

1 **COMPLEXITY OF FINITE-SUM OPTIMIZATION WITH**
2 **NONSMOOTH COMPOSITE FUNCTIONS AND NON-LIPSCHITZ**
3 **REGULARIZATION***

4 XIAO WANG[†] AND XIAOJUN CHEN[‡]

5 **Abstract.** In this paper we present complexity analysis of proximal inexact gradient methods
6 for finite-sum optimization with a nonconvex nonsmooth composite function and non-Lipschitz reg-
7 ularization. By getting access to a convex approximation to the Lipschitz function and a Lipschitz
8 continuous approximation to the non-Lipschitz regularizer, we construct a proximal subproblem at
9 each iteration without using exact function values and gradients. With certain accuracy control on
10 inexact gradients and subproblem solutions, we show that the oracle complexity in terms of total
11 number of inexact gradient evaluations is in order $\mathcal{O}(\epsilon^{-2})$ to find an (ϵ, δ) -approximate first-order
12 stationary point, ensuring that within a δ -ball centered at this point the maximum reduction of an
13 approximation model does not exceed $\epsilon\delta$. This shows that we can have the same worst-case evalua-
14 tion complexity order as [5, 12] even if we introduce the non-Lipschitz singularity and the nonconvex
15 nonsmooth composite function in the objective function. Moreover, we establish that the oracle
16 complexity regarding the total number of stochastic oracles is in order $\tilde{\mathcal{O}}(\epsilon^{-2})$ with high probability
17 for stochastic proximal inexact gradient methods. We further extend the algorithm to adjust to
18 solving stochastic problems with expectation form and derive the associated oracle complexity in
19 order $\tilde{\mathcal{O}}(\epsilon^{-10/3})$ with high probability.

20 **Key words.** nonconvexity, nonsmoothness, non-Lipschitz regularization, inexact oracle, com-
21 plexity

22 **MSC codes.** 90C30, 90C46, 65K05

23 **1. Introduction.** In this paper, we consider the following nonconvex nonsmooth
24 optimization problem:

25 (1.1)
$$\min_{x \in \mathcal{F}} Q(x) := f(x) + h(c(x)) + \|Vx\|_p^p,$$

26 where $\mathcal{F} \subseteq \mathbb{R}^n$ is nonempty, bounded, closed and convex, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $c : \mathbb{R}^n \rightarrow \mathbb{R}^r$
27 are continuously differentiable with Lipschitz continuous gradients over \mathcal{F} , $h : \mathbb{R}^r \rightarrow \mathbb{R}$
28 is Lipschitz continuous and convex but possibly nonsmooth, $V \in \mathbb{R}^{\bar{n} \times n}$ with $\bar{n} \leq n$ and
29 $p \in (0, 1)$. We assume that rows of V , denoted by v_i^T , $i = 1, \dots, \bar{n}$, are orthonormal,
30 without loss of generality. Problem (1.1) has numerous applications in data science,
31 where f is a loss function, h is a penalty function and $\|\cdot\|_p^p$ is a sparse regularization.
32 For instance, with the increasing interest of group sparsity regularization for neural
33 networks (see e.g. [4, 20, 25, 31]), the loss function f may rely on a large data set
34 and be defined in the form

35 (1.2)
$$f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x),$$

36 where $f_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 1, \dots, N$, are continuously differentiable and the sample
37 size N can be very large such that it may be time-consuming and sometimes even
38 prohibitive to access all component functions to compute the exact gradient of f at

*Submitted to the editors DATE.

Funding: This work was funded by by the Major Key Project of PCL (No. PCL2022A05), the Chinese NSF Grant 12271278 and Hong Kong Research Grants Council, grant PolyU15300219.

[†]Department of AI Computing, Peng Cheng Laboratory, Shenzhen, China (wangx07@pcl.ac.cn).

[‡]Department of Applied Mathematics, Hong Kong Polytechnic University, Kowloon, Hong Kong (maxjchen@polyu.edu.hk).

39 a query point. Moreover, constraints are often imposed to enforce specific conditions
 40 on variables. For example, constraints of the form $c(x) \leq 0$, where $c : \mathbb{R}^n \rightarrow \mathbb{R}^r$,
 41 are prevalent in a wide range of applications, including image restoration [1], film
 42 restoration [18] and SVM [21]. However, ensuring feasibility of iterates with respect
 43 to these constraints throughout the algorithmic process can be challenging. To tackle
 44 this issue, infeasible methods are commonly employed, which allow for violations of
 45 the constraints. Specifically, with the aid of a penalty function, for instance, ℓ_1 penalty
 46 function, one can remove the constraints by introducing a nonsmooth penalty term
 47 in the objective, e.g. $h(c(x)) = \rho \|(c(x))_+\|_1$ with $(c(x))_+ = \max(c(x), 0)$ and ρ being
 48 a penalty parameter. As studied in [10, 15], the resulting problem can be an exact
 49 penalty formulation of the original one to some extent, with nice properties regarding
 50 their minimizers. As is well studied in the literature, the nonconvex, nonsmooth and
 51 non-Lipschitz ℓ_p ($0 < p < 1$) regularizer has shown a good performance for sparse
 52 variable selection. However, in general the non-Lipschitz regularized problems are
 53 strongly NP-hard [9]. Challenges often arise in algorithm design and analysis. The
 54 past decade has witnessed highly productive progress on the study of ℓ_p ($0 < p < 1$)
 55 optimization and a surge of works has been proposed, to name a few but not limited
 56 to [2, 3, 13, 14, 16, 19, 21, 30].

57 Cartis *et al.* study the evaluation complexities of minimizing $f(x) + h(c(x))$ to
 58 reach the first-order critical measure within ϵ in [5] and to reach high-order approxi-
 59 mate minimizers in [7]. In recent work [8], they consider minimizing $f(x) + h(c(x))$ over
 60 a convex set and apply high-order approximation model to reach high-order approxi-
 61 mate minimizers. Gratton *et al.* in [17] propose an adaptive regularization algorithm
 62 using inexact function and gradient evaluations for minimizing $f(x) + h(c(x))$ and show
 63 that their algorithm needs at most $O(|\log(\epsilon)|\epsilon^{-2})$ evaluations of the functions and
 64 their derivatives for finding an ϵ -approximate first-order stationary point. In [11, 12],
 65 high-order algorithms for solving minimization problems with non-Lipschitzian group
 66 sparsity terms are studied, where the objective is the sum of a smooth function and
 67 a non-Lipschitz regularizer. Compared with problems studied in [5, 8, 11, 12, 17],
 68 the objective function in (1.1) has not only the nonsmooth function $h(c(x))$, but also
 69 the non-Lipschitz regularizer $\|Vx\|_p^p$, whose complexity has not been established in
 70 the literature to the best of our knowledge. We also notice that those algorithms
 71 studied in previous works [5, 8, 11, 12] rely on exact function values and gradients
 72 of f , which, however, are expensive to obtain in many scenarios with a large finite-
 73 sum structure. Inspired by above points, in this paper we will focus on complexity
 74 analysis for problem (1.1) to reach approximate first-order stationary point. We will
 75 investigate whether the absence of differentiability of the Lipschitz term together with
 76 the existence of non-Lipschitz regularization will affect the worst-case complexity,
 77 compared with existing works.

78 Problems with f in finite-sum structure (1.2) face challenges when computing
 79 exact function information, due to the large number of component functions. To alle-
 80 viate possible difficulties, stochastic oracles are normally called to approximate exact
 81 information. In the past decade, along with the development of data science, studies
 82 on stochastic approximation methods for nonlinear optimization grow rapidly in pop-
 83 ularity, ranging from convex to nonconvex problems and from smooth to nonsmooth
 84 problems. Xu *et al.* [29] study a class of optimization problems with nonconvex, non-
 85 smooth regularizer, namely minimizing $g(x) - h(x) + \Lambda(x)$, where g and h are both
 86 convex and Λ is a nonconvex and nonsmooth regularizer. Moreover, it requires in the-
 87 oretical analysis that g be smooth and h be Hölder smooth. The proposed algorithm
 88 in [29] can be also applied to unconstrained ℓ_p ($0 < p < 1$) regularized optimization

with g in the finite-sum form. The associated gradient complexity to find a nearly ϵ -critical point is in order $\mathcal{O}(\epsilon^{-4})$. Metel and Takeda [22] consider unconstrained optimization with a nonconvex but Lipschitz continuous regularizer. The proposed algorithm owns $\mathcal{O}(\epsilon^{-3})$ gradient-call complexity for finite-sum minimization when a variance reduction strategy is applied. However, the Lipschitz continuity assumption fails for ℓ_p ($0 < p < 1$) regularizer. Cheng *et al.* [14] propose an interior stochastic gradient method for nonnegative constrained optimization with ℓ_p regularizer and investigate the oracle complexities to find an approximate stationary point. Xu *et al.* [28] propose stochastic proximal gradient methods for minimizing summation of a smooth function f and a nonsmooth nonconvex regularizer and show that the $\mathcal{O}(\epsilon^{-2})$ gradient complexity can be achieved to find an ϵ -stationary point. The proposed algorithm in [28] requires the proximal mapping of the nonconvex regularizer be easy to obtain. However, these existing results cannot be applied to problem (1.1) due to the nonsmoothness and nonconvexity of $f + h$ or the non-Lipschitz continuity.

Contribution. The main contribution of this paper lies in the complexity analysis of proximal inexact gradient methods for finite-sum optimization with a nonconvex nonsmooth composite function and non-Lipschitz regularization (1.1). By getting access to a convex approximation to the Lipschitz function in the objective, together with a Lipschitz continuous approximation to the non-Lipschitz regularizer, we build a proximal subproblem at each iteration without using exact function values and gradients of f . Under certain conditions on inexact gradients and inexact subproblem solutions, we prove that the oracle complexity in terms of the total number of inexact gradient evaluations to find an approximate (ϵ, δ) -approximate first-order stationary point is in order $\mathcal{O}(\epsilon^{-2})$. This verifies that adding the nonsmooth nonconvex composite function and non-Lipschitz regularizer and using inexact gradients do not affect the worst-case oracle complexity, compared with existing results [5, 8, 11, 12]. Furthermore, we use the finite-sum structure of f and propose a stochastic variant of the algorithm through calls to stochastic first-order oracles. We show that the corresponding oracle complexity in terms of total number of stochastic first-order oracles is in order $\tilde{\mathcal{O}}(\epsilon^{-2})$ with high probability, where we use $\tilde{\mathcal{O}}$ to hide the dependence on logarithmic factor in the complexity order. Furthermore, we extend the proposed algorithm to solve stochastic problems with f in expectation form and obtain the $\tilde{\mathcal{O}}(\epsilon^{-10/3})$ -oracle complexity with high probability. We present more details on the significant differences from existing works.

(i) The related convergence and iteration complexity in [5, 8, 11, 12] are established within trust region schemes, which require accurate function values and derivatives of function f . However, those analysis cannot be applied to stochastic optimization problems, where only approximate or stochastic gradients are available. In this scenario, the behavior of the objective function can only be characterized based on inexact derivatives of f . Hence, it is imperative to modify the primary algorithmic framework to accommodate the reliance on approximate or stochastic gradients and the absence of a trust region scheme. This adaptation necessitates rigorous analysis under these altered conditions.

(ii) While our method draws inspiration from existing techniques, such as the convex approximation to the composite part and the Lipschitz continuous approximation to the non-Lipschitz regularizer, the coexistence of these two aspects brings significant challenges for the theoretical analysis, which makes it different from existing works. For instance, one particular challenge is about ensuring the existence of an inexact subproblem solution s_k that satisfies the required conditions, as presented in Lemma 2.5. Such detailed analysis, however, is not provided in [11, 12]. Another challenge

139 arises when considering the approximate criticality of the output of Algorithm 2.1.
 140 Due to the significant modifications made to adapt to the stochastic setting and the
 141 absence of a trust region scheme, the algorithm framework’s analysis differs substan-
 142 tially from existing methods. In addition to addressing these challenges, we present a
 143 unified framework by incorporating various elements and leveraging the strengths of
 144 each element. This enables our algorithm to tackle a broader range of problems.

145 (iii) When adapting the deterministic proximal inexact gradient method to stochastic
 146 settings, including the finite-sum setting and expectation setting, it causes nontrivial
 147 challenges to the theoretical analysis. In our paper, we go beyond a simple replace-
 148 ment of the deterministic gradient with a stochastic gradient, recognizing the need
 149 for careful consideration of oracle complexity analysis in the stochastic counterpart,
 150 which contributes to the value of our work. Particularly, the extension of the analysis
 151 in [25] to non-Lipschitz regularized optimization and to the expectation case proves
 152 to be a nontrivial task. Our oracle complexity analysis heavily relies on the essential
 153 property of the proposed algorithm, specifically the boundedness of $\sum_{k \in \mathcal{K}} \|s_k\|^2$ as
 154 demonstrated in Theorem 3.8.

155 **Organization.** This paper is organized as follows. In Section 2 we present a
 156 detailed algorithmic framework for proximal inexact gradient methods for (1.1). In
 157 Section 3 we explore the oracle complexity of the proposed framework to find an
 158 (ϵ, δ) -approximate first-order stationary point. In Section 4 we propose a stochastic
 159 variant of the algorithm for problems with f in finite-sum structure (1.2) and establish
 160 the oracle complexity accordingly. In Section 5 we propose an extended stochastic
 161 variant for problems in expectation case and investigate the related oracle complexity.
 162 In Section 6 we illustrate our algorithm by a numerical example. Finally, concluding
 163 remarks are drawn in Section 7.

164 **2. Algorithm description.** In this section, we will present an algorithmic
 165 framework for proximal inexact gradient methods for solving (1.1). As the objective
 166 function Q is nonconvex, nonsmooth and non-Lipschitz, it is generally intractable to
 167 approximately find a global or even a local minimizer. Thus our algorithm aims for an
 168 approximate first-order stationary point of (1.1). The core of our algorithm design is
 169 to construct a Lipschitz continuous approximation model of the objective function at
 170 each iteration. We then perform a search within a local neighborhood of the current
 171 iterate while aiming to minimize the approximation model as much as possible. The
 172 use of Lipschitz continuous approximation models helps us predict the behavior of
 173 the objective function while minimizing the impact of approximation errors in the
 174 optimization process. The proposed algorithm differs from existing works on com-
 175 plexity analysis, such as [5, 11, 12], where a trust region scheme is typically employed,
 176 requiring exact evaluations of the function value and its derivatives. In contrast, it
 177 only relies on getting access to inexact first-order derivatives of the objective function,
 178 which enables us to extend its applicability to stochastic variants. By utilizing these
 179 inexact derivatives, we can effectively navigate the search space and make progress to-
 180 wards the optimal solution without the need for precise function value and derivative
 181 evaluations. By adopting this approach, we strike a balance between computational
 182 efficiency and accuracy, making our algorithm more suitable for scenarios where exact
 183 evaluations may be costly or impractical.

184 We first define the following index sets at a point x for a given nonnegative
 185 constant ϵ :

$$186 \quad \mathcal{A}(x, \epsilon) = \{i \in [\bar{n}] : |v_i^T x| > \epsilon\}, \quad \mathcal{R}(x, \epsilon) = \bigcap_{i \in [\bar{n}] \setminus \mathcal{A}(x, \epsilon)} \ker(v_i^T),$$

where $[\bar{n}] := \{1, \dots, \bar{n}\}$. Then for any $d \in \mathcal{R}(x, \epsilon)$, it holds that $v_i^T d = 0$, $i \in [\bar{n}] \setminus \mathcal{A}(x, \epsilon)$. Define the function

$$Q_\epsilon(x) := f(x) + h(c(x)) + \sum_{i \in \mathcal{A}(x, \epsilon)} |v_i^T x|^p.$$

187 Note that $|v_i^T x|^p$, $i \in \mathcal{A}(x, \epsilon)$, is differentiable at x , and Q_ϵ is a continuous lower
188 approximation to Q . Also define

$$189 \quad (2.1) \quad \psi_Q^{\epsilon, \delta}(x) := Q_\epsilon(x) - \min_{\substack{x+d \in \mathcal{F} \\ d \in \mathcal{R}(x, \epsilon), \|d\| \leq \delta}} T_{Q_\epsilon}(x, d)$$

190

191 with

$$192 \quad T_{Q_\epsilon}(x, d) := f(x) + \nabla f(x)^T d + h(c(x) + J(x)d) + \sum_{i \in \mathcal{A}(x, \epsilon)} (|v_i^T x|^p + \nabla(|v_i^T x|^p)^T d),$$

193 where $J(x) = (\nabla c_1(x), \dots, \nabla c_r(x))^T$. Here, T_{Q_ϵ} is a convex approximation to Q_ϵ ,
194 obtained through linearization of smooth functions w.r.t. d , i.e., $f(x+d)$, $c(x+d)$
195 and $|v_i^T(x+d)|^p$, $i \in \mathcal{A}(x, \epsilon)$. The function $\psi_Q^{\epsilon, \delta}$ plays a crucial role in characterizing
196 the optimality condition of a local minimizer of (1.1). It represents the maximum
197 reduction of T_{Q_ϵ} within a neighborhood of current iterate. Intuitively, when current
198 iterate x is a local minimizer of (1.1) and $\epsilon = 0$, around x there is no feasible point
199 that can yield a greater reduction in the function value. By [8, Lemma 3.2] and [11,
200 Theorem 2.1] we obtain the following lemma.

201 **LEMMA 2.1.** *Let x_* be a local minimizer of (1.1). Then there exists $\bar{\delta} \in (0, 1]$
202 such that for any $\delta \in (0, \bar{\delta}]$, $\psi_Q^{0, \delta}(x_*) = 0$.*

Proof. As x_* is a local minimizer of (1.1), there exists $\delta_1 > 0$ such that x_* is a
global minimizer of (1.1) on $\mathcal{B}(x_*, \delta_1) \cap \mathcal{F}$. Let

$$\delta_2 = \min \left\{ 1, \min_{i \in \mathcal{A}(x_*, 0)} |v_i^T x_*| \right\}.$$

203 Obviously, $\delta_2 \in (0, 1]$. Note that there exists $\bar{\delta} \in (0, \min(\delta_1, \delta_2))$ such that for any
204 $x_* + d$ in the ball $\mathcal{B}(x_*, \bar{\delta})$,

$$205 \quad |v_i^T(x_* + d)| \geq |v_i^T x_*| - |v_i^T d| \geq \delta_2 - \bar{\delta} > 0, \quad i \in \mathcal{A}(x_*, 0).$$

206 Then $\sum_{i \in \mathcal{A}(x_*, 0)} |v_i^T x|$ is continuously differentiable in $\mathcal{B}(x_*, \bar{\delta})$. Moreover, since h is
207 Lipschitz continuous over \mathcal{F} and x_* is the global minimizer of (1.1) on $\mathcal{B}(x_*, \bar{\delta}) \cap \mathcal{F}$,
208 it holds that for any $\delta \in (0, \bar{\delta}]$,

$$\begin{aligned} 209 \quad Q(x_*) &= \min_{x_*+d \in \mathcal{F}, \|d\| \leq \delta} f(x_* + d) + h(c(x_* + d)) + \|V(x_* + d)\|_p^p \\ 210 \quad &\leq \min_{\substack{x_*+d \in \mathcal{F} \\ d \in \mathcal{R}(x_*, 0), \|d\| \leq \delta}} f(x_* + d) + h(c(x_* + d)) + \|V(x_* + d)\|_p^p \\ 211 \quad &= \min_{\substack{x_*+d \in \mathcal{F} \\ d \in \mathcal{R}(x_*, 0), \|d\| \leq \delta}} f(x_* + d) + h(c(x_* + d)) + \sum_{i \in \mathcal{A}(x_*, 0)} |v_i^T(x_* + d)|^p. \end{aligned}$$

213 Note that the equality in above relations can be reachable at $d = 0$. Thus 0 is a global
214 minimizer of the problem

$$215 \quad (2.2) \quad \min_{\substack{x_*+d \in \mathcal{F} \\ d \in \mathcal{R}(x_*, 0), \|d\| \leq \delta}} f(x_* + d) + h(c(x_* + d)) + \sum_{i \in \mathcal{A}(x_*, 0)} |v_i^T(x_* + d)|^p.$$

216 Then it yields from [8, Lemma 3.2] that $\psi_Q^{0,\delta}(x_*) = 0$ which completes the proof. \square

217 We call \bar{x} a **first-order stationary point of** (1.1), if $\psi_Q^{0,\delta}(\bar{x}) = 0$ for some
218 $\delta \in (0, 1]$.

219 *Remark 2.2.* We now show that if \bar{x} is a first-order stationary point of (1.1), i.e.
220 $\psi_Q^{0,\delta}(\bar{x}) = 0$ for some $\delta \in (0, 1]$, then \bar{x} is a limiting stationary point for a practice
221 example. The concept of a limiting stationary point for a proper lower semicontinuous
222 function has been used in the study for non-Lipschitz continuous minimization [10].
223 We recall from [24, Definition 8.3] that for a proper lower semicontinuous function Φ ,
224 the limiting subdifferential is defined as

$$225 \quad \partial\Phi(x) := \left\{ v : \exists x^k \xrightarrow{\Phi} x, v^k \rightarrow v \text{ with } \liminf_{z \rightarrow x^k} \frac{\Phi(z) - \Phi(x^k) - \langle v^k, z - x^k \rangle}{\|z - x^k\|} \geq 0, \forall k \right\},$$

226 where $x^k \xrightarrow{\Phi} x$ means both $x^k \rightarrow x$ and $\Phi(x^k) \rightarrow \Phi(x)$. In [10], a first-order stationary
227 condition using the limiting subdifferential for problem

$$228 \quad (2.3) \quad \min \Theta(x) := \lambda[(\|Ax - b\|_2^2 - \sigma^2)_+ + \|(Bx - h)_+\|_1] + \|x\|_p^p$$

229 is defined as

$$230 \quad (2.4) \quad 0 \in \partial\lambda((\|Ax - b\|_2^2 - \sigma^2)_+) + \partial\lambda\|(Bx - h)_+\|_1 + \partial\|x\|_p^p,$$

where $A \in \mathbb{R}^{r \times n}$, $B \in \mathbb{R}^{l \times n}$, $b \in \mathbb{R}^r$, $h \in \mathbb{R}^l$, $p \in (0, 1)$, $\sigma \geq 0$, and $\lambda > 0$. In
[10], a point \bar{x} is called a first-order stationary point of (2.3) if \bar{x} satisfies (2.4). Let
 $\mathcal{A}(x) = \{i : |x_i| > 0\}$. From Lemma 2.5 in [10], $\partial|t|^p = \mathbb{R}$ at $t = 0$. Hence, the
inclusion in (2.4) is trivial for $i \notin \mathcal{A}(x)$. Let

$$Q(x) = \lambda((\|Ax - b\|_2^2 - \sigma^2)_+ + \|(Bx - h)_+\|_1) + \sum_{i \in \mathcal{A}(x)} |x_i|^p,$$

$$\begin{aligned} T_Q(x, d) &= \lambda((\|Ax - b\|_2^2 - \sigma^2 + 2(Ax - b)^T Ad)_+ + \|(Bx - h + Bd)_+\|_1) \\ &\quad + \sum_{i \in \mathcal{A}(x)} (|x_i|^p + p|x_i|^{p-1} \text{sgn}(x_i) d_i) \end{aligned}$$

and

$$\psi_Q^{0,\delta}(x) = Q(x) - \min_{d \in \mathcal{R}(x), \|d\| \leq \delta} T_Q(x, d), \quad \mathcal{R}(x) = \{d \in \mathbb{R}^n : e_i^T d = 0, i \notin \mathcal{A}(x)\}.$$

Following the proof of Lemma 2.1, we can show that if \bar{x} is a local minimizer of (2.3),
then $\psi_Q^{0,\delta}(\bar{x}) = 0$ for some $\delta > 0$. Now we show that if $\psi_Q^{0,\delta}(\bar{x}) = 0$ for some $\delta > 0$,
then \bar{x} satisfies (2.4). Let $\bar{\delta} = \min_{i \in \mathcal{A}(\bar{x})} |\bar{x}_i|$. Then for any $\bar{x} + d \in \mathcal{B}(\bar{x}, \delta)$ with
 $\delta \in (0, \bar{\delta})$, we have $|\bar{x} + d|_i \geq |\bar{x}_i| - |d|_i \geq \bar{\delta} - \delta > 0, \forall i \in \mathcal{A}(\bar{x})$. Hence $\sum_{i \in \mathcal{A}(\bar{x})} |\bar{x}_i|^p$
is differentiable in $\mathcal{B}(\bar{x}, \delta)$. Moreover, we know that $(\|Ax - b\|_2^2 - \sigma^2)_+$ and $\|(Bx - h)_+\|_1$
are directionally differentiable. Therefore, Q is directionally differentiable at \bar{x} in the
direction $d \in \mathcal{R}(\bar{x})$. Additionally, the directional derivative of Q at \bar{x} in the direction
 $d \in \mathcal{R}(\bar{x})$ has the form

$$Q'(\bar{x}; d) = \lambda[2v(\bar{x})(A\bar{x} - b)^T Ad + u(\bar{x})^T Bd] + \sum_{i \in \mathcal{A}(\bar{x})} p|\bar{x}_i|^{p-1} \text{sgn}(\bar{x}_i) d_i,$$

where

$$v(\bar{x}) = \begin{cases} 1 & \text{if } \|A\bar{x} - b\|_2^2 > \sigma^2 \\ 0 & \text{if } \|A\bar{x} - b\|_2^2 < \sigma^2 \\ (\text{sgn}((A\bar{x} - b)^T Ad))_+ & \text{if } \|A\bar{x} - b\|_2^2 = \sigma^2 \end{cases}$$

and

$$u_i(\bar{x}) = \begin{cases} 1 & \text{if } (B\bar{x} - h)_i > 0 \\ 0 & \text{if } (B\bar{x} - h)_i < 0 \\ (\text{sgn}((Bd)_i))_+ & \text{if } (B\bar{x} - h)_i = 0, \quad i = 1, \dots, l. \end{cases}$$

Let $\hat{\delta} \in (0, \bar{\delta})$ such that

$$\hat{\delta} < \min \left\{ \frac{|\|A\bar{x} - b\|_2^2 - \sigma^2|}{\|(A\bar{x} - b)^T A\|_\infty}, \frac{|B\bar{x} - h|_i}{\|B\|_\infty} \right\}, \quad \text{for } \|A\bar{x} - b\|_2^2 - \sigma^2 \neq 0, (B\bar{x} - h)_i \neq 0.$$

Then it derives

$$T_Q(\bar{x}, d) = Q(\bar{x}) + Q'(\bar{x}; d), \quad \forall d \in \mathcal{R}(\bar{x}), \|d\| \leq \hat{\delta}.$$

From $\psi_Q^{0, \hat{\delta}}(\bar{x}) = 0$, we have

$$0 = Q(\bar{x}) - \min_{d \in \mathcal{R}(\bar{x}), \|d\| \leq \hat{\delta}} T_Q(\bar{x}, d) = - \min_{d \in \mathcal{R}(\bar{x}), \|d\| \leq \hat{\delta}} Q'(\bar{x}; d),$$

which implies $Q'(\bar{x}; d) \geq 0$ for any $d \in \mathcal{R}(\bar{x})$. From $\Theta(\bar{x} + d) \geq Q(\bar{x} + d)$ for $d \in \mathbb{R}^n$ and $\Theta(\bar{x}) = Q(\bar{x})$, the subderivative function $d\Theta(\bar{x})$ satisfies

$$d\Theta(\bar{x})(d) = \liminf_{\substack{t \downarrow 0 \\ d' \rightarrow d}} \frac{\Theta(\bar{x} + td') - \Theta(\bar{x})}{t} \geq \liminf_{\substack{t \downarrow 0 \\ d' \rightarrow d}} \frac{Q(\bar{x} + td') - Q(\bar{x})}{t}.$$

231 Hence $d\Theta(\bar{x})(d) \geq 0$ for $d \in \mathcal{R}(\bar{x})$ and $d\Theta(\bar{x})(d) = +\infty$ for $d \notin \mathcal{R}(\bar{x})$. By [24, Exercise
232 8.4], we find that 0 is in the regular subdifferential of Θ at \bar{x} , and thus by [24, Definition
233 8.3, Exercise 10.10], the inclusion in (2.4) holds at \bar{x} .

234 We now present the definition of an (ϵ, δ) -approximate first-order stationary point
235 of (1.1).

236 **DEFINITION 2.3.** *Given $\epsilon > 0$, we call $x \in \mathcal{F}$ an (ϵ, δ) -approximate first-order
237 stationary point of (1.1), if $\psi_Q^{\epsilon, \delta}(x) \leq \epsilon\delta$ for some $\delta \in (0, 1]$.*

238 The concept of (ϵ, δ) -approximate first-order stationary points has been used in
239 [6, 7, 11], which generalizes the concept of ϵ -approximate first-order stationary points
240 with $\delta = 1$ in some papers, e.g. [5, 12, 17]. Our definitions of first-order stationary
241 point and (ϵ, δ) -approximate first-order stationary point are based on the concepts
242 in [5, 6, 7, 8, 11, 12, 17] and related articles. In Lemma 2.1, we show that a local
243 minimizer x^* of (1.1) is a $(0, \delta)$ -approximate first-order stationary point of (1.1) for
244 some $\delta > 0$, which implies that x^* is an (ϵ, δ) -approximate first-order stationary
245 point of (1.1) for $\epsilon > 0$. Within a δ -ball centered at an (ϵ, δ) -approximate first-order
246 stationary point, the maximum reduction of the approximation model does not exceed
247 $\epsilon\delta$. In practice, the choice of (ϵ, δ) depends on the users' need for the quality of a
248 computed solution. For each k , let x_k be an (ϵ_k, δ_k) -approximate first-order stationary
249 point of (1.1) for some δ_k with $1 \geq \delta_k > 0$ and $\epsilon_k > 0$. If $\{\delta_k\}$ has a uniform positive
250 lower bound as $\epsilon_k \rightarrow 0$, following the proof of [11, Theorem 2.2] we can obtain that
251 any cluster point of $\{x_k\}$ is a first-order stationary point of (1.1).

252 In the following context, we consider $\epsilon > 0$. We now prepare for the design of the
 253 main algorithm. The main step of the algorithm is to construct a model function to
 254 predict the behavior of the objective function Q at current iterate x along a direction
 255 s . For the non-Lipschitz regularizer in the objective function, we focus on indices in
 256 $\mathcal{A}(x, \epsilon)$ and discard those close to non-Lipschitz continuity. We define the following
 257 Lipschitz continuous approximation of $|v_i^T(x+s)|^p$ in a similar approach in [11] and
 258 [12]:

$$259 \quad (2.5) \quad m_i(x, s) := |v_i^T x|^p + p|v_i^T x|^{p-1} (|v_i^T(x+s)| - |v_i^T x|), \quad i \in \mathcal{A}(x, \epsilon).$$

260 Supposing that $v_i^T(x+s) \neq 0$, $i \in \mathcal{A}(x, \epsilon)$, as analyzed in [11], m_i is the first-
 261 order Taylor's expansion of $|v_i^T x + \zeta_i \frac{v_i^T x}{|v_i^T x|}|^p$ expressed as a function of the scalar
 262 $\zeta_i := |v_i^T(x+s)| - |v_i^T x| \geq -|v_i^T s|$. Regarding the smooth function f , the calculation
 263 of exact first-order derivatives of f can be expensive sometimes even impossible in
 264 many scenarios. We can only get access to approximate gradients of f . For ease of
 265 notations, given x_k and s_k we denote g_k as an approximation to ∇f at x_k , and

$$266 \quad \mathcal{A}_k = \mathcal{A}(x_k, \epsilon), \mathcal{R}_k = \mathcal{R}(x_k, \epsilon), c_k = c(x_k), J_k = J(x_k) \text{ and } s_i^k = v_i^T s_k$$

268 for $i \in [\bar{n}]$. Due to existence of the convex but possibly nonsmooth function h , we
 269 design the following proximal type subproblem at k th iteration:

$$270 \quad (2.6) \quad \min_{\substack{x_k+s \in \mathcal{F} \\ s \in \mathcal{R}_k}} m(x_k, s) := g_k^T s + h(c_k + J_k s) + \sum_{i \in \mathcal{A}_k} m_i(x_k, s) + \frac{1}{2\eta} \|s\|^2,$$

271 where $\eta > 0$ is a proximal parameter. It is worth noting that subproblem (2.6) is
 272 a strongly convex minimization problem over a convex set, thus it admits a unique
 273 global minimizer. Note that resolution of (2.6) only involves matrix-vector products
 274 and does not affect the evaluations of (inexact) derivatives of f , thus has no impact on
 275 the iteration complexity and oracle complexity of the proposed algorithm. Moreover,
 276 when \mathcal{F} and h exhibit polyhedral structures, for example, $\mathcal{F} = [b_l, b_u] \subseteq \mathbb{R}^n$ with
 277 $-b_l, b_u \in \mathbb{R}_+^n$, and $h(\cdot) = \|(\cdot)_+\|_1$, by introducing $\bar{z} = (c_k + J_k s)_+ \in \mathbb{R}^r$, (2.6) is
 278 equivalent to the following linearly constrained convex program:

$$279 \quad \min_{s, \bar{z}} g_k^T s + e^T \bar{z} + \sum_{i \in \mathcal{A}_k} p|v_i^T x_k|^{p-1} |v_i^T(x_k+s)| + \frac{1}{2} \|s\|^2$$

$$\text{s.t. } b_l \leq x_k + s \leq b_u, \quad 0 \leq \bar{z}, c_k + J_k s \leq \bar{z}, \quad v_i^T s = 0, i \notin \mathcal{A}_k,$$

280 where $e = (1, 1, \dots, 1)^T \in \mathbb{R}^r$. Numerous state-of-the-art approaches have been ex-
 281 tensively studied for solving linearly constrained convex program in the literature.

282 In theoretical analysis, however, an inexact solution of (2.6) can be enough.
 283 Specifically, we solve (2.6) to look for s_k with $m(x_k, s_k) < m(x_k, 0)$ such that the
 284 near optimality is achieved in that

$$285 \quad (2.7) \quad \psi_m^{\epsilon, \delta}(x_k, s_k) \leq \min \left\{ \theta \epsilon, p \min_{i \in \mathcal{A}(x_k+s_k, \epsilon)} |v_i^T(x_k+s_k)| \right\} \delta, \quad \text{for some } \delta \in (0, 1],$$

286 where $\theta \in (0, 1)$ and

$$287 \quad \psi_m^{\epsilon, \delta}(x_k, s_k) := h(c_k + J_k s_k) - \min_{\substack{x_k+s_k+d \in \mathcal{F} \\ d \in \mathcal{R}(x_k+s_k, \epsilon), \|d\| \leq \delta}} T_m(x_k, s_k; d)$$

288 with $m_0(x_k, s) := g_k^T s + \frac{1}{2\eta} \|s\|^2$ and

$$289 \quad T_m(x_k, s; d) := h(c_k + J_k(s + d)) + \nabla_s m_0(x_k, s)^T d + \sum_{i \in \mathcal{A}(x_k + s, \epsilon)} \nabla_s m_i(x_k, s)^T d.$$

290 It is noteworthy that $\psi_m^{\epsilon, \delta}$ describes the potential maximum reduction of T_m within
 291 a neighborhood of s_k with radius δ . This measure is defined in a similar way to that
 292 in Definition 2.3. When the reduction is below a certain level, s_k is regarded as an
 293 inexact minimizer of (2.4). Moreover, by the definition of \mathcal{R}_k , for any $i \in [\bar{n}] \setminus \mathcal{A}_k$,
 294 $v_i^T s_k = 0$, thus $v_i^T(x_k + s_k) = v_i^T x_k$. That is, once $|v_i^T x_k| \leq \epsilon$ for some $i \in [\bar{n}]$, the
 295 value of $v_i^T(x_k + s_k)$ will be fixed and the remaining minimization will be carried out
 296 on $\mathcal{R}(x_k + s_k, \epsilon)$. Therefore, the following relations hold:

$$297 \quad (2.8) \quad \mathcal{R}_k^+ := \mathcal{R}(x_k + s_k, \epsilon) \subseteq \mathcal{R}_k, \quad \mathcal{A}_k^+ := \mathcal{A}(x_k + s_k, \epsilon) \subseteq \mathcal{A}_k.$$

298 We are now ready to present the main algorithm framework for proximal inexact
 299 gradient methods for (1.1) as Algorithm 2.1.

Algorithm 2.1

Input: $x_0 \in \mathcal{F}$, $\epsilon \in (0, 1]$, $\eta > 0$, $\bar{\beta} \in (0, \bar{w})$ with $\bar{w} \in (0, 1)$, $s_{-1} = 0$.

- 1: **for** $k = 0, 1, \dots$, **do**
 - 2: Obtain g_k from **InexactOracle**.
 - 3: Solve (2.6) to find an approximate minimizer s_k with $m(x_k, s_k) < m(x_k, 0)$ satisfying (2.7), then go to Step 5. If the solution of (2.6) is zero, then go to Step 4.
 - 4: Set $s_k = 0$. If $s_{k-1} = 0$, terminate and return x_k ; otherwise, go to Step 5.
 - 5: Set $x_{k+1} = x_k + s_k$. If $\|s_k\| + \|s_{k-1}\| \leq \bar{\beta}\epsilon$ and $\mathcal{A}_k \setminus \mathcal{A}_{k+1} = \emptyset$, terminate and return x_{k+1} .
 - 6: $k := k + 1$.
 - 7: **end for**
-

300 *Remark 2.4.* In Algorithm 2.1 two termination criteria are employed. One is
 301 $s_{k-1} = s_k = 0$ in Step 4. In this case, similar to [6, 11] we terminate the algorithm
 302 and return x_k . It will be shown in Lemma 3.1 that x_k is an (ϵ, δ) -approximate first-
 303 order stationary point of (1.1) for some $\delta \in (0, 1]$. On the other hand, if $\mathcal{A}_{k+1} = \mathcal{A}_k$
 304 (and hence $\mathcal{R}_{k+1} = \mathcal{R}_k$) and $\|s_k\| + \|s_{k-1}\|$ is sufficiently small and $\mathcal{A}_k \setminus \mathcal{A}_{k+1} = \emptyset$
 305 then there is no i s.t. $|v_i^T x_k| > \epsilon$ but $|v_i^T(x_k + s_k)| < \epsilon$, we return x_{k+1} and will prove
 306 that x_{k+1} is an approximate first-order stationary point when $s_k = 0$ in Lemma 3.1
 307 and when $s_k \neq 0$ in Lemma 3.6, respectively. In addition, as $s_k \in \mathcal{R}_k$ for any $k \geq 1$,
 308 it follows from (2.8) that $\mathcal{A}_{k+1} \setminus \mathcal{A}_k = \emptyset$ for any $k \geq 1$. Moreover, in Algorithm 2.1
 309 we obtain inexact gradient g_k through calling the subroutine **InexactOracle**, which
 310 may adopt different ways to generate an inexact gradient of f at the inquiry iterate
 311 x_k . So we simply omit the required inputs by **InexactOracle** here and specify them
 312 when necessary.

313 In the following, we denote the unique global minimizer of (2.6) by s_k^* . If $s_k^* \neq 0$,
 314 it obviously holds that $m(x_k, s_k^*) < m(x_k, 0)$. Moreover, we can guarantee that $s_k = s_k^*$
 315 satisfies (2.7) for some $\delta \in (0, 1]$, as shown in the lemma below.

316 **LEMMA 2.5.** *Suppose that $s_k^* \neq 0$. Then there exists $\underline{\mu}_k \in (0, 1]$ such that (2.7)*
 317 *holds for $s_k = s_k^*$ and any $\delta \in (0, \underline{\mu}_k]$.*

318 *Proof.* Consider the auxiliary problem
 (2.9)

$$319 \quad \min_{\substack{x_k + s_k^* + d \in \mathcal{F} \\ d \in \mathcal{R}(x_k + s_k^*, \epsilon)}} h(c_k + J_k(s_k^* + d)) + m_0(x_k, s_k^* + d) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^* + d).$$

320 Due to the strong convexity, (2.9) has a unique global minimizer, which is denoted
 321 by \bar{s}_k . As $\bar{s}_k \in \mathcal{R}(x_k + s_k^*, \epsilon) \subseteq \mathcal{R}_k$, we have $m_i(x_k, s_k^*) = m_i(x_k, s_k^* + \bar{s}_k)$ for any
 322 $i \in \mathcal{A}_k \setminus \mathcal{A}(x_k + s_k^*, \epsilon)$. Then it yields that

$$\begin{aligned} 323 \quad & m(x_k, s_k^* + \bar{s}_k) \\ 324 \quad & = h(c_k + J_k(s_k^* + \bar{s}_k)) + m_0(x_k, s_k^* + \bar{s}_k) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^* + \bar{s}_k) \\ 325 \quad & + \sum_{i \in \mathcal{A}_k \setminus \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^* + \bar{s}_k) \\ 326 \quad & \leq h(c_k + J_k s_k^*) + m_0(x_k, s_k^*) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^*) \\ 327 \quad & + \sum_{i \in \mathcal{A}_k \setminus \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^*) = m(x_k, s_k^*), \\ 328 \end{aligned}$$

329 where the inequality follows from the optimality of \bar{s}_k .

330 Due to the optimality and uniqueness of s_k^* as the global minimizer of (2.6), we
 331 obtain $\bar{s}_k = 0$. Thus 0 is the global minimizer of (2.9). Then for any $d \in \mathcal{R}(x_k + s_k^*, \epsilon)$
 332 satisfying $x_k + s_k^* + d \in \mathcal{F}$, it holds that

$$\begin{aligned} 333 \quad & g_k^T s_k^* + \frac{1}{2\eta} \|s_k^*\|^2 + h(c_k + J_k s_k^*) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^*) \\ 334 \quad & \leq g_k^T (s_k^* + d) + \frac{1}{2\eta} \|s_k^* + d\|^2 + h(c_k + J_k (s_k^* + d)) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} m_i(x_k, s_k^* + d) \\ 335 \end{aligned}$$

336 which yields

$$337 \quad h(c_k + J_k s_k^*) - h(c_k + J_k (s_k^* + d)) - g_k^T d - \frac{1}{\eta} (s_k^*)^T d - \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} \nabla m_i(x_k, s_k^*)^T d$$

(2.10)

$$338 \quad \leq \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} (m_i(x_k, s_k^* + d) - m_i(x_k, s_k^*) - \nabla_s m_i(x_k, s_k^*)^T d) + \frac{1}{2\eta} \|d\|^2.$$

340 Note that there exists $\hat{\mu}_k$ such that for any $d \in \mathcal{R}(x_k + s_k^*, \epsilon)$ with $\|d\| \leq \hat{\mu}_k$ and
 341 $x_k + s_k^* + d \in \mathcal{F}$,

(2.11)

$$342 \quad \text{sgn}(v_i^T (x_k + s_k^* + d)) = \text{sgn}(v_i^T (x_k + s_k^*)) \quad \text{and} \quad |v_i^T (x_k + s_k^* + d)| > \epsilon, \quad \forall i \in \mathcal{A}(x_k + s_k^*, \epsilon),$$

343 which together with (2.10) indicate that

$$\begin{aligned} 344 \quad & h(c_k + J_k s_k^*) - h(c_k + J_k (s_k^* + d)) - g_k^T d - \frac{1}{\eta} s_k^{*T} d - \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} \nabla_s m_i(x_k, s_k^*)^T d \\ 345 \quad & \leq \frac{1}{2\eta} \|d\|^2. \\ 346 \end{aligned}$$

347 Hence, by the definition of $\psi_m^{\epsilon, \delta}(x_k, s_k^*)$, there exists $\underline{\mu}_k \in (0, \min\{1, \hat{\mu}_k\}]$ such that for
 348 any $\delta \in (0, \underline{\mu}_k]$,

$$349 \quad (2.12) \quad \psi_m^{\epsilon, \delta}(x_k, s_k^*) \leq \frac{1}{2\eta} \delta^2 \leq \min \left\{ \theta\epsilon, p \min_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} |v_i^T(x_k + s_k^*)| \right\} \delta. \quad \square$$

350 Define the set $\mathcal{S}_k := \{s : x_k + s \in \mathcal{F}\} \cap \{s : s - s_k^* \in \mathcal{R}(x_k + s_k^*, \epsilon)\}$. Obviously
 351 $s_k^* \in \mathcal{S}_k$. Without loss of generality, we assume in the following that $\mathcal{S}_k \setminus \{s_k^*\} \neq \emptyset$.

352 LEMMA 2.6. *Suppose that $s_k^* \neq 0$. Then there exist $\tilde{\mu}_k, \bar{\mu}_k \in (0, 1]$ such that for
 353 any $\delta \in (0, \tilde{\mu}_k]$ and any $s \in \mathcal{S}_k \cap \mathcal{B}(s_k^*, \bar{\mu}_k)$, we have*

$$354 \quad (2.13) \quad m(x_k, s) < m(x_k, 0) \quad \text{and} \quad \psi_m^{\epsilon, \delta}(x_k, s) \leq \min \left\{ \theta\epsilon, p \min_{i \in \mathcal{A}(x_k + s, \epsilon)} |v_i^T(x_k + s)| \right\} \delta.$$

355 *Proof.* Note that if $s_k^* \neq 0$, there exists $\bar{\mu}_k \in (0, 1]$ such that for any $s \in \mathcal{B}(s_k^*, \bar{\mu}_k)$,
 356 $m(x_k, s) < m(x_k, 0)$ and $\mathcal{A}(x_k + s_k^*, \epsilon) \subseteq \mathcal{A}(x_k + s, \epsilon)$. Hence, for any $s \in \mathcal{S}_k \cap \mathcal{B}(s_k^*, \bar{\mu}_k)$,

$$357 \quad (2.14) \quad \mathcal{A}(x_k + s_k^*, \epsilon) = \mathcal{A}(x_k + s, \epsilon) \subseteq \mathcal{A}_k.$$

358 For any given $s \in \mathcal{S}_k \cap \mathcal{B}(s_k^*, \bar{\mu}_k)$, we define $\mathcal{F}_s := \{d : x_k + s + d \in \mathcal{F}\}$ which is
 359 obviously convex due to the convexity of \mathcal{F} . For any $d \in \mathcal{F}_s$, we denote its projection
 360 onto $\mathcal{F}_{s_k^*}$ as \bar{d} . If $d = \bar{d}$, then set $d_1 = d$. Otherwise, as $x_k + s_k^* + d \notin \mathcal{F}$, there exists
 361 $d_1 \in \mathcal{F}_{s_k^*}$ such that $x_k + s_k^* + d_1$ is the projection of $x_k + s_k^* + d$ onto \mathcal{F} . Then it follows
 362 from definition of the projection operator and $x_k + s + d \in \mathcal{F}$ that $\|d - \bar{d}\| \leq \|d - d_1\|$
 363 and

$$364 \quad \|(x_k + s_k^* + d) - (x_k + s_k^* + d_1)\| \leq \|(x_k + s_k^* + d) - (x_k + s + d)\| = \|s_k^* - s\|,$$

365 thus

$$366 \quad (2.15) \quad \|d - \bar{d}\| \leq \|s_k^* - s\|.$$

367 Then by definition of $T_m(x_k, s; d)$ and (2.14) we obtain that for any $d \in \mathcal{F}_s$,

$$\begin{aligned} 368 \quad & h(c_k + J_k s) - T_m(x_k, s; d) \\ 369 \quad & = h(c_k + J_k s_k^*) + h(c_k + J_k s) - h(c_k + J_k s_k^*) - \left[h(c_k + J_k (s_k^* + d)) \right. \\ 370 \quad & \quad \left. + h(c_k + J_k (s + d)) - h(c_k + J_k (s_k^* + d)) + d^T \nabla_s m_0(x_k, s_k^*) \right. \\ 371 \quad & \quad \left. + d^T (\nabla_s m_0(x_k, s) - \nabla_s m_0(x_k, s_k^*)) + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} d^T \nabla_s m_i(x_k, s_k^*) \right. \\ 372 \quad & \quad \left. + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} d^T \nabla_s m_i(x_k, s) - \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} d^T \nabla_s m_i(x_k, s_k^*) \right] \\ 373 \quad & = h(c_k + J_k s_k^*) - \left[h(c_k + J_k (s_k^* + d)) + d^T \nabla_s m_0(x_k, s_k^*) \right. \\ 374 \quad & \quad \left. + \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} d^T \nabla_s m_i(x_k, s_k^*) \right] + \Gamma_k \\ 375 \quad & = h(c_k + J_k s_k^*) - T_m(x_k, s_k^*; d) + \Gamma_k \\ 376 \quad & = h(c_k + J_k s_k^*) - T_m(x_k, s_k^*; \bar{d}) + \Gamma_k + T_m(x_k, s_k^*; \bar{d}) - T_m(x_k, s_k^*; d), \end{aligned}$$

378 where

$$\begin{aligned}
379 \quad \Gamma_k &= h(c_k + J_k s) - h(c_k + J_k s_k^*) - (h(c_k + J_k(s + d)) - h(c_k + J_k(s_k^* + d))) \\
380 &\quad - d^T (\nabla_s m_0(x_k, s) - \nabla_s m_0(x_k, s_k^*)) \\
381 &\quad - \sum_{i \in \mathcal{A}(x_k + s_k^*, \epsilon)} d^T (\nabla_s m_i(x_k, s) - \nabla_s m_i(x_k, s_k^*)). \\
382
\end{aligned}$$

383 Note that, on the one hand, by the definition of m_0 ,

$$384 \quad \|\nabla_s m_0(x_k, s) - \nabla_s m_0(x_k, s_k^*)\| \leq \frac{1}{\eta} \|s - s_k^*\|,$$

385 while on the other hand, by (2.5) and (2.14),

$$386 \quad \nabla_s m_i(x_k, s) = p |v_i^T x_k|^{p-1} \text{sgn}(v_i^T x_k) v_i = \nabla_s m_i(x_k, s_k^*), \quad \forall i \in \mathcal{A}(x_k + s_k^*, \epsilon).$$

387 Recall that h is Lipschitz continuous over \mathcal{F} . It together with the boundedness of J_k
388 derives

$$389 \quad (2.16) \quad \Gamma_k = \mathcal{O}(\|s - s_k^*\|).$$

390 Besides, it indicates from definition of T_m that

$$391 \quad T_m(x_k, s_k^*; \bar{d}) - T_m(x_k, s_k^*; d) = \mathcal{O}(\|d - \bar{d}\|) = \mathcal{O}(\|s - s_k^*\|).$$

392 Therefore, there exists $\tilde{\mu}_k \in (0, \min\{\underline{\mu}_k, \bar{\mu}_k\})$ such that $\tilde{\mu}_k + \tilde{\mu}_k^{(1+\varrho)} < \underline{\mu}_k$ with $\varrho > 0$,
393 and for any $\delta \in (0, \tilde{\mu}_k]$ and $s \in \mathcal{S}_k \cap \mathcal{B}(s_k^*, \delta^{1+\varrho})$, the following relations can be derived:

$$\begin{aligned}
394 \quad & h(c_k + J_k s) - \min_{\substack{x_k + s + d \in \mathcal{F} \\ d \in \mathcal{R}(x_k + s, \epsilon), \|d\| \leq \delta}} T_m(x_k, s; d) \\
395 & \leq h(c_k + J_k s_k^*) - \min_{\substack{x_k + s_k^* + \bar{d} \in \mathcal{F} \\ \bar{d} \in \mathcal{R}(x_k + s_k^*, \epsilon), \|\bar{d}\| \leq \delta + \delta^{1+\varrho}}} T_m(x_k, s_k^*; \bar{d}) + \mathcal{O}(\|s - s_k^*\|) \\
396 & \leq \psi_m^{\epsilon, \delta + \delta^{1+\varrho}}(x_k, s_k^*) + \mathcal{O}(\|s - s_k^*\|) \\
397 \quad (2.17) & \leq \frac{1}{2\eta} (\delta + \delta^{1+\varrho})^2 + \mathcal{O}(\delta^{1+\varrho}) \leq \min \left\{ \theta \epsilon, p \min_{i \in \mathcal{A}(x_k + s_k^*)} |v_i^T(x_k + s_k^*)| \right\} \delta, \\
398
\end{aligned}$$

399 where $\underline{\mu}_k$ is introduced in Lemma 2.5, the first inequality is due to $\|\bar{d}\| \leq \|d\| + \|s - s_k^*\|$
400 $\leq \delta + \delta^{1+\varrho}$ and the third inequality follows from (2.12). The proof is completed. \square

401 **3. Oracle complexity.** In this section, we will analyze the oracle complexity
402 of Algorithm 2.1 in terms of the total number of inexact gradient evaluations until
403 the algorithm terminates. In the following, we use \mathcal{K} to denote the set of all iteration
404 indices until the termination of Algorithm 2.1. Let $\{x_k\}$ be the iterate sequence gener-
405 ated during the algorithm. Since f and c are Lipschitz continuously differentiable and
406 h is Lipschitz continuous over \mathcal{F} , there exist positive constants $M_F, \kappa, L_f, L_h, L_c^0, L_c^1$
407 such that for any $x, y \in \mathcal{F}$, $\|x\| \leq M_F$ and $\|\nabla f(x)\| \leq \kappa$ and

$$\begin{aligned}
408 \quad & \|\nabla f(x) - \nabla f(y)\| \leq L_f \|x - y\|, \quad |h(x) - h(y)| \leq L_h \|x - y\|, \\
409 & \|c(x) - c(y)\| \leq L_c^0 \|x - y\|, \quad \|\nabla c(x) - \nabla c(y)\| \leq L_c^1 \|x - y\|.
\end{aligned}$$

411 To derive desired theoretical properties of Algorithm 2.1, we lay out the following
412 assumption on gradient approximations returned by **InexactOracle**.

413 ASSUMPTION 3.1. For any $k \in \mathcal{K}$, the gradient approximation g_k satisfies

$$414 \quad (3.1) \quad \|g_k - \nabla f(x_k)\| \leq \beta \max\{\bar{L} \min(\|s_{k-1}\|, D), \epsilon\},$$

415 where $\beta \in (0, \frac{1-\theta}{17})$, $D > 0$ and $\bar{L} \in (0, \bar{\beta}^{-1})$.

416 The parameter θ in Assumption 3.1 was introduced initially in (2.7). And Assumption 3.1 ensures that g_k is uniformly upper bounded, namely,

$$418 \quad (3.2) \quad \|g_k\| \leq \chi := \kappa + \beta \max\{\bar{L}D, \epsilon\} \quad \text{for any } k \in \mathcal{K}.$$

419 We will show in Lemma 3.1 that x_k is an (ϵ, δ) -approximate first-order stationary
420 point of (1.1) when Algorithm 2.1 terminates in Step 4 or in Step 5 with $s_k = 0$.
421 When Algorithm 2.1 terminates in Step 5 with $s_k \neq 0$, we will show in Lemma 3.6
422 that the output x_{k+1} is an approximate first-order stationary point of (1.1).

423 LEMMA 3.1. Suppose that Algorithm 2.1 terminates in Step 4 or in Step 5 with
424 $s_k = 0$. Then x_k is an (ϵ, δ) -approximate first-order stationary point of (1.1) for some
425 $\delta \in (0, 1]$.

426 *Proof.* Whenever Algorithm 2.1 terminates in Step 4 or in Step 5 with $s_k = 0$,
427 it holds that $s_k = 0$ and $\|s_{k-1}\| \leq \bar{\beta}\epsilon$. As $s_k = 0$, by the algorithmic framework
428 there is no step s_k with $m(x_k, s_k) < m(x_k, 0)$ satisfying (2.7). Then it follows from
429 Lemma 2.5 that $m(x_k, d) \geq m(x_k, 0)$ for any $d \in \mathcal{R}_k$ satisfying $x_k + d \in \mathcal{F}$. And by
430 the definition of T_{Q_ϵ} the following equalities hold:

$$\begin{aligned} 431 \quad & m(x_k, d) - m(x_k, 0) \\ 432 \quad &= h(c_k + J_k d) - h(c_k) + g_k^T d + \frac{1}{2\eta} \|d\|^2 + \sum_{i \in \mathcal{A}_k} [m_i(x_k, d) - m_i(x_k, 0)] \\ 433 \quad &= -\left[f(x_k) + h(c_k) + \sum_{i \in \mathcal{A}_k} |v_i^T x_k|^p - T_{Q_\epsilon}(x_k, d) \right] + (g_k - \nabla f(x_k))^T d \\ 434 \quad &+ \frac{1}{2\eta} \|d\|^2 + \sum_{i \in \mathcal{A}_k} (m_i(x_k, d) - m_i(x_k, 0) - (\nabla |v_i^T x|^p|_{x=x_k})^T d) \\ 435 \quad &= -(Q_\epsilon(x_k) - T_{Q_\epsilon}(x_k, d)) + (g_k - \nabla f(x_k))^T d + \frac{1}{2\eta} \|d\|^2 \\ 436 \quad &+ \sum_{i \in \mathcal{A}_k} (m_i(x_k, d) - m_i(x_k, 0) - (\nabla |v_i^T x|^p|_{x=x_k})^T d). \end{aligned}$$

438 Note that there exists $\bar{\delta} \in (0, 1]$ such that for all d with $\|d\| \leq \bar{\delta}$, $\text{sgn}(v_i^T(x_k + d)) =$
439 $\text{sgn}(v_i^T x_k)$ for any $i \in \mathcal{A}_k$, then

$$\begin{aligned} 440 \quad & m_i(x_k, d) - m_i(x_k, 0) = p|v_i^T x_k|^{p-1} (|v_i^T(x_k + d)| - |v_i^T x_k|) \\ 441 \quad &= \text{sgn}(v_i^T x_k) \cdot p|v_i^T x_k|^{p-1} v_i^T d \\ 442 \quad &= (\nabla (|v_i^T x|^p)|_{x=x_k})^T d \quad \text{for any } i \in \mathcal{A}_k. \end{aligned}$$

444 It follows from $\|s_{k-1}\| \leq \bar{\beta}\epsilon$ and Assumption 3.1 that $\|g_k - \nabla f(x_k)\| \leq \beta \max\{\bar{L}\bar{\beta}, 1\}\epsilon$.
445 Hence, by $\beta < (1-\theta)/(32M_F + 1)$ and $\bar{L} < \bar{\beta}^{-1}$ we can choose $\delta < \bar{\delta}$ sufficiently small
446 such that for any $d \in \mathcal{R}_k$ satisfying $x_k + d \in \mathcal{F}$ and $\|d\| \leq \delta$,

$$447 \quad Q_\epsilon(x_k) - T_{Q_\epsilon}(x_k, d) \leq (g_k - \nabla f(x_k))^T d + \frac{\|d\|^2}{2\eta} \leq \beta \max\{\bar{L}\bar{\beta}, 1\}\epsilon \|d\| + \frac{\|d\|^2}{2\eta} \leq \epsilon\delta,$$

448 which yields the conclusion by Definition 2.3. \square

450 In the following, we assume that Algorithm 2.1 does not terminate at k th iteration
 451 with $s_k = 0$. It derives from definitions of Q_ϵ and \mathcal{A}_k that for any $x \in \mathcal{F}$,

$$452 \quad Q_\epsilon(x) = Q(x) - \sum_{i \in [\bar{n}] \setminus \mathcal{A}_k} |v_i^T x|^p \geq Q^* - |[\bar{n}] \setminus \mathcal{A}_k| \epsilon^p \geq Q^* - \bar{n} \epsilon^p =: Q_\epsilon^*,$$

453 where Q^* is the lower bound of Q on \mathcal{F} . The lemma below provides an upper bound
 454 on the accumulated square of step lengths.

455 **LEMMA 3.2.** *Suppose that $\eta < (L_f + L_h L_c^1)^{-1}$. Then it holds that*

$$456 \quad (3.3) \quad \left(\frac{1}{2\eta} - \frac{L_f + L_h L_c^1}{2} \right) \sum_{k \in \mathcal{K}} \|s_k\|^2 \leq \sum_{k \in \mathcal{K}} (\nabla f(x_k) - g_k)^T s_k + Q_\epsilon(x_0) - Q_\epsilon^*.$$

457 *Proof.* It follows from (2.8) and Lipschitz continuity of ∇f , h and ∇c that

$$\begin{aligned} 458 \quad & Q_\epsilon(x_k + s_k) - Q_\epsilon(x_k) \\ 459 \quad &= f(x_k + s_k) + h(c(x_k + s_k)) - f(x_k) - h(c_k) + \sum_{i \in \mathcal{A}_k^+} |v_i^T(x_k + s_k)|^p - \sum_{i \in \mathcal{A}_k} |v_i^T x_k|^p \\ 460 \quad &\leq h(c(x_k + s_k)) - h(c_k) + (\nabla f(x_k))^T s_k + \frac{L_f}{2} \|s_k\|^2 + \sum_{i \in \mathcal{A}_k} (m_i(x_k, s_k) - m_i(x_k, 0)) \\ 461 \quad &= h(c_k + J_k s_k) - h(c_k) + g_k^T s_k + \frac{1}{2\eta} \|s_k\|^2 + \sum_{i \in \mathcal{A}_k} (m_i(x_k, s_k) - m_i(x_k, 0)) \\ 462 \quad &+ (\nabla f(x_k) - g_k)^T s_k + \left(\frac{L_f}{2} - \frac{1}{2\eta} \right) \|s_k\|^2 + h(c(x_k + s_k)) - h(c_k + J_k s_k) \\ 463 \quad &(3.4) \\ 464 \quad &\leq m(x_k, s_k) - m(x_k, 0) + (\nabla f(x_k) - g_k)^T s_k + \left(\frac{L_f + L_h L_c^1}{2} - \frac{1}{2\eta} \right) \|s_k\|^2, \end{aligned}$$

465 where the first inequality is due to $\mathcal{A}_k^+ \subseteq \mathcal{A}_k$ and [11, Lemma 3.2] which shows that
 466 $m_i(x_k, s_k) \geq |v_i^T(x_k + s_k)|^p$ for $i \in \mathcal{A}_k$. Then (3.4) indicates from $m(x_k, 0) \geq m(x_k, s_k)$
 467 that

$$468 \quad (3.5) \quad Q_\epsilon(x_k + s_k) - Q_\epsilon(x_k) \leq (\nabla f(x_k) - g_k)^T s_k + \left(\frac{L_f + L_h L_c^1}{2} - \frac{1}{2\eta} \right) \|s_k\|^2.$$

469 Hence, summing up (3.5) over $k \in \mathcal{K}$ and by $Q_\epsilon(x) \geq Q_\epsilon^*$ for all $x \in \mathcal{F}$ implies (3.3). \square

470 For a given $\mu > 0$ which is independent of ϵ , we define

$$471 \quad (3.6) \quad \mathcal{O}_{k,\mu} := \{i \in \mathcal{A}_k^+ : \min\{|v_i^T x_k|, |v_i^T(x_k + s_k)|\} \geq \mu\},$$

$$472 \quad (3.7) \quad \bar{Q}_{k,\mu}(x) := f(x) + \sum_{i \in \mathcal{O}_{k,\mu}} m_i(x, 0),$$

$$473 \quad (3.8) \quad \bar{T}_{k,\mu}(x, s) := f(x) + \nabla f(x)^T s + \sum_{i \in \mathcal{O}_{k,\mu}} m_i(x, s).$$

475 The following lemma characterizes the relation between derivatives of $\bar{Q}_{k,\mu}$ and $\bar{T}_{k,\mu}$.

476 **LEMMA 3.3.** *It holds that for any $k \geq 1$,*

$$477 \quad (3.9) \quad \|\nabla \bar{Q}_{k,\mu}(x_k + s_k) - \nabla_s \bar{T}_{k,\mu}(x_k, s_k)\| \leq L(\mu) \|s_k\|,$$

478 where $L(\mu) := L_f + \frac{p(2-p)}{1-p} \mu^{p-2}$.

479 *Proof.* By definitions of $Q_{k,\mu}$ and $T_{k,\mu}$, it is easy to obtain

$$\begin{aligned}
 & \|\nabla \bar{Q}_{k,\mu}(x_k + s_k) - \nabla_s \bar{T}_{k,\mu}(x_k, s_k)\| \\
 481 \quad (3.10) \quad & \leq \|\nabla f(x_k + s_k) - \nabla f(x_k)\| + \sum_{i \in \mathcal{O}_{k,\mu}} \|\nabla(|v_i^T x|^p)|_{x=x_k+s_k} - \nabla_s m_i(x_k, s_k)\|. \\
 & 482
 \end{aligned}$$

483 On the one hand, the Lipschitz continuity of ∇f ensures

$$484 \quad (3.11) \quad \|\nabla f(x_k + s_k) - \nabla f(x_k)\| \leq L_f \|s_k\|.$$

485 On the other hand, it follows from [11, Lemma 5.2] that

$$\begin{aligned}
 & \sum_{i \in \mathcal{O}_{k,\mu}} \|\nabla(|v_i^T x|^p)|_{x=x_k+s_k} - \nabla_s m_i(x_k, s_k)\| \leq \frac{p(2-p)}{1-p} \mu^{p-2} |v_i^T s_k| \\
 486 \quad (3.12) \quad & \leq \frac{p(2-p)}{1-p} \mu^{p-2} \|s_k\|. \\
 & 487 \\
 & 488
 \end{aligned}$$

489 Hence, plugging (3.11) and (3.12) into (3.10) leads to the conclusion. \square

490 To proceed, we assume the following assumption holds.

491 **ASSUMPTION 3.2.** *For problem (1.1), it holds that*

$$492 \quad 0 \in v_i^T \mathcal{F}, \quad \text{Proj}_{\ker(v_i^T)} \mathcal{F} \subseteq \mathcal{F}, \quad i = 1, \dots, \bar{n}.$$

493 Under Assumption 3.2, it is easy to check that for any $x \in \mathcal{F}$, $(I - v_i v_i^T)x \in \mathcal{F}$ due
 494 to $\|v_i\| = 1$, for any $i = 1, \dots, \bar{n}$. A simple example of \mathcal{F} satisfies Assumption 3.2 is
 495 that $\mathcal{F} = \{x | \underline{l} \leq Vx \leq \underline{u}\}$, where $-\underline{l}, \underline{u} \in \mathbb{R}_+^{\bar{n}}$.

496 We now set ω satisfying

$$497 \quad (3.13) \quad 0 < \omega < \min \left\{ 6^{\frac{1}{p-1}}, \left(\frac{p}{2(L_h L_c^0 + \chi + 2M_F/\eta)} \right)^{\frac{1}{1-p}} \right\}.$$

498 Next lemma characterizes properties of points that are close to singularity.

499 **LEMMA 3.4.** *Suppose $\epsilon < \omega$, $|v_i^T x_k| < \omega$ for some $i \in [\bar{n}]$ and $s_k \neq 0$. Then it*
 500 *holds that $|v_i^T(x_k + s_k)| \leq \epsilon$ or $|v_i^T(x_k + s_k)| \geq \omega$.*

501 *Proof.* It is straightforward to obtain the conclusion if $i \in [\bar{n}] \setminus \mathcal{A}_k^+$. We now
 502 assume by contradiction that $|v_i^T x_k| < \omega$ and

$$503 \quad (3.14) \quad |v_i^T(x_k + s_k)| \in (\epsilon, \omega) \quad \text{for some } i \in \mathcal{A}_k^+.$$

504 Besides, by (2.7) there exists $\delta_k \in (0, 1]$ such that $\psi_m^{\epsilon, \delta_k}(x_k, s_k) \leq p|v_i^T(x_k + s_k)|\delta_k$.
 505 As $v_i^T, i \in [\bar{n}]$ are orthogonal, by the definition of \mathcal{R}_k^+ and (2.8) we have $\mathcal{R}_{\{i\}} :=$
 506 $\text{span}\{v_i\} \subseteq \mathcal{R}_k^+$. Consider the following minimization problem:

$$507 \quad (3.15) \quad \min_{\substack{x_k + s_k + d \in \mathcal{F} \\ d \in \mathcal{R}_{\{i\}}, \|d\| \leq \delta_k}} q_k(d)$$

with

$$q_k(d) := h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k) + d^T \nabla_s m_0(x_k, s_k) + \sum_{i \in \mathcal{A}_k^+} d^T \nabla_s m_i(x_k, s_k).$$

508 It is worthy to note that $d = 0$ is a feasible point of (3.15). Then the optimal function
509 value of (3.15) must be nonpositive, thus

(3.16)

$$510 \quad \left| \min_{\substack{x_k + s_k + d \in \mathcal{F} \\ d \in \mathcal{R}_{\{i\}}, \|d\| \leq \delta_k}} q_k(d) \right| \leq \left| \min_{\substack{x_k + s_k + d \in \mathcal{F} \\ d \in \mathcal{R}_k^+, \|d\| \leq \delta_k}} q_k(d) \right| = \psi_m^{\varepsilon, \delta_k}(x_k, s_k) \leq p\delta_k |v_i^T(x_k + s_k)|.$$

511
512 Note that it follows from $\text{Proj}_{\ker(v_i^T)} \mathcal{F} \subseteq \mathcal{F}$ and $x_k + s_k \in \mathcal{F}$ that

$$513 \quad x_k + s_k - v_i^T(x_k + s_k)v_i = x_k + s_k - v_i v_i^T(x_k + s_k) = \text{Proj}_{\ker(v_i^T)}(x_k + s_k) \in \mathcal{F}.$$

514
515 Then by the convexity of \mathcal{F} and $\delta_k \in (0, 1]$ we obtain $x_k + s_k + d \in \mathcal{F}$, where
516 $d = -\delta_k(v_i^T(x_k + s_k))v_i$. Obviously, $d \in \text{span}\{v_i\} = \mathcal{R}_{\{i\}}$. And it follows from (3.14)
517 that $|v_i^T(x_k + s_k)| < \omega < 1$, thus $\|d\| = \delta_k |v_i^T(x_k + s_k)| < \delta_k$. Then d is a feasible
518 point of problem (3.15). Moreover, it holds that

$$519 \quad (3.17) \quad q_k(d) = -\delta_k(v_i^T(x_k + s_k))\bar{\mathcal{G}}_k,$$

520 where

$$521 \quad \bar{\mathcal{G}}_k = -\frac{h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k)}{\delta_k(v_i^T(x_k + s_k))} + v_i^T \left(\nabla_s m_0(x_k, s_k) + \sum_{i \in \mathcal{A}_k^+} \nabla_s m_i(x_k, s_k) \right).$$

522 We next derive a lower bound of $|\bar{\mathcal{G}}_k|$. By the definition of m_i and $s_i^k = v_i^T s_k$, we have

$$523 \quad m_i(x_k, s_k) = |v_i^T x_k|^p + p|v_i^T x_k|^{p-1} \cdot \begin{cases} s_i^k, & \text{if } v_i^T x_k > 0, v_i^T(x_k + s_k) > 0, \\ -2v_i^T x_k - s_i^k, & \text{if } v_i^T x_k > 0, v_i^T(x_k + s_k) < 0, \\ 2v_i^T x_k + s_i^k, & \text{if } v_i^T x_k < 0, v_i^T(x_k + s_k) > 0, \\ -s_i^k, & \text{if } v_i^T x_k < 0, v_i^T(x_k + s_k) < 0, \end{cases}$$

524 which implies from $|v_i^T x_k| < \omega$ that

(3.18)

$$525 \quad \text{sgn}(\nabla_{s_i} m_i(x_k, s_k)) = \text{sgn}(v_i^T(x_k + s_k)) \quad \text{and} \quad |\nabla_{s_i} m_i(x_k, s_k)| = p|v_i^T x_k|^{p-1} > p\omega^{p-1}.$$

526 As $\|x\| \leq M_F$ for any $x \in \mathcal{F}$, $\|s_k\| \leq 2M_F$. Then it indicates from (3.2) that

$$527 \quad \|\nabla_s m_0(x_k, s_k)\| = \|g_k + \frac{1}{\eta} s_k\| \leq \chi + \frac{2M_F}{\eta}.$$

528 It together with the Lipschitz continuity of h , $\|d\| = \delta_k |v_i^T(x_k + s_k)|$, (3.2), (3.18),

529 $v_i^T \sum_{i \in \mathcal{A}_k^+} \nabla_s m_i(x_k, s_k) = \nabla_{s_i} m_i(x_k, s_k)$ and $\omega < \left(\frac{p}{2(L_h L_c^0 + \chi + 2M_F/\eta)} \right)^{\frac{1}{1-p}}$ derives the

530 following lower bound:

$$531 \quad |\bar{\mathcal{G}}_k| = \left| \frac{1}{-\delta_k v_i^T(x_k + s_k)} (h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k)) + v_i^T \nabla_s m_0(x_k, s_k) \right. \\ 532 \quad \left. + \nabla_{s_i} m_i(x_k, s_k) \right| \\ 533 \quad \geq |\nabla_{s_i} m_i(x_k, s_k)| - \frac{1}{\delta_k |v_i^T(x_k + s_k)|} |h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k)| \\ 534 \quad \quad - |v_i^T \nabla_s m_0(x_k, s_k)| \\ 535 \quad \geq p|v_i^T x_k|^{p-1} - L_h L_c^0 - \|\nabla_s m_0(x_k, s_k)\| \\ 536 \quad (3.19) \quad > p\omega^{p-1} - \left(L_h L_c^0 + \chi + \frac{2M_F}{\eta} \right) \geq \frac{1}{2} p\omega^{p-1}. \\ 537$$

538 Furthermore, (3.18) indicates $\text{sgn}(\bar{\mathcal{G}}_k) = \text{sgn}(\nabla_{s_i} m_i(x_k, s_k)) = \text{sgn}(v_i^T(x_k + s_k))$, thus
 539 by (3.17), $q_k(d) = -\delta_k(v_i^T(x_k + s_k))\bar{\mathcal{G}}_k = -\delta_k|v_i^T(x_k + s_k)||\bar{\mathcal{G}}_k| < 0$.

540 We now denote by d^* the optimal solution of (3.15). Obviously, $d^* \neq 0$. As
 541 $d^* \in \mathcal{R}_{\{i\}}$, there exists $\alpha \in \mathbb{R}$ such that $d^* = \alpha v_i$, thus $\alpha \neq 0$ and $\|d^*\| = |\alpha|$. Then
 542 we obtain $q_k(d^*) = \alpha \mathcal{G}_k$, where

$$543 \quad \mathcal{G}_k = \frac{1}{\alpha}(h(c_k + J_k(s_k + d^*)) - h(c_k + J_k s_k)) + v_i^T(\nabla_s m_0(x_k, s_k) + \sum_{i \in \mathcal{A}_k^+} \nabla_s m_i(x_k, s_k)).$$

544 Again by the negativeness of the optimal function value of (3.15) it holds that

$$545 \quad (3.20) \quad \text{sgn}(\alpha) = -\text{sgn}(\mathcal{G}_k) \quad \text{and} \quad |q_k(d^*)| = |\alpha \mathcal{G}_k| = \|d^*\| |\mathcal{G}_k|.$$

546 Meanwhile, by the optimality of d^* we obtain

$$547 \quad (3.21) \quad \|d^*\| |\mathcal{G}_k| \geq \delta_k |v_i^T(x_k + s_k)| |\bar{\mathcal{G}}_k|.$$

548 We next derive a lower bound of $\|d^*\|$. From (3.19) it follows that

$$549 \quad \left| -\frac{h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k)}{\delta_k(v_i^T(x_k + s_k))} + v_i^T \nabla_s m_0(x_k, s_k) \right| \leq \frac{1}{2} |\nabla_{s_i} m_i(x_k, s_k)|.$$

550 Moreover, analogy to (3.19) we can obtain $|\mathcal{G}_k| \geq \frac{1}{2} p \omega^{p-1}$ and

$$551 \quad \left| \frac{h(c_k + J_k(s_k + d^*)) - h(c_k + J_k s_k)}{\alpha} + v_i^T \nabla_s m_0(x_k, s_k) \right| \leq \frac{1}{2} |\nabla_{s_i} m_i(x_k, s_k)|.$$

552 Then by definitions of \mathcal{G}_k and $\bar{\mathcal{G}}_k$, we have

$$553 \quad \frac{|\bar{\mathcal{G}}_k|}{|\mathcal{G}_k|} \geq \frac{|\nabla_{s_i} m_i(x_k, s_k)| - \left| -\frac{h(c_k + J_k(s_k + d)) - h(c_k + J_k s_k)}{\delta_k(v_i^T(x_k + s_k))} + v_i^T \nabla_s m_0(x_k, s_k) \right|}{\left| \frac{h(c_k + J_k(s_k + d^*)) - h(c_k + J_k s_k)}{\alpha} + v_i^T \nabla_s m_0(x_k, s_k) \right|}$$

$$554 \quad \geq \frac{\frac{1}{2} |\nabla_{s_i} m_i(x_k, s_k)|}{\frac{3}{2} |\nabla_{s_i} m_i(x_k, s_k)|} = \frac{1}{3},$$

556 which indicates from (3.21) that $\|d^*\| \geq \frac{1}{3} \delta_k |v_i^T(x_k + s_k)|$. Based on above inequality
 557 together with (3.16), (3.20) and $|\mathcal{G}_k| \geq \frac{1}{2} p \omega^{p-1}$ we obtain

$$558 \quad \frac{1}{6} p \omega^{p-1} \delta_k |v_i^T(x_k + s_k)| \leq |q_k(d^*)| \leq p \delta_k |v_i^T(x_k + s_k)|,$$

560 which, however, contradicts $\omega < 6^{\frac{1}{p-1}}$. Thus, the conclusion is proved by contradic-
 561 tion. \square

562 To analyze oracle complexity of Algorithm 2.1, we first introduce the following
 563 index sets:

$$564 \quad \mathcal{K}_u := \{k \in \mathcal{K} : x_k = x_{k+1}\}, \quad \mathcal{K}_\epsilon := \{k \in \mathcal{K} \setminus \mathcal{K}_u : \mathcal{A}_k \setminus \mathcal{A}_{k+1} \neq \emptyset\},$$

$$565 \quad \mathcal{K}_\omega := \{k \in \mathcal{K} \setminus \mathcal{K}_u : \|s_k\| \geq \frac{1}{4} \omega\}, \quad \mathcal{K}_\heartsuit := \mathcal{K} \setminus (\mathcal{K}_u \cup \mathcal{K}_\epsilon \cup \mathcal{K}_\omega).$$

567 Due to the monotonely non-increasing property of \mathcal{A}_k , it is easy to have

$$568 \quad (3.22) \quad |\mathcal{K}_\epsilon| \leq \bar{n}.$$

569 Since Algorithm 2.1 terminates when both k and $k - 1$ belong to \mathcal{K}_u , it must hold
570 that

$$571 \quad (3.23) \quad |\mathcal{K}_u| \leq |\mathcal{K} \setminus \mathcal{K}_u| + 2 \leq |\mathcal{K}_\heartsuit \cup \mathcal{K}_\omega| + \bar{n} + 2.$$

573 Define $\alpha = \frac{3}{4}\omega$, where ω satisfies (3.13). The following lemma shows properties
574 of \mathcal{A}_k and \mathcal{A}_{k+1} with $k \in \mathcal{K}_\heartsuit$ which are also discussed in [11].

575 LEMMA 3.5. *Suppose that $\epsilon < \alpha$. Then the following relations hold:*

$$576 \quad (3.24) \quad \mathcal{A}_k = \mathcal{A}_{k+1} = \mathcal{O}_{k,\alpha}, \quad k \in \mathcal{K}_\heartsuit,$$

577 where $\mathcal{O}_{k,\alpha}$ is defined in (3.6).

578 *Proof.* By (2.8) and the definition of \mathcal{K}_\heartsuit , it is easy to have $\mathcal{A}_k = \mathcal{A}_{k+1}$ for any
579 $k \in \mathcal{K}_\heartsuit$. For any $k \in \mathcal{K}_\heartsuit$, we partition \mathcal{A}_k into the following sets:

$$\begin{aligned} 580 \quad \mathcal{I}_{\heartsuit,k} &:= \{i \in \mathcal{A}_k : \min\{|v_i^T x_k|, |v_i^T(x_k + s_k)|\} \geq \alpha\}, \\ 581 \quad \mathcal{I}_{\diamond,k} &:= \{i \in \mathcal{A}_k : (|v_i^T x_k| \geq \omega, |v_i^T(x_k + s_k)| \in (\epsilon, \alpha)) \\ 582 \quad &\quad \text{or } (|v_i^T x_k| \in (\epsilon, \alpha), |v_i^T(x_k + s_k)| \geq \omega)\}, \\ 583 \quad \mathcal{I}_{\clubsuit,k} &:= \{i \in \mathcal{A}_k : |v_i^T x_k| \in (\epsilon, \omega) \text{ and } |v_i^T(x_k + s_k)| \in (\epsilon, \omega)\}. \end{aligned}$$

585 Note that for any $i \in \mathcal{I}_{\diamond,k}$,

$$586 \quad \|s_k\| \geq |v_i^T s_k| \geq \left| |v_i^T(x_k + s_k)| - |v_i^T x_k| \right| \geq \omega - \alpha = \frac{1}{4}\omega.$$

587 It then indicates $i \in \mathcal{K}_\omega$. Thus $\mathcal{I}_{\diamond,k} = \emptyset$. Meanwhile, it follows from Lemma 3.4 that
588 $\mathcal{I}_{\clubsuit,k} = \emptyset$. Thus, $\mathcal{A}_k = \mathcal{I}_{\heartsuit,k}$, namely, $\mathcal{A}_k = \{i : \min\{|v_i^T x_k|, |v_i^T(x_k + s_k)|\} \geq \alpha\}$,
589 $k \in \mathcal{K}_\heartsuit$. It then yields (3.24) by definition of $\mathcal{O}_{k,\alpha}$ in (3.6). \square

590 Motivated by Lemma 2.1, we suppose that δ_k , $k \in \mathcal{K}$ is uniformly lower bounded
591 by $\delta > 0$ which is independent of ϵ . Then by the boundedness of \mathcal{F} , there exists
592 $M > 0$ such that $\|s_k\| \leq \|x_{k+1}\| + \|x_k\| \leq 2M_F \leq M\delta \leq M\delta_k$ for any $k \in \mathcal{K}$. The
593 lemma below shows that when Algorithm 2.1 terminates at Step 5 with $s_k \neq 0$, the
594 output is an approximate first-order stationary point of (1.1), provided that input \bar{w}
595 and $\bar{\beta}$ in Algorithm 2.1 satisfy

$$596 \quad (3.25) \quad \bar{\beta} \leq \min \left\{ \frac{1}{3}\bar{w}, \frac{1 - \beta - \theta}{\max(L(\alpha) + 1/\eta + L_h L_c^1(M+1), \bar{L})} \right\}.$$

597 We would like to mention that (3.25) can ensure $\bar{\beta}\bar{L} < 1$, which meets the requirement
598 on \bar{L} in Assumption 3.1.

599 LEMMA 3.6. *Suppose that $\epsilon < \alpha$. If Algorithm 2.1 terminates at Step 5 with
600 $s_k \neq 0$ and $\bar{\beta}$ satisfies (3.25), then x_{k+1} is an (ϵ, δ) -approximate first-order stationary
601 point of (1.1).*

Proof. When Algorithm 2.1 terminates at Step 5 with $s_k \neq 0$ and $k \notin \mathcal{K}_\epsilon$, $k \notin \mathcal{K}_u$.
Besides, it follows from the algorithmic framework that

$$\|s_k\| + \|s_{k-1}\| \leq \bar{\beta}\epsilon \leq \frac{1}{3}\bar{w}\epsilon < \frac{1}{3}\alpha = \frac{1}{4}\omega,$$

602 which indicates $k \notin \mathcal{K}_\omega$, thus $k \in \mathcal{K}_\heartsuit$ and (3.24) holds. Recall that (2.7) holds with
603 $\delta = \delta_k$, for some $\delta_k \in (0, 1]$, i.e.

$$604 \quad \psi_m^{\epsilon, \delta_k}(x_k, s_k) \leq \min \left\{ \theta\epsilon, p \min_{i \in \mathcal{A}_{k+1}} |v_i^T x_{k+1}| \right\} \delta_k, \quad \text{for some } \delta_k \in (0, 1].$$

605 Note that by (2.1) and (3.24) as well as (3.7),

$$\begin{aligned}
 606 \quad & \psi_Q^{\epsilon, \delta_k}(x_{k+1}) \\
 607 \quad & = Q_\epsilon(x_{k+1}) - \min_{\substack{x_{k+1}+d \in \mathcal{F} \\ d \in \mathcal{R}_{k+1}, \|d\| \leq \delta_k}} T_{Q_\epsilon}(x_{k+1}, d) \\
 608 \quad & = h(c_{k+1}) \\
 609 \quad & \quad - \min_{\substack{x_{k+1}+d \in \mathcal{F} \\ d \in \mathcal{R}_{k+1}, \|d\| \leq \delta_k}} \left\{ h(c_{k+1} + J_{k+1}d) + d^T \nabla \left(f(x) + \sum_{i \in \mathcal{A}_{k+1}} |v_i^T x|^p \right) \Big|_{x=x_{k+1}} \right\} \\
 610 \quad (3.26) \quad & = h(c_{k+1}) - \min_{\substack{x_{k+1}+d \in \mathcal{F} \\ d \in \mathcal{R}_{k+1}, \|d\| \leq \delta_k}} \left\{ h(c_{k+1} + J_{k+1}d) + \nabla \bar{Q}_{k,\alpha}(x_{k+1})^T d \right\}. \\
 611 \quad &
 \end{aligned}$$

612 As the minimization problem in (3.26) is convex, it admits a global minimizer, which
 613 we still denote as d with a slight abuse of notation. Obviously, $\|d\| \leq \delta_k$. We next
 614 show by contradiction that $\psi_Q^{\epsilon, \delta_k}(x_{k+1}) \leq \epsilon \delta_k$. We now assume that it were not true.
 615 Then it holds that $\psi_Q^{\epsilon, \delta_k}(x_{k+1}) = h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) - \nabla \bar{Q}_{k,\alpha}(x_{k+1})^T d > \epsilon \delta_k$.
 616 It can further derive

$$\begin{aligned}
 617 \quad & \psi_Q^{\epsilon, \delta_k}(x_{k+1}) \\
 618 \quad & = -(\nabla \bar{Q}_{k,\alpha}(x_{k+1}))^T d + (\nabla_s \bar{T}_{k,\alpha}(x_k, s_k))^T d - (\nabla_s \bar{T}_{k,\alpha}(x_k, s_k))^T d \\
 619 \quad & \quad - \frac{1}{2\eta} [\nabla(\|s\|^2)|_{s=s_k}]^T d + \frac{1}{2\eta} [\nabla(\|s\|^2)|_{s=s_k}]^T d + h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) \\
 620 \quad & \leq \|\nabla \bar{Q}_{k,\alpha}(x_{k+1}) - \nabla_s \bar{T}_{k,\alpha}(x_k, s_k)\| \|d\| - \left[\nabla_s \left(\bar{T}_{k,\alpha}(x_k, s) + \frac{\|s\|^2}{2\eta} \right) \Big|_{s=s_k} \right]^T d \\
 621 \quad & \quad + \frac{1}{\eta} s_k^T d + h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) \\
 622 \quad & \leq \left(L(\alpha) + \frac{1}{\eta} + L_h L_c^1 (M+1) \right) \|s_k\| \delta_k + \|\nabla f(x_k) - g_k\| \delta_k + \theta \epsilon \delta_k, \\
 623 \quad &
 \end{aligned}$$

624 where the last inequality follows from $\|d\| \leq \delta_k$, (3.9), (3.24), and

$$\begin{aligned}
 625 \quad & h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) - \left[\nabla_s \left(\bar{T}_{k,\alpha}(x_k, s) + \frac{1}{2\eta} \|s\|^2 \right) \Big|_{s=s_k} \right]^T d \\
 626 \quad & = -(\nabla f(x_k) - g_k)^T d + h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) \\
 627 \quad & \quad - \left[\nabla_s m_0(x_k, s_k) + \sum_{i \in \mathcal{A}_k} \nabla_s m_i(x_k, s_k) \right]^T d \\
 628 \quad & \leq \|\nabla f(x_k) - g_k\| \|d\| + \max \left\{ 0, h(c_{k+1}) - h(c_{k+1} + J_{k+1}d) \right. \\
 629 \quad & \quad \left. - \left[\nabla_s m_0(x_k, s_k) + \sum_{i \in \mathcal{A}_{k+1}} \nabla_s m_i(x_k, s_k) \right]^T d \right\} \\
 630 \quad & \leq \|\nabla f(x_k) - g_k\| \|d\| + \psi_m^{\epsilon, \delta_k}(x_k, s_k) + |h(c_{k+1}) - h(c_k + J_k s_k)| \\
 631 \quad & \quad + |h(c_k + J_k(s_k + d)) - h(c_{k+1} + J_{k+1}d)| \\
 632 \quad & \leq \|\nabla f(x_k) - g_k\| \delta_k + \theta \epsilon \delta_k + L_h L_c^1 (M+1) \|s_k\| \delta_k \\
 633 \quad &
 \end{aligned}$$

634 due to $\psi_m^{\epsilon, \delta_k}(x_k, s_k) \leq \theta \epsilon \delta_k$,

$$635 \quad |h(c_{k+1}) - h(c_k + J_k s_k)| \leq \frac{L_h L_c^1}{2} \|s_k\|^2 \leq \frac{L_h L_c^1 M}{2} \delta_k \|s_k\|$$

636 and

$$\begin{aligned}
637 & |h(c_k + J_k(s_k + d)) - h(c_{k+1} + J_{k+1}d)| \\
638 & \leq L_h \|c_k + J_k(s_k + d) - c_{k+1} - J_{k+1}d\| \\
639 & \leq L_h \|c_k + J_k s_k - c_{k+1}\| + L_h \|J_k - J_{k+1}\| \|d\| \\
640 & \leq \frac{L_h L_c^1}{2} \|s_k\|^2 + L_h L_c^1 \|s_k\| \|d\| \\
641 & \leq L_h L_c^1 \left(\frac{M}{2} + 1 \right) \delta_k \|s_k\|. \\
642 &
\end{aligned}$$

643 Then it follows from $\psi_Q^{\epsilon, \delta_k}(x_{k+1}) > \epsilon \delta_k$ and Assumption 3.1 with $\beta < 1 - \theta$ that

$$\begin{aligned}
644 & \epsilon \delta_k < \left(L(\alpha) + \frac{1}{\eta} + L_h L_c^1 (M + 1) \right) \|s_k\| \delta_k + \bar{L} \|s_{k-1}\| \delta_k + (\beta + \theta) \epsilon \delta_k \\
645 &
\end{aligned}$$

646 which implies

$$\begin{aligned}
647 & (1 - \beta - \theta) \epsilon < \max \left\{ L(\alpha) + \frac{1}{\eta} + L_h L_c^1 (M + 1), \bar{L} \right\} (\|s_k\| + \|s_{k-1}\|). \\
648 &
\end{aligned}$$

649 However, this contradicts $\|s_k\| + \|s_{k-1}\| \leq \bar{\beta} \epsilon$ by the setting of $\bar{\beta}$. Therefore, x_{k+1} is
650 an (ϵ, δ) -approximate first-order stationary point of (1.1). \square

651 *Remark 3.7.* Lemmas 3.1 and 3.6 show that Algorithm 2.1 can always return an
652 approximate first-order stationary point of (1.1) when it terminates.

653 We now partition \mathcal{K}_\heartsuit into $\mathcal{K}_\heartsuit^1 \cup \mathcal{K}_\heartsuit^2$, where

$$\begin{aligned}
654 & \mathcal{K}_\heartsuit^1 := \{k \in \mathcal{K}_\heartsuit : \|s_k\| + \|s_{k-1}\| \geq \bar{\beta} \epsilon\}, \quad \mathcal{K}_\heartsuit^2 := \{k \in \mathcal{K}_\heartsuit : \|s_k\| + \|s_{k-1}\| < \bar{\beta} \epsilon\}.
\end{aligned}$$

655 By the definition of \mathcal{K}_\heartsuit , Lemma 3.6 and termination conditions of Algorithm 2.1, we
656 know that $|\mathcal{K}_\heartsuit^2| \leq 1$, thus $|\mathcal{K}_\heartsuit| \leq |\mathcal{K}_\heartsuit^1| + 1$. Then it together with (3.22) and (3.23)
657 implies that the total number of iterations until Algorithm 2.1 terminates satisfies

$$\begin{aligned}
658 & |\mathcal{K}| \leq |\mathcal{K}_u| + |\mathcal{K}_\heartsuit \cup \mathcal{K}_\omega| + |\mathcal{K}_\epsilon| \leq |\mathcal{K}_\heartsuit \cup \mathcal{K}_\omega| + \bar{n} + 2 + |\mathcal{K}_\heartsuit \cup \mathcal{K}_\omega| + \bar{n} \\
\text{\textcircled{3.27}} & \leq 2|\mathcal{K}_\heartsuit^1 \cup \mathcal{K}_\omega| + 2\bar{n} + 4.
\end{aligned}$$

661 Based on above relations, to estimate the upper bound of $|\mathcal{K}|$, it suffices to derive an
662 upper bound on $|\mathcal{K}_\heartsuit^1 \cup \mathcal{K}_\omega|$. Inspired by this, we establish the oracle complexity of
663 Algorithm 2.1 in the theorem below. In the following we assume that the positive
664 parameter η in (2.6) satisfies

$$\begin{aligned}
665 & (3.28) \quad \frac{1}{16\eta} - \frac{L_f + L_h L_c^1}{16} - \beta \bar{L} - \frac{\beta}{\bar{\beta}} \geq 1, \quad \frac{1}{4\eta} - \frac{L_f + L_h L_c^1}{4} - \beta \bar{L} - 3\beta \geq 1.
\end{aligned}$$

666 It is noteworthy that the setting of $\bar{\beta}$ in (3.25) together with (3.28) and Assumption
667 3.1 ensures the existence of desired input parameters $\bar{\omega}, \bar{\beta}, \eta, \bar{L}$ and β . We now proceed
668 under such parameter settings.

669 **THEOREM 3.8.** *Suppose that $\epsilon < \alpha$. Then there exists a positive constant $C =$
670 $\mathcal{O}(1)$ such that $\sum_{k \in \mathcal{K}} \|s_k\|^2 \leq C$. Furthermore, the maximum iteration number until
671 Algorithm 2.1 terminates is in order of $\mathcal{O}(\epsilon^{-2})$.*

672 *Proof.* It follows from Lemma 3.2, $s_{-1} = 0$ and Assumption 3.1 that

$$\begin{aligned}
 673 \quad & \left(\frac{1}{2\eta} - \frac{L_f + L_h L_c^1}{2} \right) \sum_{k \in \mathcal{K}} \|s_k\|^2 \\
 674 \quad & \leq \sum_{k \in \mathcal{K}} \|\nabla f(x_k) - g_k\| \|s_k\| + Q_\epsilon(x_0) - Q_\epsilon^* \\
 675 \quad & \leq \sum_{k \in \mathcal{K}} \beta \max\{\bar{L} \min\{\|s_{k-1}\|, D\}, \epsilon\} \|s_k\| + Q_\epsilon(x_0) - Q_\epsilon^* \\
 676 \quad & \leq \sum_{k \in \mathcal{K}} \beta (\bar{L} \|s_{k-1}\| \|s_k\| + \epsilon \|s_k\|) + Q_\epsilon(x_0) - Q_\epsilon^* \\
 677 \quad & \leq \sum_{k \in \mathcal{K}} \beta \left(\frac{\bar{L}}{2} (\|s_k\|^2 + \|s_{k-1}\|^2) + \epsilon \|s_k\| \right) + Q_\epsilon(x_0) - Q_\epsilon^* \\
 678 \quad (3.29) \quad & \leq \sum_{k \in \mathcal{K}} \beta (\bar{L} \|s_k\|^2 + \epsilon \|s_k\|) + Q_\epsilon(x_0) - Q_\epsilon^*.
 \end{aligned}$$

679 As

$$681 \quad (3.30) \quad \|s_k\| \geq \frac{1}{3}\alpha \geq \frac{1}{3}\epsilon, \quad k \in \mathcal{K}_\omega,$$

682 it indicates

$$683 \quad (3.31) \quad \bar{L} \|s_k\|^2 + \epsilon \|s_k\| \leq (\bar{L} + 3) \|s_k\|^2, \quad k \in \mathcal{K}_\omega.$$

684 Moreover, by definition of \mathcal{K}_\heartsuit^1 we have

$$685 \quad (3.32) \quad \|s_k\| + \|s_{k-1}\| \geq \bar{\beta}\epsilon, \quad k \in \mathcal{K}_\heartsuit^1,$$

687 thus

$$\begin{aligned}
 688 \quad & \bar{L} \|s_k\|^2 + \epsilon \|s_k\| \leq \bar{L} \|s_k\|^2 + \bar{\beta}^{-1} (\|s_k\| + \|s_{k-1}\|) \|s_k\| \\
 689 \quad (3.33) \quad & \leq (\bar{L} + \bar{\beta}^{-1}) (\|s_k\| + \|s_{k-1}\|) \|s_k\|, \quad k \in \mathcal{K}_\heartsuit^1.
 \end{aligned}$$

691 Since $s_k = 0$ for any $k \in \mathcal{K}_\omega$, plugging (3.31) and (3.33) into (3.29) yields

$$\begin{aligned}
 692 \quad & \left(\frac{1}{2\eta} - \frac{L_f + L_h L_c^1}{2} \right) \sum_{k \in \mathcal{K}} \|s_k\|^2 \\
 693 \quad & \leq \sum_{k \in \mathcal{K}_\heartsuit^1} \beta (\bar{L} + \bar{\beta}^{-1}) (\|s_k\| + \|s_{k-1}\|) \|s_k\| + \sum_{k \in \mathcal{K}_\omega} \beta (\bar{L} + 3) \|s_k\|^2 \\
 694 \quad (3.34) \quad & + \sum_{k \in \mathcal{K}_\epsilon \cup \mathcal{K}_\heartsuit^2} \beta (\bar{L} \|s_k\|^2 + \epsilon \|s_k\|) + Q_\epsilon(x_0) - Q_\epsilon^*.
 \end{aligned}$$

696 Recall that $\|s_k\| < \bar{\beta}\epsilon < 1$ for any $k \in \mathcal{K}_\heartsuit^2$. Besides, by the boundedness of \mathcal{F} we have
 697 $\|s_k\| \leq 2M_F$ for any $k \in \mathcal{K}_\epsilon$. Then it follows from (3.34) that

$$\begin{aligned}
 698 \quad & \left(\frac{1}{2\eta} - \frac{L_f + L_h L_c^1}{2} \right) \sum_{k \in \mathcal{K}} \|s_k\|^2 \\
 699 \quad & \leq \sum_{k \in \mathcal{K}_\heartsuit^1} \beta (\bar{L} + \bar{\beta}^{-1}) (\|s_k\| + \|s_{k-1}\|)^2 + \sum_{k \in \mathcal{K}_\omega} \beta (\bar{L} + 3) \|s_k\|^2 \\
 700 \quad (3.35) \quad & + \bar{n}\beta (4\bar{L}M_F^2 + 2\epsilon M_F) + \beta (\bar{L} + \epsilon) + Q_\epsilon(x_0) - Q_\epsilon^*,
 \end{aligned}$$

702 where the last term in above inequality uses the facts that $|\mathcal{K}_\epsilon| \leq \bar{n}$ and $|\mathcal{K}_\heartsuit^2| \leq 1$.
 703 Notice that

$$\begin{aligned}
 704 \quad \sum_{k \in \mathcal{K}} \|s_k\|^2 &= \frac{1}{2} \left(\sum_{k \in \mathcal{K}} \|s_k\|^2 + \sum_{k \in \mathcal{K}} \|s_k\|^2 \right) \\
 705 \quad &\geq \frac{1}{4} \sum_{k \in \mathcal{K}} (\|s_k\|^2 + \|s_{k-1}\|^2) + \frac{1}{2} \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2 \\
 706 \quad &\geq \frac{1}{8} \sum_{k \in \mathcal{K}} (\|s_k\| + \|s_{k-1}\|)^2 + \frac{1}{2} \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2 \\
 707 \quad &\geq \frac{1}{8} \sum_{k \in \mathcal{K}_\heartsuit^1} (\|s_k\| + \|s_{k-1}\|)^2 + \frac{1}{2} \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2, \\
 708
 \end{aligned}$$

709 which further derives

$$\begin{aligned}
 710 \quad &\left(\frac{1}{2\eta} - \frac{L_f + L_h L_c^1}{2} \right) \sum_{k \in \mathcal{K}} \|s_k\|^2 \\
 711 \quad &\geq \left(\frac{1}{16\eta} - \frac{L_f + L_h L_c^1}{16} \right) \sum_{k \in \mathcal{K}_\heartsuit^1} (\|s_k\| + \|s_{k-1}\|)^2 + \left(\frac{1}{4\eta} - \frac{L_f + L_h L_c^1}{4} \right) \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2. \\
 712
 \end{aligned}$$

713 Then it together with (3.35) and the boundedness of \mathcal{F} implies that

$$\begin{aligned}
 714 \quad &\left(\frac{1}{16\eta} - \frac{L_f + L_h L_c^1}{16} - \beta \bar{L} - \frac{\beta}{\bar{\beta}} \right) \sum_{k \in \mathcal{K}_\heartsuit^1} (\|s_k\| + \|s_{k-1}\|)^2 \\
 715 \quad &+ \left(\frac{1}{4\eta} - \frac{L_f + L_h L_c^1}{4} - \beta \bar{L} - 3\beta \right) \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2 \leq \bar{\Gamma} \\
 716
 \end{aligned}$$

717 with $\bar{\Gamma} = \bar{n}\beta(4\bar{L}M_F^2 + 2\epsilon M_F) + \beta(\bar{L} + \epsilon) + Q_\epsilon(x_0) - Q_\epsilon^*$. Furthermore, from the setting
 718 of η as in (3.28) we attain

$$719 \quad (3.36) \quad \sum_{k \in \mathcal{K}_\heartsuit^1} (\|s_k\| + \|s_{k-1}\|)^2 + \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2 \leq \bar{\Gamma}$$

720 which leads to the conclusion from (3.35) with $C = \frac{\bar{\Gamma}(1 + \beta\bar{L} + \beta \max\{\bar{\beta}^{-1}, 3\})}{1/(2\eta) - (L_f + L_h L_c^1)/2}$. Obviously,
 721 $C = \mathcal{O}(1)$.

722 Moreover, by (3.30) and (3.32) we obtain

$$723 \quad \sum_{k \in \mathcal{K}_\heartsuit^1} (\|s_k\| + \|s_{k-1}\|)^2 + \sum_{k \in \mathcal{K}_\omega} \|s_k\|^2 \geq \bar{\beta}^2 \epsilon^2 |\mathcal{K}_\heartsuit^1| + \frac{1}{9} \epsilon^2 |\mathcal{K}_\omega|.$$

724 Then it together with (3.36) implies $|\mathcal{K}_\heartsuit^1| + |\mathcal{K}_\omega| = \mathcal{O}(\epsilon^{-2})$. Hence by (3.27) the
 725 maximum iteration number until the termination of Algorithm 2.1 is in order $\mathcal{O}(\epsilon^{-2})$. \square

726 Since only one inexact gradient is evaluated at each iteration, the oracle complex-
 727 ity of Algorithm 2.1 is in order $\mathcal{O}(\epsilon^{-2})$.

728 **4. Stochastic variant.** For problem (1.1), when f owns a finite-sum structure
 729 (1.2), as the sample size N can be very large, it will be expensive to go through all
 730 component functions to compute exact gradients, thereby only approximate gradients
 731 are available. To cope with this type of problems, we propose a stochastic variant of
 732 Algorithm 2.1. The proposed algorithm follows the main framework of Algorithm 2.1,
 733 with **InexactOracle** specified in Algorithm 4.1. Here inexact gradients are computed
 734 by calling stochastic first-order oracles in a recursive way [23] and l is a positive integer.

Algorithm 4.1 InexactOracle($x_k, x_{k-1}, g_{k-1}, k, l$)

Input: Index set \mathcal{I}_k generated uniformly at random without replacement from
 $\{1, \dots, N\}$.
 1: **if** $\text{mod}(k, l) = 0$ **then**
 2: Compute $g_k = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \nabla f_i(x_k)$.
 3: **else**
 4: Compute $g_k = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} (\nabla f_i(x_k) - \nabla f_i(x_{k-1})) + g_{k-1}$.
 5: **end if**

735
 736 Note that Assumption 3.1 plays a key role in supervising the analysis in previous
 737 section. In this section, adopting a proper sampling strategy we can guarantee Assumption 3.1 with high probability. We then establish the complexity of the proposed algorithm, in terms of number of stochastic first-order oracles, to find an approximate first-order stationary point. To proceed the analysis, we first introduce a lemma regarding the concentration inequality under sampling without replacement. As this lemma is a duplicate of [27, Theorem 4], we omit its proof here.

743 **LEMMA 4.1.** *Let $\mathcal{X} = \{X_i \in \mathbb{R}^n, i = 1, \dots, N\}$. Suppose $\|X_i\| \leq \sigma$ for all $i =$
 744 $1, \dots, N$ and some $\sigma > 0$. Denote $\lambda = \frac{1}{N} \sum_{i=1}^N X_i$. Let $A_1, \dots, A_\nu, \nu < N$ be
 745 samples from \mathcal{X} under the sampling without replacement. Then, for any $\epsilon > 0$, the
 746 following bound holds:*

$$747 \quad \text{Prob}\left(\left\|\frac{1}{\nu} \sum_{i=1}^{\nu} A_i - \lambda\right\| \geq \epsilon\right) \leq 2(n+1) \exp\left(-\frac{\nu \epsilon^2}{8\sigma^2 \left(1 + \frac{1}{\nu}\right) \left(1 - \frac{\nu}{N}\right)}\right).$$

748 Given $\zeta \in (0, 1)$, following Lemma 4.1, we can achieve

$$749 \quad \text{Prob}\left(\left\|\frac{1}{\nu} \sum_{i=1}^{\nu} A_i - \lambda\right\| \leq \epsilon\right) \geq 1 - \zeta, \text{ if } \nu \geq \left[\frac{1}{N} + \frac{\epsilon^2}{16\sigma^2 \log(2(n+1)/\zeta)}\right]^{-1}.$$

750 We assume that $\nabla f_i, i = 1, \dots, N$ are Lipschitz continuously differentiable. With a
 751 slight abuse of notations, we still use L_f and κ to denote the Lipschitz constant and
 752 upper bound of $\nabla f_i, i = 1, \dots, N$ over \mathcal{F} . Then for any k with $\text{mod}(k, l) = 0$, g_k
 753 generated by Algorithm 4.1 satisfies

$$754 \quad \text{Prob}(\|g_k - \nabla f(x_k)\| \leq \beta\epsilon) \geq 1 - \zeta, \text{ if } |\mathcal{I}_k| \geq \left[\frac{1}{N} + \frac{\beta^2 \epsilon^2}{16\kappa^2 \log(2(n+1)/\zeta)}\right]^{-1}.$$

755 For those k with $\text{mod}(k, l) \neq 0$, the lemma below provides a sampling strategy such
 756 that Assumption 3.1 can be satisfied with high probability.

757 LEMMA 4.2. *Under sampling without replacement, for any k with $\text{mod}(k, l) \neq 0$,*
 758 *g_k generated by Algorithm 4.1 satisfies Assumption 3.1 with probability at least $1 - \zeta$,*
 759 *provided that*

$$760 \quad (4.1) \quad |\mathcal{I}_j| \geq \begin{cases} \left[\frac{1}{N} + \frac{\beta^2 \epsilon^2 / l^2}{256 L_f^2 \|x_j - x_{j-1}\|^2 \log(4(n+1)l/\zeta)} \right]^{-1}, & j = k, k-1, \dots, \lfloor k/l \rfloor l + 1, \\ \left[\frac{1}{N} + \frac{\beta^2 \epsilon^2}{256 \kappa^2 \log(4(n+1)/\zeta)} \right]^{-1}, & j = \lfloor k/l \rfloor l. \end{cases}$$

761 *Proof.* For any k with $\text{mod}(k, l) \neq 0$, it follows from the algorithmic framework
 762 that

$$763 \quad g_k - \nabla f(x_k) \\ 764 \quad = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} [\nabla f_i(x_k) - \nabla f_i(x_{k-1})] + g_{k-1} - \nabla f(x_k) \\ 765 \quad = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} [\nabla f_i(x_k) - \nabla f_i(x_{k-1}) - (\nabla f(x_k) - \nabla f(x_{k-1}))] + g_{k-1} - \nabla f(x_{k-1}).$$

767 We thus obtain $g_k - \nabla f(x_k) = \sum_{j=\lfloor k/l \rfloor l}^k Y_j$, where $Y_j := \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} Z_{j,i}$ with

$$768 \quad Z_{j,i} := \begin{cases} \nabla f_i(x_j) - \nabla f_i(x_{j-1}) - (\nabla f(x_j) - \nabla f(x_{j-1})), & j = k, k-1, \dots, \lfloor k/l \rfloor l + 1, \\ \nabla f_i(x_j) - \nabla f(x_j), & j = \lfloor k/l \rfloor l \end{cases}$$

770 for $i = 1, \dots, N$. Define $\bar{\epsilon}_k := \beta \max\{\bar{L} \min(\|s_{k-1}\|, D), \epsilon\}$ with \bar{L} , D and β in
 771 Assumption 3.1. For notation simplicity we denote by B_j the event $\|Y_j\| \leq \frac{\bar{\epsilon}_k}{2(k-\lfloor k/l \rfloor l)}$
 772 with $j = k, \dots, \lfloor k/l \rfloor l + 1$, and by B_j the event $\|Y_j\| \leq \frac{\bar{\epsilon}_k}{2}$ with $j = \lfloor k/l \rfloor l$. We use
 773 \bar{B}_j to denote the complement of B_j . Then

$$774 \quad \text{Prob}(\|g_k - \nabla f(x_k)\| \leq \bar{\epsilon}_k) \geq \text{Prob}\left(\bigcap_{j=\lfloor k/l \rfloor l}^k B_j\right) = 1 - \text{Prob}\left(\bigcup_{j=\lfloor k/l \rfloor l}^k \bar{B}_j\right)$$

776 which is no less than $1 - \sum_{j=\lfloor k/l \rfloor l}^k \text{Prob}(\bar{B}_j)$ by the union bound. Hence, to achieve
 777 that (3.1) holds with probability at least ζ , it suffices to require

$$778 \quad (4.2) \quad \text{Prob}(\bar{B}_j) = \begin{cases} \text{Prob}\left(\|Y_j\| > \frac{\bar{\epsilon}_k}{2(k-\lfloor k/l \rfloor l)}\right) \leq \frac{\zeta}{2(k-\lfloor k/l \rfloor l)}, & j = k, \dots, \lfloor k/l \rfloor l + 1, \\ \text{Prob}\left(\|Y_j\| > \frac{\bar{\epsilon}_k}{2}\right) \leq \frac{\zeta}{2}, & j = \lfloor k/l \rfloor l. \end{cases}$$

779 Due to the smoothness of f_i , $\|Z_{j,i}\| \leq 2L_f \|x_j - x_{j-1}\|$, $j = k, k-1, \dots, \lfloor k/l \rfloor l + 1$ and
 780 $\|Z_{\lfloor k/l \rfloor l, i}\| \leq 2\kappa$, $i = 1, \dots, N$. As $\sum_{i=1}^N Z_{j,i} = 0$ for any $j = k, \dots, \lfloor k/l \rfloor l$, by Lemma
 781 4.1 with $\lambda = 0$, $\nu = |\mathcal{I}_j|$ and $A_{i'} = Z_{j,i'}$, $i' = 1, \dots, \nu$, $i' \in [N]$, we obtain that (4.2)
 782 can be achieved provided that

$$783 \quad |\mathcal{I}_j| \geq \begin{cases} \left[\frac{1}{N} + \frac{\beta^2 \max\{\bar{L}^2 \min(\|s_{k-1}\|^2, D^2), \epsilon^2\}}{64 L_f^2 \|x_j - x_{j-1}\|^2 \log(4(n+1)(k-\lfloor k/l \rfloor l)/\zeta)} \right]^{-1}, & j = k, \dots, \lfloor k/l \rfloor l + 1, \\ \left[\frac{1}{N} + \frac{\beta^2 \max\{\bar{L}^2 \min(\|s_{k-1}\|^2, D^2), \epsilon^2\}}{64 \kappa^2 \log(4(n+1)/\zeta)} \right]^{-1}, & j = \lfloor k/l \rfloor l, \end{cases}$$

785 which can be guaranteed by (4.1) due to $\bar{\epsilon}_k \geq \beta \epsilon$ and $k - \lfloor k/l \rfloor l \leq l$. \square

786 We are now ready to present the oracle complexity in terms of total number of
 787 stochastic first-order oracles required to guarantee that Algorithm 2.1 can find an
 788 (ϵ, δ) -approximate first-order stationary point of (1.1)-(1.2).

789 **THEOREM 4.3.** *Suppose that conditions of Theorem 3.8 and Lemma 4.2 with $l =$
 790 $\mathcal{O}(N^{1/3})$ hold, and Algorithm 2.1 with Algorithm 4.1 called to compute inexact oracles
 791 terminates in finite iterations. Then for given $\rho \in (0, 1)$, with probability at least $1 - \rho$,
 792 it returns an (ϵ, δ) -approximate first-order stationary point of (1.1)-(1.2) with the
 793 oracle complexity in order $\mathcal{O}(N + N^{\frac{2}{3}}\epsilon^{-2} \log(\frac{4(n+1)N^{\frac{1}{3}}}{\epsilon^2\rho}))$. Consequently, the oracle
 794 complexity of Algorithm 2.1 with Algorithm 4.1 is in order $\tilde{\mathcal{O}}(\epsilon^{-2})$.*

795 *Proof.* We still use \mathcal{K} to denote the set of all iteration indices until termination. As
 796 can be seen from previous section, to make sure the algorithm returns an approximate
 797 stationary point with probability at least $1 - \rho$, it suffices to guarantee with probability
 798 at least $1 - \rho$ that Assumption 3.1 holds for all iterations in \mathcal{K} . To realize this,
 799 Assumption 3.1 should be satisfied at each one of the iterations with probability at
 800 least $1 - \zeta$ for some $\zeta \in [0, 1]$ such that $1 - |\mathcal{K}|\zeta \geq 1 - \rho$. We may simply set $\zeta = \frac{\rho}{|\mathcal{K}|}$.
 801 Furthermore, to achieve Assumption 3.1 with probability at least $1 - \zeta$ at j th iteration
 802 for any given $j \in \mathcal{K}$, by Lemma 4.2 the size of \mathcal{I}_j can be equal to the right side of
 803 (4.1) after rounding up. With above settings, it holds with probability at least $1 - \rho$
 804 that $|\mathcal{K}| = \mathcal{O}(\epsilon^{-2})$ and $\sum_{j \in \mathcal{K}} \|s_j\|^2 \leq C$, where $C = \mathcal{O}(1)$ by Theorem 3.8. Hence, to
 805 reach an (ϵ, δ) -approximate first-order stationary point with probability at least $1 - \rho$,
 806 the total number of stochastic first-order oracles is bounded by

$$\begin{aligned}
 807 \quad \sum_{i \in \mathcal{K}} |\mathcal{I}_i| &= \sum_{i: \text{mod}(i, l) = 0} |\mathcal{I}_i| + \sum_{i = \lfloor |\mathcal{K}|/l \rfloor l + 1}^{|\mathcal{K}|} |\mathcal{I}_i| + \sum_{i=0}^{\lfloor |\mathcal{K}|/l \rfloor - 1} \sum_{j=1}^{l-1} |\mathcal{I}_{il+j}| \\
 808 \quad &\leq \lceil \frac{|\mathcal{K}|}{l} \rceil N + \sum_{i = \lfloor |\mathcal{K}|/l \rfloor l + 1}^{|\mathcal{K}|} \left[\frac{1}{N} + \frac{\beta^2 \epsilon^2 / l^2}{256 L_f^2 \|x_i - x_{i-1}\|^2 \log(4(n+1)l/\zeta)} \right]^{-1} \\
 809 \quad &+ \sum_{i=0}^{\lfloor |\mathcal{K}|/l \rfloor - 1} \sum_{j=1}^{l-1} \left[\frac{1}{N} + \frac{\beta^2 \epsilon^2 / l^2}{256 L_f^2 \|x_{il+j} - x_{il+j-1}\|^2 \log(4(n+1)l/\zeta)} \right]^{-1} + |\mathcal{K}| \\
 810 \quad &\leq \lceil \frac{|\mathcal{K}|}{l} \rceil N + \sum_{i = \lfloor |\mathcal{K}|/l \rfloor l + 1}^{|\mathcal{K}|} \frac{256 L_f^2 \|x_i - x_{i-1}\|^2 \log(4(n+1)l/\zeta)}{\beta^2 \epsilon^2 / l^2} \\
 811 \quad &+ \sum_{i=0}^{\lfloor |\mathcal{K}|/l \rfloor - 1} \sum_{j=1}^{l-1} \frac{256 L_f^2 \|x_{il+j} - x_{il+j-1}\|^2 \log(4(n+1)l/\zeta)}{\beta^2 \epsilon^2 / l^2} + |\mathcal{K}| \\
 812 \quad &\leq \lceil \frac{|\mathcal{K}|}{l} \rceil N + 256 C l^2 L_f^2 \frac{1}{\beta^2 \epsilon^2} \log(4(n+1)l/\zeta) + |\mathcal{K}| \\
 813
 \end{aligned}$$

814 which derives the oracle complexity order by the setting of l . □

815 **5. Extension to expectation case.** In this section, we focus on solving the
 816 problem with f in the expectation form, given by:

$$817 \quad (5.1) \quad \min_{x \in \mathcal{F}} Q(x) := f(x) + h(c(x)) + \|Vx\|_p^p \quad \text{with} \quad f(x) := \mathbb{E}[\mathbf{F}(x, \xi)].$$

818 Here, $\xi \in \Xi$ represents a random variable following the probability function \mathbb{P} , and
 819 $\mathbf{F} : \mathbb{R}^n \times \Xi \rightarrow \mathbb{R}$ is continuously differentiable with respect to $x \in \mathcal{F}$ for almost

every $\xi \in \Xi$. To address the challenges posed by problems in the expectation form, where the sample set can be infinite, we propose a modification to Algorithm 4.1 by randomly generating a subset of samples from Ξ , presented in Algorithm 5.1.

Algorithm 5.1 InexactOracle($x_k, x_{k-1}, g_{k-1}, k, l$)

Input: Generate a sample subset ξ_k uniformly at random from Ξ .

- 1: **if** $\text{mod}(k, l) = 0$ **then**
 - 2: Compute $g_k = \frac{1}{|\xi_k|} \sum_{\xi \in \xi_k} \nabla_x \mathbf{F}(x_k, \xi)$.
 - 3: **else**
 - 4: Compute $g_k = \frac{1}{|\xi_k|} \sum_{\xi \in \xi_k} (\nabla_x \mathbf{F}(x_k, \xi) - \nabla_x \mathbf{F}(x_{k-1}, \xi)) + g_{k-1}$.
 - 5: **end if**
-

The aim of this section is to investigate the oracle complexity of Algorithm 2.1 with Algorithm 5.1 called to compute stochastic first-order oracles. Before delving into the analysis, we introduce an assumption that stochastic oracles satisfy.

ASSUMPTION 5.1. *There exist $\Delta, L_f > 0$ such that for all $x \in \mathcal{F}$,*

$$\mathbb{E}[\nabla_x \mathbf{F}(x, \xi)] = \nabla f(x), \quad \|\nabla_x \mathbf{F}(x, \xi) - \nabla f(x)\| \leq \Delta \text{ almost surely,}$$

and for any $x, y \in \mathcal{F}$, $\|\nabla_x \mathbf{F}(x, \xi) - \nabla_x \mathbf{F}(y, \xi)\| \leq L_f \|x - y\|$ almost surely.

The following lemma presents the matrix Bernstein inequality [26].

LEMMA 5.1. *Let X_1, \dots, X_ν be i.i.d. random vectors in \mathbb{R}^n , and satisfy $\mathbb{E}[X_i] = 0$ and $\|X_i\| \leq \sigma$ almost surely for some $\sigma > 0$ and any $i = 1, \dots, \nu$. Define $M := \max(\|\sum_{i=1}^\nu \mathbb{E}[X_i X_i^T]\|, \|\sum_{i=1}^\nu \mathbb{E}[X_i^T X_i]\|)$. Then for any $t \geq 0$,*

$$\text{Prob} \left(\left\| \sum_{i=1}^\nu X_i \right\| \geq t \right) \leq (n+1) \cdot \exp \left(\frac{-t^2/2}{M + \sigma t/3} \right).$$

Note that $M \leq \sum_{i=1}^\nu \mathbb{E}[\|X_i\|^2] \leq \nu \sigma^2$. By Lemma 5.1, we obtain that for any $\epsilon > 0$,

$$\text{Prob} \left(\frac{1}{\nu} \left\| \sum_{i=1}^\nu X_i \right\| \geq \epsilon \right) \leq (n+1) \cdot \exp \left(\frac{-\nu \epsilon^2/2}{\sigma^2 + \sigma \epsilon/3} \right).$$

Then for g_k , generated by Algorithm 5.1 with k s.t. $\text{mod}(k, l) = 0$, under Assumption 5.1 and by Lemma 5.1 we attain

$$\text{Prob}(\|g_k - \nabla f(x_k)\| \leq \beta \epsilon) \geq 1 - \zeta, \text{ if } |\xi_k| \geq \left(\frac{2\Delta^2}{\beta^2 \epsilon^2} + \frac{2\Delta}{3\beta \epsilon} \right) \log \left(\frac{n+1}{\zeta} \right).$$

For those k with $\text{mod}(k, l) \neq 0$, similar to Lemma 4.2, we can provide a sampling strategy such that Assumption 3.1 holds with high probability.

LEMMA 5.2. *Let g_k be generated by Algorithm 5.1. For any k with $\text{mod}(k, l) \neq 0$, Assumption 3.1 holds at k th iteration with probability at least $1 - \zeta$, provided that*

$$(5.2) \quad |\xi_j| \geq \begin{cases} \left(\frac{32L_f^2 \|x_j - x_{j-1}\|^2 l^2}{\beta^2 \epsilon^2} + \frac{8L_f \|x_j - x_{j-1}\| l}{3\beta \epsilon} \right) \log \left(\frac{2(n+1)l}{\zeta} \right), & j = k, \dots, \lfloor k/l \rfloor l + 1, \\ \left(\frac{8\Delta^2}{\beta^2 \epsilon^2} + \frac{4\Delta}{3\beta \epsilon} \right) \log \left(\frac{2(n+1)}{\zeta} \right), & j = \lfloor k/l \rfloor l. \end{cases}$$

846 *Proof.* By the computation of g_k in Algorithm 5.1 and $Y_j := \frac{1}{|\xi_j|} \sum_{\xi \in \xi_j} Z_j(\xi)$
 847 with

$$848 \quad Z_j(\xi) = \begin{cases} \nabla_x \mathbf{F}(x_j, \xi) - \nabla_x \mathbf{F}(x_{j-1}, \xi) - \nabla f(x_j) + \nabla f(x_{j-1}), & j = k, \dots, \lfloor k/l \rfloor l + 1, \\ \nabla_x \mathbf{F}(x_j, \xi) - \nabla f(x_j), & j = \lfloor k/l \rfloor l, \end{cases}$$

850 we obtain $g_k - \nabla f(x_k) = \sum_{j=\lfloor k/l \rfloor l}^k Y_j$. Under Assumption 5.1 and due to the sm-
 851oothness of f , $\|Z_j(\xi)\| \leq 2L_f \|x_j - x_{j-1}\|$, $j = k, \dots, \lfloor k/l \rfloor l + 1$ and $\|Z_{\lfloor k/l \rfloor l}(\xi)\| \leq \Delta$.
 852 Similar to the analysis of Lemma 5.2, the remainder is to ensure (4.2). It follows from
 853 Lemma 5.1 that to achieve (4.2) it suffices to require

$$854 \quad |\xi_j| \geq \begin{cases} \left(\frac{32L_f^2 \|x_j - x_{j-1}\|^2 (k - \lfloor k/l \rfloor l)^2}{\bar{\epsilon}_k^2} + \frac{8L_f \|x_j - x_{j-1}\| (k - \lfloor k/l \rfloor l)}{3\bar{\epsilon}_k} \right) \log \left(\frac{2(n+1)(k - \lfloor k/l \rfloor l)}{\zeta} \right), & j = k, \dots, \lfloor k/l \rfloor l + 1, \\ \left(\frac{8\Delta^2}{\bar{\epsilon}_k^2} + \frac{4\Delta}{3\bar{\epsilon}_k} \right) \log \left(\frac{2(n+1)}{\zeta} \right), & j = \lfloor k/l \rfloor l, \end{cases}$$

856 which can be guaranteed by (5.2) and $\bar{\epsilon}_k \geq \beta\epsilon$. \square

857 We slightly abuse the notation and continue to use \mathcal{K} to represent all the iteration
 858 indices until Algorithm 2.1 terminates, with Algorithm 5.1 being called to compute
 859 inexact oracles. According to Lemma 5.2, in order to achieve an (ϵ, δ) -approximate
 860 first-order stationary point with a probability at least $1 - \rho$, where $\rho \in (0, 1)$, Assump-
 861 tion 3.1 must hold at each iteration with probability at least $1 - \zeta$ for $\zeta \in (0, 1)$ such
 862 that $1 - |\mathcal{K}|\zeta \geq 1 - \rho$. Therefore, we set $\zeta = \frac{\rho}{|\mathcal{K}|}$. Consequently, by applying Theorem
 863 3.8 and setting $l = \mathcal{O}(|\mathcal{K}|^{1/3})$ we can conclude that the total number of stochastic
 864 first-order oracles is bounded by:

$$\begin{aligned} 865 \quad \sum_{i \in \mathcal{K}} |\xi_i| &= \sum_{i: \text{mod}(i, l) = 0} |\xi_i| + \sum_{i = \lfloor |\mathcal{K}|/l \rfloor l + 1}^{|\mathcal{K}|} |\xi_i| + \sum_{i=0}^{\lfloor |\mathcal{K}|/l \rfloor - 1} \sum_{j=1}^{l-1} |\xi_{il+j}| \\ 866 \quad &\leq \lceil \frac{|\mathcal{K}|}{l} \rceil \left(\frac{8\Delta^2}{\beta^2 \epsilon^2} + \frac{4\Delta}{3\beta\epsilon} \right) \log \left(\frac{2(n+1)}{\zeta} \right) \\ 867 \quad &\quad + \sum_{i = \lfloor |\mathcal{K}|/l \rfloor l + 1}^{|\mathcal{K}|} \left(\frac{32L_f^2 \|x_j - x_{j-1}\|^2 l^2}{\beta^2 \epsilon^2} + \frac{8L_f \|x_j - x_{j-1}\| l}{3\beta\epsilon} \right) \log \left(\frac{2(n+1)l}{\zeta} \right) \\ 868 \quad &\quad + \sum_{i=0}^{\lfloor |\mathcal{K}|/l \rfloor - 1} \sum_{j=1}^{l-1} \left(\frac{32L_f^2 \|x_j - x_{j-1}\|^2 l^2}{\beta^2 \epsilon^2} + \frac{8L_f \|x_j - x_{j-1}\| l}{3\beta\epsilon} \right) \log \left(\frac{2(n+1)l}{\zeta} \right) + |\mathcal{K}| \\ 869 \quad &= \mathcal{O} \left(\frac{|\mathcal{K}|}{l} \left(\frac{1}{\epsilon^2} + \frac{1}{\epsilon} \right) + \frac{l^2}{\epsilon^2} + |\mathcal{K}|^{1/2} \frac{l}{\epsilon} \right) \log \left(\frac{2(n+1)l}{\zeta} \right) + |\mathcal{K}| \\ 870 \quad &= \mathcal{O} \left(\epsilon^{-10/3} \log \left(\frac{1}{\rho\epsilon} \right) \right). \end{aligned}$$

872 We summarize above analysis into the following theorem.

873 **THEOREM 5.3.** *Suppose that conditions of Theorem 3.8 and Lemma 5.2 hold, with*
 874 *$l = \mathcal{O}(|\mathcal{K}|^{1/3})$, and Algorithm 2.1 with Algorithm 5.1 called to compute inexact oracles*
 875 *terminates in finite iterations. Then for given $\rho \in (0, 1)$, with probability at least $1 - \rho$,*
 876 *the algorithm returns an (ϵ, δ) -approximate first-order stationary point of (5.1) with*
 877 *the oracle complexity in order $\mathcal{O}(\epsilon^{-10/3} \log(1/(\rho\epsilon)))$, i.e., $\tilde{\mathcal{O}}(\epsilon^{-10/3})$.*

878 **6. Numerical simulation.** In this section, we consider the problem

$$879 \quad (6.1) \quad \min_{x \in \mathcal{F}} f(x) + \|Vx\|_p^p, \quad \text{s.t.} \quad Bx \leq b,$$

880 where $\mathcal{F} = \{x \in \mathbb{R}^n : b_l \leq x \leq b_u\}$, $B \in \mathbb{R}^{r \times n}$, $b \in \mathbb{R}^r$, and $f(x) = \frac{1}{N} \sum_{i=1}^N ((A_i x -$
881 $c_i)_+)^2$ with $A_i^T \in \mathbb{R}^n$ and $c_i \in \mathbb{R}$. By penalizing the constraints of (6.1) with τ being
882 a penalty parameter, we obtain the penalty approximation problem in the form of
883 (1.1)-(1.2):

$$884 \quad (6.2) \quad \min_{x \in [b_l, b_u]} \frac{1}{N} \sum_{i=1}^N ((A_i x - c_i)_+)^2 + \tau \|(Bx - b)_+\|_1 + \|Vx\|_p^p.$$

885 We apply Algorithm 2.1 to solve (6.2) by calling Algorithm 4.1 at k th iteration to
886 compute inexact first-order oracle g_k , $k \geq 0$. Following (2.6), the subproblem at k th
887 iteration is defined as

$$888 \quad \min_{s \in \mathbb{R}^n} g_k^T s + \tau \|(Bx_k + Bs - b)_+\|_1 + \sum_{i \in \mathcal{A}_k} p |v_i^T x_k|^{p-1} |v_i^T (x_k + s)| + \frac{1}{2} \|s\|^2$$

$$889 \quad \text{s.t.} \quad b_l \leq x_k + s \leq b_u, \quad v_i^T s = 0, \quad i \notin \mathcal{A}_k.$$

891 By introducing $\bar{z} = (Bx_k + Bs - b)_+ \in \mathbb{R}^r$, and $\hat{z} = (\hat{z}_i, i \in \mathcal{A}_k)^T \in \mathbb{R}^{|\mathcal{A}_k|}$ with
892 $\hat{z}_i = |v_i^T (x_k + s)|$, $i \in \mathcal{A}_k$, we obtain the following linearly constrained quadratic
893 program:

$$894 \quad \min_{s, \bar{z}, \hat{z}} g_k^T s + \tau e^T \bar{z} + \sum_{i \in \mathcal{A}_k} p |v_i^T x_k|^{p-1} \hat{z}_i + \frac{1}{2} \|s\|^2$$

$$895 \quad \text{s.t.} \quad b_l \leq x_k + s \leq b_u, \quad v_i^T s = 0, \quad i \notin \mathcal{A}_k,$$

$$896 \quad 0 \leq \bar{z}, \quad Bx_k + Bs - b \leq \bar{z}, \quad -\hat{z}_i \leq v_i^T (x_k + s) \leq \hat{z}_i, \quad i \in \mathcal{A}_k.$$

898 The numerical implementation was conducted in MATLAB R2022a on a PC
899 with Intel I7-12700H 2.3GHZ CPU processor, 16GB RAM memory and a Windows
900 operating system. We use Matlab default solver **quadprog** to solve each quadratic
901 program. We generate the optimal solution x^* with $\|x^*\|_0 = K$ and set V , b_l , b_u , B ,
902 b , A , c as follows.

```
903 IndexK = randperm(n);  $x_0$  = randn(n, 1);  $x^*$  = zeros(n, 1);
904  $x^*(\text{IndexK}(1 : K)) = 2 * (\text{randn}(K, 1) > 0.5) - 1$ ;  $V = 0.1 * \text{eye}(n)$ ;
905  $b_l = -100 * \text{ones}(n, 1)$ ;  $b_u = 100 * \text{ones}(n, 1)$ ;  $B = \text{rand}(n, n)$ ;  $B = \text{orth}(B')$ ;
906  $b = B * x^*$ ,  $A = \text{randn}(N, n)$ ;  $c = \max(A * x^* + 0.01 * \text{randn}(N, 1), 0)$ ;
```

908 In particular, we set parameters $n = 100$, $N = 10^5$, $K = 10$, $\epsilon = 10^{-4}$, $\bar{\beta} = 0.2$, $\tau =$
909 200 , $\eta = 0.01$, $l = 10$ and the batch size as 1000. In Figure 1, we report the perfor-
910 mances of the proposed algorithm. Specifically, Figures 1(a)-(d) showcase the behav-
911 ior of different metrics, including the function value error $f(x_k) - f(x^*)$, the relative
912 error between the iterate and x^* given by $\frac{\|x_k - x^*\|}{\|x_k\|}$, the number of nonzero entries in
913 the iterate denoted as $\|x_k\|_0$, and the comparison between the nonzero entries of the
914 output and x^* , respectively.

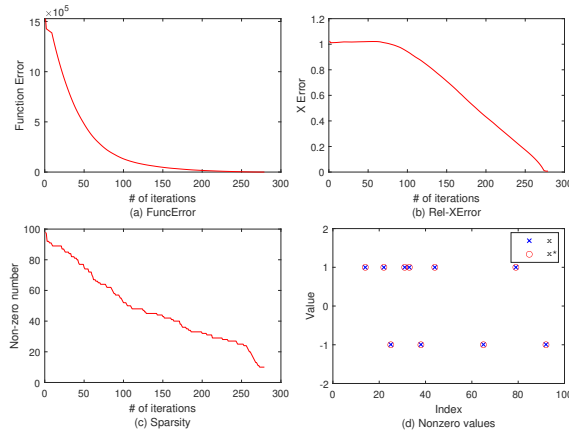


FIG. 1. Numerical profiles on test problem (6.2)

915 **7. Conclusions.** We present complexity analysis of proximal inexact gradient
 916 methods for finite-sum optimization with nonsmooth composite functions and a non-
 917 Lipschitz regularizer (1.1). Existence of the nonsmooth function h and non-Lipschitz
 918 term makes it inadequate to build an approximation model simply based on Taylor
 919 expansion as in [5, 8, 11, 12]. Moreover, those algorithms in [5, 8, 11, 12] rely on exact
 920 function values and gradients of f , which have difficulties in computation of problem
 921 (1.1) with the large scale finite-sum of f . In our Algorithm 2.1, we solve a strongly
 922 convex proximal subproblem (2.6) at each iteration without computing the function
 923 values and exact gradients of f , based on convex approximation to $f(x) + h(c(x))$ and
 924 a Lipschitz continuous approximation to $\|Vx\|_p^p$. By controlling inexactness of inexact
 925 gradients as well as subproblem solutions, we establish $\mathcal{O}(\epsilon^{-2})$ oracle complexity to
 926 find an (ϵ, δ) -approximate first-order stationary point of problem (1.1). This verifies
 927 that the worst-case oracle complexity still keeps the same with the absence of the
 928 differentiability of the Lipschitz term compared to [11, 12] and with the existence
 929 of non-Lipschitz regularizer in contrast to [5, 8]. Moreover, we propose a stochastic
 930 variant of Algorithm 2.1, by calling stochastic first-order oracles in a recursive way
 931 and applying a proper sampling strategy. We establish that the oracle complexity is
 932 in order $\tilde{\mathcal{O}}(\epsilon^{-2})$ to find an (ϵ, δ) -approximate first-order stationary point with high
 933 probability. We further extend the stochastic variant of algorithm to solve problems
 934 in the expectation form and derive the oracle complexity in order $\tilde{\mathcal{O}}(\epsilon^{-10/3})$ with high
 935 probability. Numerical performances of the proposed algorithm are also reported on
 936 a test problem.

937

REFERENCES

938 [1] W. BIAN AND X. CHEN, *Linearly constrained non-Lipschitz optimization for image restoration*,
 939 SIAM J. Imaging Sci., 8 (2015), pp. 2294–2322.
 940 [2] W. BIAN AND X. CHEN, *Optimality and complexity for constrained optimization problems with*
 941 *nonconvex regularization*, Math. Oper. Res., 42 (2017), pp. 1063–1084.
 942 [3] W. BIAN, X. CHEN, AND Y. YE, *Complexity analysis of interior point algorithms for non-*
 943 *Lipschitz and nonconvex minimization*, Math. Program., 149 (2015), pp. 301–327.
 944 [4] K. BUI, F. PARK, S. ZHANG, Y. QI, AND J. XIN, *Structured sparsity of convolutional neural*
 945 *networks via nonconvex sparse group regularization*, Front. Appl. Math. Stat., (2021),
 946 <https://doi.org/10.3389/fams.2020.529564>.
 947 [5] C. CARTIS, N. GOULD, AND P. TOINT, *On the evaluation complexity of composite function*

- 948 *minimization with applications to nonconvex nonlinear programming*, SIAM J. Optim., 21
 949 (2011), pp. 1721–1739.
- 950 [6] C. CARTIS, N. GOULD, AND P. TOINT, *Sharp worst-case evaluation complexity bounds for*
 951 *arbitrary-order nonconvex optimization with inexpensive constraints*, SIAM J. Optim., 30
 952 (2020), pp. 513–541.
- 953 [7] C. CARTIS, N. GOULD, AND P. TOINT, *Strong evaluation complexity of an inex-*
 954 *act trust-region algorithm for arbitrary-order unconstrained nonconvex optimization*,
 955 <https://doi.org/10.48550/arXiv.2011.00854>, (2021).
- 956 [8] C. CARTIS, N. GOULD, AND P. TOINT, *Strong evaluation complexity bounds for*
 957 *arbitrary-order optimization of nonconvex nonsmooth composite functions*, Proceed-
 958 *ings of the International Congress of Mathematicians (ICM 2022)*. EMS Press.
 959 <https://doi.org/10.4171/ICM2022/95>, (2022).
- 960 [9] X. CHEN, D. GE, Z. WANG, AND Y. YE, *Complexity of unconstrained ℓ_2 - ℓ_p minimization*,
 961 *Math. Program.*, 143 (2014), pp. 371–383.
- 962 [10] X. CHEN, Z. LU, AND T. PONG, *Penalty methods for a class of non-Lipschitz optimization*
 963 *problems*, SIAM J. Optim., 26 (2016), pp. 1465–1492.
- 964 [11] X. CHEN AND P. TOINT, *High-order evaluation complexity for convexly-constrained optimiza-*
 965 *tion with non-Lipschitzian group sparsity terms*, *Math. Program.*, 187 (2020), pp. 47–78.
- 966 [12] X. CHEN, P. TOINT, AND H. WANG, *Complexity of partially-separable convexly-constrained*
 967 *optimization with non-Lipschitzian singularities*, SIAM J. Optim., 29 (2019), pp. 874–903.
- 968 [13] X. CHEN, F. XU, AND Y. YE, *Lower bound theory of nonzero entries in solutions of ℓ_2 - ℓ_p*
 969 *minimization*, SIAM J. Sci. Comput., 32 (2010), pp. 2832–2852.
- 970 [14] W. CHENG, X. WANG, AND X. CHEN, *An interior stochastic gradient method for a class of*
 971 *non-Lipschitz optimization problems*, *J. Sci. Comput.*, 92 (2022).
- 972 [15] Y. CUI, Z. HE, AND J.-S. PANG, *Multicomposite nonconvex optimization for training deep*
 973 *neural networks*, SIAM J. Optim., 30 (2020), pp. 1693–1723.
- 974 [16] D. GE, R. HE, AND S. HE, *An improved algorithm for the L_2 - L_p minimization problem*, *Math.*
 975 *Program.*, 166 (2017), pp. 131–158.
- 976 [17] S. GRATTON, E. SIMON, AND P. TOINT, *An algorithm for the minimization of nonsmooth non-*
 977 *convex functions using inexact evaluations and its worst-case complexity*, *Math. Program.*,
 978 187 (2021), pp. 1–24.
- 979 [18] A. JAIN, *Fundamentals of Digital Image Processing*, Upper Saddle River, NJ: Prentice-Hall,
 980 1989.
- 981 [19] D. LI, Z. SUN, AND X. ZHANG, *A constrained optimization reformulation and a feasible descent*
 982 *direction method for $L_{1/2}$ regularization*, *Comput. Optim. Appl.*, 59 (2014), pp. 263–284.
- 983 [20] Y. LI, S. GU, AND R. T. C. MAYER, L. V. GOOL, *Group sparsity: The hinge between filter*
 984 *pruning and decomposition for network compression*, CVPR, (2020).
- 985 [21] Y. LIU, S. MA, Y. DAI, AND S. ZHANG, *A smoothing SQP framework for a class of composite*
 986 *L_q minimization over polyhedron*, *Math. Program.*, 158 (2016), pp. 467–500.
- 987 [22] M. METEL AND A. TAKEDA, *Simple stochastic gradient methods for non-smooth non-convex*
 988 *regularized optimization*, ICML, (2019), pp. 4537–4545.
- 989 [23] L. NGUYEN, J. LIU, K. SCHEINBERG, AND M. TAKAC, *SARAH: a novel method for machine*
 990 *learning problems using stochastic recursive gradient*, ICML, (2017), pp. 2613–2621.
- 991 [24] R. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Springer Verlag, Heidelberg,
 992 Berlin, New York, 1998.
- 993 [25] S. SCARDAPANE, D. COMMINELO, A. HUSSAIN, AND A. UNCINI, *Group sparse regularization*
 994 *for deep neural networks*, *Neurocomputing*, 241 (2017), pp. 81–89.
- 995 [26] J. A. TROPP, *User-friendly tail bounds for sums of random matrices*, *Found Comput Math*,
 996 12 (2012), pp. 389–434.
- 997 [27] Z. WANG, Y. ZHOU, Y. LIANG, AND G. LAN, *Stochastic variance-reduced cubic regularization*
 998 *for nonconvex optimization*, Proceedings of the Twenty-Second International Conference
 999 on Artificial Intelligence and Statistics, PMLR, 89 (2019), pp. 2731–2740.
- 1000 [28] Y. XU, R. JIN, AND T. YANG, *Non-asymptotic analysis of stochastic methods for non-smooth*
 1001 *non-convex regularized problems*, in NeurIPS, H. Wallach, H. Larochelle, A. Beygelzimer,
 1002 F. Alché-Buc, E. Fox, and R. Garnett, eds., vol. 32, Curran Associates, Inc., 2019.
- 1003 [29] Y. XU, Q. QI, Q. LIN, R. JIN, AND T. YANG, *Stochastic optimization for DC functions and non-*
 1004 *smooth non-convex regularizers with non-asymptotic convergence*, ICML, (2019), pp. 6942–
 1005 6951.
- 1006 [30] Z. XU, X. CHANG, F. XU, AND H. ZHANG, *$L_{1/2}$ regularization: A thresholding representation*
 1007 *theory and a fast solver*, *IEEE T. Neur. Net. Lear.*, 23 (2012), pp. 1013–1027.
- 1008 [31] J. YOON AND S. HWANG, *Combined group and exclusive sparsity for deep neural networks*,
 1009 ICML, (2017), pp. 3958–3966.