

UNIFIED ACCELERATION OF HIGH-ORDER ALGORITHMS UNDER GENERAL HÖLDER CONTINUITY*

CHAOBING SONG[†], YONG JIANG[‡], AND YI MA[§]

Abstract. In this paper, through an intuitive *vanilla proximal method* perspective, we derive a concise *unified acceleration framework (UAF)* for minimizing a convex function that has Hölder continuous derivatives with respect to general (non-Euclidean) norms. The UAF reconciles the two different high-order acceleration approaches, one by Nesterov and Baes [27, 3, 31] and one by Monteiro and Svaiter [23]. As a result, the UAF unifies the high-order acceleration instances [27, 3, 31, 15, 16, 23, 18, 6, 14] of the two approaches by only two problem-related parameters and two additional parameters for framework design. Meanwhile, the UAF (and its analysis) is *the first approach* to make high-order methods applicable for high-order smoothness conditions with respect to non-Euclidean norms. Furthermore, the UAF is the first approach that can match the existing lower bound of the iteration complexity for minimizing a convex function with Hölder continuous derivatives [16]. For practical implementation, we introduce a new and effective heuristic that significantly simplifies the binary search procedure required by the framework. We use experiments to verify the effectiveness of the heuristic and demonstrate clear and consistent advantages of high-order acceleration methods over first-order ones, in terms of run-time complexity. Finally, the UAF is proposed directly in the general composite convex setting, thus show that the existing high-order algorithms [27, 3, 31, 16, 6, 14] can be naturally extended to the general composite convex setting.

Key words. High-order algorithms, Nesterov’s acceleration, proximal method, non-Euclidean norm.

AMS subject classifications. 49M15, 49M37, 65K05, 68Q25, 90C25, 90C30

1. Introduction. In optimization, people often consider the problem of minimizing a convex function:

$$(1.1) \quad \min_{x \in \mathbb{R}^d} f(x).$$

A typical assumption is that $f(x)$ has L -Lipschitz continuous gradients with respect to (*w.r.t.*) the Euclidean norm $\|\cdot\|_2$, *i.e.*,

$$(1.2) \quad \|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2,$$

where $L > 0$ is the Lipschitz constant. For this problem, to find an ϵ -accurate solution x such that $f(x) - f(x^*) \leq \epsilon$, the classic gradient descent method:

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

with $\eta \leq 1/L$ takes $O(\epsilon^{-1})$ iterations. Nevertheless, it is known that from [28], for a convex function $f(x)$ with L -Lipschitz continuous gradients, a lower-bound for the number of iterations for any first-order algorithms is known to be

$$(1.3) \quad O(\epsilon^{-1/2}), \quad (L\text{-Lipschitz continuous gradients}).$$

In the seminal work [26], Nesterov has introduced an acceleration technique, the so-called *accelerated gradient descent (AGD)* algorithm, that achieves this optimal lower bound. This algorithm dramatically improves the convergence rate of smooth convex optimization with negligible per-iteration cost. Besides the smooth convex problem (1.1) under the Euclidean norm setting (1.2), AGD can also be generalized to solve the composite convex problem [5, 1, 13], in which the objective function may contain a second possibly non-smooth but simple convex term (see (2.6)). Meanwhile, AGD can be extended to the more general (non-Euclidean) norm settings [12, 1, 21], also achieving the optimal rate (1.3).

*Submitted to the editors on October 09, 2019.

Funding: This work was done during Chaobing Song’s visit to Professor Yi Ma’s group at UC Berkeley. The work is partially supported by the TBSI program and EECS Startup fund of Professor Yi Ma.

[†]Tsinghua-Berkeley Shenzhen Institute (TBSI), Tsinghua University (songcb16@mails.tsinghua.edu.cn).

[‡]Tsinghua-Berkeley Shenzhen Institute (TBSI), Tsinghua University (jiangy@sz.tsinghua.edu.cn)

[§]EECS Department, University of California, Berkeley (yima@eecs.berkeley.edu).

38 **1.1. High-order Acceleration Methods with Lipschitz Continuity.** To hope for a bet-
 39 ter iteration complexity beyond $O(\epsilon^{-1/2})$, $f(x)$ needs to be smooth for its high-order deriva-
 40 tives. A common assumption is that $f(x)$ has (p, ν, L) -Hölder continuous derivatives:

$$41 \quad (1.4) \quad \frac{1}{(p-1)!} \|\nabla^p f(x) - \nabla^p f(y)\|_2 \leq L \|x - y\|_2^\nu,$$

42 for some $\nu \in [0, 1]$, $p \in \mathbb{Z}_+$. Notice that for $p = 1$ and $\nu = 1$, this condition becomes the first
 43 order L -Lipschitz continuous gradient (1.2) above. Here, for $p \geq 2$, the $\|\cdot\|_2$ norm of a p -th
 44 order tensor denotes its operator norm [31] *w.r.t.* the vector 2-norm $\|\cdot\|_2$. Sometimes, when
 45 $\nu = 1$, the function is said to have (p, L) -Lipschitz continuous derivatives:

$$46 \quad (1.5) \quad \frac{1}{(p-1)!} \|\nabla^p f(x) - \nabla^p f(y)\|_2 \leq L \|x - y\|_2.$$

47 If a convex function $f(x)$ satisfies (1.5), the recent work [2] has given a lower-bound on the
 48 complexity: any deterministic algorithm would need at least

$$49 \quad (1.6) \quad O(\epsilon^{-\frac{2}{3p+1}}), \quad ((p, L)\text{-Lipschitz continuous derivatives})$$

50 iterations to find an ϵ -accurate solution. For the special case $p = 2$, [27] has proposed an
 51 “*accelerated cubic regularized Newton method*” (ACNM) that achieves an iteration complexity
 52 of $O(\epsilon^{-\frac{1}{3}})$. From a different approach, after proposing an accelerated hybrid proximal
 53 extragradient (A-HPE) framework, [23] has implemented an “*accelerated Newton proximal*
 54 *extragradient*” (A-NPE) instance of the A-HPE framework that has achieved the optimal
 55 complexity $O(\epsilon^{-\frac{2}{3}})$ ¹ for $p = 2$, although each iteration requires a nontrivial binary search
 56 procedure.

57 To achieve better complexity results and also being encouraged by the fact that third-order
 58 methods can often be implemented as efficiently as second-order methods [31], there is an
 59 increasing interest to extend ACNM and implement the A-HPE framework to even higher-order
 60 smoothness settings ($p \in \{3, 4, \dots\}$) [3, 31, 18, 14, 6]. In particular, by extending ACNM,
 61 [3] and [31] have proposed accelerated tensor methods with $O(\epsilon^{-\frac{1}{p+1}})$ iteration complexity
 62 for $p \in \{2, 3, \dots\}$. Meanwhile, by implementing A-HPE, [23, 18, 14, 6] have proposed
 63 accelerated methods that achieve the optimal $O(\epsilon^{-\frac{2}{3p+1}})$ iteration complexity, although just
 64 like A-NPE, all these methods need the nontrivial binary search procedure.

65 Hence the current situation seems to be: methods [27, 3, 31, 15, 16] by extending ACNM
 66 have advantages with simpler implementation, while methods [23, 18, 14, 6] by implementing
 67 A-HPE can in theory achieve the optimal rate $O(\epsilon^{-\frac{2}{3p+1}})$. However, it remains somewhat
 68 mysterious how we could reconcile the differences between these two approaches. In addition,
 69 the A-HPE framework is somewhat abstract so implementing it in the high-order setting
 70 requires rather nontrivial techniques [23, 18, 14, 6]. It remains unclear how to propose
 71 a concise but equivalently powerful alternative to the A-HPE framework and obtain these
 72 different instances of A-HPE in a unified way. Furthermore, although AGD can be generalized
 73 to general non-Euclidean norm settings, up to now, it is not known whether high-order
 74 methods can have a similar generalization. Finally, both the ACNM and A-HPE approaches
 75 do not directly address the composite convex setting (see (2.6)) at the framework level,
 76 hence obtaining high-order algorithms in this setting is highly desired and seems nontrivial
 77 [23, 15, 18].

¹When talking about iteration complexity, we mean the complexity in terms of outer iteration without concerning about the inner implementation of subproblems.

78 **1.2. Acceleration under Hölder Continuity and Our Results.** Besides the Lipschitz
 79 continuous setting, the more general Hölder continuous setting (1.4) is also of increased
 80 interest, partly for designing universal optimization schemes [29, 37, 16, 10]. If $f(x)$ has
 81 $(1, \nu, L)$ -Hölder continuous gradients, a lower bound for the iteration complexity is known to
 82 be [25]:

$$83 \quad (1.7) \quad O\left(\epsilon^{-\frac{2}{1+3\nu}}\right), \quad ((1, \nu, L)\text{-Hölder continuous gradients}).$$

84 An algorithm that can achieve this lower bound has been proposed in [24].

85 For the more general setting of (p, ν, L) -Hölder continuous derivatives, during the prepara-
 86 tion of this paper, [16] has given a lower bound of iteration complexity

$$87 \quad (1.8) \quad O\left(\epsilon^{-\frac{2}{3(p+\nu)-2}}\right), \quad ((p, \nu, L)\text{-Hölder continuous derivatives}).$$

88 By extending Nesterov's method in [31], [16] has proposed a method that achieves the iteration
 89 complexity $O(\epsilon^{-\frac{1}{p+\nu}})$. To the best of our knowledge, methods that can achieve the lower
 90 bound $O(\epsilon^{-\frac{2}{3(p+\nu)-2}})$ are still unknown.

91 In this paper, for the minimization of convex functions with (p, ν, L) -Hölder continuous
 92 derivatives, we propose a *unified acceleration framework* (UAF), see Algorithm 5.1, that
 93 achieves the iteration complexity of $O(\epsilon^{-\frac{2}{3(p+\nu)-2}})$ ($p \in \{1, 2, \dots\}, \nu \in [0, 1]$ with $p + \nu \geq 2$,
 94 $L > 0$), which matches the lower bound [16]. To be more precise, if a convex function $f(x)$
 95 has (p, ν, L) -Hölder continuous derivatives, our algorithm can find an ϵ -accurate solution with

$$96 \quad (1.9) \quad O\left(\epsilon^{-\frac{q}{(q+1)(p+\nu)-q}}\right)$$

97 iterations, where q is a tunable parameter² such that $2 \leq q \leq p + \nu$. Notice that our result and
 98 algorithm unify previously known results as (important) special cases:

- 99 • For the case of L -Lipschitz continuous gradients [28] where $p = \nu = 1$ and $q = 2$,
 100 the rate (1.9) of the proposed algorithm achieves the lower bound $O(\epsilon^{-\frac{1}{2}})$ of (1.3).
 - 101 • For the more general setting of (p, ν, L) -Hölder continuous derivatives: when $p \in$
 102 $\{2, 3, \dots\}, q = p + \nu$, it recovers the complexity $O(\epsilon^{-\frac{1}{p+\nu}})$ of the method in [16].
- 103 Meanwhile, by setting $q = 2$, the rate (1.9) of the UAF is *the first* convergence result that
 104 matches the lower bound $O(\epsilon^{-\frac{2}{3(p+\nu)-2}})$ of (1.8) [16] under the Hölder continuous setting.

105 Besides the unified convergence rate (1.9), the UAF has several significant improvements
 106 over the ACNM approach and the A-HPE framework. First, the UAF provides a continuous
 107 transition from the ACNM approach to the A-HPE framework by choosing q from $p + \nu$ to
 108 2. Second, as we will soon see, the UAF can be conveniently instantiated by only specifying
 109 two problem-related parameters and two adjustable parameters for framework design, and
 110 thus recover the high-order acceleration algorithms [27, 3, 31, 15, 16, 23, 18, 6, 14] without
 111 extra effort. Third, we provide *the first* and also a unified convergence rate analysis for both
 112 the Euclidean and non-Euclidean norm settings, and thus opens the possibility of applying
 113 high-order methods in the non-Euclidean norm setting.³ Fourth, the UAF is proposed and
 114 analyzed directly under the composite convex setting (see (2.6)), hence our results imply that
 115 all existing high-order algorithms [27, 3, 31, 16, 6, 14] can be naturally extended to the general
 116 composite convex setting.

²As we will later see, q is the order of the uniform convexity of the proxy-function for framework design. [16] has used a uniformly convex proxy-function with $q = (p + \nu)$ -th order, while [18, 14, 6] have used a uniformly convex proxy-function with $q = 2$ -nd order.

³which is pertinent to many important practical problems such as logistic regression loss in machine learning, see Example 2.4.

117 In terms of implementation for high-order acceleration algorithms, to obtain the optimal
 118 rate that matches the lower bound [2], we must employ a binary search procedure to find a
 119 suitable coupling coefficient in each iteration, which may substantially slow down the practical
 120 performance [31]. Therefore, in addition to the above theoretical results, we introduce a
 121 simple heuristic for finding the coupling coefficient, suggested by our analysis, so that the
 122 resulting implementation does not need a binary search procedure required by the optimal
 123 acceleration method. Our experiments show that this simple heuristic is extremely effective
 124 and can easily ensure the conditions needed to achieve the optimal rate. This leads to a very
 125 practical implementation of the optimal acceleration algorithms without extra implementation
 126 cost, alleviating concerns raised by [31]. Last but not the least, with a general *restart scheme*,
 127 our analysis for the general convex setting extends to the uniformly convex setting. The
 128 resulting algorithm complexity can match the existing lower bounds [2] in most important
 129 cases⁴.

130 **1.3. Our Approach.** In this paper, instead of directly designing an algorithm and then
 131 analyzing its iteration complexity, we consider a different paradigm to make our approach and
 132 algorithm more intuitive and explainable. The paradigm is inspired by the unified theory for
 133 first-order algorithms [13] and the continuous-time interpretations of Nesterov’s acceleration
 134 [35, 21, 22, 36]. Our approach to the algorithmic design is based on an idealized but impractical
 135 algorithm called *vanilla proximal method (VPM)*, introduced in Section 3. The VPM aims to
 136 solve a regularized program of the original one with an arbitrary convergence rate depending
 137 on parameters of our choice. However, the VPM serves more as an ideal target and is itself
 138 computationally infeasible to realize.

139 We show that, in Section 4, to overcome the computational hurdle, one can instead
 140 solve a continuous-time *convex approximation* to the VPM. Then an accelerated continuous-
 141 time dynamics can be derived simply as sufficient conditions to ensure that solution to the
 142 approximate convex program achieves the same convergence rate as the original VPM. Such
 143 point of view unifies the existing continuous-time accelerated dynamics introduced in [35],
 144 [21] and [36] and serves as an arguably better guideline for the design of practical algorithms
 145 in the discrete setting.

146 In practice, to realize the desired accelerated dynamics, we need to know how to implement
 147 them in the discrete setting as an iterative algorithm. To this end, we consider a discrete-time
 148 *convex approximation* to the VPM. However, as we will see in Section 5, in order for the
 149 discrete-time approximation to achieve the same convergence rate as VPM, we must solve a
 150 fixed-point problem which itself is computationally infeasible (if not impossible) in practice.
 151 To circumvent this difficulty, we propose to solve the fixed-point problem *approximately* by
 152 solving a *smooth approximation* to the VPM which becomes a tractable problem. Finally, by
 153 combing the convex approximation and the smooth approximation to the VPM, we propose
 154 the implementable discrete-time unified acceleration framework which achieves the optimal
 155 iteration complexity given in (1.9) for the minimization of convex functions with (p, ν, L) -
 156 Hölder continuous derivatives (for $p \in \{1, 2, \dots\}, \nu \in [0, 1]$ with $p + \nu \geq 2$ and $L > 0$).

157 **2. Preliminaries.** Before we proceed, we first introduce some notations. Let $:=$ denote
 158 a definition. Let $[n]$ denote the set $\{1, 2, \dots, n\}$. For $p = \{1, 2, \dots\}$, let $p! := 1 \times 2 \times \dots \times p$
 159 with $0! := 1$. Let $\|\cdot\|$ denote a norm of vectors and $\|\cdot\|_*$ denote the dual norm of $\|\cdot\|$. For
 160 $x \in \mathbb{R}^d$ and $q \geq 1$, let $\|x\|_q := (\sum_{i=1}^d |x_i|^q)^{\frac{1}{q}}$. For a matrix $B \in \mathbb{R}^{d \times d}$ and $p, q \geq 1$, denote
 161 the operator norm $\|B\|_{p,q} := \max_{x \in \mathbb{R}^d} \{\|Bx\|_q : \|x\|_p \leq 1\}$. By a little abuse of notation,
 162 for a convex function $f(x)$ defined on \mathbb{R}^d , let $\nabla f(x)$ denote the gradient at x or one point
 163 in the subgradient set $\partial f(x)$. For a function $f(x; y)$, x denotes the variable of $f(x; y)$, y the

⁴Due to limit of space, one may refer to details in the longer arXiv version: <https://arxiv.org/pdf/1906.00582.pdf>.

164 parameter of $f(x; y)$, and $\nabla f(x; y)$ is the gradient or one point in the subgradient set $\partial f(x; y)$
 165 *w.r.t.* x .

166 Similar to the notations in [31], for $p \in \{1, 2, \dots\}$, we use $\nabla^p f(x)[y_1, y_2, \dots, y_p]$ to
 167 denote the directional derivative of a function f at x along the directions $y_i \in \mathbb{R}^d, i =$
 168 $1, 2, \dots, p$. Then $\nabla^p f(x)[\cdot]$ is a symmetric p -linear form and its operator norm *w.r.t.* a norm
 169 $\|\cdot\|$ is defined as

$$170 \quad (2.1) \quad \|\nabla^p f(x)\|_* := \max_{y_1, y_2, \dots, y_p} \{\nabla^p f(x)[y_1, \dots, y_p] : \|y_i\| \leq 1, i = 1, 2, \dots, p\}.$$

171 **DEFINITION 2.1** (Strictly, Uniformly, or Strongly Convex). *We say a continuous function*
 172 *$f(x)$ is convex on \mathbb{R}^d , if $\forall x, y \in \mathbb{R}^d$, one has*

$$173 \quad (2.2) \quad f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle;$$

174 *$f(x)$ is strictly convex on \mathbb{R}^d , if the equality sign in (2.2) holds if and only if $x = y$;*

175 *$f(x)$ is (s, σ) -uniformly convex on \mathbb{R}^d *w.r.t.* a norm $\|\cdot\|$, if $\forall x, y \in \mathbb{R}^d$, one has*

$$176 \quad (2.3) \quad f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\sigma}{s} \|x - y\|^s,$$

177 *where $s \geq 2$ is the order of uniform convexity and $\sigma \geq 0$ the constant of uniform convexity;*

178 *$f(x)$ is σ -strongly convex on \mathbb{R}^d *w.r.t.* $\|\cdot\|$, if $f(x)$ is $(2, \sigma)$ -uniformly convex on \mathbb{R}^d
 179 *w.r.t.* $\|\cdot\|$.*

180 In Definition 2.1, uniform convexity can be viewed as an extension to the better known concept
 181 of strong convexity. Example 2.2 gives two cases of uniform convexity.

182 **Example 2.2** (Uniform Convexity). $\frac{1}{2}\|x\|_q^2 (1 < q \leq 2)$ is $(2, q - 1)$ -uniformly convex
 183 on \mathbb{R}^d *w.r.t.* $\|\cdot\|_q$ [4]; $\frac{1}{q}\|x\|_2^q (q \geq 2)$ is $(q, 2^{2-q})$ -uniformly convex on \mathbb{R}^d *w.r.t.* $\|\cdot\|_2$ [27].

184 Starting from the work of [29], there is an increasing interest to replace the Lipschitz
 185 continuity assumption by the Hölder continuity assumption [37, 34, 30, 16] and to propose
 186 universal algorithms in the sense that the convergence of algorithms can optimally adapt to
 187 the Hölder parameter. [29, 37] have considered first-order algorithms with Hölder continuous
 188 gradients *w.r.t.* $\|\cdot\|_2$; [15] has proposed cubic regularized Newton methods for minimizing
 189 functions with Hölder continuous Hessians *w.r.t.* $\|\cdot\|_2$; [10, 16] have considered tensor
 190 methods for minimizing convex functions with p -th Hölder continuous derivatives *w.r.t.* $\|\cdot\|_2$
 191 ($p \in \{2, 3, \dots\}$). In this paper, we extend the definition of Hölder continuous derivatives *w.r.t.*
 192 any norm $\|\cdot\|$, including non-Euclidean norms. Our analysis and results will be applicable to
 193 this general setting (Example 2.4 shows some important cases in machine learning.).

194 **DEFINITION 2.3** (Hölder Continuous Derivative). *We say a function $f(x)$ on \mathbb{R}^d has*
 195 *(p, ν, L) -Hölder continuous derivatives *w.r.t.* $\|\cdot\|$, if $\forall x, y \in \mathbb{R}^d$, one has*

$$196 \quad (2.4) \quad \frac{1}{(p-1)!} \|\nabla^p f(x) - \nabla^p f(y)\|_* \leq L \|x - y\|^\nu,$$

197 *where $p \in \{1, 2, 3, \dots\}$ denotes the order of derivative, $0 \leq \nu \leq 1$ denotes the Hölder*
 198 *parameter and $L > 0$ is the constant of smoothness.*

199 *$f(x)$ is said to have (p, L) -Lipschitz continuous derivatives on \mathbb{R}^d *w.r.t.* $\|\cdot\|$ if $f(x)$ has
 200 $(p, 1, L)$ -Hölder continuous derivatives on \mathbb{R}^d *w.r.t.* $\|\cdot\|$.*

201 In Definition 2.3, for $p = 1$, $\|\cdot\|_*$ denotes the dual norm of $\|\cdot\|$; for $p \in \{2, 3, \dots\}$, $\|\cdot\|_*$
 202 denotes the operator norm of tensor of p -th order *w.r.t.* $\|\cdot\|$, which is defined by (2.1).

203 **Example 2.4** (Non-Euclidean High-order Smoothness). Consider the objective function
 204 $f(x) := \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-\bar{b}_j \bar{a}_j^T x))$ for logistic regression, where $j \in [n], \bar{a}_j \in \mathbb{R}^d, \bar{b}_j \in$

205 $\{1, -1\}$. Denote $B := \frac{1}{n} \sum_{j=1}^n \bar{a}_j \bar{a}_j^T$. For $1 \leq p \leq 2$ and q satisfying $1/p + 1/q = 1$, let
 206 $\|\nabla^s f(x)\|_q$ denote the operator norm of $\nabla^s f(x)$ ($s = 2, 3$) in (2.1) w.r.t. the vector norm
 207 $\|\cdot\|_p$. Then we have

$$208 \quad (2.5) \quad \|\nabla^2 f(x) - \nabla^2 f(y)\|_q \leq \|B\|_{p,q} \max_{j \in [n]} \|\bar{a}_j\|_q^\nu \cdot \|x - y\|_p^\nu.$$

209 *Proof.* See Section A.1. ■

210 In the paper, we consider the following composite convex optimization problem:

$$211 \quad (2.6) \quad \min_{x \in \mathbb{R}^d} f(x) := g(x) + l(x),$$

212 where $g(x)$ is a closed proper convex function and $l(x)$ is a simple convex but maybe non-
 213 smooth function. We consider the case when $g(x)$ has (p, ν, L) -Hölder continuous derivatives
 214 for all $x \in \mathbb{R}^d$. Then we can define the following two auxiliary functions:

$$215 \quad (2.7) \quad \hat{f}(x; y) := g(y) + \langle \nabla g(y), x - y \rangle + l(x),$$

$$216 \quad (2.8) \quad \tilde{f}(x; y) := g(y) + \sum_{i=1}^p \frac{1}{i!} \nabla^i g(y) [x - y]^i + l(x),$$

217 where we do not linearize the term $l(x)$ which may be nonsmooth. Formally, we have:

218 **LEMMA 2.5.** *If $g(x)$ and $l(x)$ are convex, and $g(x)$ has (p, ν, L) -Hölder continuous*
 219 *derivatives, then we have: for all $x, y \in \mathbb{R}^d$,*

$$220 \quad (2.9) \quad \hat{f}(x; y) \leq f(x),$$

$$221 \quad (2.10) \quad |f(x) - \tilde{f}(x; y)| \leq \frac{L}{p} \|x - y\|^{p+\nu},$$

$$222 \quad (2.11) \quad \|\nabla f(x) - \nabla \tilde{f}(x; y)\|_* \leq L \|x - y\|^{p+\nu-1}.$$

223 *Proof.* See Section A.2. ■

224 Because of (2.9), in this paper, $\hat{f}(x; y)$ is viewed as a lower-bound convex approximation
 225 to $f(x)$ for any parameter $y \in \mathbb{R}^d$. $\tilde{f}(x; y)$ satisfies (2.10) and (2.11), and gives a high-order
 226 smooth approximation to $f(x)$ for any parameter $y \in \mathbb{R}^d$. In our analysis, the convexity and
 227 smoothness assumptions are used only by the two inequalities (2.9) and (2.11), which allow
 228 a unified treatment for the smooth and the composite convex settings (with or without the
 229 term $l(x)$). Meanwhile, because we only need the property (2.11) of high-order smoothness, it
 230 implies that in convex optimization, the high-order smoothness is mainly used to give a more
 231 accurate estimation of the ‘‘implicit gradient’’ $\nabla f(x)$.

232 Finally, we give two inequalities in Lemma 2.6 which will be used in our analysis.

233 **LEMMA 2.6.** *Given a sequence $\{b_k\}_{k \geq 0}$ with $b_0 = 0$ and $b_k > 0$ ($k \geq 1$). One has:*

- 234 • For $\rho \geq 1$ and $C > 0$, if $\forall k \geq 1, (b_k - b_{k-1})^\rho \geq C b_k^{\rho-1}$, then $b_k \geq C \left(\frac{k}{\rho}\right)^\rho$;
- 235 • For $\rho \geq 1, \delta > 0$ and $C > 0$, if $\forall k \geq 1, \sum_{i=1}^k \left(\frac{b_i^{\rho-1}}{(b_i - b_{i-1})^\rho}\right)^\delta \leq C$, then $b_k \geq$
 236 $C^{-\frac{1}{\delta}} \left(\frac{k}{\rho}\right)^{\rho + \frac{1}{\delta}}$.

237 *Proof.* See Section A.3. ■

238 **3. Vanilla Proximal Method.** Let us start our study by considering the composite convex
 239 optimization problem in (2.6). In the following discussion, we assume that x^* is a minimizer
 240 of $f(x)$ on \mathbb{R}^d . To design an acceleration algorithm to minimize $f(x)$, we first introduce a

Algorithm 3.1 Vanilla Proximal Method (VPM)

- 1: **Input:** an initialized point $x_0 \in \mathbb{R}^d$, a scalar $A > 0$, a convex function $h(x; x_0)$.
- 2: Find a $z \in \mathbb{R}^d$ that satisfies

$$(3.1) \quad z := \operatorname{argmin}_{x \in \mathbb{R}^d} \{ \psi^{\text{vpm}}(x) := Af(x) + h(x; x_0) \}.$$

- 3: **return** z .

241 so-called *vanilla proximal method (VPM)*, that considers to minimize an auxiliary function
 242 $\psi^{\text{vpm}}(x)$ as in Algorithm 3.1.

243 In $\psi^{\text{vpm}}(x)$, the convex term $h(x; x_0)$ should satisfy the non-negative property:

244 *Assumption 3.1.* $\forall x, x_0 \in \mathbb{R}^d, h(x; x_0) \geq 0$ with $h(x; x_0) = 0$ if and only if $x = x_0$.

245 In the VPM, (3.1) is a convex program and thus there exists a minimizer z . By using only the
 246 optimality condition of (3.1) and Assumption 3.1, we can characterize the ‘‘convergence rate’’
 247 of the VPM as below.

248 **THEOREM 3.2.** *The solution z generated by Algorithm 3.1 satisfies*

$$249 \quad (3.2) \quad f(z) - f(x^*) \leq \frac{h(x^*; x_0)}{A}.$$

250 *Proof.* By the definition of $\psi^{\text{vpm}}(x)$ in (3.1), one has

$$251 \quad (3.3) \quad \min_{x \in \mathbb{R}^d} \psi^{\text{vpm}}(x) \leq Af(x^*) + h(x^*; x_0).$$

252 Then by the optimality condition of z and the nonnegativity of $h(x; x_0)$, one has

$$254 \quad (3.4) \quad Af(z) \leq Af(z) + h(z; x_0) = \min_{x \in \mathbb{R}^d} \psi^{\text{vpm}}(x).$$

255 By the upper bound of $\min_{x \in \mathbb{R}^d} \psi^{\text{vpm}}(x)$ in (3.3) and lower bound of $\min_{x \in \mathbb{R}^d} \psi^{\text{vpm}}(x)$ in
 256 (3.4), after a simple rearrangement, Theorem 3.2 is proved. ■

257 By Theorem 3.2, the VPM may converge with any convergence rate if A is chosen to
 258 a large enough value. Although solving the subproblem (3.1) is impractical in general, it
 259 provides us a good starting point to design practical algorithms: by making certain assumptions
 260 on the objective function $f(x)$ and the proxy function $h(x; x_0)$, it is possible to achieve or
 261 approach the convergence rate of the VPM by solving a tractable approximation to (3.1).
 262

263 **4. Continuous-time Accelerated Descent Dynamics.** The subproblem (3.1) in the
 264 VPM is merely conceptual as it is almost as difficult as minimizing the original function.
 265 Nevertheless, if $f(x)$ is convex, one can always seek more tractable approximations. From
 266 an acceleration perspective, the convex approximation $\hat{f}(x; y)$ in Lemma 2.5 gives a lower
 267 bound for $f(x)$ at the state y . The minimizer of $\hat{f}(x; y)$ would suggest an aggressive direction
 268 and step for the next iterate to go to. However, for such iterates not to diverge too far from
 269 the landscape of $f(x)$, we also need a good upper bound. A basic idea is that up to time t , we
 270 have already traversed a path $x_\tau, \tau \in [0, t)$ over the landscape of $f(x)$. We could potentially
 271 use all the lower-bounds $\hat{f}(x; x_\tau)$ of $f(x)$ to construct a good upper bound to guide the next
 272 step. The simplest possible form for such an upper bound we could consider is a superposition
 273 (or an integral) of these lower bounds to guide the descent trajectory as follows.

$$274 \quad (4.1) \quad z_t := \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \psi_t^{\text{cont}}(x) := \int_0^t a_\tau \hat{f}(x; x_\tau) d\tau + h(x; x_0) \right\},$$

275 where $\forall 0 < \tau \leq t, a_\tau > 0$ and satisfies $\int_0^t a_\tau d\tau = A_t$ with $a_0 = A_0 = 0$ and $\{x_\tau\}_{0 \leq \tau \leq t}$ is
 276 the optimization path and its relationship with z_t will be determined soon.

277 In this section, our main goal is to show that the widely studied continuous-time ac-
 278 celerated dynamics arise from a sufficient condition that allows (4.1) to achieve the same
 279 convergence rate as the original VPM. First, the upper bound (3.3) of ψ^{VPM} is extended to
 280 ψ_t^{cont} as follows.

281 LEMMA 4.1. $\forall t \geq 0$, we have $\min_{x \in \mathbb{R}^d} \psi_t^{\text{cont}}(x) \leq A_t f(x^*) + h(x^*; x_0)$.

282 *Proof.* Lemma 4.1 can be easily proven by using (2.9). ■

283 In other words, Lemma 4.1 provides a lower bound of $A_t f(x^*)$. Second, the lower bound
 284 (3.4) of ψ^{VPM} can be extended to ψ_t^{cont} as follows, at least approximately.

285 LEMMA 4.2. $\forall t \geq 0$, we have $A_t f(x_t) \leq \min_{x \in \mathbb{R}^d} \psi_t^{\text{cont}}(x) + \int_0^t \langle \nabla f(x_\tau), A_\tau \dot{x}_\tau -$
 286 $a_\tau(z_\tau - x_\tau) \rangle d\tau$.

287 *Proof.* See Section B.1 ■

288 We would like to make this approximation as close as possible and establish $\min_{x \in \mathbb{R}^d} \psi_t^{\text{cont}}(x)$
 289 as an upper bound of $A_t f(x_t)$, at least along certain path by our choice. To this end, based on
 290 Lemmas 4.1 and 4.2, we have the following theorem.

291 THEOREM 4.3 (Continuous-Time VPM). *If the continuous-time trajectories $\{x_t\}_{t \geq 0}$
 292 and $\{z_t\}_{t \geq 0}$ are evolved according to the dynamics:*

$$293 \quad (4.2) \quad \begin{cases} A_t \dot{x}_t = a_t(z_t - x_t), \\ z_t = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \int_0^t a_\tau \hat{f}(x; x_\tau) d\tau + h(x; x_0) \right\}, \end{cases}$$

294 where $\forall 0 < \tau \leq t, a_\tau > 0, \int_0^t a_\tau d\tau = A_t$, and $a_0 = A_0 = 0$, then for all $t > 0$, one has

$$295 \quad (4.3) \quad f(x_t) - f(x^*) \leq \frac{h(x^*; x_0)}{A_t}.$$

296 *Proof.* If $A_t \dot{x}_t = a_t(z_t - x_t)$, from Lemma 4.2, one has $A_t f(x_t) \leq \min_{x \in \mathbb{R}^d} \psi_t^{\text{cont}}(x)$.
 298 Combining Lemma 4.1, we have (4.3). ■

299 In Theorem 4.3, (4.2) does not specify any concrete values or forms for a_t and A_t , except
 300 the condition $a_\tau > 0, \int_0^t a_\tau d\tau = A_t$ ($0 < \tau \leq t$) and $a_0 = A_0 = 0$;⁵ meanwhile it does not
 301 specify any concrete form for $h(x; x_0)$. As result, by instantiating the dynamical system (4.2)
 302 for different choices of a_t, A_t , and h , one may obtain all ODEs previously introduced and
 303 studied for algorithm acceleration in the literature [35, 21, 22, 36], respectively.⁶

304 *Remark 4.4.* Although we have derived the dynamics (4.2) from a different perspective,
 305 it should be noted that the dynamical system (4.2) is an extension and refinement to the ODE
 306 derived by the ‘‘approximate duality gap technique (ADGT)’’ [13]. The main difference is
 307 that instead of giving an upper bound of $f(x_t)$ and a lower bound of $f(x^*)$, we give an upper
 308 bound of $A_t f(x_t)$ and a lower bound of $A_t f(x^*)$. This modification allows us to set $A_0 = 0$
 309 rather than $A_0 > 0$, and thus the initialization expression about A_0 can be removed. Such a
 310 modification simplifies future derivation and analysis greatly.

311 **5. Unified Acceleration Framework.** To achieve the same convergence rate of the VPM,
 312 the continuous-time approximation needs the extra ODE condition in (4.2), which is reasonable
 313 to assume in the continuous setting. In the discrete-time setting, if all other conditions remain

⁵Theoretically A_t should be chosen such that the differential equation has a unique solution.

⁶See examples given in the long arXiv version: <https://arxiv.org/pdf/1906.00582.pdf> for details.

314 unchanged, except that we replace the weighted continuous-time approximation (4.1) by a
 315 weighted discrete-time counterpart, one may see that the ODE will be replaced by a condition
 316 that requires us to find a solution to a *fixed-point* problem (which will be clear in Lemma
 317 5.3). Unfortunately, directly solving this fixed-point problem is computationally infeasible in
 318 practice. To remedy this difficulty, we need stronger conditions for the proxy function $h(x; x_0)$,
 319 the associated norm $\|\cdot\|$, and the smooth component $g(x)$ of $f(x)$ as given in Assumption 5.1
 320 below.

321 *Assumption 5.1.* $\forall x, x_0 \in \mathbb{R}^d, p \in \{1, 2, \dots\}, \nu \in [0, 1]$ with $p + \nu \geq 2, q \in [2, p + \nu]$,
 322 $\gamma > 0$ and $L > 0$, we have

- 323 1. $h(x; x_0)$ satisfies Assumption 3.1 and is (q, γ) -uniformly convex *w.r.t.* a norm $\|\cdot\|$.
- 324 2. $\frac{1}{q}\|x\|^q$ is (q, β) -uniformly convex *w.r.t.* the norm $\|\cdot\|$.
- 325 3. $g(x)$ has (p, ν, L) -Hölder continuous derivatives *w.r.t.* the norm $\|\cdot\|$.

326 Based on Assumption 5.1, we consider a weighted discrete-time convex approximation of
 327 (3.1): for $k \geq 0$,

$$328 \quad (5.1) \quad z_k := \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \psi_k^{\text{dis}}(x) := \sum_{i=1}^k a_i \hat{f}(x; x_i) + h(x; x_0) \right\},$$

329 where we assume that $\forall 1 \leq i \leq k, a_i > 0, A_i := \sum_{j=1}^i a_j$ and $a_0 = A_0 = 0, h(x; x_0)$
 330 satisfies Assumption 5.1, and $\hat{f}(x; x_i)$ is defined in Lemma 2.5. Meanwhile, in (5.1), when
 331 $k = 0$, we let $\psi_0^{\text{dis}}(x) = h(x; x_0)$ and thus $z_0 = \operatorname{argmin}_{x \in \mathbb{R}^d} h(x; x_0) = x_0$. We now derive
 332 the unified acceleration framework by analyzing the conditions needed to emulate the same
 333 rate of the VPM.

334 First, the upper bound (3.3) of ψ^{VPM} can be extended to the discrete case ψ_k^{dis} trivially.

335 LEMMA 5.2. $\forall k \geq 0$, one has $\min_{x \in \mathbb{R}^d} \psi_k^{\text{dis}}(x) \leq A_k f(x^*) + h(x^*; x_0)$.

336 *Proof.* Lemma 5.2 can be easily proven by using (2.9). ■

337 Then in Lemma 5.3 below, we show how the lower bound (3.4) of ψ^{VPM} can be extended
 338 to the discrete case ψ_k^{dis} with some extra terms.

339 LEMMA 5.3. $\forall i \geq 1$, let $E_i := A_i \left\langle \nabla f(x_i), x_i - \frac{a_i}{A_i} z_i - \frac{A_{i-1}}{A_i} x_{i-1} \right\rangle - \frac{\gamma}{q} \|z_i - z_{i-1}\|^q$.

340 Then $\forall k \geq 1$, one has $A_k f(x_k) - \psi_k^{\text{dis}}(z_k) \leq \sum_{i=1}^k E_i$.

341 *Proof.* See Section C.1. ■

342 In Lemma 5.3, the extra negative term $-\frac{\gamma}{q} \|z_i - z_{i-1}\|^q$ in E_i is from the uniform convexity
 343 of $h(x; x_0)$. If $h(x; x_0)$ is only convex (i.e. $\gamma = 0$), this negative term does not exist and thus
 344 a sufficient condition for $E_i \leq 0$ is:

$$345 \quad (5.2) \quad x_i = \frac{a_i}{A_i} z_i + \frac{A_{i-1}}{A_i} x_{i-1}, \quad \forall 1 \leq i \leq k.$$

346 By (5.1), z_i is a function of x_i . Therefore finding x_i to satisfy (5.2) is reduced to a fixed-point
 347 problem (so is it for z_i). It is computationally infeasible (if not impossible) to find an exact
 348 solution to this problem in general. Nevertheless, if $h(x; x_0)$ satisfies Assumption 5.1, the term
 349 E_i contains a negative term $-\frac{\gamma}{q} \|z_i - z_{i-1}\|^q$. So there is hope that an approximate solution to
 350 the fixed-point problem (5.2) can still make $E_i \leq 0$.

351 To approximately solve the fixed-point problem, for convenient analysis, inspired by
 352 [17, 12], we define a pair $(\hat{x}_{i-1}, \hat{z}_i)$ such that

$$353 \quad (5.3) \quad \hat{x}_{i-1} := \frac{a_i}{A_i} z_{i-1} + \frac{A_{i-1}}{A_i} x_{i-1}, \quad \hat{z}_i := \frac{A_i}{a_i} x_i - \frac{A_{i-1}}{a_i} x_{i-1}.$$

354 By the definition of \hat{z}_i in (5.3), we have $x_i = \frac{\alpha_i}{A_i} \hat{z}_i + \frac{A_{i-1}}{A_i} x_{i-1}$. Therefore (5.3) can be viewed
 355 as two-step fixed-point iterations for x_i based on \hat{x}_{i-1} and \hat{z}_i . Here \hat{z}_i can be viewed as the
 356 best estimate of the desired fixed point z_i based on the calculated x_i in our algorithm. It is
 357 defined for convenience and will only be used in our analysis but not in the algorithm.

358 Based on the definition of $(\hat{x}_{i-1}, \hat{z}_i)$, Assumption 5.1, and the definition of E_i in Lemma
 359 5.3, we have the following result.

360 LEMMA 5.4. For $i \geq 1$ and any $\gamma'_i \in (0, \gamma]$, we have

$$361 \quad E_i \leq a_i \left\langle \nabla f(x_i) + \frac{\gamma'_i A_i^{q-1}}{a_i^q} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q, \hat{z}_i - z_i \right\rangle$$

$$362 \quad (5.4) \quad -\gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right).$$

363 *Proof.* See Section C.2. ■

364 In Lemma 5.4, we purposely introduce a new parameter γ'_i , which as we will soon show,
 365 helps unify the four high-order instances [23, 18, 6, 18] of the A-HPE framework. Meanwhile,
 366 because of the uniform convexity of $\frac{1}{q} \|\cdot\|^q$, the negative term $-\frac{\beta}{q} \|\hat{z}_i - z_i\|^q$ is reduced to
 367 two negative terms and an inner product.

368 By Lemma 5.4, if we can find x_i such that

$$369 \quad (5.5) \quad \nabla f(x_i) + \frac{\gamma'_i A_i^{q-1}}{a_i^q} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q = 0,$$

370 then we can ensure $E_i \leq 0$. However the problem of finding x_i that satisfies (5.5) is equivalent
 371 to solving the VPM problem exactly in (3.1) with the settings $x_0 := \hat{x}_{i-1}$, $h(x; x_0) :=$
 372 $\frac{1}{q} \|x - \hat{x}_{i-1}\|^q$, $A := \frac{\alpha_i^q}{\gamma'_i A_i^{q-1}}$, which is computationally infeasible in general. Fortunately,
 373 the two negative terms in (5.4) may dominate small errors if we can solve the VPM problem
 374 (5.5) approximately. Hence we approximate the intermediate VPM problem (5.5) by a smooth
 375 approximation $\tilde{f}(x_i; \hat{x}_{i-1})$ using the fact

$$376 \quad (5.6) \quad \|\nabla f(x_i) - \nabla \tilde{f}(x_i; \hat{x}_{i-1})\|_* \leq L \|x_i - \hat{x}_{i-1}\|^{p+\nu-1},$$

377 from Lemma 2.5. Then by Lemmas 2.5 and 5.4, we have Lemma 5.5.

378 LEMMA 5.5. Denote $c_q := (\beta(q-1)^{1-q})^{\frac{1}{q}}$ and $\lambda'_i := \frac{\alpha_i^q}{c_q \gamma'_i A_i^{q-1}}$. For $i \geq 1$, one has

$$379 \quad E_i \leq \left((L \lambda'_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q})^{\frac{q}{q-1}} - 1 \right) \frac{\gamma'_i A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q$$

$$380 \quad (5.7) \quad + a_i \left\langle \nabla \tilde{f}(x_i; \hat{x}_{i-1}) + \frac{\gamma'_i A_i^{q-1}}{a_i^q} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q, \hat{z}_i - z_i \right\rangle.$$

381 *Proof.* See Section C.3. ■

382 From Lemma 5.5, to ensure $E_i \leq 0$, the VPM problem (5.5) can be reduced to an easier
 383 smooth approximation problem

$$384 \quad (5.8) \quad \nabla \tilde{f}(x_i; \hat{x}_{i-1}) + \frac{\gamma'_i A_i^{q-1}}{a_i^q} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q = 0,$$

385 and we also need the condition

$$386 \quad (5.9) \quad L \lambda'_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 \leq 1$$

387 to hold, where $\theta_2 \in (0, 1]$ is a constant.

388 We here discuss the role of the parameter γ'_i . So far our derivation works for any
 389 $\gamma'_i \in (0, \gamma]$. A simple choice of γ'_i would be $\gamma'_i := \gamma$. Nevertheless, under the condition (5.9),
 390 for any $\alpha \in [0, 1]$, we could choose γ'_i to satisfy:

$$391 \quad (5.10) \quad \gamma'_i = \left(\frac{L\lambda'_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q}}{\theta_2} \right)^{\frac{\alpha}{1-\alpha}} \gamma,$$

392 where for $\alpha = 1$, we set $\frac{\alpha}{1-\alpha} = \lim_{\alpha \rightarrow 1^-} \frac{\alpha}{1-\alpha} = +\infty$. This would still ensure $\gamma'_i \in (0, \gamma]$.
 393 But notice that λ'_i in the RHS depends on γ'_i . To sort out an explicit expression for so-defined
 394 γ'_i , we denote

$$395 \quad (5.11) \quad \lambda_i := \frac{a_i^q}{c_q \gamma A_i^{q-1}}.$$

396 Then by the definition of λ'_i in Lemma 5.5, with (5.10) and (5.11), we can write γ'_i in the form:

$$397 \quad (5.12) \quad \gamma'_i = \left(\frac{L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q}}{\theta_2} \right)^\alpha \gamma.$$

398 Then by the fact for all $s \geq 0, t \geq 2, x \in \mathbb{R}^d$,

$$399 \quad (5.13) \quad \|x\|^s \nabla \frac{1}{t} \|x\|^t = \nabla \frac{1}{s+t} \|x\|^{t+s},$$

400 and combing (5.11) and (5.12), it follows that (5.8) is equivalent to

$$401 \quad (5.14) \quad \nabla \tilde{f}(x_i; \hat{x}_{i-1}) + \frac{L^\alpha}{c_q \lambda_i^{(1-\alpha)} \theta_2^\alpha} \nabla \frac{1}{\alpha(p+\nu) + (1-\alpha)q} \|x_i - \hat{x}_{i-1}\|^{\alpha(p+\nu) + (1-\alpha)q} = 0.$$

402 Or equivalently, let $\varsigma := \alpha(p+\nu) + (1-\alpha)q$, and then x_i is the solution to the following
 403 minimization problem:

$$404 \quad (5.15) \quad x_i := \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \tilde{f}(x; \hat{x}_{i-1}) + \frac{L^\alpha}{c_q \lambda_i^{(1-\alpha)} \theta_2^\alpha \varsigma} \|x - \hat{x}_{i-1}\|^\varsigma \right\}.$$

406 In (5.15), because $\alpha \in [0, 1]$, the power of the norm $\|x - \hat{x}_{i-1}\|$ ranges from $p+\nu$ to q freely,
 407 which unifies the choice $\alpha = 0$ in [23] and $\alpha = 1$ in [6]. Meanwhile, [18, 14] has used a
 408 mixture of both $\alpha = 0$ and $\alpha = 1$ in their formulations, which are also equivalent to (5.14) by
 409 (5.13). A surprising phenomenon is that, as our analysis shows, the choice of α in (5.15) does
 410 not affect the convergence rate (except the constant).

411 Meanwhile, by (5.12), (5.9) is equivalent to

$$412 \quad (5.16) \quad \omega_i := L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 \leq 1,$$

413 where we call ω_i as a *convergence indicator* in the sense that if for all $1 \leq i \leq k$, $\omega_i \leq 1$, then
 414 the iterate x_k will converge according to the following theorem; otherwise, the convergence of
 415 x_k is not guaranteed. Based on the equivalence relationship between (5.8) and (5.15), (5.9)
 416 and (5.16), we have Theorem 5.6.

417 **THEOREM 5.6 (Discrete-Time VPM).** *If Assumption 5.1 is true, and in (5.1), $\forall i \geq 1$,
 418 the sequences $\{a_i\}, \{A_i\}$ satisfy $a_i > 0, A_i = A_{i-1} + a_i$ with $a_0 = A_0 = 0, \{x_i\}$ satisfies
 419 (5.15), and $\{\lambda_i\}$ defined in (5.11) satisfies (5.16), then for $k \geq 1$, one has*

$$420 \quad (5.17) \quad f(x_k) - f(x^*) \leq \frac{h(x^*; x_0)}{A_k}.$$

421 *Proof.* See Section C.4. ■

422 As we see, Theorem 5.6 is very much like the discrete-time approximation version of Theorem
423 3.2. To accurately characterize the convergence rate from (5.17), we need to have a good
424 lower-bound for A_k . In Theorem 5.6, by $A_i = A_{i-1} + a_i$, the definition of λ_i (5.11) and the
425 condition (5.16), it follows that A_i must satisfy the condition

$$426 \quad (5.18) \quad \frac{L(A_i - A_{i-1})^q}{c_q \gamma A_i^{q-1}} \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 \leq 1.$$

427 Therefore A_i cannot be chosen as an arbitrarily large value as in the continuous-time setting.
428 Except the basic condition $A_0 = 0$ and for $i \geq 1$, $A_i > 0$, (5.18) is the only condition A_i
429 needs to satisfy, therefore one may expect that the tightest bound of A_i should be obtained if
430 $\frac{L(A_i - A_{i-1})^q}{c_q \gamma A_i^{q-1}} \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} = O(1)$. In other words, we hope that

$$431 \quad (5.19) \quad 0 < \theta_1 \leq L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} = \frac{L(A_i - A_{i-1})^q}{c_q \gamma A_i^{q-1}} \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 \leq 1,$$

432 where θ_1 and θ_2 are $O(1)$ constants. To verify this point of view, we discuss below the two
433 settings $q = p + \nu$ and $q < p + \nu$, respectively.

434 When $q = p + \nu$, we have $\lambda_i = \frac{(A_i - A_{i-1})^q}{c_q \gamma A_i^{q-1}}$ and $\|x_i - \hat{x}_{i-1}\|^{p+\nu-q} = 1$. Taking A_i as a
435 variable, then for all $A_i > A_{i-1}$, by the fact $q \geq 2$ and

$$436 \quad (5.20) \quad \frac{d \log \lambda_i}{d A_i} = \frac{(q-1)A_{i-1} + A_i}{A_i(A_i - A_{i-1})} > 0,$$

438 we have λ_i is a strictly monotonically increasing function *w.r.t.* A_i , which is an one to one
439 mapping. Therefore determining the lower bound of A_i is equivalent to determining the lower
440 bound of λ_i . To ensure $E_i \leq 0$, by the condition (5.16), when $q = p + \nu$, $L\lambda_i$ is upper
441 bounded by the constant $\theta_2 \leq 1$. Therefore the tightest lower bound for λ_i is obtained if $L\lambda_i$
442 is lower bounded by a constant $\theta_1 \in (0, \theta_2]$. Then by Lemma 2.6 and Theorem 5.6, we obtain
443 Theorem 5.7.

444 **THEOREM 5.7** (Convergence Rate for the Case $q = p + \nu$). *If Assumption 5.1 is*
445 *true, c_q is defined in Lemma 5.5, and in (5.1), $\forall i \geq 1$, the sequences $\{a_i\}, \{A_i\}$ satisfy*
446 *$a_i > 0, A_i = A_{i-1} + a_i$ with $a_0 = A_0 = 0$, $\{x_i\}$ satisfies (5.15), and $\{\lambda_i\}$ defined in (5.11)*
447 *satisfies*

$$448 \quad (5.21) \quad 0 < \theta_1 \leq L\lambda_i \leq \theta_2 \leq 1,$$

449 then for $k \geq 1$, we have

$$450 \quad (5.22) \quad A_k \geq \frac{\theta_1 c_q \gamma}{L} \left(\frac{k}{p + \nu} \right)^{p+\nu},$$

451 and

$$452 \quad (5.23) \quad f(x_k) - f(x^*) \leq \frac{h(x^*; x_0)}{A_k} \leq \frac{L}{\theta_1 c_q \gamma} h(x^*; x_0) \left(\frac{p + \nu}{k} \right)^{p+\nu}.$$

453 *Proof.* See Section C.5. ■

454 When $q < p + \nu$, because the condition of λ_i to ensure $E_i \leq 0$ involves the unknown x_i ,
455 the situation seems to be more complicated. Nevertheless, under the conditions (5.15) and
456 (5.16), and combining Lemmas 5.2 and 5.3, we can obtain a condition as in Lemma 5.8 below
457 that leads to a good lower bound for A_k .

458 LEMMA 5.8. Assume $\{x_i\}$ satisfies (5.15) and $\{\omega_i\}$ satisfies (5.16). Then if $2 \leq q <$
 459 $p + \nu$, we have

(5.24)

$$460 \sum_{i=1}^k \omega_i^{\frac{\varsigma}{p+\nu-q}} \left(\frac{A_i^{p+\nu-1}}{(A_i - A_{i-1})^{p+\nu}} \right)^{\frac{q}{p+\nu-q}} \leq q\theta_2^\alpha (1 - \theta_2^{\frac{q}{q-1}})^{-1} \gamma^{-\frac{p+\nu}{p+\nu-q}} \left(\frac{L}{c_q} \right)^{\frac{q}{p+\nu-q}} h(x^*; x_0).$$

462 *Proof.* See Section C.6. ■

463 In Lemma 5.8, if $\theta_2 \in (0, 1)$, then the RHS of (5.24) will be a positive constant. Therefore
 464 if ω_i on the LHS of (5.24) is lower bounded by a constant $\theta_1 \in (0, \theta_2]$, then we use Lemma
 465 2.6 to give a lower bound about A_i . Based on the above analysis, and combining Lemma
 466 2.6, Theorem 5.6 and Lemma 5.8, we can characterize the convergence rate of the proposed
 467 iteration when $2 \leq q < p + \nu$.

468 THEOREM 5.9 (Convergence Rate for the Case $2 \leq q < p + \nu$). If Assumption 5.1
 469 is true, c_q is defined in Lemma 5.5, and in (5.1), $\forall i \geq 1$, the sequences $\{a_i\}, \{A_i\}$ satisfy
 470 $a_i > 0, A_i = A_{i-1} + a_i$ with $a_0 = A_0 = 0$, $\{x_i\}$ satisfies (5.15), and $\{\lambda_i\}$ defined in (5.11)
 471 satisfies

$$472 (5.25) \quad 0 < \theta_1 \leq \omega_i = L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 < 1,$$

473 then by defining $C_0 := \left(q\theta_2^\alpha (1 - \theta_2^{\frac{q}{q-1}})^{-1} \right)^{-\frac{p+\nu-q}{q}} \theta_1^{\frac{\alpha(p+\nu)+(1-\alpha)q}{q}} \gamma^{\frac{p+\nu}{q}} c_q$, we have

$$474 (5.26) \quad A_k \geq \frac{C_0}{L} (h(x^*; x_0))^{-\frac{p+\nu-q}{q}} \left(\frac{k}{p+\nu} \right)^{\frac{(q+1)(p+\nu)-q}{q}}$$

475 and

$$476 (5.27) \quad f(x_k) - f(x^*) \leq \frac{h(x^*; x_0)}{A_k} \leq \frac{L}{C_0} (h(x^*; x_0))^{\frac{p+\nu}{q}} \left(\frac{p+\nu}{k} \right)^{\frac{(q+1)(p+\nu)-q}{q}}.$$

477 *Proof.* See Section C.7. ■

In Theorems 5.7 and 5.9, if we do not consider the constants, in both $q = p + \nu$ and
 $2 \leq q < p + \nu$ settings, we can find an ϵ -accurate solution x such that $f(x) - f(x^*) \leq \epsilon$ with

$$O\left(\epsilon^{-\frac{q}{(q+1)(p+\nu)-q}}\right)$$

478 iterations, where $q \in [2, p + \nu]$. It is easy to find that the rate will be the best as $O\left(\epsilon^{-\frac{2}{3(p+\nu)-2}}\right)$
 479 if we set $q = 2$. In fact $O\left(\epsilon^{-\frac{2}{3(p+\nu)-2}}\right)$ matches the lower bound of iteration complexity [16]
 480 for all the settings of $p \in \{1, 2, \dots\}$ and $\nu \in [0, 1]$ as long as $p + \nu \geq 2$. As q becomes large,
 481 the rate $O\left(\epsilon^{-\frac{q}{(q+1)(p+\nu)-q}}\right)$ will become worse. However, particularly, when $q = p + \nu$, λ_i
 482 can be determined trivially and thus the setting $q = p + \nu$ is suboptimal but has the advantage
 483 of algorithmic implementation, as we will elaborate on later.

484 Regarding the other two parameters θ_1, θ_2 , when $q = p + \nu$, based on Theorem 5.7,
 485 to minimize the bound in (5.23), the optimal choice will be $\theta_1 = 1$ and thus $\theta_2 = 1$ by
 486 $\theta_1 \leq \theta_2 \leq 1$. When $q < p + \nu$, based on Theorem 5.9, one can optimize the choice of θ_1, θ_2
 487 by minimizing the bound in (5.27) under the constraint $0 < \theta_1 \leq \theta_2 < 1$.

488 As we have noted before, by varying the parameter α from 0 to 1 in (5.15), the range of
 489 the power of $\|x - \hat{x}_{i-1}\|$ changes from q to $p + \nu$. For $q = p + \nu$, as Theorem 5.7 indicates,
 490 choice of α has no influence on the convergence rate; for $2 \leq q < p + \nu$, as Theorem 5.9

TABLE 1
Algorithmic Instances of the Unified Acceleration Framework with $h(x; x_0) := \frac{1}{q} \|x - x_0\|_2^q$.

Instances	p	ν	q	α
[3, 31]	$\{2, 3, \dots\}$	1	$p + 1$	1
[16]	$\{2, 3, \dots\}$	$[0, 1]$	$p + \nu$	1
[23]	2	1	2	0
[6]	$\{2, 3, \dots\}$	1	2	1
[18, 14]	$\{2, 3, \dots\}$	1	2	a mixture of 0 and 1

491 shows, α only has a minor influence on the constant in the bound. Therefore, our result shows
 492 that α can be chosen according to implementation convenience without worrying about the
 493 convergence rate.

494 Compared with the existing results for high-order optimization [32, 27, 31, 16] and
 495 [23, 6, 18, 14], our convergence results are given under the Hölder continuous assumption
 496 *w.r.t.* a general norm $\|\cdot\|$ that satisfies Assumption 5.1. Such general norms include the
 497 Euclidean norm $\|x\|_2$ and the generalized Euclidean norm $\sqrt{x^T B x}$ as special cases, where B
 498 is any positive definite matrix. To the best of our knowledge, this is the first convergence result
 499 for high-order optimization that can be applied to the high-order non-Euclidean smoothness
 500 setting. To this end, we have adopted a new proof paradigm inspired by the intuitive proof
 501 techniques for the accelerated extra-gradient descent (AXGD) algorithm [12] for first order
 502 methods.

503 Summarizing the above results, we obtain a *unified acceleration framework* (UAF) shown
 504 in Algorithm 5.1. In the algorithm, the parameters p, ν are from the problem setting and the
 505 parameters q, α and the proxy function $h(x; x_0)$ are for framework design. These parameters
 506 can vary in their entire feasible ranges. By specifying p, ν, q, α and $h(x; x_0)$, we obtain
 507 algorithmic instances of UAF. As results, Algorithm 5.1 recovers many existing algorithms.
 508 We give a few examples in Table 1⁷.

509 Meanwhile, Algorithm 5.1 also includes several new interesting instances. First, if we set
 510 $p = \nu = 1, q = 2, \alpha \in [0, 1]$, Algorithm 5.1 defines a new variant of AGD with an $O(1/k^2)$
 511 convergence rate. Such variant is similar to the variant AXGD of AGD. One advantage of
 512 this variant is that Algorithm 5.1 allows $h(x; x_0)$ to be any strongly convex function *w.r.t.*
 513 $\|\cdot\|$, while AXGD assumes that $h(x; x_0)$ is the Bregman divergence of a strongly convex
 514 function *w.r.t.* $\|\cdot\|$. Second, if we set $p \in \{2, 3, \dots\}, \nu \in [0, 1], q = 2, \alpha \in [0, 1]$, then we
 515 obtain the first kind of high-order algorithms that can attain the optimal rate $O(\epsilon^{-\frac{2}{3(p+\nu)-2}})$ for
 516 the composite minimization problem (2.6) with the smooth component $g(x)$ having (p, ν, L) -
 517 Hölder continuous derivatives *w.r.t.* $\|\cdot\|$.

518 For the loop from Step 4 to 7 in Algorithm 5.1, we need solve two subproblems:

- 519 • The first one is about finding λ_i such that the minimizer x_i of the objective (5.15),
 520 together with λ_i , satisfy the conditions (5.28) and (5.29).
- 521 • The second one is about finding the solution z_i of a discrete-time convex approxima-
 522 tion problem of the VPM in Step 6. Because in our setting the convex approximation
 523 $\hat{f}(x; y)$ defined in Lemma 2.5 is a linear function plus a simple convex function $l(x)$,
 524 the subproblem of finding z_i can be solved efficiently.

525 When $p = \nu = 1$ and $q = 2$, the subproblem associated with Step 5, namely (5.15), is
 526 reduced to a proximal gradient decent step [33], which can be solved efficiently. However,

⁷As we have mentioned, [18, 14] have used a mixture of 0 and 1 for α , which is also equivalent to (5.14) by (5.13). To simplify presentation and practical implementation, we usually only consider choosing a single α in Algorithm 5.1.

Algorithm 5.1 Unified Acceleration Framework (UAF)

-
- 1: **Input:** $\hat{f}(x; y)$ and $\tilde{f}(x; y)$ in (2.7) and (2.8), $h(x; x_0)$, $\|\cdot\|$ in Assumption 5.1.
 - 2: Set the constant $\theta_2 \in (0, 1]$ if $q = p + \nu$, $\theta_2 \in (0, 1)$ if $2 \leq q < p + \nu$; set $\theta_1 \in (0, \theta_2]$.
 - 3: Set $\alpha \in [0, 1]$, $c_q = (\beta(q-1)^{1-q})^{\frac{1}{q}}$, $\varsigma = \alpha(p + \nu) + (1 - \alpha)q$.
 - 4: Set $A_0 = 0$, $x_0 = z_0 \in \mathbb{R}^d$.
 - 5: **for** $i = 1$ **to** k **do**
 - 6: $x_i = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \tilde{f}(x; \hat{x}_{i-1}) + \frac{L^\alpha}{c_q \lambda_i^{(1-\alpha)} \theta_2^\alpha \varsigma} \|x - \hat{x}_{i-1}\|^\varsigma \right\}$,
where we find a $\lambda_i > 0$ such that a_i, A_i, λ_i and $\hat{x}_{i-1} \in \mathbb{R}^d$ satisfy
$$(5.28) \quad A_i = A_{i-1} + a_i, \quad \lambda_i = \frac{a_i^q}{c_q \gamma A_i^{q-1}}, \quad \hat{x}_{i-1} = \frac{A_{i-1}}{A_i} x_{i-1} + \frac{a_i}{A_i} z_{i-1},$$

$$(5.29) \quad \theta_1 \leq L \lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2.$$
 - 7: Update $z_i = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \sum_{j=1}^i a_j \hat{f}(x; x_j) + h(x; x_0) \right\}$.
 - 8: **end for**
 - 9: **return** x_k .
-

527 in the setting of high-order optimization, *i.e.*, $p \in \{2, 3, \dots\}$, (5.15) is nontrivial and will
528 dominate the per-iteration cost in general. Finding a general efficient procedure to solve
529 this subproblem remains active research. Nevertheless, for some special important cases,
530 there already exist efficient algorithms. For example, if $p = 2, \nu = 1, \alpha = 1$ and the maybe
531 nonsmooth part $l(x) = 0$, (5.15) is reduced to an iteration of cubic regularized Newton method
532 (CNM), which can be solved efficiently by the Lanczos method [9]; if $p = 3, \nu = 1, \alpha = 1$
533 and $l(x) = 0$, (5.15) is reduced to a third-order convex multivariate polynomial and can be
534 solved as efficiently as the iteration of CNM in many cases [31, 7].

535 Notice that, in Step 5, for the setting $q = p + \nu$, λ_i can be determined easily as it does not
536 depend on x_i and thus A_i, a_i can be solved efficiently by solving a simple one-dimensional
537 equation with Newton method. However, for the setting $2 \leq q < p + \nu$, the condition (5.29)
538 depends on the solution x_i and can not be determined so trivially. In fact, as of now, when
539 $2 \leq q < p + \nu$, we do not even know whether we can find such a pair (x_i, λ_i) that satisfies
540 all the conditions simultaneously. However, as nearly a trivial extension to [6], the following
541 Proposition 5.10 ensures such a pair always exists until we attain the minimizer.

542 **PROPOSITION 5.10.** *Let $A \geq 0, \lambda \geq 0, x, y \in \mathbb{R}^d$ such that $f(x) \neq f(x^*)$. Assume that*
543 *$a(\lambda)$ is implicitly defined by*

$$544 \quad (5.30) \quad \lambda = \frac{(a(\lambda))^q}{c_q \gamma (A + a(\lambda))^{q-1}}, \quad \text{and} \quad x(\lambda) = \frac{a(\lambda)}{A + a(\lambda)} x + \frac{A}{A + a(\lambda)} y,$$

$$545 \quad (5.31) \quad w(v) = \operatorname{argmin}_{z \in \mathbb{R}^d} \left\{ \tilde{f}(z; v) + \frac{L^\alpha}{c_q \lambda^{(1-\alpha)} \theta_2^\alpha \varsigma} \|z - v\|^\varsigma \right\},$$

$$546 \quad (5.32) \quad \chi(\lambda) = L \lambda \|w(x(\lambda)) - x(\lambda)\|^{p+\nu-q},$$

547 where the constants $p, q, \nu, \alpha, c_q, \gamma, L, \varsigma$ and θ_2 are given in Algorithm 5.1. Then $\chi(\lambda)$ is a
550 continuous function with $\chi(0) = 0$ and $\chi(+\infty) = +\infty$.

551 *Proof.* See Section C.8 ■

552 By Proposition 5.10, with the setting $A := A_{i-1}, x := z_{i-1}, y := x_{i-1}$, we can always use
553 a binary search procedure to find a pair (x_i, λ_i) such that $\chi(\lambda_i) = L \lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q}$

554 satisfies the condition (5.29). For the case with $\alpha = 0, q = 2$ and $\|\cdot\| := \|\cdot\|_2$, a
 555 complexity analysis for a binary search procedure can be found in [18]; for the case with
 556 $\alpha = 1, p \in \{2, 3, \dots\}, \nu = 1$ and $\|\cdot\| := \|\cdot\|_2$, a complexity analysis for a binary search
 557 procedure can be found in [6]. Although it is possible to give a complexity analysis of binary
 558 search for the general setting in (5.31), in this paper we consider another perspective.

In the Discussion section of [31], Nesterov claims that from the view of practical efficiency, the Algorithm 5.1 with the suboptimal setting $q = p + \nu$ may be better than the Algorithm 5.1 with the optimal setting $q = 2$, where “optimal” is in the sense of iteration complexity. If we do not consider the implementation cost in the Step 5 of Algorithm 5.1 and ignore the difference of constants in the bound of Theorems 5.7 and 5.9, to attain an ϵ -accurate solution such that $f(x) - f(x^*) \leq \epsilon$, the ratio from the number of iterations of the suboptimal algorithm with $q = p + \nu$ to that of the optimal algorithm with $q = 2$ is

$$O\left(\left(\frac{1}{\epsilon}\right)^{\frac{1}{p+\nu} - \frac{2}{3(p+\nu)-2}}\right) = O\left(\left(\frac{1}{\epsilon}\right)^{\frac{p+\nu-2}{(p+\nu)(3(p+\nu)-2)}}\right).$$

559 If $p = 2, \nu = 1$, *i.e.*, the commonly second-order setting, the ratio will be $O\left(\left(\frac{1}{\epsilon}\right)^{\frac{1}{21}}\right)$, which
 560 implies that when we pursue an accuracy $\epsilon = 2^{-21} \approx 10^{-6}$, if the per-iteration cost of the
 561 optimal setting $q = 2$ (or the settings $2 \leq q < p + \nu$) is twice larger than the suboptimal
 562 setting $q = p + \nu$, then the small advantage of the optimal setting will be removed by the
 563 additional implementation complexity. Because of this effect, a binary search procedure
 564 which involves $O(\log \frac{1}{\epsilon})$ calls to the subprocedure for finding x_i may be rather unrealistic in
 565 practice. Therefore in this paper, instead of binary search, we propose a simple heuristic to
 566 find a pair (x_i, λ_i) that satisfy the condition (5.29). The proposed heuristic only needs call the
 567 subprocedure once for finding x_i , which will be explained in Section 6.

568 *Remark 5.11.* The idea that two-step fixed-point iterations lead to acceleration is first
 569 introduced in [12], which has proposed the variant AXGD of AGD. In this paper, such point
 570 of view motivates us to simplify the analysis by defining an intermediate variable \hat{z}_i in (5.3),
 571 whereas the main strategy leading to acceleration in this paper is to use a combination of a
 572 convex approximation (5.1) to the original VPM problem (3.1) and a smooth approximation
 573 (5.15) to the intermediate VPM problem (5.5).

574 *Remark 5.12.* Similar to [16], it is also possible to give a universal version of Algorithm
 575 5.1 in the sense that, by modifying Algorithm 5.1 according to the paradigm of [16], we can
 576 obtain a near-optimal rate even if the Hölder parameter ν is unknown. Such an improvement is
 577 interesting, however it goes beyond the scope of this paper and will be left for further research.

6. Implementation Details and Experimental Validation. High-order optimization is a very different situation from first-order optimization in that the optimal acceleration method (*e.g.*, the UAF with $q = 2$) requires certain conditions to be met (in each iteration). Those conditions sometimes are not so trivial to be satisfied. In fact, in the UAF Algorithm 5.1, for $2 \leq q < p + \nu$, it is not trivial to find a pair (λ_i, x_i) that satisfies the condition (5.29). In (5.16), we have defined $\omega_i = L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q}$ as a convergence indicator in the sense that $\forall 2 \leq q \leq p + \nu$, if

$$\omega_i = L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \leq \theta_2 \leq 1,$$

578 Algorithm 5.1 will converge according to Theorem 5.6; otherwise, the convergence behavior
 579 of Algorithm 5.1 cannot be guaranteed. More specifically, when $q = p + \nu$ if ω_i satisfies
 580 (5.21), then Algorithm 5.1 converges according to Theorem 5.7; when $2 \leq q < p + \nu$, if ω_i
 581 satisfies (5.25), then Algorithm 5.1 converges according to Theorem 5.9. When $q = p + \nu$, we

582 can easily find $0 < \theta_1 \leq \omega_i \leq \theta_2 \leq 1$ to satisfy (5.21); while when $2 \leq q < p + \nu$, because
 583 ω_i involves the variable x_i , it is nontrivial to find a $0 < \theta_1 \leq \omega_i \leq \theta_2 < 1$ to satisfy (5.25). A
 584 standard technique to ensure that the value of the convergence indicator ω_i stays in $[\theta_1, \theta_2]$ is
 585 through a binary search procedure [23, 18, 6]. However, as per our discussion at the end of
 586 Section 5, the cost of the binary search procedure could substantially reduce the advantage of
 587 convergence rate of the optimal method in practice.

588 **6.1. A Good Heuristic for Practical Implementation.** In this section, inspired by the
 589 analysis of Theorem 5.9, for the Algorithm 5.1 with $2 \leq q < p + \nu$, instead of using a binary
 590 search, we introduce a simple heuristic: in the i -th iteration of Algorithm 5.1, A_i is set as its
 591 lower bound such that

$$592 \quad (6.1) \quad A_i = \frac{C_0}{L} (h(x^*; x_0))^{-\frac{p+\nu-q}{q}} \left(\frac{i}{p+\nu} \right)^{\frac{(q+1)(p+\nu)-q}{q}},$$

593 where all the constants are from Theorem 5.9. With so assigned A_i , λ_i and a_i can be easily
 594 determined by (5.28). Therefore the per-iteration cost under the setting $2 \leq q < p + \nu$ will
 595 remain the same as the setting $q = p + \nu$.

596 However, if we use the heuristic (6.1) of A_i for $2 \leq q < p + \nu$, there is no theoretical
 597 guarantee for convergence of the algorithm. In this section, we conduct experiments to show
 598 that this heuristic (6.1) is surprisingly effective: the values of the convergence indicator (5.16)
 599 will always remain within the range $(0, 1)$, hence Algorithm 5.1 converges according to
 600 Theorem 5.9.

601 To be more precise, we consider the commonly second-order (*i.e.*, $p = 2$) setting with
 602 Euclidean Lipschitz smoothness Hessians (*i.e.*, $\nu = 1$), and set $h(x; x_0) := \frac{1}{q} \|x - x_0\|_2^q$, where
 603 q is chosen as $q \in \{2, 2.5, 3\} \subset [2, p + \nu]$. Meanwhile, as shown in Theorems 5.7 and 5.9,
 604 the parameter α of Algorithm 5.1 has only a minor influence on performance. Therefore to
 605 simplify our implementation, we always set $\alpha = 1$. By setting $\alpha = 1$, when $p = 2, \nu = 1$,
 606 given \hat{x}_{i-1} and λ_i , the subproblem of finding x_i in the Step 5 of UAF is a standard cubic
 607 regularized Newton step [8]. We solve this subproblem to high accuracy by an implementaion
 608 [20]⁸ of the Lanczos method [9]. Furthermore, in the heuristic (6.1) for A_i , C_0 is determined
 609 by the parameters $p, \nu, q, \theta_1, \theta_2, \beta$ and γ , while we already set the values of p, ν . By the
 610 uniformly convexity of $h(x; x_0) = \frac{1}{q} \|x - x_0\|_2^q$ ($q \geq 2$) [27], we have $\gamma = \beta = 2^{2-q}$.
 611 We simply choose $\theta_1 = 0.5, \theta_2 = 0.67$. The Lipschitz smoothness constant L is tuned in
 612 $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3\}$ to optimize the convergence speed in terms of run time,
 613 while the value of $h(x^*; x_0) = \frac{1}{q} \|x^* - x_0\|_2^q$ is determined by setting $x_0 = 0$ and using an
 614 approximation of x^* to replace x^* .

615 Under the above setting, three instances of the UAF Algorithm 5.1 with $q = 2, 2.5, 3$
 616 respectively will be tested. The instance with $q = 3$ is equivalent to the accelerated cubic
 617 regularized Newton method (ACNM) [27]. For the instance with $q = 2$ or 2.5, we always use
 618 the heuristic (6.1) to determine the values of A_i, a_i and λ_i in each iteration.

619 **6.2. Experiments on Large-Scale Classification Datasets.** To verify the performance
 620 of the proposed UAF and the effectiveness of the heuristic (6.1) in all three instances, we
 621 consider large-scale optimization associated with the logistic regression problem as follows

$$622 \quad (6.2) \quad \min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-\bar{b}_j \bar{a}_j^T x)) \right\},$$

⁸The GitHub URL: https://github.com/dalab/subsampled_cubic_regularization

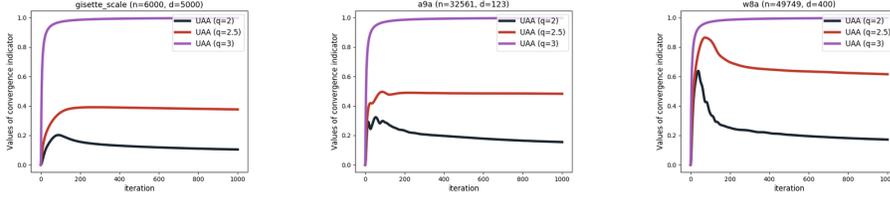


FIG. 1. The values of the convergence indicator (5.16) versus the number of iterations, for the 3 datasets “gisette_scale”, “a9a”, “w8a,” respectively.

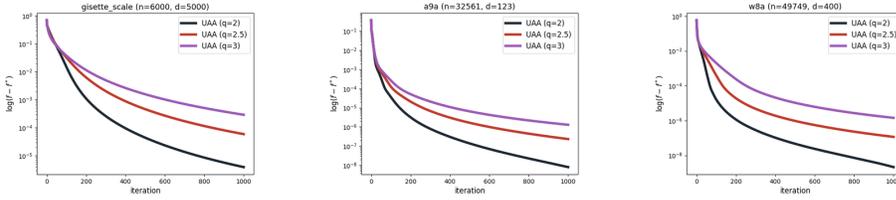


FIG. 2. Accuracy of the objective function (6.2) versus the number of iterations, for the 3 datasets “gisette_scale”, “a9a”, “w8a,” respectively.

623 where $\{(\bar{a}_j, \bar{b}_j)\}_{j=1}^n$ denotes a dataset. (For $j \in [n]$, $\bar{a}_j \in \mathbb{R}^d$ denotes the j -th sample and
 624 $\bar{b}_j \in \{1, -1\}$ denotes the corresponding label of \bar{a}_j .) In our experiments, we choose the
 625 three datasets “gisette_scale”, “a9a” and “w8a” from the LIBSVM library [11] to validate the
 626 performance of our algorithm.

627 In Figure 1, we show the values of the convergence indicator (5.16) of UAF along the
 628 iterations. It is interesting (and somewhat surprising) to see that after several initial steps, the
 629 convergence indicator will approach to a constant in $[0, 1]$. For the case with $q = 3$, *i.e.*, the
 630 ACNM [27], the value of the indicator will approach to 1, which matches the condition (5.21)
 631 with the optimal choice $\theta_1 = \theta_2 = 1$. For the cases with $q = 2$ and 2.5, the values of the
 632 indicator will stay stable around a constant in $(0, 1)$.

633 Because the values of the indicators satisfy the condition (5.21) when $q = 3$ and the
 634 condition (5.25) when $q = 2$ and 2.5, the UAF algorithm will converge according to the
 635 rates in Theorems 5.7 and 5.9 respectively, which is shown in Figure 2. In Figure 2, with the
 636 heuristic (6.1), then the UAF with $q = 2$ has the fastest convergence speed, which matches
 637 the theoretical result that the setting $q = 2$ gives us the best possible iteration complexity
 638 $O\left(k^{-\frac{3(p+\nu)-2}{2}}\right)$.

639 An interesting phenomenon is that the speed edge for the cases $q = 2$ and 2.5 is beyond
 640 our expectation based on the bound (5.27). In the k -th iteration, from the theoretical bound
 641 in Theorems 5.7 and 5.9, the error ratio from the setting $q = 3$ to the setting $q \in [2, p + \nu)$
 642 should be

$$643 \quad (6.3) \quad O\left(\frac{C_0}{\theta_1 c_q \gamma} \left(\frac{k}{(p + \nu)h(x^*; x_0)}\right)^{\frac{p+\nu-q}{q}}\right).$$

644 In the experiments on all the 3 datasets, we found empirically that $h(x^*; x_0) > 1$. Meanwhile,
 645 by simple calculation, we also know that $\frac{\theta_1 c_q \gamma}{C_0} > 1$. Therefore in the 1000-th iteration, by the
 646 theoretical bound (6.3), the error ratio from $q = 3$ to 2 should not go beyond $(\frac{1000}{3})^{\frac{1}{2}} < 20$.
 647 However, in practice the ratio is beyond 100. A possible explanation for this phenomenon
 648 is that even we do not add any strongly convex regularizer in (6.2), the problem itself may
 649 have some kind of local strong convexity around the minimum point (also known as implicit
 650 regularization). Such implicit strong convexity makes the algorithms converges faster as the

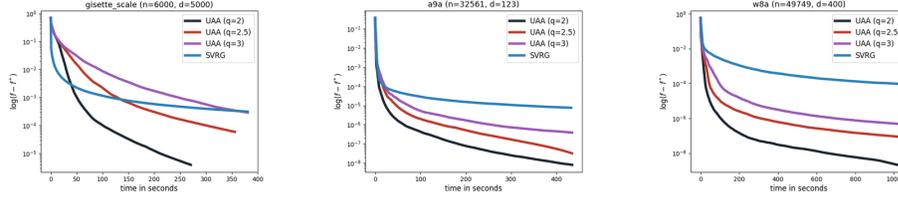


FIG. 3. Accuracy of the objective function (6.2) versus algorithm run time, for the 3 datasets “gisette_scale”, “a9a”, “w8a,” respectively.

651 iterate approaches the minimizer.

652 In Figure 3, we show the performance comparison measured by error versus run time. Here
 653 we add a *stochastic variance reduction gradient* (SVRG) [19] method to show the practical
 654 efficiency of the proposed UAF algorithm. SVRG is a representative first-order algorithm
 655 for finite-sum stochastic convex optimization. The implementation of SVRG is also from the
 656 GitHub project of [20] and the learning rate of SVRG is tuned in $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1},$
 657 $10, 10^2\}$.

658 As shown in Figure 3, SVRG can effectively exploit the finite-sum structure of the
 659 objective (6.2) and shows advantage in obtaining a low-accurate solution quickly. However,
 660 when further pursuing a high-accuracy solution, the high-order UAFs demonstrate clear edges
 661 of their faster convergence rates. In particular, with the effective heuristic (6.1), the UAF with
 662 $q = 2$ demonstrates consistent and superior performance in terms of run time behaviors.

663 **7. Conclusions.** In this paper, inspired by recent work on high-order acceleration meth-
 664 ods, we have introduced a rather unified framework towards developing and understanding
 665 high-order acceleration algorithms for convex optimization. We show how various ideas, tech-
 666 niques, results, and algorithms can be derived from a simple vanilla proximal method (VPM).
 667 Based on this framework, through careful analysis, we are able to derive a unified acceleration
 668 framework (UAF) that achieves the optimal lower bounds for functions that have Hölder
 669 continuous derivatives. Our analysis and results also seem to unify many results known for the
 670 first order and high order methods, as well as results previously obtained through two separate
 671 approaches, namely the ACNM [27] and A-HPE [23] approaches. Meanwhile, the UAF is
 672 the first high-order acceleration approach that can be used in general (non-Euclidean) norm
 673 settings. Furthermore, for practical implementation of the proposed algorithm, through a new
 674 heuristic inspired from our analysis, our experiments show how the binary search procedure
 675 required by the optimal acceleration methods can be significantly simplified or forgone. This
 676 helps alleviate concerns about practical efficiency of optimal high-order acceleration methods
 677 versus suboptimal ones [31]. Finally, combined with a general *restart scheme* similar to that
 678 in [27], our analysis for the general convex setting can be easily extended to the uniformly
 679 convex setting. The resulted complexity results can match the existing lower bounds [2] in
 680 most important cases. Particularly, we shaved off the logarithmic factor of the upper bound in
 681 [2] so that matching the lower bound [2] for the σ -strongly convex minimization problem with
 682 $(2, L)$ -Lipschitz continuous derivatives.⁹

683 Appendix A. Proofs for Section 2.

684 **A.1. Proof of Example 2.4.** By direct computation, for $x, h \in \mathbb{R}^d$, we have

$$685 \quad (\text{A.1}) \quad \langle \nabla^2 f(x)h, h \rangle = \frac{1}{n} \sum_{j=1}^n \frac{\exp(-\bar{b}_j \bar{a}_j^T x)}{(1 + \exp(-\bar{b}_j \bar{a}_j^T x))^2} \langle \bar{a}_j, h \rangle^2 \leq \frac{1}{n} \sum_{j=1}^n \langle \bar{a}_j, h \rangle^2 = h^T B h.$$

⁹See details in the longer arXiv version: <https://arxiv.org/pdf/1906.00582.pdf>.

686 Meanwhile, we have

$$\begin{aligned}
687 \quad \nabla^3 f(x)[h, h, h] &= \frac{1}{n} \sum_{j=1}^n \frac{\exp(-\bar{b}_j \bar{a}_j^T x) (1 - \exp(-\bar{b}_j \bar{a}_j^T x))}{(1 + \exp(-\bar{b}_j \bar{a}_j^T x))^3} \langle \bar{b}_j \bar{a}_j, h \rangle^3 \leq \frac{1}{n} \sum_{j=1}^n |\langle \bar{a}_j, h \rangle|^3 \\
688 \quad (\text{A.2}) \quad &\leq \frac{1}{n} \left(\sum_{j=1}^n \langle \bar{a}_j, h \rangle^2 \right) \max_{j \in [n]} |\langle \bar{a}_j, h \rangle| = h^T B h \cdot \max_{j \in [n]} |\langle \bar{a}_j, h \rangle|.
\end{aligned}$$

689 Therefore,

$$\begin{aligned}
690 \quad \|\nabla^2 f(x)\|_q &= \max_{h \in \mathbb{R}^d: \|h\|_p \leq 1} \langle \nabla^2 f(x) h, h \rangle \leq \max_{h \in \mathbb{R}^d: \|h\|_p \leq 1} \|B h\|_q \leq \|B\|_{p,q}, \\
691 \quad \|\nabla^3 f(x)\|_q &= \max_{h \in \mathbb{R}^d: \|h\|_p \leq 1} \nabla^3 f(x)[h, h, h] \leq \|B\|_{p,q} \max_{j \in [n]} \|\bar{a}_j\|_q.
\end{aligned}$$

693 Let $L(\nu) := \sup_{x, y \in \mathbb{R}^d, x \neq y} \frac{\|\nabla^2 f(x) - \nabla^2 f(y)\|_q}{\|x - y\|_p^\nu}$, $\nu \in [0, 1]$. Then $L(0) = \|\nabla^2 f(x)\|_q$, $L(1) =$

694 $\|\nabla^3 f(x)\|_q$. Note that $L(\nu)$ is log-convex, therefore we have

$$695 \quad (\text{A.3}) \quad L(\nu) \leq L^{1-\nu}(0) L^\nu(1) \leq \|B\|_{p,q} \max_{j \in [n]} \|\bar{a}_j\|_q^\nu.$$

696 Example 2.4 is proved.

697 **A.2. Proof of Lemma 2.5.** By the convexity of $g(x)$, (2.9) holds trivially.

698 If $g(x)$ has p -th derivatives, for $i \in \{0, 1, 2, \dots, p-1\}$, we define a sequence

$$699 \quad (\text{A.4}) \quad C_i := \frac{1}{i!} \int_0^1 (1-\tau)^i \nabla^{i+1} g(y + \tau(x-y)) [x-y]^{i+1} d\tau.$$

700 Then one has

$$\begin{aligned}
701 \quad C_0 &= \int_0^1 \nabla g(y + \tau(x-y)) [x-y] d\tau = \int_0^1 \langle \nabla g(y + \tau(x-y)), x-y \rangle d\tau \\
702 \quad (\text{A.5}) \quad &= g(y + \tau(x-y)) \Big|_{\tau=0}^1 = g(x) - g(y).
\end{aligned}$$

703 Meanwhile,

$$\begin{aligned}
704 \quad C_i &= \frac{1}{i!} \int_0^1 (1-\tau)^i d(\nabla^i g(y + \tau(x-y)) [x-y]^i) \\
705 \quad &= \frac{1}{i!} (\nabla^i g(y + \tau(x-y)) [x-y]^i) \Big|_{\tau=0}^1 \\
706 \quad &\quad - \frac{1}{i!} \int_0^1 (\nabla^i g(y + \tau(x-y)) [x-y]^i) d(1-\tau)^i \\
707 \quad &= -\frac{1}{i!} \nabla^i g(y) [x-y]^i + \frac{1}{(i-1)!} \int_0^1 (1-\tau)^{i-1} (\nabla^i g(y + \tau(x-y)) [x-y]^i) d\tau \\
708 \quad (\text{A.6}) \quad &= -\frac{1}{i!} \nabla^i g(y) [x-y]^i + C_{i-1}.
\end{aligned}$$

709 Therefore by (A.5) and (A.6), one has

$$\begin{aligned}
710 \quad C_{p-1} &= \sum_{i=1}^{p-1} (C_i - C_{i-1}) + C_0 = \sum_{i=1}^{p-1} -\frac{1}{i!} \nabla^i g(y) [x-y]^i + g(x) - g(y) \\
711 \quad &= f(x) - \tilde{f}(x; y) + \frac{1}{p!} \nabla^p g(y) [x-y]^p \\
712 \quad (\text{A.7}) \quad &= f(x) - \tilde{f}(x; y) + \frac{1}{(p-1)!} \nabla^p g(y) [x-y]^p \int_0^1 (1-\tau)^{p-1} d\tau.
\end{aligned}$$

713 Then by (A.7), it follows that

$$\begin{aligned}
714 \quad & |f(x) - \tilde{f}(x; y)| = \left| C_{p-1} - \frac{1}{(p-1)!} \nabla^p g(y) [x-y]^p \int_0^1 (1-\tau)^{p-1} d\tau \right| \\
715 \quad & = \frac{1}{(p-1)!} \left| \int_0^1 (1-\tau)^{p-1} (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [x-y]^p d\tau \right| \\
716 \quad & \leq \frac{1}{(p-1)!} \int_0^1 (1-\tau)^{p-1} d\tau \max_{\tau \in [0,1]} \left| (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [x-y]^p \right| \\
717 \quad & \leq \frac{1}{(p-1)!} \int_0^1 (1-\tau)^{p-1} d\tau \max_{\tau \in [0,1]} \|\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)\|_* \|x-y\|^p \\
718 \quad & \leq \frac{1}{(p-1)!} \frac{1}{p} \max_{\tau \in [0,1]} ((p-1)! L \|\tau(x-y)\|^\nu) \|x-y\|^p \\
719 \quad & \leq \frac{L}{p} \|x-y\|^{p+\nu},
\end{aligned}$$

720 Therefore (2.10) holds. Then by (A.7), by taking gradient *w.r.t.* x , one has

$$\begin{aligned}
721 \quad & \nabla C_{p-1} = \nabla f(x) - \nabla \tilde{f}(x; y) + \frac{1}{(p-1)!} \nabla^p g(y) [x-y]^{p-1} \\
722 \quad (A.8) \quad & = \nabla f(x) - \nabla \tilde{f}(x; y) + \frac{p}{(p-1)!} \nabla^p g(y) [x-y]^{p-1} \int_0^1 (1-\tau)^{p-1} d\tau,
\end{aligned}$$

723 while by (A.4), one also has

$$724 \quad (A.9) \quad \nabla C_{p-1} = \frac{p}{(p-1)!} \int_0^1 (1-\tau)^{p-1} \nabla^p g(y + \tau(x-y)) [x-y]^{p-1} d\tau.$$

725 By (A.8) and (A.9), it follows that

$$\begin{aligned}
726 \quad & \|\nabla f(x) - \nabla \tilde{f}(x; y)\|_* = \left\| \nabla C_{p-1} - \frac{p}{(p-1)!} \nabla^p g(y) [x-y]^{p-1} \int_0^1 (1-\tau)^{p-1} d\tau \right\|_* \\
727 \quad & = \left\| \frac{p}{(p-1)!} \int_0^1 (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [x-y]^{p-1} (1-\tau)^{p-1} d\tau \right\|_* \\
728 \quad & = \max_{v: \|v\| \leq 1} \frac{p}{(p-1)!} \int_0^1 (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [v] [x-y]^{p-1} (1-\tau)^{p-1} d\tau \\
729 \quad & \leq \frac{p}{(p-1)!} \int_0^1 \max_{v: \|v\| \leq 1} (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [v] [x-y]^{p-1} (1-\tau)^{p-1} d\tau \\
730 \quad & \leq \frac{p}{(p-1)!} \int_0^1 (1-\tau)^{p-1} d\tau \max_{\tau \in [0,1]} \max_{v: \|v\| \leq 1} (\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)) [v] [x-y]^{p-1} \\
731 \quad & \leq \frac{p}{(p-1)!} \cdot \frac{1}{p} \cdot \max_{\tau \in [0,1]} \max_{v: \|v\| \leq 1} \|\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)\|_* \cdot \|v\| \cdot \|x-y\|^{p-1} \\
732 \quad & = \frac{p}{(p-1)!} \cdot \frac{1}{p} \cdot \max_{\tau \in [0,1]} \|\nabla^p g(y + \tau(x-y)) - \nabla^p g(y)\|_* \cdot \|x-y\|^{p-1} \\
733 \quad & \leq \frac{p}{(p-1)!} \cdot \frac{1}{p} \cdot \max_{\tau \in [0,1]} (p-1)! L \|\tau(x-y)\|^\nu \cdot \|x-y\|^{p-1} \\
734 \quad &
\end{aligned}$$

735 Then $\|\nabla f(x) - \nabla \tilde{f}(x; y)\|_* \leq L \|x-y\|^{p+\nu-1}$, *i.e.*, (2.11) holds. Lemma 2.5 is proved.

736 **A.3. Proof of Lemma 2.6.** For the first statement, by the condition we have $b_k - b_{k-1} \geq$
 737 $C^{\frac{1}{\rho}} b_k^{\frac{\rho-1}{\rho}}$. Then by $b_0 = 0$,

$$738 \quad b_k = \sum_{i=1}^k (b_i - b_{i-1}) \geq C^{\frac{1}{\rho}} \sum_{i=1}^k b_i^{\frac{\rho-1}{\rho}}.$$

739 Then in [6, Lemma 12], for $i \geq 1$, by setting $B_i := b_i^{\frac{\rho-1}{\rho}}$, $\alpha := \frac{\rho}{\rho-1}$, $c := C^{\frac{1}{\rho}}$, then one has
 740 $b_k^{\frac{\rho-1}{\rho}} = B_k \geq \left(\frac{1}{\rho} C^{\frac{1}{\rho}} k\right)^{\rho-1}$.

741 Then after a simple rearrangement, we obtain the first statement.

742 For the second statement, by the reverse Hölder inequality, $\|fg\|_1 \geq \|f\|_{\frac{1}{t}} \|g\|_{-\frac{1}{t-1}}$ for
 743 $t \geq 1$ and invoking this with $t = \rho\delta + 1$ and by $b_0 = 0$, then

$$744 \quad \sum_{i=1}^k \left(\frac{b_i^{\rho-1}}{(b_i - b_{i-1})^\rho}\right)^\delta = \sum_{i=1}^k b_i^{(\rho-1)\delta} (b_i - b_{i-1})^{-\rho\delta}$$

$$745 \quad \geq \left(\sum_{i=1}^k b_i^{(\rho-1)\delta \cdot \frac{1}{t}}\right)^t \left(\sum_{i=1}^k (b_i - b_{i-1})^{-\rho\delta \cdot \frac{1}{t-1}}\right)^{-(t-1)}$$

$$746 \quad = \left(\sum_{i=1}^k b_i^{\frac{(\rho-1)\delta}{\rho\delta+1}}\right)^{\rho\delta+1} \left(\sum_{i=1}^k (b_i - b_{i-1})\right)^{-\rho\delta} = \left(\sum_{i=1}^k b_i^{\frac{(\rho-1)\delta}{\rho\delta+1}}\right)^{\rho\delta+1} b_k^{-\rho\delta}.$$

747 Then by the corresponding condition, we have $b_k^{\frac{\rho\delta}{\rho\delta+1}} \geq C^{-\frac{1}{\rho\delta+1}} \left(\sum_{i=1}^k b_i^{\frac{(\rho-1)\delta}{\rho\delta+1}}\right)$. Then in [6,
 748 Lemma 12], for $i \geq 1$, by setting $B_i := b_i^{\frac{(\rho-1)\delta}{\rho\delta+1}}$, $\alpha := \frac{\rho}{\rho-1}$, $c := C^{-\frac{1}{\rho\delta+1}}$, one has

$$749 \quad b_k^{\frac{(\rho-1)\delta}{\rho\delta+1}} = B_k \geq \left(\frac{1}{\rho} C^{-\frac{1}{\rho\delta+1}} k\right)^{\rho-1}.$$

750 Then after a simple rearrangement, we obtain the second statement. Lemma 2.6 is proved.

751 Appendix B. Proofs for Section 4.

752 **B.1. Proof of Lemma 4.2.** By the optimality condition of $z_t := \operatorname{argmin}_{x \in \mathbb{R}^d} \psi_t^{\operatorname{cont}}(x)$,
 753 one has $\left\langle \int_0^t a_\tau \nabla \hat{f}(z_t; x_\tau) d\tau + \nabla h(z_t; x_0), \dot{z}_t \right\rangle \geq 0$. It follows that

$$754 \quad \frac{d}{dt} \left(\min_{x \in \mathbb{R}^d} \psi_t^{\operatorname{cont}}(x) \right) = \frac{d}{dt} \psi_t^{\operatorname{cont}}(z_t) = \frac{d}{dt} \left(\int_0^t a_\tau \hat{f}(z_t; x_\tau) d\tau + h(z_t; x_0) \right)$$

$$755 \quad = a_t \hat{f}(z_t; x_t) + \left\langle \int_0^t a_\tau \nabla \hat{f}(z_t; x_\tau) d\tau + \nabla h(z_t; x_0), \dot{z}_t \right\rangle$$

$$756 \quad \geq a_t (\hat{f}(x_t; x_t) + \langle \nabla \hat{f}(x_t; x_t), z_t - x_t \rangle) + \left\langle \int_0^t a_\tau \nabla \hat{f}(z_t; x_\tau) d\tau + \nabla h(z_t; x_0), \dot{z}_t \right\rangle$$

$$757 \quad \stackrel{\textcircled{1}}{=} a_t (f(x_t) + \langle \nabla f(x_t), z_t - x_t \rangle) + \left\langle \int_0^t a_\tau \nabla \hat{f}(z_t; x_\tau) d\tau + \nabla h(z_t; x_0), \dot{z}_t \right\rangle$$

$$758 \quad \text{(B.1)} \geq a_t (f(x_t) + \langle \nabla f(x_t), z_t - x_t \rangle),$$

759 where $\textcircled{1}$ is by the definition of $\hat{f}(x; y)$ in (2.7). Furthermore, one has

$$760 \quad \text{(B.2)} \quad \frac{d(A_t f(x_t))}{dt} = a_t f(x_t) + A_t \langle \nabla f(x_t), \dot{x}_t \rangle.$$

761

762 By Combing (B.1) and (B.2), one has

$$763 \quad (\text{B.3}) \quad \frac{d}{dt} \left(A_t f(x_t) - \min_{x \in \mathbb{R}^d} \psi_t^{\text{cont}}(x) \right) \leq \langle \nabla f(x_t), A_t \dot{x}_t - a_t(z_t - x_t) \rangle.$$

764 Meanwhile by $A_0 = 0$ and $\min_{x \in \mathbb{R}^d} \psi_0(x) = 0$, one has $A_0 f(x_0) - \min_{x \in \mathbb{R}^d} \psi_0^{\text{cont}}(x) = 0$.
765 Taking integral from $\tau = 0$ to t for (B.3), then Lemma 4.2 is proved.

766 Appendix C. Proofs for Section 5.

767 **C.1. Proof of Lemma 5.3.** First, in (5.1), by $A_0 = 0$ and $z_0 = x_0$, we have

$$768 \quad (\text{C.1}) \quad A_0 f(x_0) - \psi_0^{\text{dis}}(z_0) = 0.$$

769 Then by our assumption, $\hat{f}(x; x_i)$ is convex *w.r.t.* $\|\cdot\|$ and $h(x; x_0)$ is (q, γ) -uniformly convex
770 *w.r.t.* $\|\cdot\|$. Therefore for all $x, y \in \mathbb{R}^d$, it follows that

$$771 \quad (\text{C.2}) \quad \psi_i^{\text{dis}}(x) \geq \psi_i^{\text{dis}}(y) + \langle \nabla \psi_i^{\text{dis}}(y), x - y \rangle + \frac{\gamma}{q} \|x - y\|^q.$$

773 Then by the optimality condition of z_i , it follows that for all $x \in \mathbb{R}^d$, $\langle \nabla \psi_i^{\text{dis}}(z_i), x - z_i \rangle \geq 0$.
774 Therefore, it follows that $\psi_i^{\text{dis}}(x) \geq \psi_i^{\text{dis}}(z_i) + \frac{\gamma}{q} \|x - z_i\|^q$. Therefore we have,

$$775 \quad (\text{C.3}) \quad \psi_i^{\text{dis}}(x) = \psi_{i-1}^{\text{dis}}(x) + a_i \hat{f}(x; x_i) \geq \psi_{i-1}^{\text{dis}}(z_{i-1}) + \frac{\gamma}{q} \|x - z_{i-1}\|^q + a_i \hat{f}(x; x_i).$$

776 Meanwhile, we can lower bound the last term of RHS of (C.3).

$$\begin{aligned} 777 \quad a_i \hat{f}(x; x_i) &\stackrel{\textcircled{1}}{\geq} a_i (\hat{f}(x_i; x_i) + \langle \nabla \hat{f}(x_i; x_i), x - x_i \rangle) \stackrel{\textcircled{2}}{=} a_i (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) \\ 778 &\stackrel{\textcircled{3}}{=} A_i (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) - A_{i-1} (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) \\ 779 &= A_i \left(f(x_i) + \left\langle \nabla f(x_i), \frac{a_i}{A_i} x + \frac{A_{i-1}}{A_i} x_{i-1} - x_i \right\rangle \right) \\ 780 &\quad - A_{i-1} (f(x_i) + \langle \nabla f(x_i), x_{i-1} - x_i \rangle) \\ 781 &\stackrel{\textcircled{4}}{\geq} A_i \left(f(x_i) + \left\langle \nabla f(x_i), \frac{a_i}{A_i} x + \frac{A_{i-1}}{A_i} x_{i-1} - x_i \right\rangle \right) - A_{i-1} f(x_{i-1}) \\ 782 &= A_i f(x_i) - A_{i-1} f(x_{i-1}) + A_i \left\langle \nabla f(x_i), \frac{a_i}{A_i} x + \frac{A_{i-1}}{A_i} x_{i-1} - x_i \right\rangle, \end{aligned}$$

783 where $\textcircled{1}$ is by the convexity of $\hat{f}(x; y)$ *w.r.t.* x , $\textcircled{2}$ is by the definition of $\hat{f}(x; y)$ in (2.7), $\textcircled{3}$
784 is by the identity $a_i = A_i - A_{i-1}$, and $\textcircled{4}$ is by the convexity of $f(x)$.

785 Therefore it follows that

$$\begin{aligned} 786 \quad \psi_i^{\text{dis}}(x) &\geq \psi_{i-1}^{\text{dis}}(z_{i-1}) + \frac{\gamma}{q} \|x - z_{i-1}\|^q + A_i f(x_i) - A_{i-1} f(x_{i-1}) \\ 787 \quad (\text{C.4}) &\quad + A_i \left\langle \nabla f(x_i), \frac{a_i}{A_i} x + \frac{A_{i-1}}{A_i} x_{i-1} - x_i \right\rangle. \end{aligned}$$

788 By setting $x := z_i$ and a simple arrangement of (C.4), we have

$$\begin{aligned} 789 &\quad (A_i f(x_i) - \psi_i^{\text{dis}}(z_i)) - (A_{i-1} f(x_{i-1}) - \psi_{i-1}^{\text{dis}}(z_{i-1})) \\ 790 \quad (\text{C.5}) &\leq A_i \left\langle \nabla f(x_i), x_i - \frac{a_i}{A_i} z_i - \frac{A_{i-1}}{A_i} x_{i-1} \right\rangle - \frac{\gamma}{q} \|z_i - z_{i-1}\|^q \end{aligned}$$

791 Summing (C.5) from $i = 0$ to $k - 1$ and by (C.1), it follows that

$$\begin{aligned}
792 \quad A_k f(x_k) - \psi_k^{\text{dis}}(z_k) &\leq A_0 f(x_0) - \psi_0^{\text{dis}}(z_0) \\
793 \quad &+ \sum_{i=1}^k \left(A_i \left\langle \nabla f(x_i), x_i - \frac{a_i}{A_i} z_i - \frac{A_{i-1}}{A_i} x_{i-1} \right\rangle - \frac{\gamma}{q} \|z_i - z_{i-1}\|^q \right) \\
794 \quad &= \sum_{i=1}^k \left(A_i \left\langle \nabla f(x_i), x_i - \frac{a_i}{A_i} z_i - \frac{A_{i-1}}{A_i} x_{i-1} \right\rangle - \frac{\gamma}{q} \|z_i - z_{i-1}\|^q \right).
\end{aligned}$$

795 Then by the definition of E_i , Lemma 5.3 is proved.

796 **C.2. Proof of Lemma 5.4.** By the definition of E_i , one has

$$\begin{aligned}
797 \quad E_i &\stackrel{\textcircled{1}}{\leq} a_i \langle \nabla f(x_i), \hat{z}_i - z_i \rangle - \frac{\gamma}{q} \|z_i - z_{i-1}\|^q \stackrel{\textcircled{2}}{\leq} a_i \langle \nabla f(x_i), \hat{z}_i - z_i \rangle - \frac{\gamma'_i}{q} \|z_i - z_{i-1}\|^q \\
798 \quad &\stackrel{\textcircled{3}}{\leq} \left\langle a_i \nabla f(x_i) + \gamma'_i \nabla \frac{1}{q} \|\hat{z}_i - z_{i-1}\|^q, \hat{z}_i - z_i \right\rangle - \gamma'_i \left(\frac{1}{q} \|\hat{z}_i - z_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right) \\
799 \quad &\stackrel{\textcircled{4}}{=} \left\langle a_i \nabla f(x_i) + \frac{\gamma'_i A_i^{q-1}}{a_i^{q-1}} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q, \hat{z}_i - z_i \right\rangle - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right)
\end{aligned}$$

800 where $\textcircled{1}$ is (5.3), $\textcircled{2}$ is by $\gamma \geq \gamma'_i$, $\textcircled{3}$ is by Assumption 5.1, $\textcircled{4}$ is (5.3) such that $\nabla \|\hat{z}_i - z_i\|^q =$
801 $\frac{A_i^{q-1}}{a_i^{q-1}} \nabla \|x_i - \hat{x}_{i-1}\|^q$. Therefore Lemma 5.4 is proved.

802 **C.3. Proof of Lemma 5.5.** By Lemma 5.4, one has

$$\begin{aligned}
803 \quad E_i &\leq a_i \left\langle \nabla f(x_i) + \frac{\gamma'_i A_i^{q-1}}{a_i^{q-1}} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q, \hat{z}_i - z_i \right\rangle - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right) \\
804 \quad &\leq a_i \langle \nabla f(x_i) - \nabla \tilde{f}(x_i; \hat{x}_{i-1}), \hat{z}_i - z_i \rangle + a_i \left\langle \nabla \tilde{f}(x_i; \hat{x}_{i-1}) + \frac{\gamma'_i A_i^{q-1}}{a_i^{q-1}} \nabla \frac{1}{q} \|x_i - \hat{x}_{i-1}\|^q, \hat{z}_i - z_i \right\rangle \\
805 \quad &\quad \text{(C.6)} \quad - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right).
\end{aligned}$$

806 Meanwhile, it follows that

$$\begin{aligned}
807 \quad &a_i \langle \nabla f(x_i) - \nabla \tilde{f}(x_i; \hat{x}_{i-1}), \hat{z}_i - z_i \rangle - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right) \\
808 \quad &\leq a_i \|\nabla f(x_i) - \nabla \tilde{f}(x_i; \hat{x}_{i-1})\|_* \|\hat{z}_i - z_i\| - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right) \\
809 \quad &\stackrel{\textcircled{1}}{\leq} a_i L \|x_i - \hat{x}_{i-1}\|^{p+\nu-1} \|\hat{z}_i - z_i\| - \gamma'_i \left(\frac{A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q + \frac{\beta}{q} \|\hat{z}_i - z_i\|^q \right) \\
810 \quad &\stackrel{\textcircled{2}}{\leq} \frac{q-1}{q} (\beta \gamma'_i)^{-\frac{1}{q-1}} (a_i L)^{\frac{q}{q-1}} \|x_i - \hat{x}_{i-1}\|^{\frac{q(p+\nu-1)}{q-1}} - \frac{\gamma'_i A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q \\
811 \quad &\stackrel{\textcircled{3}}{=} \left(\left(L \frac{a_i^q}{c_q \gamma'_i A_i^{q-1}} \|x_i - \hat{x}_{i-1}\|^{p+\nu-q} \right)^{\frac{q}{q-1}} - 1 \right) \frac{\gamma'_i A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q \\
812 \quad \text{(C.7)} \quad &\stackrel{\textcircled{4}}{=} \left((L \lambda'_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q})^{\frac{q}{q-1}} - 1 \right) \frac{\gamma'_i A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q,
\end{aligned}$$

813 where $\textcircled{1}$ is by (5.6), $\textcircled{2}$ is by [27, Lemma 1.3], $\textcircled{3}$ is by a simple rearrangement and the
814 definition of c_q in Lemma 5.5, and $\textcircled{4}$ is by the definition of λ'_i .

815 Combing (C.6) and (C.7), Lemma 5.5 is proved.

816 **C.4. Proof of Theorem 5.6.** First, for $i \geq 1$, if the conditions (5.15) and (5.16) are true,
 817 then one can know that (5.8) and (5.9) are true and thus for $i \geq 1$, $E_i \leq 0$. Then by Lemma
 818 5.3, one has $A_k f(x_k) - \psi_k^{\text{dis}}(z_k) \leq \sum_{i=1}^k E_i \leq 0$. Then combing Lemma 5.2, one has

$$819 \quad (\text{C.8}) \quad A_k f(x_k) \leq \psi_k^{\text{dis}}(z_k) \leq A_k f(x^*) + h(x^*; x_0).$$

820 Theorem 5.6 is proved.

821 **C.5. Proof of Theorem 5.7.** First, by our assumption, $\{\lambda_i\}$ defined in (5.11) satisfies
 822 (5.21), therefore $\{\lambda_i\}$ satisfies (5.16); meanwhile $\{x_i\}$ satisfies (5.15). Therefore Theorem
 823 5.6 holds, *i.e.*, $f(x_k) - f(x^*) \leq \frac{h(x^*; x_0)}{A_k}$. Then by (5.21), because $L\lambda_i = \frac{L a_i^q}{c_q \gamma A_i^{q-1}} =$
 824 $\frac{L(A_i - A_{i-1})^q}{c_q \gamma A_i^{q-1}} \geq \theta_1$, in Lemma 2.6, by setting $b_i := A_i, \rho := p + \nu, C := \frac{\theta_1 c_q \gamma}{L}$, we can
 825 obtain the lower bound $A_k \geq \frac{\theta_1 c_q \gamma}{L} \left(\frac{k}{p+\nu}\right)^{p+\nu}$. Finally, (5.23) is obtained and Theorem 5.7
 826 is proved.

827 **C.6. Proof of Lemma 5.8.** When (5.15) and (5.16) are satisfied, by Lemma 5.5, we have

$$\begin{aligned} 828 \quad E_i &\stackrel{\textcircled{1}}{\leq} (\theta_2^{\frac{q}{q-1}} - 1) \frac{\gamma'_i A_i^q}{q a_i^q} \|x_i - \hat{x}_{i-1}\|^q \stackrel{\textcircled{2}}{\leq} (\theta_2^{\frac{q}{q-1}} - 1) \frac{A_i^q}{q a_i^q} \left(\frac{\omega_i}{\theta_2}\right)^\alpha \gamma \|x_i - \hat{x}_{i-1}\|^q \\ 829 \quad &\stackrel{\textcircled{3}}{=} (\theta_2^{\frac{q}{q-1}} - 1) \frac{A_i^q}{q a_i^q} \left(\frac{\omega_i}{\theta_2}\right)^\alpha \gamma (L\lambda_i \|x_i - \hat{x}_{i-1}\|^{p+\nu-q})^{\frac{q}{p+\nu-q}} (L\lambda_i)^{-\frac{q}{p+\nu-q}} \\ 830 \quad &\stackrel{\textcircled{4}}{=} \frac{1}{q\theta_2^\alpha} (\theta_2^{\frac{q}{q-1}} - 1) \frac{A_i^q}{a_i^q} \omega_i^{\frac{\varsigma}{p+\nu-q}} \gamma (L\lambda_i)^{-\frac{q}{p+\nu-q}} \\ 831 \quad (\text{C.9}) \quad &\stackrel{\textcircled{5}}{=} \frac{1}{q\theta_2^\alpha} (\theta_2^{\frac{q}{q-1}} - 1) \omega_i^{\frac{\varsigma}{p+\nu-q}} \gamma \left(\frac{c_q \gamma A_i^{p+\nu-1}}{L(A_i - A_{i-1})^{p+\nu}}\right)^{\frac{q}{p+\nu-q}}, \end{aligned}$$

832 where $\textcircled{1}$ is by (5.15) and (5.16), $\textcircled{2}$ is by the value of γ'_i in (5.12) and the definition of ω_i in
 833 Lemma 5.8, $\textcircled{3}$ is by a simple rearrangement, $\textcircled{4}$ is by definition of ω_i and $\varsigma = \alpha(p + \nu) +$
 834 $(1 - \alpha)q$, $\textcircled{5}$ is by definition of λ_i in (5.11) and the fact $a_i = A_i - A_{i-1}$.

835 Then by combing Lemmas 5.2 and 5.3, it follows that

$$836 \quad (\text{C.10}) \quad A_k f(x_k) \leq \psi_k^{\text{dis}}(z_k) + \sum_{i=1}^k E_i \leq A_k f(x^*) + h(x^*; x_0) + \sum_{i=1}^k E_i.$$

837 Then by combing (C.9) and (C.10), and $f(x_k) \geq f(x^*)$, $A_k \geq 0$, one has

$$838 \quad \frac{1}{q\theta_2^\alpha} (1 - \theta_2^{\frac{q}{q-1}}) \gamma \sum_{i=1}^k \omega_i^{\frac{\varsigma}{p+\nu-q}} \left(\frac{c_q \gamma A_i^{p+\nu-1}}{L(A_i - A_{i-1})^{p+\nu}}\right)^{\frac{q}{p+\nu-q}} \leq \sum_{i=1}^k -E_i \leq h(x^*; x_0).$$

839 Then after a simple rearrangement, we have Lemma 5.8.

840 **C.7. Proof of Theorem 5.9.** First, by our assumption, $\{\lambda_i\}$ defined in (5.11) satisfies
 841 (5.25), therefore $\{\lambda_i\}$ satisfies (5.16); meanwhile $\{x_i\}$ satisfies (5.15). Therefore Theorem
 842 5.6 holds, *i.e.*,

$$843 \quad (\text{C.11}) \quad f(x_k) - f(x^*) \leq \frac{h(x^*; x_0)}{A_k}.$$

844

845 Then by Lemma 5.8 and the assumption that $\omega_i \geq \theta_1$, we have

$$846 \quad (C.12) \quad \sum_{i=1}^k \left(\frac{A_i^{p+\nu-1}}{(A_i - A_{i-1})^{p+\nu}} \right)^{\frac{q}{p+\nu-q}} \leq (C_0^{-1}L)^{\frac{q}{p+\nu-q}} h(x^*; x_0),$$

847 where C_0 is defined in Theorem 5.9.

848 In Lemma 2.6, for $1 \leq i \leq k$, by setting $b_i := A_i$, $\rho := p + \nu$, $\delta := \frac{q}{p+\nu-q}$, $C :=$
849 $(C_0^{-1}L)^{\frac{q}{p+\nu-q}} h(x^*; x_0)$, then we obtain the lower bound

$$850 \quad (C.13) \quad A_k \geq \frac{C_0}{L} (h(x^*; x_0))^{-\frac{p+\nu-q}{q}} \left(\frac{k}{p+\nu} \right)^{\frac{(q+1)(p+\nu)-q}{q}}.$$

851 Then combing (C.11), we obtain (5.27).

852 **C.8. Proof of Proposition 5.10.** First by our assumption about $\|\cdot\|$ and $\tilde{f}(x; y)$, (5.31)
853 is a strictly convex function, therefore $w(v)$ is a continuous function of v . Meanwhile $x(\lambda)$ is
854 continuous about λ . Therefore $\chi(\lambda)$ is continuous *w.r.t.* λ .

855 Next by the fact $\varsigma = \alpha(p + \nu) + (1 - \alpha)q \in [q, p + \nu]$ and

$$856 \quad (C.14) \quad \tilde{f}(z; v) + \frac{L^\alpha}{c_q \lambda^{(1-\alpha)\theta_2^\alpha} \varsigma} \|z - v\|^\varsigma \leq \tilde{f}(v; v) = f(v) < +\infty,$$

857 as $\lambda \rightarrow 0$, $\|z - v\| \rightarrow 0$ if $\varsigma \in [q, p + \nu]$ or $\|z - v\|$ is a finite value if $\varsigma = p + \nu$. In both
858 cases, we have $\chi(0) = 0$. Then since $f(v) \neq f(x^*)$, we will also have as $\lambda \rightarrow +\infty$, it is
859 easy to find that $\frac{\alpha(\lambda)}{A + \alpha(\lambda)} \rightarrow 1$ and thus $x(\lambda) = x$. Since $f(x) \neq f(x^*)$, we have $\omega(x) \neq x$.
860 Therefore $\chi(+\infty) = +\infty$.

861

REFERENCES

- 862 [1] Z. ALLEN-ZHU, *Katyusha: The first direct acceleration of stochastic gradient methods*, The Journal of
863 Machine Learning Research, 18 (2017), pp. 8194–8244.
- 864 [2] Y. ARJEVANI, O. SHAMIR, AND R. SHIFF, *Oracle complexity of second-order methods for smooth convex*
865 *optimization*, Math. Program., (2017), pp. 1–34.
- 866 [3] M. BAESE, *Estimate sequence methods: extensions and approximations*, manuscript, (2009).
- 867 [4] K. BALL, E. A. CARLEN, AND E. H. LIEB, *Sharp uniform convexity and smoothness inequalities for trace*
868 *norms*, Inventiones mathematicae, 115 (1994), pp. 463–482.
- 869 [5] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*,
870 SIAM journal on imaging sciences, 2 (2009), pp. 183–202.
- 871 [6] S. BUBECK, Q. JIANG, Y. T. LEE, Y. LI, AND A. SIDFORD, *Near-optimal method for highly smooth convex*
872 *optimization*, in Proceedings of the Thirty-Second Conference on Learning Theory, vol. 99 of Proc. Mach.
873 Learn. Res, Phoenix, USA, 2019, pp. 492–507, <http://proceedings.mlr.press/v99/bubeck19a.html>.
- 874 [7] B. BULLINS, *Fast minimization of structured convex quartics*, arXiv preprint arXiv:1812.10349, (2018).
- 875 [8] Y. CARMON AND J. DUCHI, *Gradient descent finds the cubic-regularized nonconvex newton step*, SIAM J.
876 Optim., 29 (2019), pp. 2146–2178, <https://doi.org/10.1137/17M1113898>.
- 877 [9] Y. CARMON AND J. C. DUCHI, *Analysis of krylov subspace solutions of regularized nonconvex quadratic*
878 *problems*, in Proceedings of the 32Nd International Conference on Neural Information Processing Sys-
879 tems (NeurIPS), Montr al, Canada, 2018, pp. 10728–10738, [http://dl.acm.org/citation.cfm?id=3327546.](http://dl.acm.org/citation.cfm?id=3327546.3327730)
880 [3327730](http://dl.acm.org/citation.cfm?id=3327546.3327730).
- 881 [10] C. CARTIS, N. I. GOULD, AND P. L. TOINT, *Universal regularization methods: varying the power, the*
882 *smoothness and the accuracy*, SIAM J. Optim., 29 (2019), pp. 595–615.
- 883 [11] C.-C. CHANG AND C.-J. LIN, *LIBSVM: A library for support vector machines*, ACM T. Intel. Syst. Tec., 2
884 (2011), pp. 27:1–27:27. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- 885 [12] J. DIAKONIKOLAS AND L. ORECCHIA, *Accelerated extra-gradient descent: A novel accelerated first-order*
886 *method*, in 9th Innovations in Theoretical Computer Science Conference (ITCS 2018), Schloss Dagstuhl-
887 Leibniz-Zentrum fuer Informatik, 2018.

- 888 [13] J. DIAKONIKOLAS AND L. ORECCHIA, *The approximate duality gap technique: A unified theory of first-order*
889 *methods*, SIAM J. Optim., 29 (2019), pp. 660–689.
- 890 [14] A. GASNIKOV, P. DVURECHENSKY, E. GORBUNOV, E. VORONTSOVA, D. SELIKHANOVYCH, AND C. A.
891 URIBE, *Optimal tensor methods in smooth convex and uniformly convex optimization*, in Proceedings of
892 the Thirty-Second Conference on Learning Theory, vol. 99 of Proc. Mach. Learn. Res., Phoenix, USA,
893 2019, PMLR, pp. 1374–1391.
- 894 [15] G. N. GRAPIGLIA AND Y. NESTEROV, *Accelerated regularized newton methods for minimizing composite*
895 *convex functions*, SIAM J. Optim., 29 (2019), pp. 77–99.
- 896 [16] G. N. GRAPIGLIA AND Y. NESTEROV, *Tensor methods for minimizing functions with Hölder continuous*
897 *higher-order derivatives*, arXiv preprint arXiv:1904.12559, (2019).
- 898 [17] HAIRER AND PETERS, *Solving ordinary differential equations I*, Springer Berlin Heidelberg, 1987.
- 899 [18] B. JIANG, H. WANG, AND S. ZHANG, *An optimal high-order tensor method for convex optimization*, in
900 Proceedings of the Thirty-Second Conference on Learning Theory, vol. 99 of Proc. Mach. Learn. Res.,
901 Phoenix, USA, 2019, PMLR, pp. 1799–1801, <http://proceedings.mlr.press/v99/jiang19a.html>.
- 902 [19] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*, in
903 Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1,
904 Lake Tahoe, USA, 2013, pp. 315–323, <http://dl.acm.org/citation.cfm?id=2999611.2999647>.
- 905 [20] J. M. KOHLER AND A. LUCCHI, *Sub-sampled cubic regularization for non-convex optimization*, in Proceedings
906 of the 34th International Conference on Machine Learning (ICML), vol. 70, 2017, pp. 1895–1904.
- 907 [21] W. KRICHENE, A. BAYEN, AND P. L. BARTLETT, *Accelerated mirror descent in continuous and discrete*
908 *time*, in Advances in neural information processing systems, 2015, pp. 2845–2853.
- 909 [22] W. KRICHENE, A. BAYEN, AND P. L. BARTLETT, *Adaptive averaging in accelerated descent dynamics*, in
910 Advances in Neural Information Processing Systems (NeurIPS), 2016, pp. 2991–2999.
- 911 [23] R. D. MONTEIRO AND B. F. SVAITER, *An accelerated hybrid proximal extragradient method for convex*
912 *optimization and its implications to second-order methods*, SIAM J. Optim., 23 (2013), pp. 1092–1125.
- 913 [24] A. S. NEMIROVSKII AND Y. NESTEROV, *Optimal methods of smooth convex minimization*, USSR Computa-
914 tional Mathematics and Mathematical Physics, 25 (1985), pp. 21–30.
- 915 [25] A. S. NEMIROVSKY AND D. B. YUDIN, *Problem complexity and method efficiency in optimization.*, (1983).
- 916 [26] Y. NESTEROV, *A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$*
917 *(in Doklady AN USSR, vol. 269, 1983, pp. 543–547.*
- 918 [27] Y. NESTEROV, *Accelerating the cubic regularization of Newton’s method on convex problems*, Math. Program.,
919 112 (2008), pp. 159–181.
- 920 [28] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Springer Publishing Company,
921 Incorporated, 1 ed., 2014.
- 922 [29] Y. NESTEROV, *Universal gradient methods for convex optimization problems*, Math. Program., 152 (2015),
923 pp. 381–404.
- 924 [30] Y. NESTEROV, *Complexity bounds for primal-dual methods minimizing the model of objective function*, Math.
925 Program., 171 (2018), pp. 311–330.
- 926 [31] Y. NESTEROV, *Implementable tensor methods in unconstrained convex optimization*, Universite catholique de
927 Louvain, Center for Operations Research and Econometrics (CORE), (2018).
- 928 [32] Y. NESTEROV AND B. T. POLYAK, *Cubic regularization of Newton method and its global performance*, Math.
929 Program., 108 (2006), pp. 177–205.
- 930 [33] N. PARIKH, S. BOYD, ET AL., *Proximal algorithms*, Foundations and Trends® in Optimization, 1 (2014),
931 pp. 127–239.
- 932 [34] V. ROULET AND A. D’ASPREMONT, *Sharpness, restart and acceleration*, in Proceedings of the 31st Inter-
933 national Conference on Neural Information Processing Systems (NeurIPS), Long Beach, USA, 2017,
934 pp. 1119–1129, <http://dl.acm.org/citation.cfm?id=3294771.3294878>.
- 935 [35] W. SU, S. BOYD, AND E. J. CANDÈS, *A differential equation for modeling nesterov’s accelerated gradient*
936 *method: Theory and insights*, in Proceedings of the 27th International Conference on Neural Information
937 Processing Systems (NeurIPS), Montreal, Canada, 2014, pp. 2510–2518, [http://dl.acm.org/citation.cfm?](http://dl.acm.org/citation.cfm?id=2969033.2969107)
938 [id=2969033.2969107](http://dl.acm.org/citation.cfm?id=2969033.2969107).
- 939 [36] A. WIBISONO, A. C. WILSON, AND M. I. JORDAN, *A variational perspective on accelerated methods in*
940 *optimization*, Proc. Natl. Acad. Sci., 113 (2016), pp. E7351–E7358.
- 941 [37] A. YURTSEVER, Q. TRAN-DINH, AND V. CEVHER, *A universal primal-dual convex optimization framework*,
942 in Proceedings of the 28th International Conference on Neural Information Processing Systems (NeurIPS),
943 Montreal, Canada, 2015, pp. 3150–3158, <http://dl.acm.org/citation.cfm?id=2969442.2969591>.