proximal AL method [60, Algorithm 3] and compare their  
589 performance. For a fair comparison, we use a minimum eigenvalue oracle that returns  
590 a deterministic output for them so that they both certainly output an approximate  
591 second-order stationary point. Specifically, we use the Matlab subroutine  $[v, \lambda] =$   
592  $\text{eigs}(H, 1, 'smallestreal')$  as the minimum eigenvalue oracle to find the minimum eigenvalue  
593  $\lambda$  and its associated unit eigenvector  $v$  of a real symmetric matrix  $H$ . In addition,  
594 for both methods, we choose the initial point as  $z^0 = (1/\sqrt{n}, \dots, 1/\sqrt{n})^T$ , the initial  
595 Lagrangian multiplier as  $\lambda^0 = 0$ , and the other parameters as

$n$	$m$	$\mu$	Objective value		Feasibility violation ( $\times 10^{-4}$ )		Total inner iterations		CPU time (seconds)	
			Algorithm 2	Prox-AL	Algorithm 2	Prox-AL	Algorithm 2	Prox-AL	Algorithm 2	Prox-AL
100	10	1	7.1	7.1	0.18	0.27	40.9	97.3	0.73	2.2
100	50	1	46.6	46.6	0.21	0.30	37.0	86.3	0.78	1.7
100	90	1	87.0	87.0	0.12	0.40	39.5	68.6	1.1	1.9
500	50	5	44.4	44.4	0.40	0.68	59.0	343.4	11.4	134.9
500	250	5	244.3	244.3	0.37	0.47	59.0	543.3	11.7	178.2
500	450	5	444.0	444.0	0.27	0.53	66.7	634.1	17.1	158.2
1000	100	10	92.8	92.8	0.28	0.42	95.0	2054.6	46.3	1516.8
1000	500	10	491.9	491.9	0.22	0.72	68.3	756.2	39.5	558.6
1000	900	10	893.4	893.4	0.19	0.37	81.8	1281.4	57.7	1099.6

TABLE 3  
Numerical results for problem (5.2)

•  $\Lambda = 100$ ,  $\rho_0 = 10$ ,  $\alpha = 0.25$ , and  $r = 10$  for Algorithm 4.1;  
 •  $\eta = 1$ ,  $q = 10$  and  $T_0 = 2$  for the proximal AL method ([60]).

The computational results of Algorithm 4.1 and the proximal AL method in [60] (abbreviated as Prox-AL) for solving problem (5.2) for the instances randomly generated above are presented in Table 3. In detail, the value of  $n$ ,  $m$ , and  $\mu$  is listed in the first three columns, respectively. For each triple  $(n, m, \mu)$ , the average CPU time (in seconds), the average total number of inner iterations, the average final objective value, and the average final feasibility violation over 10 random instances are given in the rest columns. One can observe that both methods output an approximate solution of similar quality in terms of objective value and feasibility violation, while our Algorithm 4.1 vastly outperforms the proximal AL method in [60] in terms of CPU time. This corroborates our theoretical finding that Algorithm 4.1 achieves a significantly better operation complexity than the proximal AL method in [60] for finding an approximate SOSP.

610 **6. Proof of the main results.** We provide proofs of our main results in Sections  
 611 3 and 4, including Theorem 3.2, Lemma 4.4, and Theorems 4.6, 4.8, 4.10 and 4.14.

612 **6.1. Proof of the main results in Section 3.** In this subsection we first  
 613 establish several technical lemmas and then use them to prove Theorem 3.2.

614 One can observe from Assumption 3.1(b) that for all  $x$  and  $y \in \Omega$ ,

$$(6.1) \|\nabla F(y) - \nabla F(x) - \nabla^2 F(x)(y - x)\| \leq L_H^F \|y - x\|^2/2,$$

$$(6.2) F(y) \leq F(x) + \nabla F(x)^T(y - x) + (y - x)^T \nabla^2 F(x)(y - x)/2 + L_H^F \|y - x\|^3/6.$$

617 The next lemma provides useful properties of the output of Algorithm A.1, whose  
 618 proof is similar to the ones in [56, Lemma 3] and [54, Lemma 7] and thus omitted here.

619 **LEMMA 6.1.** *Suppose that Assumption 3.1 holds and the direction  $d^t$  results from  
 620 the output  $d$  of Algorithm A.1 with a type specified in  $d\_type$  at some iteration  $t$  of  
 621 Algorithm 3.1. Then the following statements hold.*

622 (i) *If  $d\_type=$ SOL, then  $d^t$  satisfies*

$$(6.3) \epsilon_H \|d^t\|^2 \leq (d^t)^T (\nabla^2 F(x^t) + 2\epsilon_H I) d^t,$$

$$(6.4) \|d^t\| \leq 1.1\epsilon_H^{-1} \|\nabla F(x^t)\|,$$

$$(6.5) (d^t)^T \nabla F(x^t) = -(d^t)^T (\nabla^2 F(x^t) + 2\epsilon_H I) d^t,$$

$$(6.6) \|(\nabla^2 F(x^t) + 2\epsilon_H I)d^t + \nabla F(x^t)\| \leq \epsilon_H \zeta \|d^t\|/2.$$

627 (ii) *If  $d\_type=$ NC, then  $d^t$  satisfies  $(d^t)^T \nabla F(x^t) \leq 0$  and*

$$(6.7) (d^t)^T \nabla^2 F(x^t) d^t / \|d^t\|^2 = -\|d^t\| \leq -\epsilon_H.$$

629 The next lemma shows that when the search direction  $d^t$  in Algorithm 3.1 is of  
 630 type ‘SOL’, the line search step results in a sufficient reduction on  $F$ .

631 LEMMA 6.2. *Suppose that Assumption 3.1 holds and the direction  $d^t$  results from  
 632 the output  $d$  of Algorithm A.1 with  $d\_type=SOL$  at some iteration  $t$  of Algorithm 3.1.  
 633 Let  $U_g^F$  and  $c_{\text{sol}}$  be given in (3.3) and (3.13), respectively. Then the following statements  
 634 hold.*

635 (i) *The step length  $\alpha_t$  is well-defined, and moreover,*

$$636 \quad (6.8) \quad \alpha_t \geq \min \left\{ 1, \sqrt{\frac{\min\{6(1-\eta), 2\}}{1.1L_H^F U_g^F}} \theta \epsilon_H \right\}.$$

637 (ii) *The next iterate  $x^{t+1} = x^t + \alpha_t d^t$  satisfies*

$$638 \quad (6.9) \quad F(x^t) - F(x^{t+1}) \geq c_{\text{sol}} \min\{\|\nabla F(x^{t+1})\|^2 \epsilon_H^{-1}, \epsilon_H^3\}.$$

639 *Proof.* One can observe that  $F$  is descent along the iterates (whenever well-defined)  
 640 generated by Algorithm 3.1, which together with  $x^0 = u^0$  implies that  $F(x^t) \leq F(u^0)$   
 641 and hence  $\|\nabla F(x^t)\| \leq U_g^F$  due to (3.3). In addition, since  $d^t$  results from the output  
 642  $d$  of Algorithm A.1 with  $d\_type=SOL$ , one can see that  $\|\nabla F(x^t)\| > \epsilon_g$  and (6.3)-(6.6)  
 643 hold for  $d^t$ . Moreover, by  $\|\nabla F(x^t)\| > \epsilon_g$  and (6.6), one can conclude that  $d^t \neq 0$ .

644 We first prove statement (i). If (3.10) holds for  $j = 0$ , then  $\alpha_t = 1$ , which clearly  
 645 implies that (6.8) holds. We now suppose that (3.10) fails for  $j = 0$ . Claim that for  
 646 all  $j \geq 0$  that violate (3.10), it holds that

$$647 \quad (6.10) \quad \theta^{2j} \geq \min\{6(1-\eta), 2\} \epsilon_H (L_H^F)^{-1} \|d^t\|^{-1}.$$

648 Indeed, suppose that (3.10) is violated by some  $j \geq 0$ . We now show that (6.10) holds  
 649 for such  $j$  by considering two separate cases below.

650 Case 1)  $F(x^t + \theta^j d^t) > F(x^t)$ . Let  $\phi(\alpha) = F(x^t + \alpha d^t)$ . Then  $\phi(\theta^j) > \phi(0)$ . Also,  
 651 since  $d^t \neq 0$ , by (6.3) and (6.5), one has  $\phi'(0) = \nabla F(x^t)^T d^t = -(d^t)^T (\nabla^2 F(x^t) +  
 652 2\epsilon_H I) d^t \leq -\epsilon_H \|d^t\|^2 < 0$ . Using these, we can observe that there exists a local  
 653 minimizer  $\alpha^* \in (0, \theta^j)$  of  $\phi$  such that  $\phi'(\alpha^*) = \nabla F(x^t + \alpha^* d^t)^T d^t = 0$  and  $\phi(\alpha^*) < \phi(0)$ ,  
 654 which implies that  $F(x^t + \alpha^* d^t) < F(x^t) \leq F(u^0)$ . Hence, (6.1) holds for  $x = x^t$  and  
 655  $y = x^t + \alpha^* d^t$ . Using this,  $0 < \alpha^* < \theta^j \leq 1$  and  $\nabla F(x^t + \alpha^* d^t)^T d^t = 0$ , we obtain

$$\begin{aligned} 656 \quad & \frac{(\alpha^*)^2 L_H^F}{2} \|d^t\|^3 \stackrel{(6.1)}{\geq} \|d^t\| \|\nabla F(x^t + \alpha^* d^t) - \nabla F(x^t) - \alpha^* \nabla^2 F(x^t) d^t\| \\ & \geq (d^t)^T (\nabla F(x^t + \alpha^* d^t) - \nabla F(x^t) - \alpha^* \nabla^2 F(x^t) d^t) \\ & = -(d^t)^T \nabla F(x^t) - \alpha^* (d^t)^T \nabla^2 F(x^t) d^t \\ & \stackrel{(6.5)}{=} (1 - \alpha^*) (d^t)^T (\nabla^2 F(x^t) + 2\epsilon_H I) d^t + 2\alpha^* \epsilon_H \|d^t\|^2 \\ & \stackrel{(6.3)}{\geq} (1 + \alpha^*) \epsilon_H \|d^t\|^2 \geq \epsilon_H \|d^t\|^2, \end{aligned}$$

657 which along with  $d^t \neq 0$  implies that  $(\alpha^*)^2 \geq 2\epsilon_H (L_H^F)^{-1} \|d^t\|^{-1}$ . Using this and  
 658  $\theta^j > \alpha^*$ , we conclude that (6.10) holds in this case.

659 Case 2)  $F(x^t + \theta^j d^t) \leq F(x^t)$ . This together with  $F(x^t) \leq F(u^0)$  implies that  
 660 (6.2) holds for  $x = x^t$  and  $y = x^t + \theta^j d^t$ . Then, because  $j$  violates (3.10), we obtain

$$\begin{aligned} 661 \quad & -\eta \epsilon_H \theta^{2j} \|d^t\|^2 \leq F(x^t + \theta^j d^t) - F(x^t) \\ 662 \quad & \leq \theta^j \nabla F(x^t)^T d^t + \frac{\theta^{2j}}{2} (d^t)^T \nabla^2 F(x^t) d^t + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3 \\ 663 \quad & \stackrel{(6.5)}{=} -\theta^j (d^t)^T (\nabla^2 F(x^t) + 2\epsilon_H I) d^t + \frac{\theta^{2j}}{2} (d^t)^T \nabla^2 F(x^t) d^t + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3 \\ 664 \quad & = -\theta^j \left( 1 - \frac{\theta^j}{2} \right) (d^t)^T (\nabla^2 F(x^t) + 2\epsilon_H I) d^t - \theta^{2j} \epsilon_H \|d^t\|^2 + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3 \end{aligned}$$

$$\begin{aligned}
665 \quad & \stackrel{(6.3)}{\leq} -\theta^j \left(1 - \frac{\theta^j}{2}\right) \epsilon_H \|d^t\|^2 - \theta^{2j} \epsilon_H \|d^t\|^2 + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3 \\
666 \quad (6.11) \quad & \leq -\theta^j \epsilon_H \|d^t\|^2 + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3.
\end{aligned}$$

667 Recall that  $d^t \neq 0$ . Dividing both sides of (6.11) by  $L_H^F \theta^j \|d^t\|^3 / 6$  and using  $\eta, \theta \in (0, 1)$ ,  
668 we obtain that  $\theta^{2j} \geq 6(1 - \theta^j \eta) \epsilon_H (L_H^F)^{-1} \|d^t\|^{-1} \geq 6(1 - \eta) \epsilon_H (L_H^F)^{-1} \|d^t\|^{-1}$ . Hence,  
669 (6.10) also holds in this case.

670 Combining the above two cases, we conclude that (6.10) holds for any  $j \geq 0$  that  
671 violates (3.10). By this and  $\theta \in (0, 1)$ , one can see that all  $j \geq 0$  that violate (3.10)  
672 must be bounded above. It then follows that the step length  $\alpha_t$  associated with (3.10)  
673 is well-defined. We next prove (6.8). Observe from the definition of  $j_t$  in Algorithm 3.1  
674 that  $j = j_t - 1$  violates (3.10) and hence (6.10) holds for  $j = j_t - 1$ . Then, by (6.10)  
675 with  $j = j_t - 1$  and  $\alpha_t = \theta^{j_t}$ , one has

$$676 \quad (6.12) \quad \alpha_t = \theta^{j_t} \geq \sqrt{\min\{6(1 - \eta), 2\} \epsilon_H (L_H^F)^{-1}} \theta \|d^t\|^{-1/2},$$

677 which, along with (6.4) and  $\|\nabla F(x^t)\| \leq U_g^F$ , implies (6.8). This proves statement (i).

678 We next prove statement (ii) by considering two separate cases below.

679 Case 1)  $\alpha_t = 1$ . By this, one knows that (3.10) holds for  $j = 0$ . It then follows that  
680  $F(x^t + d^t) \leq F(x^t) \leq F(u^0)$ , which implies that (6.1) holds for  $x = x^t$  and  $y = x^t + d^t$ .  
681 By this and (6.6), one has

$$\begin{aligned}
682 \quad \|\nabla F(x^{t+1})\| = \|\nabla F(x^t + d^t)\| & \leq \|\nabla F(x^t + d^t) - \nabla F(x^t) - \nabla^2 F(x^t) d^t\| \\
& \quad + \|(\nabla^2 F(x^t) + 2\epsilon_H I) d^t + \nabla F(x^t)\| + 2\epsilon_H \|d^t\| \\
& \leq \frac{L_H^F}{2} \|d^t\|^2 + \frac{4+\zeta}{2} \epsilon_H \|d^t\|,
\end{aligned}$$

683 where the last inequality follows from (6.1) and (6.6). Solving the above inequality for  
684  $\|d^t\|$  and using the fact that  $\|d^t\| > 0$ , we obtain that

$$\begin{aligned}
685 \quad \|d^t\| & \geq \frac{-(4+\zeta)\epsilon_H + \sqrt{(4+\zeta)^2 \epsilon_H^2 + 8L_H^F \|\nabla F(x^{t+1})\|}}{2L_H^F} \\
& \geq \frac{-(4+\zeta)\epsilon_H + \sqrt{(4+\zeta)^2 \epsilon_H^2 + 8L_H^F \epsilon_H^2}}{2L_H^F} \min\{\|\nabla F(x^{t+1})\|/\epsilon_H^2, 1\} \\
& = \frac{4}{4+\zeta + \sqrt{(4+\zeta)^2 + 8L_H^F}} \min\{\|\nabla F(x^{t+1})\|/\epsilon_H, \epsilon_H\},
\end{aligned}$$

686 where the second inequality follows from the inequality  $-a + \sqrt{a^2 + bs} \geq (-a + \sqrt{a^2 + b}) \min\{s, 1\}$  for all  $a, b, s \geq 0$ , which can be verified by performing a rationalization to the terms  $-a + \sqrt{a^2 + b}$  and  $-a + \sqrt{a^2 + bs}$ , respectively. By this,  $\alpha_t = 1$ ,  
687 (3.10) and (3.13), one can see that (6.9) holds.

688 Case 2)  $\alpha_t < 1$ . It then follows that  $j = 0$  violates (3.10) and hence (6.10) holds  
689 for  $j = 0$ . Now, letting  $j = 0$  in (6.10), we obtain that  $\|d^t\| \geq \min\{6(1 - \eta), 2\} \epsilon_H / L_H^F$ ,  
690 which together with (3.10) and (6.12) implies that

$$693 \quad F(x^t) - F(x^{t+1}) \geq \eta \epsilon_H \theta^{2j_t} \|d^t\|^2 \geq \eta \frac{\min\{6(1 - \eta), 2\} \epsilon_H^2}{L_H^F} \theta^2 \|d^t\|^2 \geq \eta \left[ \frac{\min\{6(1 - \eta), 2\} \theta}{L_H^F} \right]^2 \epsilon_H^3.$$

694 By this and (3.13), one can see that (6.9) also holds in this case.  $\square$

695 The following lemma shows that when the search direction  $d^t$  in Algorithm 3.1 is  
696 of type 'NC', the line search step results in a sufficient reduction on  $F$  as well.

697 LEMMA 6.3. *Suppose that Assumption 3.1 holds and the direction  $d^t$  results from  
698 either the output  $d$  of Algorithm A.1 with  $d\_type=NC$  or the output  $v$  of Algorithm B.1  
699 at some iteration  $t$  of Algorithm 3.1. Let  $c_{nc}$  be defined in (3.14). Then the following  
700 statements hold.*

701 (i) The step length  $\alpha_t$  is well-defined, and  $\alpha_t \geq \min\{1, \theta/L_H^F, 3(1-\eta)\theta/L_H^F\}$ .  
 702 (ii) The next iterate  $x^{t+1} = x^t + \alpha_t d^t$  satisfies  $F(x^t) - F(x^{t+1}) \geq c_{\text{nc}} \epsilon_H^3$ .

703 *Proof.* Observe that  $F$  is descent along the iterates (whenever well-defined) generated by Algorithm 3.1. Using this and  $x^0 = u^0$ , we have  $F(x^t) \leq F(u^0)$ . By the assumption on  $d^t$ , one can see from Algorithm 3.1 that  $d^t$  is a negative curvature direction given in (3.7) or (3.9). Also, notice that the vector  $v$  returned from Algorithm B.1 satisfies  $\|v\| = 1$ . By these, Lemma 6.1(ii), (3.7) and (3.9), one can observe that

708 (6.13) 
$$\nabla F(x^t)^T d^t \leq 0, \quad (d^t)^T \nabla^2 F(x^t) d^t = -\|d^t\|^3 < 0.$$

709 We first prove statement (i). If (3.11) holds for  $j = 0$ , then  $\alpha_t = 1$ , which clearly  
 710 implies that  $\alpha_t \geq \min\{1, \theta/L_H^F, 3(1-\eta)\theta/L_H^F\}$ . We now suppose that (3.11) fails for  
 711  $j = 0$ . Claim that for all  $j \geq 0$  that violate (3.11), it holds that

712 (6.14) 
$$\theta^j \geq \min\{1/L_H^F, 3(1-\eta)/L_H^F\}.$$

713 Indeed, suppose that (3.11) is violated by some  $j \geq 0$ . We now show that (6.14) holds  
 714 for such  $j$  by considering two separate cases below.

715 Case 1)  $F(x^t + \theta^j d^t) > F(x^t)$ . Let  $\phi(\alpha) = F(x^t + \alpha d^t)$ . Then  $\phi(\theta^j) > \phi(0)$ .  
 716 Also, by (6.13), one has  $\phi'(0) = \nabla F(x^t)^T d^t \leq 0$  and  $\phi''(0) = (d^t)^T \nabla^2 F(x^t) d^t < 0$ .  
 717 Using these, we can observe that there exists a local minimizer  $\alpha^* \in (0, \theta^j)$  of  
 718  $\phi$  such that  $\phi(\alpha^*) < \phi(0)$ , namely,  $F(x^t + \alpha^* d^t) < F(x^t)$ . By the second-order  
 719 optimality condition of  $\phi$  at  $\alpha^*$ , one has  $\phi''(\alpha^*) = (d^t)^T \nabla^2 F(x^t + \alpha^* d^t) d^t \geq 0$ . Since  
 720  $F(x^t + \alpha^* d^t) < F(x^t) \leq F(u^0)$ , it follows that (3.2) holds for  $x = x^t$  and  $y = x^t + \alpha^* d^t$ .  
 721 Using this, the second relation in (6.13) and  $(d^t)^T \nabla^2 F(x^t + \alpha^* d^t) d^t \geq 0$ , we obtain  
 722 that

723 
$$L_H^F \alpha^* \|d^t\|^3 \stackrel{(3.2)}{\geq} \|d^t\|^2 \|\nabla^2 F(x^t + \alpha^* d^t) - \nabla^2 F(x^t)\|$$
  
 724 (6.15) 
$$\geq (d^t)^T (\nabla^2 F(x^t + \alpha^* d^t) - \nabla^2 F(x^t)) d^t \geq -(d^t)^T \nabla^2 F(x^t) d^t = \|d^t\|^3.$$

725 Recall from (6.13) that  $d^t \neq 0$ . It then follows from (6.15) that  $\alpha^* \geq 1/L_H^F$ , which  
 726 along with  $\theta^j > \alpha^*$  implies that  $\theta^j > 1/L_H^F$ . Hence, (6.14) holds in this case.

727 Case 2)  $F(x^t + \theta^j d^t) \leq F(x^t)$ . It follows from this and  $F(x^t) \leq F(u^0)$  that (6.2)  
 728 holds for  $x = x^t$  and  $y = x^t + \theta^j d^t$ . By this and the fact that  $j$  violates (3.11), one has

729 
$$-\frac{\eta}{2} \theta^{2j} \|d^t\|^3 \leq F(x^t + \theta^j d^t) - F(x^t) \stackrel{(6.2)}{\leq} \theta^j \nabla F(x^t)^T d^t + \frac{\theta^{2j}}{2} (d^t)^T \nabla^2 F(x^t) d^t + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3$$
  

$$\stackrel{(6.13)}{\leq} -\frac{\theta^{2j}}{2} \|d^t\|^3 + \frac{L_H^F}{6} \theta^{3j} \|d^t\|^3,$$

730 which together with  $d^t \neq 0$  implies that  $\theta^j \geq 3(1-\eta)/L_H^F$ . Hence, (6.14) also holds in  
 731 this case.

732 Combining the above two cases, we conclude that (6.14) holds for any  $j \geq 0$  that  
 733 violates (3.11). By this and  $\theta \in (0, 1)$ , one can see that all  $j \geq 0$  that violate (3.11)  
 734 must be bounded above. It then follows that the step length  $\alpha_t$  associated with (3.11)  
 735 is well-defined. We next derive a lower bound for  $\alpha_t$ . Notice from the definition of  $j_t$  in  
 736 Algorithm 3.1 that  $j = j_t - 1$  violates (3.11) and hence (6.14) holds for  $j = j_t - 1$ . Then,  
 737 by (6.14) with  $j = j_t - 1$  and  $\alpha_t = \theta^{j_t}$ , one has  $\alpha_t = \theta^{j_t} \geq \min\{\theta/L_H^F, 3(1-\eta)\theta/L_H^F\}$ ,  
 738 which immediately yields  $\alpha_t \geq \min\{1, \theta/L_H^F, 3(1-\eta)\theta/L_H^F\}$  as desired.

739 We next prove statement (ii) by considering two separate cases below.

740 Case 1)  $d^t$  results from the output  $d$  of Algorithm A.1 with  $\text{d\_type}=\text{NC}$ . It then  
 741 follows from (6.7) that  $\|d^t\| \geq \epsilon_H$ . This together with (3.11) and statement (i) implies  
 742 that statement (ii) holds.

743 Case 2)  $d^t$  results from the output  $v$  of Algorithm B.1. Notice from Algorithm B.1  
 744 that  $\|v\| = 1$  and  $v^T \nabla^2 F(x^t) v \leq -\epsilon_H/2$ , which along with (3.9) yields  $\|d^t\| \geq \epsilon_H/2$ .  
 745 By this, (3.11) and statement (i), one can see that statement (ii) again holds.  $\square$

746 *Proof of Theorem 3.2.* For notational convenience, we let  $\{x^t\}_{t \in \mathbb{T}}$  denote all the  
 747 iterates generated by Algorithm 3.1, where  $\mathbb{T}$  is a set of consecutive nonnegative  
 748 integers starting from 0. Notice that  $F$  is descent along the iterates generated by  
 749 Algorithm 3.1, which together with  $x^0 = u^0$  implies that  $x^t \in \{x : F(x) \leq F(u^0)\}$ . It  
 750 then follows from (3.3) that  $\|\nabla^2 F(x^t)\| \leq U_H^F$  holds for all  $t \in \mathbb{T}$ .

751 (i) Suppose for contradiction that the total number of calls of Algorithm B.1  
 752 in Algorithm 3.1 is more than  $T_2$ . Notice from Algorithm 3.1 and Lemma 6.3(ii)  
 753 that each of these calls, except the last one, returns a sufficiently negative curvature  
 754 direction, and each of them results in a reduction on  $F$  of at least  $c_{nc}\epsilon_H^3$ . Hence,  
 755  $T_2 c_{nc}\epsilon_H^3 \leq \sum_{t \in \mathbb{T}} [F(x^t) - F(x^{t+1})] \leq F(x^0) - F_{\text{low}} = F_{\text{hi}} - F_{\text{low}}$ , which contradicts  
 756 the definition of  $T_2$  given in (3.12). Hence, statement (i) of Theorem 3.2 holds.

757 (ii) Suppose for contradiction that the total number of calls of Algorithm A.1  
 758 in Algorithm 3.1 is more than  $T_1$ . Observe that if Algorithm A.1 is called at some  
 759 iteration  $t$  and generates the next iterate  $x^{t+1}$  satisfying  $\|\nabla F(x^{t+1})\| \leq \epsilon_g$ , then  
 760 Algorithm B.1 must be called at the next iteration  $t+1$ . In view of this and statement  
 761 (i) of Theorem 3.2, we see that the total number of such iterations  $t$  is at most  $T_2$ .  
 762 Hence, the total number of iterations  $t$  of Algorithm 3.1 at which Algorithm A.1  
 763 is called and generates the next iterate  $x^{t+1}$  satisfying  $\|\nabla F(x^{t+1})\| > \epsilon_g$  is at least  
 764  $T_1 - T_2 + 1$ . Moreover, for each of such iterations  $t$ , we observe from Lemmas 6.2(ii)  
 765 and 6.3(ii) that  $F(x^t) - F(x^{t+1}) \geq \min\{c_{\text{sol}}, c_{\text{nc}}\} \min\{\epsilon_g^2 \epsilon_H^{-1}, \epsilon_H^3\}$ . It then follows that  
 766  $(T_1 - T_2 + 1) \min\{c_{\text{sol}}, c_{\text{nc}}\} \min\{\epsilon_g^2 \epsilon_H^{-1}, \epsilon_H^3\} \leq \sum_{t \in \mathbb{T}} [F(x^t) - F(x^{t+1})] \leq F_{\text{hi}} - F_{\text{low}}$ ,  
 767 which contradicts the definition of  $T_1$  and  $T_2$  given in (3.12). Hence, statement (ii) of  
 768 Theorem 3.2 holds.

769 (iii) Notice that either Algorithm A.1 or B.1 is called at each iteration of Algo-  
 770 rithm 3.1. It follows from this and statements (i) and (ii) of Theorem 3.2 that the total  
 771 number of iterations of Algorithm 3.1 is at most  $T_1 + T_2$ . In addition, the relation  
 772 (3.15) follows from (3.13), (3.14) and (3.12). One can also observe that the output  $x^t$   
 773 of Algorithm 3.1 satisfies  $\|\nabla F(x^t)\| \leq \epsilon_g$  deterministically and  $\lambda_{\min}(\nabla^2 F(x^t)) \geq -\epsilon_H$   
 774 with probability at least  $1 - \delta$  for some  $0 \leq t \leq T_1 + T_2$ , where the latter part is due  
 775 to Algorithm B.1. This completes the proof of statement (ii) of Theorem 3.2.

776 (iv) By Theorem A.1 with  $(H, \epsilon) = (\nabla^2 F(x^t), \epsilon_H)$  and the fact that  $\|\nabla^2 F(x^t)\| \leq$   
 777  $U_H^F$ , one can observe that the number of Hessian-vector products required by each call  
 778 of Algorithm A.1 with input  $U = 0$  is at most  $\tilde{\mathcal{O}}(\min\{n, (U_H^F/\epsilon_H)^{1/2}\})$ . In addition,  
 779 by Theorem B.1 with  $(H, \epsilon) = (\nabla^2 F(x^t), \epsilon_H)$ ,  $\|\nabla^2 F(x^t)\| \leq U_H^F$ , and the fact that  
 780 each iteration of the Lanczos method requires only one matrix-vector product, one  
 781 can observe that the number of Hessian-vector products required by each call of  
 782 Algorithm B.1 is also at most  $\tilde{\mathcal{O}}(\min\{n, (U_H^F/\epsilon_H)^{1/2}\})$ . Based on these observations  
 783 and statement (iii) of Theorem 3.2, we see that statement (iv) of this theorem holds.  $\square$

784 **6.2. Proof of the main results in Section 4.** Recall from Assumption 4.1(a)  
 785 that  $\|c(z_{\epsilon_1})\| \leq \epsilon_1/2 < 1$ . By virtue of this, (4.2) and the definition of  $\tilde{c}$  in (4.4), we  
 786 obtain that

$$787 (6.16) \quad f(x) + \gamma\|\tilde{c}(x)\|^2 \geq f(x) + \gamma\|c(x)\|^2/2 - \gamma\|c(z_{\epsilon_1})\|^2 \geq f_{\text{low}} - \gamma, \quad \forall x \in \mathbb{R}^n.$$

788 We now prove the following auxiliary lemma that will be used frequently later.

789 LEMMA 6.4. *Suppose that Assumption 4.1 holds. Let  $\gamma$ ,  $f_{\text{hi}}$  and  $f_{\text{low}}$  be given in  
 790 Assumption 4.1. Assume that  $\rho > 2\gamma$ ,  $\lambda \in \mathbb{R}^m$ , and  $x \in \mathbb{R}^n$  satisfy*

791 (6.17)  $\tilde{\mathcal{L}}(x, \lambda; \rho) \leq f_{\text{hi}},$

792 where  $\tilde{\mathcal{L}}$  is defined in (4.5). Then the following statements hold.

793 (i)  $f(x) \leq f_{\text{hi}} + \|\lambda\|^2/(2\rho).$

794 (ii)  $\|\tilde{c}(x)\| \leq \sqrt{2(f_{\text{hi}} - f_{\text{low}} + \gamma)/(\rho - 2\gamma)} + \|\lambda\|^2/(\rho - 2\gamma)^2 + \|\lambda\|/(\rho - 2\gamma).$

795 (iii) If  $\rho \geq \|\lambda\|^2/(2\tilde{\delta}_f)$  for some  $\tilde{\delta}_f > 0$ , then  $f(x) \leq f_{\text{hi}} + \tilde{\delta}_f$ .

796 (iv) If

797 (6.18)  $\rho \geq 2(f_{\text{hi}} - f_{\text{low}} + \gamma)\tilde{\delta}_c^{-2} + 2\|\lambda\|\tilde{\delta}_c^{-1} + 2\gamma$

798 for some  $\tilde{\delta}_c > 0$ , then  $\|\tilde{c}(x)\| \leq \tilde{\delta}_c$ .

799 *Proof.* (i) It follows from (6.17) and the definition of  $\tilde{\mathcal{L}}$  in (4.5) that

800  $f_{\text{hi}} \geq f(x) + \lambda^T \tilde{c}(x) + \frac{\rho}{2} \|\tilde{c}(x)\|^2 = f(x) + \frac{\rho}{2} \left\| \tilde{c}(x) + \frac{\lambda}{\rho} \right\|^2 - \frac{\|\lambda\|^2}{2\rho} \geq f(x) - \frac{\|\lambda\|^2}{2\rho}.$

801 Hence, statement (i) holds.

802 (ii) In view of (6.16) and (6.17), one has

803 
$$\begin{aligned} f_{\text{hi}} &\stackrel{(6.17)}{\geq} f(x) + \lambda^T \tilde{c}(x) + \frac{\rho}{2} \|\tilde{c}(x)\|^2 = f(x) + \gamma \|\tilde{c}(x)\|^2 + \frac{\rho-2\gamma}{2} \left\| \tilde{c}(x) + \frac{\lambda}{\rho-2\gamma} \right\|^2 - \frac{\|\lambda\|^2}{2(\rho-2\gamma)} \\ &\stackrel{(6.16)}{\geq} f_{\text{low}} - \gamma + \frac{\rho-2\gamma}{2} \left\| \tilde{c}(x) + \frac{\lambda}{\rho-2\gamma} \right\|^2 - \frac{\|\lambda\|^2}{2(\rho-2\gamma)}. \end{aligned}$$

804 It then follows that  $\left\| \tilde{c}(x) + \frac{\lambda}{\rho-2\gamma} \right\| \leq \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho-2\gamma} + \frac{\|\lambda\|^2}{(\rho-2\gamma)^2}}$ , which implies that statement (ii) holds.

805 (iii) Statement (iii) immediately follows from statement (i) and  $\rho \geq \|\lambda\|^2/(2\tilde{\delta}_f)$ .

806 (iv) Suppose that (6.18) holds. Multiplying both sides of (6.18) by  $\tilde{\delta}_c^2$  and rearranging the terms, we have  $(\rho - 2\gamma)\tilde{\delta}_c^2 - 2\|\lambda\|\tilde{\delta}_c - 2(f_{\text{hi}} - f_{\text{low}} + \gamma) \geq 0$ . Recall that  $\rho > 2\gamma$  and  $\tilde{\delta}_c > 0$ . Solving this inequality for  $\tilde{\delta}_c$  yields

807  $\tilde{\delta}_c \geq \sqrt{2(f_{\text{hi}} - f_{\text{low}} + \gamma)/(\rho - 2\gamma) + \|\lambda\|^2/(\rho - 2\gamma)^2} + \|\lambda\|/(\rho - 2\gamma),$

808 which along with statement (ii) implies that  $\|\tilde{c}(x)\| \leq \tilde{\delta}_c$ . Hence, statement (iv) holds.  $\square$

809 *Proof of Lemma 4.4.* (i) Let  $x$  be any point such that  $\tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)$ .

810 It then follows from (4.9) that  $\tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq f_{\text{hi}}$ . By this,  $\|\lambda^k\| \leq \Lambda$ ,  $\rho_k \geq \rho_0 > 2\gamma$ ,  $\delta_{f,1} \leq \delta_f$ ,  $\delta_{c,1} \leq \delta_c$ , and Lemma 6.4 with  $(\lambda, \rho) = (\lambda^k, \rho_k)$ , one has  $f(x) \leq f_{\text{hi}} + \|\lambda^k\|^2/(2\rho_k) \leq f_{\text{hi}} + \Lambda^2/(2\rho_0) = f_{\text{hi}} + \delta_{f,1} \leq f_{\text{hi}} + \delta_f$  and

811 (6.19) 
$$\begin{aligned} \|\tilde{c}(x)\| &\leq \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho_k - 2\gamma} + \frac{\|\lambda^k\|^2}{(\rho_k - 2\gamma)^2}} + \frac{\|\lambda^k\|}{\rho_k - 2\gamma} \\ &\leq \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho_0 - 2\gamma} + \frac{\Lambda^2}{(\rho_0 - 2\gamma)^2}} + \frac{\Lambda}{\rho_0 - 2\gamma} = \delta_{c,1} \leq \delta_c. \end{aligned}$$

812 Also, recall from the definition of  $\tilde{c}$  in (4.4) and  $\|c(z_{\epsilon_1})\| \leq 1$  that  $\|c(x)\| \leq 1 + \|\tilde{c}(x)\|$ . This together with the above inequalities and (4.3) implies  $x \in \mathcal{S}(\delta_f, \delta_c)$ . Hence, statement (i) of Lemma 4.4 holds.

813 (ii) Note that  $\inf_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) = \inf_{x \in \mathbb{R}^n} \{\tilde{\mathcal{L}}(x, \lambda^k; \rho_k) : \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)\}$ .

814 Consequently, to prove statement (ii) of Lemma 4.4, it suffices to show that

815 (6.20)  $\inf_{x \in \mathbb{R}^n} \{\tilde{\mathcal{L}}(x, \lambda^k; \rho_k) : \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)\} \geq f_{\text{low}} - \gamma - \Lambda\delta_c.$

816 To this end, let  $x$  be any point satisfying  $\tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)$ . We then know from (6.19) that  $\|\tilde{c}(x)\| \leq \delta_c$ . By this,  $\|\lambda^k\| \leq \Lambda$ ,  $\rho_k > 2\gamma$ , and (6.16), one has

$$\begin{aligned}
& \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) = f(x) + \gamma \|\tilde{c}(x)\|^2 + (\lambda^k)^T \tilde{c}(x) + \frac{\rho_k - 2\gamma}{2} \|\tilde{c}(x)\|^2 \\
& \geq f(x) + \gamma \|\tilde{c}(x)\|^2 - \Lambda \|\tilde{c}(x)\| \geq f_{\text{low}} - \gamma - \Lambda \delta_c,
\end{aligned}$$

826 and hence (6.20) holds as desired.  $\square$

827 *Proof of Theorem 4.6.* Suppose that Algorithm 4.1 terminates at some iteration  
828  $k$ , that is,  $\tau_k^g \leq \epsilon_1$ ,  $\tau_k^H \leq \epsilon_2$ , and  $\|c(x^{k+1})\| \leq \epsilon_1$  hold. Then, by  $\tau_k^g \leq \epsilon_1$ ,  $\tilde{\lambda}^{k+1} =$   
829  $\lambda^k + \rho_k \tilde{c}(x^{k+1})$ ,  $\nabla \tilde{c} = \nabla c$  and the second relation in (4.6), one has  $\|\nabla f(x^{k+1}) +$   
830  $\nabla c(x^{k+1}) \tilde{\lambda}^{k+1}\| = \|\nabla f(x^{k+1}) + \nabla \tilde{c}(x^{k+1})(\lambda^k + \rho_k \tilde{c}(x^{k+1}))\| = \|\nabla_x \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)\| \leq$   
831  $\tau_k^g \leq \epsilon_1$ . Hence,  $(x^{k+1}, \tilde{\lambda}^{k+1})$  satisfies the first relation in (2.4). In addition, by (4.7)  
832 and  $\tau_k^H \leq \epsilon_2$ , one can show that  $\lambda_{\min}(\nabla_{xx}^2 \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)) \geq -\epsilon_2$  with probability  
833 at least  $1 - \delta$ , which leads to  $d^T \nabla_{xx}^2 \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k) d \geq -\epsilon_2 \|d\|^2$  for all  $d \in \mathbb{R}^n$   
834 with probability at least  $1 - \delta$ . Using this,  $\tilde{\lambda}^{k+1} = \lambda^k + \rho_k \tilde{c}(x^{k+1})$ ,  $\nabla \tilde{c} = \nabla c$ , and  
835  $\nabla^2 \tilde{c}_i = \nabla^2 c_i$  for  $1 \leq i \leq m$ , we see that with probability at least  $1 - \delta$ , it holds that  
836  $d^T (\nabla^2 f(x^{k+1}) + \sum_{i=1}^m \tilde{\lambda}_i^{k+1} \nabla^2 c_i(x^{k+1}) + \rho_k \nabla c(x^{k+1}) \nabla c(x^{k+1})^T) d \geq -\epsilon_2 \|d\|^2$  for all  
837  $d \in \mathbb{R}^n$ , which implies  $d^T (\nabla^2 f(x^{k+1}) + \sum_{i=1}^m \tilde{\lambda}_i^{k+1} \nabla^2 c_i(x^{k+1})) d \geq -\epsilon_2 \|d\|^2$  for all  
838  $d \in \mathcal{C}(x^{k+1})$ , where  $\mathcal{C}(\cdot)$  is defined in (2.3). Hence,  $(x^{k+1}, \tilde{\lambda}^{k+1})$  satisfies (2.5) with  
839 probability at least  $1 - \delta$ . Combining these with  $\|c(x^{k+1})\| \leq \epsilon_1$ , we conclude that  
840  $x^{k+1}$  is a deterministic  $\epsilon_1$ -FOSP of (1.1) and an  $(\epsilon_1, \epsilon_2)$ -SOSP of (1.1) with probability  
841 at least  $1 - \delta$ . Hence, Theorem 4.6 holds.  $\square$

842 *Proof of Theorem 4.8.* It follows from (4.14) that  $\rho_{\epsilon_1} \geq 2\rho_0$ . By this, one has

$$(6.21) \quad K_{\epsilon_1} \stackrel{(4.12)}{=} \lceil \log \epsilon_1 / \log \omega_1 \rceil \stackrel{(4.11)}{=} \lceil \log 2 / \log r \rceil \leq \log(\rho_{\epsilon_1} \rho_0^{-1}) / \log r + 1.$$

844 Notice that  $\{\rho_k\}$  is either unchanged or increased by a ratio  $r$  as  $k$  increases. By this  
845 fact and (6.21), we see that

$$(6.22) \quad \max_{0 \leq k \leq K_{\epsilon_1}} \rho_k \leq r^{K_{\epsilon_1}} \rho_0 \stackrel{(6.21)}{\leq} r^{\frac{\log(\rho_{\epsilon_1} \rho_0^{-1})}{\log r} + 1} \rho_0 = r \rho_{\epsilon_1}.$$

847 In addition, notice that  $\rho_k > 2\gamma$  and  $\|\lambda^k\| \leq \Lambda$ . Using these, (4.1), the first relation in  
848 (4.6), and Lemma 6.4(ii) with  $(x, \lambda, \rho) = (x^{k+1}, \lambda^k, \rho_k)$ , we obtain that

$$(6.23) \quad \|\tilde{c}(x^{k+1})\| \leq \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho_k - 2\gamma} + \frac{\|\lambda^k\|^2}{(\rho_k - 2\gamma)^2}} + \frac{\|\lambda^k\|}{\rho_k - 2\gamma} \leq \sqrt{\frac{2(f_{\text{hi}} - f_{\text{low}} + \gamma)}{\rho_k - 2\gamma} + \frac{\Lambda^2}{(\rho_k - 2\gamma)^2}} + \frac{\Lambda}{\rho_k - 2\gamma}.$$

850 Also, we observe from  $\|c(z_{\epsilon_1})\| \leq \epsilon_1/2$  and the definition of  $\tilde{c}$  in (4.4) that

$$(6.24) \quad \|c(x^{k+1})\| \leq \|\tilde{c}(x^{k+1})\| + \|c(z_{\epsilon_1})\| \leq \|\tilde{c}(x^{k+1})\| + \epsilon_1/2.$$

852 We now prove that  $\overline{K}_{\epsilon_1}$  is finite. Suppose for contradiction that  $\overline{K}_{\epsilon_1}$  is infinite.  
853 It then follows from this and (4.15) that  $\|c(x^{k+1})\| > \epsilon_1$  for all  $k \geq K_{\epsilon_1}$ , which  
854 along with (6.24) implies that  $\|\tilde{c}(x^{k+1})\| > \epsilon_1/2$  for all  $k \geq K_{\epsilon_1}$ . It then follows that  
855  $\|\tilde{c}(x^{k+1})\| > \alpha \|\tilde{c}(x^k)\|$  must hold for infinitely many  $k$ 's. Using this and the update  
856 scheme on  $\{\rho_k\}$ , we deduce that  $\rho_{k+1} = r \rho_k$  holds for infinitely many  $k$ 's, which  
857 together with the monotonicity of  $\{\rho_k\}$  implies that  $\rho_k \rightarrow \infty$  as  $k \rightarrow \infty$ . By this and  
858 (6.23), one can see that  $\|\tilde{c}(x^{k+1})\| \rightarrow 0$  as  $k \rightarrow \infty$ , which contradicts the fact that  
859  $\|\tilde{c}(x^{k+1})\| > \epsilon_1/2$  holds for all  $k \geq K_{\epsilon_1}$ . Hence,  $\overline{K}_{\epsilon_1}$  is finite. In addition, notice from  
860 (4.11), (4.12) and (4.13) that  $(\tau_k^g, \tau_k^H) = (\epsilon_1, \epsilon_2)$  for all  $k \geq K_{\epsilon_1}$ . This along with the  
861 termination criterion of Algorithm 4.1 and the definition of  $\overline{K}_{\epsilon_1}$  implies that Algorithm  
862 4.1 must terminate at iteration  $\overline{K}_{\epsilon_1}$ .

863 We next show that (4.16) and  $\rho_k \leq r\rho_{\epsilon_1}$  hold for  $0 \leq k \leq \bar{K}_{\epsilon_1}$  by considering two  
 864 separate cases below.

865 Case 1)  $\|c(x^{K_{\epsilon_1}+1})\| \leq \epsilon_1$ . By this and (4.15), one can see that  $\bar{K}_{\epsilon_1} = K_{\epsilon_1}$ , which  
 866 together with (6.21) and (6.22) implies that (4.16) and  $\rho_k \leq r\rho_{\epsilon_1}$  hold for  $0 \leq k \leq \bar{K}_{\epsilon_1}$ .

867 Case 2)  $\|c(x^{K_{\epsilon_1}+1})\| > \epsilon_1$ . By this and (4.15), one can observe that  $\bar{K}_{\epsilon_1} > K_{\epsilon_1}$   
 868 and also  $\|c(x^{k+1})\| > \epsilon_1$  for all  $K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1$ , which together with (6.24) implies

$$869 \quad (6.25) \quad \|\tilde{c}(x^{k+1})\| > \epsilon_1/2, \quad \forall K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1.$$

870 It then follows from  $\|\lambda^k\| \leq \Lambda$ , (4.1), the first relation in (4.6), and Lemma 6.4(iv)  
 871 with  $(x, \lambda, \rho, \tilde{\delta}_c) = (x^{k+1}, \lambda^k, \rho_k, \epsilon_1/2)$  that

$$872 \quad (6.26) \quad \begin{aligned} \rho_k &< 8(f_{hi} - f_{low} + \gamma)\epsilon_1^{-2} + 4\|\lambda^k\|\epsilon_1^{-1} + 2\gamma \\ &\leq 8(f_{hi} - f_{low} + \gamma)\epsilon_1^{-2} + 4\Lambda\epsilon_1^{-1} + 2\gamma \stackrel{(4.14)}{\leq} \rho_{\epsilon_1}, \quad \forall K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1. \end{aligned}$$

873 Combining this relation, (6.22), and the fact  $\rho_{\bar{K}_{\epsilon_1}} \leq r\rho_{\bar{K}_{\epsilon_1}-1}$ , we conclude that  
 874  $\rho_k \leq r\rho_{\epsilon_1}$  holds for  $0 \leq k \leq \bar{K}_{\epsilon_1}$ . It remains to show that (4.16) holds. To this  
 875 end, let  $\mathbb{K} = \{k : \rho_{k+1} = r\rho_k, K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 2\}$ . It follows from (6.26) and the  
 876 update scheme of  $\rho_k$  that  $r^{|\mathbb{K}|}\rho_{K_{\epsilon_1}} = \max_{K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1}-1}\{\rho_k\} \leq \rho_{\epsilon_1}$ , which together  
 877 with  $\rho_{K_{\epsilon_1}} \geq \rho_0$  implies that

$$878 \quad (6.27) \quad |\mathbb{K}| \leq \log(\rho_{\epsilon_1}\rho_{K_{\epsilon_1}}^{-1})/\log r \leq \log(\rho_{\epsilon_1}\rho_0^{-1})/\log r.$$

879 Let  $\{k_1, k_2, \dots, k_{|\mathbb{K}|}\}$  denote all the elements of  $\mathbb{K}$  arranged in ascending order, and  
 880 let  $k_0 = K_{\epsilon_1}$  and  $k_{|\mathbb{K}|+1} = \bar{K}_{\epsilon_1} - 1$ . We next derive an upper bound for  $k_{j+1} - k_j$   
 881 for  $j = 0, 1, \dots, |\mathbb{K}|$ . By the definition of  $\mathbb{K}$ , one can observe that  $\rho_k = \rho_{k'}$  for  
 882  $k_j < k, k' \leq k_{j+1}$ . Using this and the update scheme of  $\rho_k$ , we deduce that

$$883 \quad (6.28) \quad \|\tilde{c}(x^{k+1})\| \leq \alpha\|\tilde{c}(x^k)\|, \quad \forall k_j < k < k_{j+1}.$$

884 On the other hand, by (4.10), (6.23) and  $\rho_k \geq \rho_0$ , one has  $\|\tilde{c}(x^{k+1})\| \leq \delta_{c,1}$  for  
 885  $0 \leq k \leq \bar{K}_{\epsilon_1}$ . By this and (6.25), one can see that

$$886 \quad (6.29) \quad \epsilon_1/2 < \|\tilde{c}(x^{k+1})\| \leq \delta_{c,1}, \quad \forall K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1.$$

887 Now, note that either  $k_{j+1} - k_j = 1$  or  $k_{j+1} - k_j > 1$ . In the latter case, we can apply  
 888 (6.28) with  $k = k_{j+1} - 1, \dots, k_j + 1$  together with (6.29) to deduce that

$$889 \quad \epsilon_1/2 < \|\tilde{c}(x^{k_{j+1}})\| \leq \alpha\|\tilde{c}(x^{k_{j+1}-1})\| \leq \dots \leq \alpha^{k_{j+1}-k_j-1}\|\tilde{c}(x^{k_j+1})\| \leq \alpha^{k_{j+1}-k_j-1}\delta_{c,1}$$

890 for all  $j = 0, 1, \dots, |\mathbb{K}|$ . Combining these two cases, we have

$$891 \quad (6.30) \quad k_{j+1} - k_j \leq |\log(\epsilon_1(2\delta_{c,1})^{-1})/\log \alpha| + 1, \quad \forall j = 0, 1, \dots, |\mathbb{K}|.$$

892 Summing up these inequalities, and using (6.21), (6.27),  $k_0 = K_{\epsilon_1}$  and  $k_{|\mathbb{K}|+1} = \bar{K}_{\epsilon_1} - 1$ ,  
 893 we have

$$894 \quad \begin{aligned} \bar{K}_{\epsilon_1} &= 1 + k_{|\mathbb{K}|+1} = 1 + k_0 + \sum_{j=0}^{|\mathbb{K}|} (k_{j+1} - k_j) \\ 895 &\stackrel{(6.30)}{\leq} 1 + K_{\epsilon_1} + (|\mathbb{K}| + 1) \left( \left| \frac{\log(\epsilon_1(2\delta_{c,1})^{-1})}{\log \alpha} \right| + 1 \right) \\ 896 \quad (6.31) \quad &\leq 2 + \frac{\log(\rho_{\epsilon_1}\rho_0^{-1})}{\log r} + \left( \frac{\log(\rho_{\epsilon_1}\rho_0^{-1})}{\log r} + 1 \right) \left( \left| \frac{\log(\epsilon_1(2\delta_{c,1})^{-1})}{\log \alpha} \right| + 1 \right) \\ 897 &= 1 + \left( \frac{\log(\rho_{\epsilon_1}\rho_0^{-1})}{\log r} + 1 \right) \left( \left| \frac{\log(\epsilon_1(2\delta_{c,1})^{-1})}{\log \alpha} \right| + 2 \right), \end{aligned}$$

898 where the second inequality is due to (6.21) and (6.27). Hence, (4.16) also holds in  
 899 this case.  $\square$

900 We next prove Theorem 4.10. Before proceeding, we introduce some notation that  
901 will be used shortly. Let  $L_{k,H}$  denote the Lipschitz constant of  $\nabla_{xx}^2 \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  on the  
902 convex open neighborhood  $\Omega(\delta_f, \delta_c)$  of  $\mathcal{S}(\delta_f, \delta_c)$ , where  $\mathcal{S}(\delta_f, \delta_c)$  is defined in (4.3),  
903 and let  $U_{k,H} = \sup_{x \in \mathcal{S}(\delta_f, \delta_c)} \|\nabla_{xx}^2 \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)\|$ . Notice from (4.4) and (4.5) that

904 (6.32) 
$$\nabla_{xx}^2 \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) = \nabla^2 f(x) + \sum_{i=1}^m \lambda_i^k \nabla^2 c_i(x) + \rho_k \left( \nabla c(x) \nabla c(x)^T + \sum_{i=1}^m \tilde{c}_i(x) \nabla^2 c_i(x) \right).$$

905 By this,  $\|\lambda^k\| \leq \Lambda$ , the definition of  $\tilde{c}$ , and the Lipschitz continuity of  $\nabla^2 f$  and  $\nabla^2 c_i$ 's  
906 (see Assumption 4.1(c)), one can observe that there exist some constants  $L_1, L_2, U_1$   
907 and  $U_2$ , depending only on  $f, c, \Lambda, \delta_f$  and  $\delta_c$ , such that

908 (6.33) 
$$L_{k,H} \leq L_1 + \rho_k L_2, \quad U_{k,H} \leq U_1 + \rho_k U_2.$$

909 *Proof of Theorem 4.10.* Let  $T_k$  and  $N_k$  denote the number of iterations and matrix-  
910 vector products performed by Algorithm 3.1 at the outer iteration  $k$  of Algorithm 4.1,  
911 respectively. It then follows from Theorem 4.8 that the total number of iterations and  
912 matrix-vector products performed by Algorithm 3.1 in Algorithm 4.1 are  $\sum_{k=0}^{\bar{K}_{\epsilon_1}} T_k$   
913 and  $\sum_{k=0}^{\bar{K}_{\epsilon_1}} N_k$ , respectively. In addition, notice from (4.14) and Theorem 4.8 that  
914  $\rho_{\epsilon_1} = \mathcal{O}(\epsilon_1^{-2})$  and  $\rho_k \leq r \rho_{\epsilon_1}$ , which yield  $\rho_k = \mathcal{O}(\epsilon_1^{-2})$ .

915 We first claim that  $(\tau_k^g)^2 / \tau_k^H \geq \min\{\epsilon_1^2 / \epsilon_2, \epsilon_2^3\}$  holds for any  $k \geq 0$ . Indeed, let  
916  $\bar{t} = \log \epsilon_1 / \log \omega_1$  and  $\psi(t) = \max\{\epsilon_1, \omega_1^t\}^2 / \max\{\epsilon_2, \omega_2^t\}$  for all  $t \in \mathbb{R}$ . It then follows  
917 from (4.13) that  $\omega_1^{\bar{t}} = \epsilon_1$  and  $\omega_2^{\bar{t}} = \epsilon_2$ . By this and  $\omega_1, \omega_2 \in (0, 1)$ , one can observe  
918 that  $\psi(t) = (\omega_1^2 / \omega_2)^t$  if  $t \leq \bar{t}$  and  $\psi(t) = \epsilon_1^2 / \epsilon_2$  otherwise. This along with  $\epsilon_2 \in (0, 1)$   
919 implies that  $\min_{t \in [0, \infty)} \psi(t) = \min\{\psi(0), \psi(\bar{t})\} = \min\{1, \epsilon_1^2 / \epsilon_2\} \geq \min\{\epsilon_1^2 / \epsilon_2, \epsilon_2^3\}$ ,  
920 which together with (4.11) yields  $(\tau_k^g)^2 / \tau_k^H = \psi(k) \geq \min\{\epsilon_1^2 / \epsilon_2, \epsilon_2^3\}$  for all  $k \geq 0$ .

921 (i) From Lemma 4.4(i) and the definitions of  $\Omega(\delta_f, \delta_c)$  and  $L_{k,H}$ , we see that  
922  $L_{k,H}$  is a Lipschitz constant of  $\nabla_{xx}^2 \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$  on a convex open neighborhood of  $\{x : \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)\}$ . Also, recall from Lemma 4.4(ii) that  $\inf_{x \in \mathbb{R}^n} \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)$   
923  $\geq f_{\text{low}} - \gamma - \Lambda \delta_c$ . By these,  $\tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k) \leq f_{\text{hi}}$  (see (4.9)) and Theorem 3.2(iii) with  
924  $(F_{\text{hi}}, F_{\text{low}}, L_H^F, \epsilon_g, \epsilon_H) = (\tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k), f_{\text{low}} - \gamma - \Lambda \delta_c, L_{k,H}, \tau_k^g, \tau_k^H)$ , one has

925 (6.34) 
$$\begin{aligned} T_k &= \mathcal{O}((f_{\text{hi}} - f_{\text{low}} + \gamma + \Lambda \delta_c) L_{k,H}^2 \max\{(\tau_k^g)^{-2} \tau_k^H, (\tau_k^H)^{-3}\}) \\ &\stackrel{(6.33)}{=} \mathcal{O}(\rho_k^2 \max\{(\tau_k^g)^{-2} \tau_k^H, (\tau_k^H)^{-3}\}) = \mathcal{O}(\epsilon_1^{-4} \max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\}), \end{aligned}$$

926 where the last equality is from  $(\tau_k^g)^2 / \tau_k^H \geq \min\{\epsilon_1^2 / \epsilon_2, \epsilon_2^3\}$ ,  $\tau_k^H \geq \epsilon_2$ , and  $\rho_k = \mathcal{O}(\epsilon_1^{-2})$ .

927 Next, if  $c(x) = Ax - b$  for some  $A \in \mathbb{R}^{m \times n}$  and  $b \in \mathbb{R}^m$ , then  $\nabla c(x) = A^T$  and  
928  $\nabla^2 c_i(x) = 0$  for  $1 \leq i \leq m$ . By these and (6.32), one has  $L_{k,H} = \mathcal{O}(1)$ . Using this  
929 and similar arguments as for (6.34), we obtain that  $T_k = \mathcal{O}(\max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\})$ . By  
930 this, (6.34) and  $\bar{K}_{\epsilon_1} = \mathcal{O}(|\log \epsilon_1|^2)$  (see Remark 4.9), we conclude that statement (i)  
931 of Theorem 4.10 holds.

932 (ii) In view of Lemma 4.4(i) and the definition of  $U_{k,H}$ , one can see that  $U_{k,H} \geq$   
933  $\sup_{x \in \mathbb{R}^n} \{\|\nabla_{xx}^2 \tilde{\mathcal{L}}(x, \lambda^k; \rho_k)\| : \tilde{\mathcal{L}}(x, \lambda^k; \rho_k) \leq \tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)\}$ . Using this,  $\tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k)$   
934  $\leq f_{\text{hi}}$  and Theorem 3.2(iv) with  $(F_{\text{hi}}, F_{\text{low}}, L_H^F, U_H^F, \epsilon_g, \epsilon_H) = (\tilde{\mathcal{L}}(x_{\text{init}}^k, \lambda^k; \rho_k), f_{\text{low}} -$   
935  $\gamma - \Lambda \delta_c, L_{k,H}, U_{k,H}, \tau_k^g, \tau_k^H)$ , we obtain that

936 (6.35) 
$$\begin{aligned} N_k &= \tilde{\mathcal{O}}((f_{\text{hi}} - f_{\text{low}} + \gamma + \Lambda \delta_c) L_{k,H}^2 \max\{(\tau_k^g)^{-2} \tau_k^H, (\tau_k^H)^{-3}\} \min\{n, (U_{k,H} / \tau_k^H)^{1/2}\}) \\ &\stackrel{(6.33)}{=} \tilde{\mathcal{O}}(\rho_k^2 \max\{(\tau_k^g)^{-2} \tau_k^H, (\tau_k^H)^{-3}\} \min\{n, (\rho_k / \tau_k^H)^{1/2}\}) \\ &= \tilde{\mathcal{O}}(\epsilon_1^{-4} \max\{\epsilon_1^{-2} \epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1} \epsilon_2^{-1/2}\}), \end{aligned}$$

938 where the last equality is from  $(\tau_k^g)^2/\tau_k^H \geq \min\{\epsilon_1^2/\epsilon_2, \epsilon_2^3\}$ ,  $\tau_k^H \geq \epsilon_2$ , and  $\rho_k = \mathcal{O}(\epsilon_1^{-2})$ .

939 On the other hand, if  $c$  is assumed to be affine, it follows from the above discussion  
940 that  $L_{k,H} = \mathcal{O}(1)$ . Using this,  $U_{k,H} \leq U_1 + \rho_k U_2$ , and similar arguments as for  
941 (6.35), we obtain that  $N_k = \tilde{\mathcal{O}}(\max\{\epsilon_1^{-2}\epsilon_2, \epsilon_2^{-3}\} \min\{n, \epsilon_1^{-1}\epsilon_2^{-1/2}\})$ . By this, (6.35)  
942 and  $\bar{K}_{\epsilon_1} = \mathcal{O}(|\log \epsilon_1|^2)$  (see Remark 4.9), we conclude that statement (ii) of Theorem  
943 4.10 holds.  $\square$

944 Next, we provide a proof of Theorem 4.14. To proceed, we first observe from  
945 Assumptions 4.1(c) and 4.12 that there exist  $U_g^f > 0$ ,  $U_g^c > 0$  and  $\sigma > 0$  such that

$$946 \quad (6.36) \quad \|\nabla f(x)\| \leq U_g^f, \|\nabla c(x)\| \leq U_g^c, \lambda_{\min}(\nabla c(x)^T \nabla c(x)) \geq \sigma^2, \forall x \in \mathcal{S}(\delta_f, \delta_c).$$

947 We next establish several technical lemmas that will be used shortly.

948 LEMMA 6.5. *Suppose that Assumptions 4.1 and 4.12 hold and that  $\rho_0$  is sufficiently  
949 large such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). Let  
950  $\{(x^k, \lambda^k, \rho_k)\}$  be generated by Algorithm 4.1. Suppose that*

951 (6.37)  $\rho_k \geq \max\{\Lambda^2(2\delta_f)^{-1}, 2(f_{\text{hi}} - f_{\text{low}} + \gamma)\delta_c^{-2} + 2\Lambda\delta_c^{-1} + 2\gamma, 2(U_g^f + U_g^c\Lambda + 1)(\sigma\epsilon_1)^{-1}\}$   
952 for some  $k \geq 0$ , where  $\gamma$ ,  $f_{\text{hi}}$ ,  $f_{\text{low}}$ ,  $\delta_f$  and  $\delta_c$  are given in Assumption 4.1, and  $U_g^f$ ,  
953  $U_g^c$  and  $\sigma$  are given in (6.36). Then it holds that  $\|c(x^{k+1})\| \leq \epsilon_1$ .

954 *Proof.* By (6.37) and  $\|\lambda^k\| \leq \Lambda$  (see step 6 of Algorithm 4.1), one can see that  $\rho_k \geq$   
955  $\max\{\|\lambda^k\|^2(2\delta_f)^{-1}, 2(f_{\text{hi}} - f_{\text{low}} + \gamma)\delta_c^{-2} + 2\|\lambda^k\|\delta_c^{-1} + 2\gamma\}$ . Using this, (4.1), the first re-  
956 lation in (4.6), and Lemma 6.4(iii) and (iv) with  $(x, \lambda, \rho, \tilde{\delta}_f, \tilde{\delta}_c) = (x^{k+1}, \lambda^k, \rho_k, \delta_f, \delta_c)$ ,  
957 we obtain that  $f(x^{k+1}) \leq f_{\text{hi}} + \delta_f$  and  $\|\tilde{c}(x^{k+1})\| \leq \delta_c$ . In addition, recall from  
958  $\|c(z_{\epsilon_1})\| \leq 1$  and the definition of  $\tilde{c}$  in (4.4) that  $\|c(x^{k+1})\| \leq 1 + \|\tilde{c}(x^{k+1})\|$ . These  
959 together with (4.3) show that  $x^{k+1} \in \mathcal{S}(\delta_f, \delta_c)$ . It then follows from (6.36) that  
960  $\|\nabla f(x^{k+1})\| \leq U_g^f$ ,  $\|\nabla c(x^{k+1})\| \leq U_g^c$ , and  $\lambda_{\min}(\nabla c(x^{k+1})^T \nabla c(x^{k+1})) \geq \sigma^2$ . By  
961  $\|\nabla f(x^{k+1})\| \leq U_g^f$ ,  $\|\nabla c(x^{k+1})\| \leq U_g^c$ ,  $\tau_k^g \leq 1$ ,  $\|\lambda^k\| \leq \Lambda$ , (4.4) and (4.6), one has

$$962 \quad \rho_k \|\nabla c(x^{k+1}) \tilde{c}(x^{k+1})\| \leq \|\nabla f(x^{k+1}) + \nabla c(x^{k+1})\lambda^k\| + \|\nabla_x \tilde{\mathcal{L}}(x^{k+1}, \lambda^k; \rho_k)\| \\ 963 \quad (6.38) \quad \stackrel{(4.6)}{\leq} \|\nabla f(x^{k+1})\| + \|\nabla c(x^{k+1})\| \|\lambda^k\| + \tau_k^g \leq U_g^f + U_g^c\Lambda + 1.$$

965 In addition, note that  $\lambda_{\min}(\nabla c(x^{k+1})^T \nabla c(x^{k+1})) \geq \sigma^2$  implies that  $\nabla c(x^{k+1})^T \nabla c(x^{k+1})$   
966 is invertible. Using this fact and (6.38), we obtain

$$967 \quad \|\tilde{c}(x^{k+1})\| \leq \|(\nabla c(x^{k+1})^T \nabla c(x^{k+1}))^{-1} \nabla c(x^{k+1})^T\| \|\nabla c(x^{k+1}) \tilde{c}(x^{k+1})\| \\ 968 \quad (6.39) \quad = \lambda_{\min}(\nabla c(x^{k+1})^T \nabla c(x^{k+1}))^{-\frac{1}{2}} \|\nabla c(x^{k+1}) \tilde{c}(x^{k+1})\| \stackrel{(6.38)}{\leq} \frac{U_g^f + U_g^c\Lambda + 1}{\sigma\rho_k}.$$

970 We also observe from (6.37) that  $\rho_k \geq 2(U_g^f + U_g^c\Lambda + 1)(\sigma\epsilon_1)^{-1}$ , which along with  
971 (6.39) proves  $\|\tilde{c}(x^{k+1})\| \leq \epsilon_1/2$ . Combining this with the definition of  $\tilde{c}$  in (4.4) and  
972  $\|c(z_{\epsilon_1})\| \leq \epsilon_1/2$ , we conclude that  $\|c(x^{k+1})\| \leq \epsilon_1$  holds as desired.  $\square$

973 The next lemma provides a stronger upper bound for  $\{\rho_k\}$  than the one in Theorem  
974 4.8.

975 LEMMA 6.6. *Suppose that Assumptions 4.1 and 4.12 hold and that  $\rho_0$  is sufficiently  
976 large such that  $\delta_{f,1} \leq \delta_f$  and  $\delta_{c,1} \leq \delta_c$ , where  $\delta_{f,1}$  and  $\delta_{c,1}$  are defined in (4.10). Let  
977  $\{\rho_k\}$  be generated by Algorithm 4.1 and*

978 (6.40)  $\tilde{\rho}_{\epsilon_1} := \max\{\Lambda^2(2\delta_f)^{-1}, 2(f_{\text{hi}} - f_{\text{low}} + \gamma)\delta_c^{-2} + 2\Lambda\delta_c^{-1} + 2\gamma, 2(U_g^f + U_g^c\Lambda + 1)(\sigma\epsilon_1)^{-1}, 2\rho_0\}$ ,  
979 where  $\gamma$ ,  $f_{\text{hi}}$ ,  $f_{\text{low}}$ ,  $\delta_f$  and  $\delta_c$  are given in Assumption 4.1, and  $U_g^f$ ,  $U_g^c$  and  $\sigma$  are given  
980 in (6.36). Then  $\rho_k \leq r\tilde{\rho}_{\epsilon_1}$  holds for  $0 \leq k \leq \bar{K}_{\epsilon_1}$ , where  $\bar{K}_{\epsilon_1}$  is defined in (4.15).

981 *Proof.* It follows from (6.40) that  $\tilde{\rho}_{\epsilon_1} \geq 2\rho_0$ . By this and similar arguments as  
 982 for (6.21), one has  $K_{\epsilon_1} \leq \log(\tilde{\rho}_{\epsilon_1}\rho_0^{-1})/\log r + 1$ , where  $K_{\epsilon_1}$  is defined in (4.12). Using  
 983 this, the update scheme for  $\{\rho_k\}$ , and similar arguments as for (6.22), we obtain

984 (6.41) 
$$\max_{0 \leq k \leq K_{\epsilon_1}} \rho_k \leq r\tilde{\rho}_{\epsilon_1}.$$

985 If  $\|c(x^{K_{\epsilon_1}+1})\| \leq \epsilon_1$ , it follows from (4.15) that  $\bar{K}_{\epsilon_1} = K_{\epsilon_1}$ , which together with (6.41)  
 986 implies that  $\rho_k \leq r\tilde{\rho}_{\epsilon_1}$  holds for  $0 \leq k \leq \bar{K}_{\epsilon_1}$ . On the other hand, if  $\|c(x^{K_{\epsilon_1}+1})\| > \epsilon_1$ ,  
 987 it follows from (4.15) that  $\|c(x^{k+1})\| > \epsilon_1$  for  $K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1$ . This together with  
 988 Lemma 6.5 and (6.40) implies that for all  $K_{\epsilon_1} \leq k \leq \bar{K}_{\epsilon_1} - 1$ ,

989 
$$\rho_k < \max\{\Lambda^2(2\delta_f)^{-1}, 2(f_{\text{hi}} - f_{\text{low}} + \gamma)\delta_c^{-2} + 2\Lambda\delta_c^{-1} + 2\gamma, 2(U_g^f + U_g^c\Lambda + 1)(\sigma\epsilon_1)^{-1}\} \stackrel{(6.40)}{\leq} \tilde{\rho}_{\epsilon_1}.$$

990 By this, (6.41), and  $\rho_{\bar{K}_{\epsilon_1}} \leq r\rho_{\bar{K}_{\epsilon_1}-1}$ , we also see that  $\rho_k \leq r\tilde{\rho}_{\epsilon_1}$  holds for  $0 \leq k \leq \bar{K}_{\epsilon_1}$ .  $\square$

991 *Proof of Theorem 4.14.* Notice from (6.40) and Lemma 6.6 that  $\tilde{\rho}_{\epsilon_1} = \mathcal{O}(\epsilon_1^{-1})$  and  
 992  $\rho_k \leq r\tilde{\rho}_{\epsilon_1}$ , which yield  $\rho_k = \mathcal{O}(\epsilon_1^{-1})$ . The conclusion of Theorem 4.14 then follows  
 993 from this and the same arguments as for the proof of Theorem 4.10 with  $\rho_k = \mathcal{O}(\epsilon_1^{-2})$   
 994 replaced by  $\rho_k = \mathcal{O}(\epsilon_1^{-1})$ .  $\square$

995 **7. Future work.** There are several possible future studies on this work. First,  
 996 it would be interesting to extend our AL method to seek an approximate SOSOP  
 997 of nonconvex optimization with inequality or more general constraints. Indeed, for  
 998 nonconvex optimization with inequality constraints, one can reformulate it as an  
 999 equality constrained problem using squared slack variables (e.g., see [7]). It can be  
 1000 shown that an SOSOP of the latter problem induces a weak SOSOP of the original problem  
 1001 and also linear independence constraint qualification holds for the latter problem if  
 1002 it holds for the original problem. As a result, it is promising to find an approximate  
 1003 weak SOSOP of an inequality constrained problem by applying our AL method to the  
 1004 equivalent equality constrained problem. Second, it is worth studying whether the  
 1005 enhanced complexity results in Section 4.3 can be derived under weaker constraint  
 1006 qualification (e.g., see [5]). Third, the development of our AL method is based on a  
 1007 strong assumption that a nearly feasible solution of the problem is known. It would  
 1008 make the method applicable to a broader class of problems if such an assumption  
 1009 could be removed by modifying the method possibly through the use of infeasibility  
 1010 detection techniques (e.g., see [19]). Lastly, more numerical studies would be helpful  
 1011 to further improve our AL method from a practical perspective.

1013 [1] N. AGARWAL, Z. ALLEN-ZHU, B. BULLINS, E. HAZAN, AND T. MA, *Finding approximate local*  
 1014 *minima faster than gradient descent*, in Proceedings of the 49th Annual ACM SIGACT  
 1015 *Symposium on Theory of Computing*, 2017, pp. 1195–1199.

1016 [2] R. ANDREANI, E. G. BIRGIN, J. M. MARTÍNEZ, AND M. L. SCHUVERDT, *On augmented Lagrangian*  
 1017 *methods with general lower-level constraints*, SIAM J. Optim., 18 (2008), pp. 1286–1309.

1018 [3] R. ANDREANI, G. HAESER, AND J. M. MARTÍNEZ, *On sequential optimality conditions for smooth*  
 1019 *constrained optimization*, Optim., 60 (2011), pp. 627–641.

1020 [4] R. ANDREANI, G. HAESER, A. RAMOS, AND P. J. SILVA, *A second-order sequential optimality*  
 1021 *condition associated to the convergence of optimization algorithms*, IMA J. Numer. Anal.,  
 1022 37 (2017), pp. 1902–1929.

1023 [5] R. ANDREANI, G. HAESER, M. L. SCHUVERDT, AND P. J. SILVA, *Two new weak constraint*  
 1024 *qualifications and applications*, SIAM J. Optim., 22 (2012), pp. 1109–1135.

1025 [6] P. ARMAND AND N. N. TRAN, *An augmented Lagrangian method for equality constrained*  
 1026 *optimization with rapid infeasibility detection capabilities*, J. Optim. Theory Appl., 181  
 1027 (2019), pp. 197–215.

1028 [7] D. P. BERTSEKAS, *Nonlinear Programming*, Athena Scientific, 1999.

1029 [8] W. BIAN, X. CHEN, AND Y. YE, *Complexity analysis of interior point algorithms for non-*  
1030 *Lipschitz and nonconvex minimization*, Math. Program., 149 (2015), pp. 301–327.

1031 [9] E. G. BIRGIN, J. GARDENGHI, J. M. MARTÍNEZ, S. A. SANTOS, AND P. L. TOINT, *Evaluation*  
1032 *complexity for nonlinear constrained optimization using unscaled KKT conditions and*  
1033 *high-order models*, SIAM J. Optim., 26 (2016), pp. 951–967.

1034 [10] E. G. BIRGIN, G. HAESER, AND A. RAMOS, *Augmented Lagrangians with constrained subproblems*  
1035 *and convergence to second-order stationary points*, Comput. Optim. Appl., 69 (2018), pp. 51–  
1036 75.

1037 [11] E. G. BIRGIN AND J. M. MARTÍNEZ, *Practical Augmented Lagrangian Methods for Constrained*  
1038 *Optimization*, SIAM, 2014.

1039 [12] E. G. BIRGIN AND J. M. MARTÍNEZ, *The use of quadratic regularization with a cubic descent*  
1040 *condition for unconstrained optimization*, SIAM J. Optim., 27 (2017), pp. 1049–1074.

1041 [13] E. G. BIRGIN AND J. M. MARTÍNEZ, *Complexity and performance of an augmented Lagrangian*  
1042 *algorithm*, Optim. Methods and Softw., 35 (2020), pp. 885–920.

1043 [14] J. F. BONNANS AND G. LAUNAY, *Sequential quadratic programming with penalization of the*  
1044 *displacement*, SIAM J. Optim., 5 (1995), pp. 792–812.

1045 [15] N. BOUMAL, V. VORONINSKI, AND A. S. BANDEIRA, *The non-convex Burer-Monteiro approach*  
1046 *works on smooth semidefinite programs*, in Advances in Neural Information Processing  
1047 Systems, vol. 29, 2016, pp. 2757–2765.

1048 [16] L. F. BUENO AND J. M. MARTÍNEZ, *On the complexity of an inexact restoration method for*  
1049 *constrained optimization*, SIAM J. Optim., 30 (2020), pp. 80–101.

1050 [17] S. BURER AND R. D. C. MONTEIRO, *A nonlinear programming algorithm for solving semidefinite*  
1051 *programs via low-rank factorization*, Math. Program., 95 (2003), pp. 329–357.

1052 [18] S. BURER AND R. D. C. MONTEIRO, *Local minima and convergence in low-rank semidefinite*  
1053 *programming*, Math. Program., 103 (2005), pp. 427–444.

1054 [19] J. V. BURKE, F. E. CURTIS, AND H. WANG, *A sequential quadratic optimization algorithm with*  
1055 *rapid infeasibility detection*, SIAM J. Optim., 24 (2014), pp. 839–872.

1056 [20] R. H. BYRD, F. E. CURTIS, AND J. NOCEDAL, *Infeasibility detection and SQP methods for*  
1057 *nonlinear optimization*, SIAM J. Optim., 20 (2010), pp. 2281–2299.

1058 [21] R. H. BYRD, R. B. SCHNABEL, AND G. A. SHULTZ, *A trust region algorithm for nonlinearly*  
1059 *constrained optimization*, SIAM J. Numer. Anal., 24 (1987), pp. 1152–1170.

1060 [22] Y. CARMON AND J. C. DUCHI, *Gradient descent finds the cubic-regularized nonconvex Newton*  
1061 *step*, SIAM J. Optim., 29 (2019), pp. 2146–2178.

1062 [23] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, “*Convex until proven guilty*”:  
1063 *Dimension-free acceleration of gradient descent on non-convex functions*, in International  
1064 Conference on Machine Learning, PMLR, 2017, pp. 654–663.

1065 [24] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, *Accelerated methods for nonconvex*  
1066 *optimization*, SIAM J. Optim., 28 (2018), pp. 1751–1772.

1067 [25] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Adaptive cubic regularisation methods for uncon-*  
1068 *strained optimization. Part II: worst-case function-and derivative-evaluation complexity*,  
1069 Math. Program., 130 (2011), pp. 295–319.

1070 [26] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *On the evaluation complexity of cubic regularization*  
1071 *methods for potentially rank-deficient nonlinear least-squares problems and its relevance to*  
1072 *constrained nonlinear optimization*, SIAM J. Optim., 23 (2013), pp. 1553–1574.

1073 [27] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *On the complexity of finding first-order critical*  
1074 *points in constrained nonlinear optimization*, Math. Program., 144 (2014), pp. 93–106.

1075 [28] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *On the evaluation complexity of constrained*  
1076 *nonlinear least-squares and general constrained nonlinear optimization using second-order*  
1077 *methods*, SIAM J. Numer. Anal., 53 (2015), pp. 836–851.

1078 [29] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Evaluation complexity bounds for smooth constrained*  
1079 *nonlinear optimization using scaled KKT conditions, high-order models and the criticality*  
1080 *measure  $\chi$* , in Approximation and Optimization, Springer, 2019, pp. 5–26.

1081 [30] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Optimality of orders one to three and beyond: char-*  
1082 *acterization and evaluation complexity in constrained nonconvex optimization*, J. Complex.,  
1083 53 (2019), pp. 68–94.

1084 [31] X. CHEN, L. GUO, Z. LU, AND J. J. YE, *An augmented Lagrangian method for non-Lipschitz*  
1085 *nonconvex programming*, SIAM J. Numer. Anal., 55 (2017), pp. 168–193.

1086 [32] D. CIFUENTES AND A. MOITRA, *Polynomial time guarantees for the Burer-Monteiro method*,  
1087 arXiv preprint arXiv:1912.01745, (2019).

1088 [33] T. F. COLEMAN, J. LIU, AND W. YUAN, *A new trust-region algorithm for equality constrained*  
1089 *optimization*, Comput. Optim. Appl., 21 (2002), pp. 177–199.

1090 [34] F. E. CURTIS, D. P. ROBINSON, C. W. ROYER, AND S. J. WRIGHT, *Trust-region Newton-CG*  
 1091 *with strong second-order complexity guarantees for nonconvex optimization*, SIAM J Optim.,  
 1092 31 (2021), pp. 518–544.

1093 [35] F. E. CURTIS, D. P. ROBINSON, AND M. SAMADI, *A trust region algorithm with a worst-case*  
 1094 *iteration complexity of  $\mathcal{O}(\epsilon^{-3/2})$  for nonconvex optimization*, Math. Program., 162 (2017),  
 1095 pp. 1–32.

1096 [36] F. E. CURTIS, D. P. ROBINSON, AND M. SAMADI, *Complexity analysis of a trust funnel algorithm*  
 1097 *for equality constrained optimization*, SIAM J. Optim., 28 (2018), pp. 1533–1563.

1098 [37] G. N. GRAPIGLIA AND Y. YUAN, *On the complexity of an augmented Lagrangian method for*  
 1099 *nonconvex optimization*, IMA J. Numer. Anal., 41 (2021), pp. 1508–1530.

1100 [38] G. HAESER, H. LIU, AND Y. YE, *Optimality condition and complexity analysis for linearly-*  
 1101 *constrained optimization without differentiability on the boundary*, Math. Program., (2019),  
 1102 pp. 1–37.

1103 [39] M. R. HESTENES, *Multiplier and gradient methods*, J. Optim. Theory Appl., 4 (1969), pp. 303–  
 1104 320.

1105 [40] M. HONG, D. HAJINEZHAD, AND M.-M. ZHAO, *Prox-PDA: The proximal primal-dual algorithm*  
 1106 *for fast distributed nonconvex optimization and learning over networks*, in International  
 1107 Conference on Machine Learning, PMLR, 2017, pp. 1529–1538.

1108 [41] C. JIN, R. GE, P. NETRAPALLI, S. M. KAKADE, AND M. I. JORDAN, *How to escape saddle points*  
 1109 *efficiently*, in International Conference on Machine Learning, PMLR, 2017, pp. 1724–1732.

1110 [42] C. KANZOW AND D. STECK, *An example comparing the standard and safeguarded augmented*  
 1111 *Lagrangian methods*, Oper. Res. Lett., 45 (2017), pp. 598–603.

1112 [43] W. KONG, J. G. MELO, AND R. D. C. MONTEIRO, *Complexity of a quadratic penalty accelerated*  
 1113 *inexact proximal point method for solving linearly constrained nonconvex composite*  
 1114 *programs*, SIAM J. Optim., 29 (2019), pp. 2566–2593.

1115 [44] J. KUCZYŃSKI AND H. WOŹNIAKOWSKI, *Estimating the largest eigenvalue by the power and*  
 1116 *Lanczos algorithms with a random start*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1094–  
 1117 1122.

1118 [45] Z. LI, P.-Y. CHEN, S. LIU, S. LU, AND Y. XU, *Rate-improved inexact augmented lagrangian*  
 1119 *method for constrained nonconvex optimization*, in International Conference on Artificial  
 1120 Intelligence and Statistics, PMLR, 2021, pp. 2170–2178.

1121 [46] S. LU, *A single-loop gradient descent and perturbed ascent algorithm for nonconvex functional*  
 1122 *constrained optimization*, in International Conference on Machine Learning, PMLR, 2022,  
 1123 pp. 14315–14357.

1124 [47] S. LU, M. RAZAVIYAYN, B. YANG, K. HUANG, AND M. HONG, *Finding second-order stationary*  
 1125 *points efficiently in smooth nonconvex linearly constrained optimization problems*, Advances  
 1126 in Neural Information Processing Systems, 33 (2020), pp. 2811–2822.

1127 [48] Z. LU AND X. LI, *Sparse recovery via partial regularization: models, theory, and algorithms*,  
 1128 Math. Oper. Res., 43 (2018), pp. 1290–1316.

1129 [49] Z. LU AND Y. ZHANG, *An augmented Lagrangian approach for sparse principal component*  
 1130 *analysis*, Math. Program., 135 (2012), pp. 149–193.

1131 [50] J. M. MARTÍNEZ AND M. RAYDAN, *Cubic-regularization counterpart of a variable-norm trust-*  
 1132 *region method for unconstrained minimization*, J. Glob. Optim., 68 (2017), pp. 367–385.

1133 [51] J. G. MELO, R. D. MONTEIRO, AND W. KONG, *Iteration-complexity of an inner accelerated*  
 1134 *inexact proximal augmented Lagrangian method based on the classical Lagrangian function*  
 1135 *and a full Lagrange multiplier update*, arXiv preprint arXiv:2008.00562, (2020).

1136 [52] Y. NESTEROV AND B. T. POLYAK, *Cubic regularization of Newton method and its global*  
 1137 *performance*, Math. Program., 108 (2006), pp. 177–205.

1138 [53] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, 2nd ed., 2006.

1139 [54] M. O’NEILL AND S. J. WRIGHT, *A log-barrier Newton-CG method for bound constrained*  
 1140 *optimization with complexity guarantees*, IMA J. Numer. Anal., 41 (2021), pp. 84–121.

1141 [55] R. T. ROCKAFELLAR, *Lagrange multipliers and optimality*, SIAM review, 35 (1993), pp. 183–238.

1142 [56] C. W. ROYER, M. O’NEILL, AND S. J. WRIGHT, *A Newton-CG algorithm with complexity*  
 1143 *guarantees for smooth unconstrained optimization*, Math. Program., 180 (2020), pp. 451–  
 1144 488.

1145 [57] C. W. ROYER AND S. J. WRIGHT, *Complexity analysis of second-order line-search algorithms*  
 1146 *for smooth nonconvex optimization*, SIAM J. Optim., 28 (2018), pp. 1448–1477.

1147 [58] M. F. SAHIN, A. EFTEKHARI, A. ALACAOGLU, F. LATORRE, AND V. CEVHER, *An inexact*  
 1148 *augmented Lagrangian framework for nonconvex optimization with nonlinear constraints*,  
 1149 *Advances in Neural Information Processing Systems*, 32 (2019).

1150 [59] Y. XIE AND S. J. WRIGHT, *Complexity of projected Newton methods for bound-constrained*  
 1151 *optimization*, arXiv preprint arXiv:2103.15989, (2021).

1152 [60] Y. XIE AND S. J. WRIGHT, *Complexity of proximal augmented Lagrangian for nonconvex*  
 1153 *optimization with nonlinear equality constraints*, J. Sci. Comput., 86 (2021), pp. 1–30.  
 1154 [61] L. YANG, D. SUN, AND K. C. TOH, *SDPNAL+*: *A majorized semismooth Newton-CG augmented*  
 1155 *Lagrangian method for semidefinite programming with nonnegative constraints*, Math.  
 1156 Program. Comput., 7 (2015), pp. 331–366.  
 1157 [62] X. ZHAO, D. SUN, AND K. C. TOH, *A Newton-CG augmented Lagrangian method for semidefinite*  
 1158 *programming*, SIAM J. Optim., 20 (2010), pp. 1737–1765.

1159 **Appendix A. A capped conjugate gradient method.** In this part we  
 1160 present the capped CG method proposed in [56, Algorithm 1] for finding either an  
 1161 approximate solution to the linear system (3.6) or a sufficiently negative curvature  
 1162 direction of the associated matrix  $H$ , which has been briefly discussed in Section 3.1.  
 Its details can be found in [56, Section 3.1].

---

**Algorithm A.1** A capped conjugate gradient method

---

*Inputs:* symmetric matrix  $H \in \mathbb{R}^{n \times n}$ , vector  $g \neq 0$ , damping parameter  $\varepsilon \in (0, 1)$ , desired relative accuracy  $\zeta \in (0, 1)$ .

*Optional input:* scalar  $U \geq 0$  (set to 0 if not provided).

*Outputs:* d\_type,  $d$ .

*Secondary outputs:* final values of  $U, \kappa, \hat{\zeta}, \tau$ , and  $T$ .

Set

$$\bar{H} := H + 2\varepsilon I, \quad \kappa := \frac{U+2\varepsilon}{\varepsilon}, \quad \hat{\zeta} := \frac{\zeta}{3\kappa}, \quad \tau := \frac{\sqrt{\kappa}}{\sqrt{\kappa+1}}, \quad T := \frac{4\kappa^4}{(1-\sqrt{\tau})^2},$$

$$y^0 \leftarrow 0, r^0 \leftarrow g, p^0 \leftarrow -g, j \leftarrow 0.$$

if  $(p^0)^T \bar{H} p^0 < \varepsilon \|p^0\|^2$  then

Set  $d \leftarrow p^0$  and terminate with d\_type = NC;

else if  $\|H p^0\| > U \|p^0\|$  then

Set  $U \leftarrow \|H p^0\| / \|p^0\|$  and update  $\kappa, \hat{\zeta}, \tau, T$  accordingly;

end if

while TRUE do

$\alpha_j \leftarrow (r^j)^T r^j / (p^j)^T \bar{H} p^j$ ; {Begin Standard CG Operations}

$y^{j+1} \leftarrow y^j + \alpha_j p^j$ ;

$r^{j+1} \leftarrow r^j + \alpha_j \bar{H} p^j$ ;

$\beta_{j+1} \leftarrow \|r^{j+1}\|^2 / \|r^j\|^2$ ;

$p^{j+1} \leftarrow -r^{j+1} + \beta_{j+1} p^j$ ; {End Standard CG Operations}

$j \leftarrow j + 1$ ;

if  $\|H p^j\| > U \|p^j\|$  then

Set  $U \leftarrow \|H p^j\| / \|p^j\|$  and update  $\kappa, \hat{\zeta}, \tau, T$  accordingly;

end if

if  $\|H y^j\| > U \|y^j\|$  then

Set  $U \leftarrow \|H y^j\| / \|y^j\|$  and update  $\kappa, \hat{\zeta}, \tau, T$  accordingly;

end if

if  $\|H r^j\| > U \|r^j\|$  then

Set  $U \leftarrow \|H r^j\| / \|r^j\|$  and update  $\kappa, \hat{\zeta}, \tau, T$  accordingly;

end if

if  $(y^j)^T \bar{H} y^j < \varepsilon \|y^j\|^2$  then

Set  $d \leftarrow y^j$  and terminate with d\_type = NC;

else if  $\|r^j\| \leq \hat{\zeta} \|r^0\|$  then

Set  $d \leftarrow y^j$  and terminate with d\_type = SOL;

else if  $(p^j)^T \bar{H} p^j < \varepsilon \|p^j\|^2$  then

Set  $d \leftarrow p^j$  and terminate with d\_type = NC;

else if  $\|r^j\| > \sqrt{T} \tau^{j/2} \|r^0\|$  then

Compute  $\alpha_j, y^{j+1}$  as in the main loop above;

Find  $i \in \{0, \dots, j-1\}$  such that

$$(y^{j+1} - y^i)^T \bar{H} (y^{j+1} - y^i) < \varepsilon \|y^{j+1} - y^i\|^2;$$

Set  $d \leftarrow y^{j+1} - y^i$  and terminate with d\_type = NC;

end if

end while

---

1163 The following theorem presents the iteration complexity of Algorithm A.1.

1164 **THEOREM A.1 (iteration complexity of Algorithm A.1).** *Consider applying*  
 1165 *Algorithm A.1 with input  $U = 0$  to the linear system (3.6) with  $g \neq 0$ ,  $\varepsilon > 0$ , and  $H$*   
 1166 *being an  $n \times n$  symmetric matrix. Then the number of iterations of Algorithm A.1 is*  
 1167  $\tilde{\mathcal{O}}(\min\{n, \sqrt{\|H\|/\varepsilon}\})$ .

1168 *Proof.* From [56, Lemma 1], we know that the number of iterations of Algorithm  
 1169 A.1 is bounded by  $\min\{n, J(U, \varepsilon, \zeta)\}$ , where  $J(U, \varepsilon, \zeta)$  is the smallest integer  $J$  such that  
 1170  $\sqrt{T}\tau^{J/2} \leq \hat{\zeta}$ , with  $U, \hat{\zeta}, T$  and  $\tau$  being the values returned by Algorithm A.1. In addition,  
 1171 it was shown in [56, Section 3.1] that  $J(U, \varepsilon, \zeta) \leq \left\lceil \left( \sqrt{\kappa} + \frac{1}{2} \right) \ln \left( \frac{144(\sqrt{\kappa}+1)^2 \kappa^6}{\zeta^2} \right) \right\rceil$ ,  
 1172 where  $\kappa = \mathcal{O}(U/\varepsilon)$  is an output by Algorithm A.1. Then one can see that  $J(U, \varepsilon, \zeta) =$   
 1173  $\tilde{\mathcal{O}}(\sqrt{U/\varepsilon})$ . Notice from Algorithm A.1 that the output  $U \leq \|H\|$ . Combining these,  
 1174 we obtain the conclusion as desired.  $\square$

1175 **Appendix B. A randomized Lanczos based minimum eigenvalue oracle.**

1176 In this part we present the randomized Lanczos method proposed in [56, Section 3.2],  
 1177 which can be used as a minimum eigenvalue oracle for Algorithm 3.1. As briefly  
 1178 discussed in Section 3.1, this oracle outputs either a sufficiently negative curvature  
 1179 direction of  $H$  or a certificate that  $H$  is nearly positive semidefinite with high probability.  
 1180 More detailed motivation and explanation of it can be found in [56, Section 3.2].

---

**Algorithm B.1** A randomized Lanczos based minimum eigenvalue oracle

*Input:* symmetric matrix  $H \in \mathbb{R}^{n \times n}$ , tolerance  $\varepsilon > 0$ , and probability parameter  $\delta \in (0, 1)$ .  
*Output:* a sufficiently negative curvature direction  $v$  satisfying  $v^T Hv \leq -\varepsilon/2$  and  $\|v\| = 1$ ; or  
 a certificate that  $\lambda_{\min}(H) \geq -\varepsilon$  with probability at least  $1 - \delta$ .

Apply the Lanczos method [44] to estimate  $\lambda_{\min}(H)$  starting with a random vector uniformly  
 generated on the unit sphere, and run it for at most

$$(B.1) \quad N(\varepsilon, \delta) := \min \left\{ n, 1 + \left\lceil \frac{\ln(2.75n/\delta^2)}{2} \sqrt{\frac{\|H\|}{\varepsilon}} \right\rceil \right\}$$

iterations. If a unit vector  $v$  with  $v^T Hv \leq -\varepsilon/2$  is found at some iteration, terminate  
 immediately and return  $v$ .

---

1181 The following theorem justifies that Algorithm B.1 is a suitable minimum eigenvalue  
 1182 oracle for Algorithm 3.1. Its proof is identical to that of [56, Lemma 2] and thus  
 1183 omitted.

1184 **THEOREM B.1 (iteration complexity of Algorithm B.1).** *Consider Algo-  
 1185 rithm B.1 with tolerance  $\varepsilon > 0$ , probability parameter  $\delta \in (0, 1)$ , and symmetric matrix*  
 1186  *$H \in \mathbb{R}^{n \times n}$  as its input. Then it either finds a sufficiently negative curvature direction*  
 1187  *$v$  satisfying  $v^T Hv \leq -\varepsilon/2$  and  $\|v\| = 1$  or certifies that  $\lambda_{\min}(H) \geq -\varepsilon$  holds with*  
 1188 *probability at least  $1 - \delta$  in at most  $N(\varepsilon, \delta)$  iterations, where  $N(\varepsilon, \delta)$  is defined in*  
 1189 *(B.1).*

1190 Notice that  $\|H\|$  is required in Algorithm B.1. In general, computing  $\|H\|$  may  
 1191 not be cheap when  $n$  is large. Nevertheless,  $\|H\|$  can be efficiently estimated via a  
 1192 randomization scheme with high confidence (e.g., see the discussion in [56, Appen-  
 1193 dix B3]).