
Practical Schemes for Finding Near-Stationary Points of Convex Finite-Sums

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The problem of finding near-stationary points in convex optimization has not been
2 adequately studied yet, unlike other optimality measures such as the function
3 value. Even in the deterministic case, the optimal method (OGM-G, due to Kim
4 and Fessler [33]) has just been discovered recently. In this work, we conduct a
5 systematic study of algorithmic techniques for finding near-stationary points of
6 convex finite-sums. Our main contributions are several algorithmic discoveries:
7 (1) we discover a memory-saving variant of OGM-G based on the performance
8 estimation problem approach [19]; (2) we design a new accelerated SVRG variant
9 that can simultaneously achieve fast rates for minimizing both the gradient norm
10 and function value; (3) we propose an adaptively regularized accelerated SVRG
11 variant, which does not require the knowledge of some unknown initial constants
12 and achieves near-optimal complexities. We put an emphasis on the simplicity and
13 practicality of the new schemes, which could facilitate future developments.

14 1 Introduction

15 Classic convex optimization usually focuses on providing guarantees for minimizing function value.
16 For this task, the optimal (up to constant factors) Nesterov’s accelerated gradient method (NAG)
17 [40, 41] has been known for decades, and there are even methods that can exactly match the lower
18 complexity bounds [30, 17, 55, 18]. On the other hand, in general non-convex optimization, near-
19 stationarity is the typical optimality measure, and there has been a flurry of recent research devoted to
20 this topic [25, 26, 23, 28, 21, 60]. Recently, there has been growing interest on devising fast schemes
21 for finding near-stationary points in convex optimization [42, 2, 22, 7, 31, 32, 33, 27, 15, 14]. This
22 line of research is basically driven by the following facts.

- 23 • Nesterov [42] studied the problem with a linear constraint: $f(x^*) = \min_{x \in Q} \{f(x) : Ax = b\}$,
24 where Q is a convex set and f is strongly convex. Assuming that Q and f are simple, we can focus
25 on the dual problem $\phi(y^*) = \max_y \{\phi(y) \triangleq \min_{x \in Q} \{f(x) + \langle y, b - Ax \rangle\}\}$. Clearly, the dual
26 objective $-\phi(y)$ is smooth convex. Letting x_y be the unique solution to the inner problem, we have
27 $\nabla \phi(y) = b - Ax_y$. Note that $f(x_y) - f(x^*) = \phi(y) - \langle y, \nabla \phi(y) \rangle - \phi(y^*) \leq \|y\| \|\nabla \phi(y)\|$.
28 Thus, in this problem, the quantity $\|\nabla \phi(y)\|$ serves as a measure of both primal optimality
29 $f(x_y) - f(x^*)$ and feasibility $\|b - Ax_y\|$, which is better than just measuring the function value.
- 30 • Matrix scaling [50] is a convex problem and its goal is to find near-stationary points [4, 9].
- 31 • Gradient norm is readily available, unlike other optimality measures ($f(x) - f(x^*)$ and $\|x - x^*\|$),
32 and is thus usable as a stopping criterion. This fact motivates the design of several parameter-free
33 algorithms [43, 39, 27], and their guarantees are established on the gradient norm.
- 34 • Designing schemes for minimizing the gradient norm can inspire new non-convex optimization
35 methods. For example, SARAH [46] was designed for convex finite-sums with gradient-norm mea-
36 sure, but was later discovered to be the near-optimal method for non-convex finite-sums [21, 47].

Table 1: Finding near-stationary points $\|\nabla f(x)\| \leq \epsilon$ of convex finite-sums.

	Algorithm	Complexity	Remark
I F C	GD [33]	$O(\frac{n}{\epsilon^2})$	
	Regularized NAG* [7]	$O(\frac{n}{\epsilon} \log \frac{1}{\epsilon})$	
	OGM-G [33]	$O(\frac{n}{\epsilon})$	$O(\frac{1}{\epsilon} + d)$ memory, optimal in ϵ
	M-OGM-G [Section 3.1]	$O(\frac{n}{\epsilon})$	$O(d)$ memory, optimal in ϵ
	L2S [37]	$O(n + \frac{\sqrt{n}}{\epsilon^2})$	Loopless variant of SARAH [46]
	Regularized Katyusha* [2]	$O((n + \frac{\sqrt{n}}{\epsilon}) \log \frac{1}{\epsilon})$	Requires the knowledge of Δ_0
	R-Acc-SVRG-G* [Section 5]	$O((n \log \frac{1}{\epsilon} + \frac{\sqrt{n}}{\epsilon}) \log \frac{1}{\epsilon})$	Without the knowledge of Δ_0
I D C	GD [42, 54]	$O(\frac{n}{\epsilon})$	
	NAG / NAG + GD [32] / [42]	$O(\frac{n}{\epsilon^{2/3}})$	
	Regularized NAG* [42, 27]	$O(\frac{n}{\sqrt{\epsilon}} \log \frac{1}{\epsilon})$	
	NAG + OGM-G [45]	$O(\frac{n}{\sqrt{\epsilon}})$	$O(\frac{1}{\sqrt{\epsilon}} + d)$ memory, optimal in ϵ
	NAG + M-OGM-G [Section 3.1]	$O(\frac{n}{\sqrt{\epsilon}})$	$O(d)$ memory, optimal in ϵ
	Katyusha + L2S [Appendix E]	$O(n \log \frac{1}{\epsilon} + \frac{\sqrt{n}}{\epsilon^{2/3}})$	
	Acc-SVRG-G [Section 4]	$O(n \log \frac{1}{\epsilon} + \frac{n^{2/3}}{\epsilon^{2/3}})^1$	$O(n \log \frac{1}{\epsilon} + \sqrt{\frac{n}{\epsilon}})$ for function at the same time, simple and elegant
	Regularized Katyusha* [2]	$O((n + \sqrt{\frac{n}{\epsilon}}) \log \frac{1}{\epsilon})$	Requires the knowledge of R_0
R-Acc-SVRG-G* [Section 5]	$O((n \log \frac{1}{\epsilon} + \sqrt{\frac{n}{\epsilon}}) \log \frac{1}{\epsilon})$	Without the knowledge of R_0	

* Indirect methods (using regularization).

37 Moreover, finding near-stationary points is a harder task than minimizing function value, because
 38 NAG has the optimal guarantee for $f(x) - f(x^*)$ but is only suboptimal for minimizing $\|\nabla f(x)\|$.

39 In this work, we consider the problem $\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$, where each f_i is L -smooth
 40 and convex. We focus on finding an ϵ -stationary point of this objective, i.e., a point with $\|\nabla f(x)\| \leq \epsilon$.
 41 We use \mathcal{X}^* to denote the set of optimal solutions, which is assumed to be nonempty. There are two
 42 different assumptions on the initial point x_0 , namely, the Initial bounded-Function Condition (**IFC**):
 43 $f(x_0) - f(x^*) \leq \Delta_0$, and the Initial bounded-Distance Condition (**IDC**): $\|x_0 - x^*\| \leq R_0$ for some
 44 $x^* \in \mathcal{X}^*$. This subtlety results in drastically different best achievable rates as studied in [7, 22].
 45 Below we categorize existing algorithmic techniques into three classes (relating to Table 1).

- 46 (i) “**IDC + IFC**”. Nesterov [42] showed that we can combine the guarantees of a method
 47 minimizing function value under IDC and a method finding near-stationary points under IFC
 48 to produce a faster one for minimizing gradient norm under IDC. For example, NAG produces
 49 $f(x_{K_1}) - f(x^*) = O(\frac{LR_0^2}{K_1^2})$ [40] and GD produces $\|\nabla f(x_{K_2})\|^2 = O(\frac{L(f(x_0) - f(x^*))}{K_2})$ [33]
 50 under IFC. Letting $x_0 = x_{K_1}$ and $K = K_1 + K_2$, by balancing the ratio of K_1 and K_2 , we
 51 obtain the guarantee $\|\nabla f(x_K)\|^2 = O(\frac{L^2 R_0^2}{K^3})$ for “NAG + GD”. We point out that we can use
 52 this technique to combine the guarantees of Katyusha [1] and SARAH² [46]; see Appendix E.
- 53 (ii) *Regularization*. Nesterov [42] used NAG (strongly convex variant) to solve the regularized
 54 objective, and showed that it achieves near-optimal complexity (optimal up to logarithmic
 55 factors). Inspired by this technique, Allen-Zhu [2] proposed recursive regularization for
 56 stochastic approximation algorithms, which also achieves near-optimal complexities [22].

¹Table 1 shows that Katyusha+L2S has a slightly better dependence on n than Acc-SVRG-G. It is due to the adoption of n -dependent step size in L2S. As studied in [37], despite having a better complexity, n -dependent step size boosts numerical performance only when n is *extremely large*. If the practically fast n -independent step size is used for L2S, Katyusha+L2S and Acc-SVRG-G have the same complexity. See also Appendix A.

²We adopt the loopless variant of SARAH in [37], which has a refined analysis for general convex objectives.

57 (iii) *Direct methods.* Due to the lack of insight, existing direct methods are mostly derived or
 58 analyzed with the help of computer-aided tools [31, 32, 54, 33]. The computer-aided approach
 59 was pioneered by Drori and Teboulle [19], who introduced the performance estimation
 60 problem (PEP). The only known optimal method OGM-G [33] was designed based on the
 61 PEP approach.

62 Observe that since $f(x) - f(x^*) \leq \|\nabla f(x)\| \|x - x^*\|$, the lower bound for finding near-stationary
 63 points must be of the same order as for minimizing function value [44]. Thus, under IDC, the lower
 64 bound is $\Omega(n + \sqrt{\frac{n}{\epsilon}})$ due to [58]. Under IFC, we can establish an $\Omega(n + \frac{\sqrt{n}}{\epsilon})$ lower bound using
 65 the techniques in [7, 58]. The main contributions of this work are three new algorithmic schemes that
 66 improve the practicalities of existing methods as summarized below (highlighted in Table 1).

- 67 • (Section 3) We propose a memory-saving variant of OGM-G for the deterministic case ($n = 1$),
 68 which does not require a pre-computed and stored parameter sequence. The derivation of the new
 69 variant is inspired by the numerical solution to a PEP problem.
- 70 • (Section 4) We propose a new accelerated SVRG [29, 59] variant that can *simultaneously*
 71 achieve fast convergence rates for minimizing both the gradient norm and function value, that is,
 72 $O(n \log \frac{1}{\epsilon} + \frac{n^{2/3}}{\epsilon^{2/3}})$ complexity for gradient norm and $O(n \log \frac{1}{\epsilon} + \sqrt{\frac{n}{\epsilon}})$ complexity for function
 73 value. Note that other stochastic approaches in Table 1 do not have this property.
- 74 • (Section 5) We propose an adaptively regularized accelerated SVRG variant, which does not
 75 require the knowledge of R_0 or Δ_0 and achieves a near-optimal complexity under IDC or IFC.

76 We put in extra efforts to make the proposed schemes as simple and elegant as possible. We believe
 77 that the simplicity makes the extensions of the new schemes easier.

78 2 Preliminaries

79 Throughout this paper, we use $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ to denote the inner product and the Euclidean norm,
 80 respectively. We let $[n]$ denote the set $\{1, 2, \dots, n\}$, \mathbb{E} denote the total expectation and \mathbb{E}_{i_k} denote
 81 the expectation with respect to a random sample i_k . We say that a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *L-smooth*
 82 if it has *L-Lipschitz* continuous gradients, i.e.,

$$\forall x, y \in \mathbb{R}^d, \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|.$$

83 A continuously differentiable f is called *μ -strongly convex* if

$$\forall x, y \in \mathbb{R}^d, f(x) - f(y) - \langle \nabla f(y), x - y \rangle \geq \frac{\mu}{2} \|x - y\|^2.$$

84 Other equivalent definitions of these two assumptions can be found in the textbook [44]. The
 85 following is an important consequence of a function f being *L-smooth* and convex:

$$\forall x, y \in \mathbb{R}^d, f(x) - f(y) - \langle \nabla f(y), x - y \rangle \geq \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2. \quad (1)$$

86 We call (1) the *interpolation condition* at (x, y) following [56]. If f is both *L-smooth* and *μ -strongly*
 87 convex, we can define a “shifted” function $h(x) = f(x) - f(x^*) - \frac{\mu}{2} \|x - x^*\|^2$ following [63]. It
 88 can be easily verified that h is $(L - \mu)$ -smooth and convex, and thus from (1),

$$\forall x, y \in \mathbb{R}^d, h(x) - h(y) - \langle \nabla h(y), x - y \rangle \geq \frac{1}{2(L - \mu)} \|\nabla h(x) - \nabla h(y)\|^2, \quad (2)$$

89 which is equivalent to the *strongly convex interpolation condition* discovered in [56].

90 Oracle complexity (or simply complexity) refers to the required number of stochastic gradient ∇f_i
 91 computations to find an ϵ -accurate solution.

92 3 OGM-G: “Momentum” Reformulation and a Memory-Saving Variant

93 In this section, we focus on the IFC case, i.e., $f(x_0) - f(x^*) \leq \Delta_0$. We use N to denote the total
 94 iteration number to prevent confusion (in other sections, we use K). Proofs in this section are given in

Algorithm 1 OGM-G: “Momentum” reformulation

Input: initial guess $x_0 \in \mathbb{R}^d$, total iteration number N .

Initialize: vector $v_0 = \mathbf{0}$, scalars $\theta_N = 1$ and $\theta_k^2 - \theta_k = \theta_{k+1}^2$, for $k = 0 \dots N - 1$.

1: **for** $k = 0, \dots, N - 1$ **do**

2: $v_{k+1} = v_k + \frac{1}{L\theta_k\theta_{k+1}^2} \nabla f(x_k)$.

3: $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k) - (2\theta_{k+1}^3 - \theta_{k+1}^2)v_{k+1}$.

4: **end for**

Output: x_N .

95 Appendix B. Recall that OGM-G has the following updates [33]. Let $y_0 = x_0$. For $k = 0, \dots, N - 1$,
96

$$\begin{aligned} y_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k), \\ x_{k+1} &= y_{k+1} + \frac{(\theta_k - 1)(2\theta_{k+1} - 1)}{\theta_k(2\theta_k - 1)}(y_{k+1} - y_k) + \frac{2\theta_{k+1} - 1}{2\theta_k - 1}(y_{k+1} - x_k), \end{aligned} \quad (3)$$

97 where $\{\theta_k\}$ is recursively defined: $\theta_N = 1$ and $\begin{cases} \theta_k^2 - \theta_k = \theta_{k+1}^2 & k = 1 \dots N - 1, \\ \theta_0^2 - \theta_0 = 2\theta_1^2 & \text{otherwise.} \end{cases}$

98 OGM-G was discovered from the numerical solution to an SDP problem and its analysis is to show
99 that the step coefficients in (3) specify a feasible solution to the SDP problem. While this analysis is
100 natural for the PEP approach, it is hard to understand how each coefficient affects the rate, especially
101 if one wants to generalize the scheme. Here we provide a simple algebraic analysis for OGM-G.

102 We start with a reformulation³ of OGM-G in Algorithm 1, which aims to simplify the proof. We
103 adopt a consistent $\{\theta_k\}$: $\theta_N = 1$ and $\theta_k^2 - \theta_k = \theta_{k+1}^2$, $k = 0 \dots N - 1$, which only costs a constant
104 factor.⁴ Interestingly, the reformulated scheme resembles the heavy-ball momentum method [49].
105 However, it can be shown that Algorithm 1 is not covered by the heavy-ball momentum scheme.
106 Defining $\theta_{N+1}^2 = \theta_N^2 - \theta_N = 0$, we provide the one-iteration analysis in the following proposition:

107 **Proposition 3.1.** *In Algorithm 1, the following holds at any iteration $k \in \{0, \dots, N - 1\}$:*

$$\begin{aligned} A_k + B_{k+1} + C_{k+1} + E_{k+1} &\leq A_{k+1} + B_k + C_k + E_k - \theta_{k+1} \langle \nabla f(x_{k+1}), v_{k+1} \rangle \\ &\quad + \sum_{i=k+1}^N \frac{\theta_i}{L\theta_k\theta_{k+1}^2} \langle \nabla f(x_k), \nabla f(x_i) \rangle, \end{aligned} \quad (4)$$

108 with $A_k \triangleq \frac{1}{\theta_k^2} (f(x_N) - f(x^*) - \frac{1}{2L} \|\nabla f(x_N)\|^2)$, $B_k \triangleq \frac{1}{\theta_k^2} (f(x_k) - f(x^*))$, $C_k \triangleq \frac{1}{2L\theta_k^2} \|\nabla f(x_k)\|^2$,

109 $E_k \triangleq \frac{\theta_{k+1}^2}{\theta_k} \langle \nabla f(x_k), v_k \rangle$.

110 **Remark 3.1.1.** *A recent work [15] also conducted an algebraic analysis of OGM-G under a potential
111 function framework. Their potential function decrease can be directly obtained from Proposition 3.1
112 by summing up (4). By contrast, our “momentum” vector $\{v_k\}$ naturally merges into the analysis,
113 which significantly simplifies the analysis. Moreover, it provides a better interpretation on how
114 OGM-G utilizes the past gradients to achieve acceleration.*

115 From (4), we see that only the last two terms do not telescope. Note that the “momentum” vector is a
116 weighted sum of the past gradients, i.e., $v_{k+1} = \sum_{i=0}^k \frac{1}{L\theta_i\theta_{i+1}^2} \nabla f(x_i)$. If we sum the terms up from
117 $k = 0, \dots, N - 1$, it can be verified that they exactly sum up to 0. The presence of these special
118 terms prevents OGM-G to have a usual potential function (e.g., those in [6]). Then, by telescoping
119 the remaining terms, we obtain the final convergence guarantee.

120 **Theorem 3.1.** *The output of Algorithm 1 satisfies $\|\nabla f(x_N)\|^2 \leq \frac{8L\Delta_0}{(N+2)^2}$.*

121 We observe two drawbacks of OGM-G (same as the algorithm description in [15]): (1) it requires
122 storing a pre-computed parameter sequence, which costs $O(\frac{1}{\epsilon})$ floats; (2) except for the last iterate,

³It can be verified that this scheme is equivalent to the original one (3) through $v_k = \frac{1}{(2\theta_k - 1)\theta_k^2} (y_k - x_k)$.

⁴The original guarantee of OGM-G can be recovered if we set $\theta_0^2 - \theta_0 = 2\theta_1^2$.

Algorithm 2 M-OGM-G: Memory-saving OGM-G

Input: initial guess $x_0 \in \mathbb{R}^d$, total iteration number N .

Initialize: vector $v_0 = \mathbf{0}$.

1: **for** $k = 0, \dots, N - 1$ **do**

2: $v_{k+1} = v_k + \frac{12}{L(N-k+1)(N-k+2)(N-k+3)} \nabla f(x_k)$.

3: $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k) - \frac{(N-k)(N-k+1)(N-k+2)}{6} v_{k+1}$.

4: **end for**

Output: x_N or $\arg \min_{x \in \{x_0, \dots, x_N\}} \|\nabla f(x)\|$.

123 all other iterates are not known to have guarantees. We resolve these issues by proposing another
 124 parameterization of Algorithm 1 in the next subsection.

125 3.1 Memory-Saving OGM-G

126 A straightforward idea to resolve the aforementioned issues is to generalize Algorithm 1. However,
 127 we find it rather difficult since the parameters in the analysis are rather strict (despite that the proof is
 128 already simple). We choose to rely on computer-aided techniques [19]. The derivation of this variant
 129 (Algorithm 2) is based on the following numerical experiment.

130 **Numerical experiment.** OGM-G was discovered when considering the relaxed PEP problem [33]:

$$\begin{aligned} & \max_{\substack{\nabla f(x_0), \dots, \nabla f(x_N) \in \mathbb{R}^d \\ f(x_0), \dots, f(x_N), f(x^*) \in \mathbb{R}}} \|\nabla f(x_N)\|^2 \\ \text{subject to } & \begin{cases} \text{interpolation condition (1) at } (x_k, x_{k+1}), & k = 0, \dots, N - 1, \\ \text{interpolation condition (1) at } (x_N, x_k), & k = 0, \dots, N - 1, \\ \text{interpolation condition (1) at } (x_N, x^*), & f(x_0) - f(x^*) \leq \Delta_0, \end{cases} \end{aligned} \quad (\mathbf{P})$$

131 where the sequence $\{x_k\}$ is defined as $x_{k+1} = x_k - \frac{1}{L} \sum_{i=0}^k h_{k+1,i} \nabla f(x_i)$, $k = 0, \dots, N - 1$ for
 132 some step coefficients $h \in \mathbb{R}^{N(N+1)/2}$. Given N , the step coefficients of OGM-G correspond to
 133 a numerical solution to the problem: $\arg \min_h \{\text{Lagrangian dual of (P)}\}$, which is denoted as (HD).
 134 Conceptually, solving problem (HD) would give us the fastest possible step coefficients under the
 135 constraints.⁵ We expect there to be some constant-time slower schemes, which are neglected when
 136 solving (HD). To identify such schemes, we relax a set of interpolation conditions in problem (P):

$$f(x_N) - f(x_k) - \langle \nabla f(x_k), x_N - x_k \rangle \geq \frac{1}{2L} \|\nabla f(x_N) - \nabla f(x_k)\|^2 - \rho \|\nabla f(x_k)\|^2,$$

137 for $k = 0, \dots, N - 1$ and some $\rho > 0$. After this relaxation, solving (HD) will no longer give us the
 138 step coefficients of OGM-G. By trying different ρ and checking the dependence on N , we discover
 139 Algorithm 2 when $\rho = \frac{1}{2L}$. Similar to our analysis of OGM-G, we provide a simple algebraic analysis
 140 for the new variant in the following theorem.

141 **Theorem 3.2.** Define $\delta_{k+1} \triangleq \frac{12}{(N-k+1)(N-k+2)(N-k+3)}$, $k = 0, \dots, N$. In Algorithm 2, it holds that
 142

$$\sum_{k=0}^N \frac{\delta_{k+1}}{2} \|\nabla f(x_k)\|^2 \leq \frac{12L\Delta_0}{(N+2)(N+3)}. \quad (5)$$

143 **Remark 3.2.1.** Algorithm 2 converges optimally on the last iterate (note that $\delta_{N+1} = 2$) and the
 144 minimum gradient since

$$\min_{k \in \{0, \dots, N\}} \|\nabla f(x_k)\|^2 \leq \frac{1}{\sum_{k=0}^N \frac{\delta_{k+1}}{2}} \sum_{k=0}^N \frac{\delta_{k+1}}{2} \|\nabla f(x_k)\|^2 \leq \frac{8L\Delta_0}{(N+2)(N+3) - 2}.$$

145 Clearly, the parameters of this variant can be computed on the fly and from (5), each iterate has a
 146 guarantee (although the guarantee degenerates quickly as $k \rightarrow 0$ since $1/\delta_{k+1} = \Omega((N-k)^3)$).
 147 Moreover, we can extend the benefits into the IDC case using the ideas in [42] as summarized below.

⁵However, since problem (HD) is non-convex, we can only obtain approximate solutions.

Algorithm 3 Acc-SVRG-G: Accelerated SVRG for Gradient minimization

Input: parameters $\{\tau_k\}$, $\{p_k\}$, initial guess $x_0 \in \mathbb{R}^d$, total iteration number K .

Initialize: vectors $z_0 = \tilde{x}_0 = x_0$ and scalars $\alpha_k = \frac{L\tau_k}{1-\tau_k}$, $\forall k$ and $\tilde{\tau} = \sum_{k=0}^{K-1} \tau_k^{-2}$.

1: **for** $k = 0, \dots, K - 1$ **do**

2: $y_k = \tau_k z_k + (1 - \tau_k) \left(\tilde{x}_k - \frac{1}{L} \nabla f(\tilde{x}_k) \right)$.

3: $z_{k+1} = \arg \min_x \left\{ \langle \mathcal{G}_k, x \rangle + (\alpha_k/2) \|x - z_k\|^2 \right\}$.

4: $\mathcal{G}_k \triangleq \nabla f_{i_k}(y_k) - \nabla f_{i_k}(\tilde{x}_k) + \nabla f(\tilde{x}_k)$, where i_k is sampled uniformly in $[n]$.

5: $\tilde{x}_{k+1} = \begin{cases} y_k & \text{with probability } p_k, \\ \tilde{x}_k & \text{with probability } 1 - p_k. \end{cases}$

6: **end for**

Output (for gradient): x_{out} is sampled from $\left\{ \text{Prob}\{x_{\text{out}} = \tilde{x}_k\} = \frac{\tau_k^{-2}}{\tilde{\tau}} \mid k \in \{0, \dots, K - 1\} \right\}$.

Output (for function value): \tilde{x}_K .

148 **Corollary 3.2.1** (IDC case). *If we first run $N/2$ iterations of NAG and then continue with $N/2$*
 149 *iterations of Algorithm 2, we obtain an output satisfying $\|\nabla f(x_N)\| = O(\frac{LR_0}{N^2})$.*

150 4 Accelerated SVRG: Fast Rates for Both Gradient Norm and Objective

151 In this section, we focus on the IDC case, i.e., $\|x_0 - x^*\| \leq R_0$ for some $x^* \in \mathcal{X}^*$. From the
 152 development in the previous section, it is natural to ask whether we can use the PEP approach to
 153 motivate new stochastic schemes. However, due to the exponential growth of the number of possible
 154 states (i_0, i_1, \dots) , we cannot directly adopt this approach. A feasible alternative is to first fix an
 155 algorithmic framework and a family of potential functions, and then use the potential-based PEP
 156 approach in [54]. However, this approach is much more restrictive. For example, it cannot identify
 157 special constructions like (4) in OGM-G. Fortunately, as we will see, we can get some inspiration
 158 from the recent development of deterministic methods. Proofs in this section are given in Appendix C.

159 Our proposed scheme is given in Algorithm 3. We adopt the elegant loopless design of SVRG in
 160 [34]. Note that the full gradient $\nabla f(\tilde{x}_k)$ is computed and stored only when $\tilde{x}_{k+1} = y_k$ at Step 5. We
 161 summarize our main technical novelty as follows.

162 **Main algorithmic novelty.** The design of stochastic accelerated methods is largely inspired by
 163 NAG. To make it clear, by setting $n = 1$, we see that Katyusha [1], MiG [61], SSNM [62], Varag [36],
 164 VRADA [52], ANITA [38], the acceleration framework in [16] and AC-SA [35, 24] all reduce to one
 165 of the following variants of NAG. We say that these methods are under the NAG framework.

$$\begin{array}{cc} \left\{ \begin{array}{l} x_k = \tau_k z_k + (1 - \tau_k) y_k, \\ z_{k+1} = z_k - \alpha_k \nabla f(x_k), \\ y_{k+1} = \tau_k z_{k+1} + (1 - \tau_k) y_k. \end{array} \right. & \left\{ \begin{array}{l} x_k = \tau_k z_k + (1 - \tau_k) y_k, \\ z_{k+1} = z_k - \alpha_k \nabla f(x_k), \\ y_{k+1} = x_k - \eta_k \nabla f(x_k). \end{array} \right. \\ \text{Auslender and Teboulle [5]} & \text{Linear Coupling [64]} \end{array}$$

166 See [57, 12] for other variants of NAG. When $n = 1$, Algorithm 3 reduces to the following scheme:

$$\begin{cases} y_k = \tau_k z_k + (1 - \tau_k) \left(y_{k-1} - \frac{1}{L} \nabla f(y_{k-1}) \right), \\ z_{k+1} = z_k - \frac{1}{\alpha_k} \nabla f(y_k). \end{cases}$$

Optimized Gradient Method (OGM) [19, 30]

167 Algorithm 3 reduces to the scheme of OGM when $n = 1$ (this point is clearer in the formulation of
 168 ITEM in [55]). OGM has a constant-time faster worst-case rate than NAG, which exactly matches
 169 the lower complexity bound established in [17]. In the following proposition, we show that the OGM
 170 framework helps us conduct a tight one-iteration analysis, which gives room for achieving our goal.

171 **Proposition 4.1.** *In Algorithm 3, the following holds at any iteration $k \geq 0$ and $\forall x^* \in \mathcal{X}^*$:*

$$\begin{aligned} & \left(\frac{1 - \tau_k}{\tau_k^2 p_k} \mathbb{E} [f(\tilde{x}_{k+1}) - f(x^*)] + \frac{L}{2} \mathbb{E} [\|z_{k+1} - x^*\|^2] \right) + \frac{(1 - \tau_k)^2}{2L\tau_k^2} \mathbb{E} [\|\nabla f(\tilde{x}_k)\|^2] \\ & \leq \left(\frac{(1 - \tau_k p_k)(1 - \tau_k)}{\tau_k^2 p_k} \mathbb{E} [f(\tilde{x}_k) - f(x^*)] + \frac{L}{2} \mathbb{E} [\|z_k - x^*\|^2] \right). \end{aligned} \quad (6)$$

172 The terms inside the parentheses form the commonly used potential function of SVRG variants. The
173 additional $\mathbb{E}[\|\nabla f(\tilde{x}_k)\|^2]$ term is created by adopting the OGM framework. In other words, we use
174 the following potential function for Algorithm 3 ($a_k, b_k, c_k \geq 0$):

$$T_k = a_k \mathbb{E} [f(\tilde{x}_k) - f(x^*)] + b_k \mathbb{E} [\|z_k - x^*\|^2] + \sum_{i=0}^{k-1} c_i \mathbb{E} [\|\nabla f(\tilde{x}_i)\|^2].$$

175 We first provide a simple parameter choice, which leads to a simple and clean analysis.

176 **Theorem 4.1** (Single-stage parameter choice). *In Algorithm 3, if we choose $p_k \equiv \frac{1}{n}$, $\tau_k = \frac{3}{k/n+6}$,
177 then the following holds at the outputs:*

$$\begin{aligned} \mathbb{E} [\|\nabla f(x_{\text{out}})\|^2] &= O \left(\frac{n^3 L (f(x_0) - f(x^*)) + n^2 L^2 R_0^2}{K^3} \right), \\ \mathbb{E} [f(\tilde{x}_K) - f(x^*)] &= O \left(\frac{n^2 (f(x_0) - f(x^*)) + n L R_0^2}{K^2} \right). \end{aligned} \quad (7)$$

178 *In other words, to guarantee that $\mathbb{E} [\|\nabla f(x_{\text{out}})\|] \leq \epsilon_g$ and $\mathbb{E} [f(\tilde{x}_K) - f(x^*)] \leq \epsilon_f$, the oracle com-
179 plexities are $O \left(\frac{n(L(f(x_0) - f(x^*)))^{1/3}}{\epsilon_g^{2/3}} + \frac{(nLR_0)^{2/3}}{\epsilon_g^{2/3}} \right)$ and $O \left(n \sqrt{\frac{f(x_0) - f(x^*)}{\epsilon_f}} + \frac{\sqrt{nLR_0}}{\sqrt{\epsilon_f}} \right)$, respectively.*

180 From (7), we see that Algorithm 3 achieves fast $O(\frac{1}{K^{1.5}})$ and $O(\frac{1}{K^2})$ rates for minimizing the
181 gradient norm and function value at the same time. However, despite being a simple choice, the oracle
182 complexities are not better than the deterministic methods in Table 1. Below we provide a two-stage
183 parameter choice, which is inspired by the idea of including a ‘‘warm-up phase’’ in [3, 36, 52, 38].
184 This theorem corresponds to the reported result in Table 1.

185 **Theorem 4.2** (Two-stage parameter choice). *In Algorithm 3, let $p_k = \max\{\frac{6}{k+8}, \frac{1}{n}\}$, $\tau_k = \frac{3}{p_k(k+8)}$.
186 The oracle complexities needed to guarantee $\mathbb{E} [\|\nabla f(x_{\text{out}})\|] \leq \epsilon_g$ and $\mathbb{E} [f(\tilde{x}_K) - f(x^*)] \leq \epsilon_f$ are*

$$O \left(n \min \left\{ \log \frac{LR_0}{\epsilon_g}, \log n \right\} + \frac{(nLR_0)^{2/3}}{\epsilon_g^{2/3}} \right) \text{ and } O \left(n \min \left\{ \log \frac{LR_0^2}{\epsilon_f}, \log n \right\} + \frac{\sqrt{nLR_0}}{\sqrt{\epsilon_f}} \right),$$

187 *respectively.*

188 If ϵ is large or n is very large, the recently proposed ANITA [38] achieves an $O(n)$ complexity, which
189 matches the lower complexity bound $\Omega(n)$ in this case [58]. Since ANITA uses the NAG framework,
190 we show that similar results can be derived under the OGM framework in the following theorem:

191 **Theorem 4.3** (Low accuracy parameter choice). *In Algorithm 3, let iteration N be the first time
192 Step 5 updates $\tilde{x}_{k+1} = y_k$. If we choose $p_k \equiv \frac{1}{n}$, $\tau_k \equiv 1 - \frac{1}{\sqrt{n+1}}$ and terminate Algorithm 3 at
193 iteration N , then the following holds at \tilde{x}_{N+1} :*

$$\mathbb{E} [\|\nabla f(\tilde{x}_{N+1})\|^2] \leq \frac{8L^2 R_0^2}{5(\sqrt{n+1}+1)} \text{ and } \mathbb{E} [f(\tilde{x}_{N+1}) - f(x^*)] \leq \frac{LR_0^2}{\sqrt{n+1}+1}.$$

194 *In particular, if the required accuracies are low (or n is very large), i.e., $\epsilon_g^2 \geq \frac{8L^2 R_0^2}{5(\sqrt{n+1}+1)}$ and
195 $\epsilon_f \geq \frac{LR_0^2}{\sqrt{n+1}+1}$, then Algorithm 3 only has an $O(n)$ oracle complexity.*

196 In the low accuracy region (specified above), the choice in Theorem 4.3 removes the $O(\log \frac{1}{\epsilon})$ factor
197 in the complexity of Theorem 4.2. We include some numerical justifications of Algorithm 3 in
198 Appendix A. We believe that the potential-based PEP approach in [54] can help us identify better
199 parameter choices of Algorithm 3, which we leave for future work.

Algorithm 4 R-Acc-SVRG-G

Input: accuracy $\epsilon > 0$, parameters $\delta_0 = L, \beta > 1$, initial guess $x_0 \in \mathbb{R}^d$.

- 1: **for** $t = 0, 1, 2, \dots$ **do**
- 2: Define $f^{\delta_t}(x) = (1/n) \sum_{i=1}^n f_i^{\delta_t}(x)$, where $f_i^{\delta_t}(x) = f_i(x) + (\delta_t/2) \|x - x_0\|^2$.
- 3: Initialize vectors $z_0 = \tilde{x}_0 = x_0$ and set $\tau_x, \tau_z, \alpha, p, C_{\text{IDC}}, C_{\text{IFC}}$ according to Proposition 5.1.
- 4: **for** $k = 0, 1, 2, \dots$ **do**
- 5: $y_k = \tau_x z_k + (1 - \tau_x) \tilde{x}_k + \tau_z (\delta_t (\tilde{x}_k - z_k) - \nabla f^{\delta_t}(\tilde{x}_k))$.
- 6: $z_{k+1} = \arg \min_x \left\{ \left\langle \mathcal{G}_k^{\delta_t}, x \right\rangle + (\alpha/2) \|x - z_k\|^2 + (\delta_t/2) \|x - y_k\|^2 \right\}$.
- 7: $\| \mathcal{G}_k^{\delta_t} \triangleq \nabla f_{i_k}^{\delta_t}(y_k) - \nabla f_{i_k}^{\delta_t}(\tilde{x}_k) + \nabla f^{\delta_t}(\tilde{x}_k)$, where i_k is sampled uniformly in $[n]$.
- 8: $\tilde{x}_{k+1} = \begin{cases} y_k & \text{with probability } p, \\ \tilde{x}_k & \text{with probability } 1 - p. \end{cases}$
- 9: **if** $\|\nabla f(\tilde{x}_k)\| \leq \epsilon$ **then** output \tilde{x}_k and terminate the algorithm.
- 10: **if** under IDC and $(1 + \frac{\delta_t}{\alpha})^k \geq \sqrt{C_{\text{IDC}}}/\delta_t$ **then** break the inner loop.
- 11: **if** under IFC and $(1 + \frac{\delta_t}{\alpha})^k \geq \sqrt{C_{\text{IFC}}}/2\delta_t$ **then** break the inner loop.
- 12: **end for**
- 13: $\delta_{t+1} = \delta_t/\beta$.
- 14: **end for**

200 5 Near-Optimal Accelerated SVRG with Adaptive Regularization

201 Currently, there is no known stochastic method that directly achieves the optimal rate in ϵ . To get near-
 202 optimal rates, the existing strategy is to use a carefully designed regularization technique [42, 2] with
 203 a method that solves strongly convex problems; see, e.g., [42, 2, 22, 11]. However, the regularization
 204 parameter requires the knowledge of R_0 or Δ_0 , which significantly limits its practicality.

205 Inspired by the recently proposed adaptive regularization technique [27], we develop a near-optimal
 206 accelerated SVRG variant (Algorithm 4) that does not require the knowledge of R_0 or Δ_0 . Note
 207 that this technique was originally proposed for NAG under the IDC assumption. Our development
 208 extends this technique to the stochastic setting, which brings an $O(\sqrt{n})$ rate improvement. Moreover,
 209 we consider both IFC and IDC cases. Proofs in this section are provided in Appendix D.

210 **Detailed design.** Algorithm 4 has a “guess-and-check” framework. In the outer loop, we first
 211 define the regularized objective f^{δ_t} using the current estimate of regularization parameter δ_t , and
 212 then we initialize an accelerated SVRG method (the inner loop) to solve the δ_t -strongly convex f^{δ_t} .
 213 If the inner loop breaks at Step 10 or 11, indicating the poor quality of the current estimate δ_t , δ_t will
 214 be divided by a fixed β . Thus, conceptually, we can adopt any method that solves strongly convex
 215 finite-sums at the optimal rate as the inner loop. However, since the constructions of Step 10 or 11
 216 require some algorithm-dependent constants, we have to fix one method as the inner loop.

217 The inner loop we adopted is a loopless variant of BS-SVRG [63]. This is because (i) BS-SVRG is
 218 the fastest known accelerated SVRG variant (for ill-conditioned problems) and (ii) it has a simple
 219 scheme, especially after using the loopless construction [34]. However, its original guarantee is built
 220 upon $\{z_k\}$. Clearly, we cannot implement the stopping criterion (Step 9) on $\|\nabla f(z_k)\|$. Interestingly,
 221 we discover that its sequence $\{\tilde{x}_k\}$ works perfectly in our regularization framework, even if we can
 222 neither establish convergence on $f(\tilde{x}_k) - f(x^*)$ nor on $\|\tilde{x}_k - x^*\|^2$.⁷ Moreover, we find that the
 223 loopless construction significantly simplifies the parameter constraints of BS-SVRG, which originally
 224 involves $\Theta(n)$ -th-order inequality. We provide the detailed parameter choice as follows:

225 **Proposition 5.1** (Parameter choice). *In Algorithm 4, we set $\tau_x = \frac{\alpha + \delta_t}{\alpha + L + \delta_t}$, $\tau_z = \frac{\tau_x}{\delta_t} - \frac{\alpha(1 - \tau_x)}{\delta_t L}$ and
 226 $p = \frac{1}{n}$. We set α as the (unique) positive root of the cubic equation $\left(1 - \frac{p(\alpha + \delta_t)}{\alpha + L + \delta_t}\right) \left(1 + \frac{\delta_t}{\alpha}\right)^2 = 1$
 227 and specify $C_{\text{IDC}} = L^2 + \frac{L\alpha^2 p}{L + (1-p)(\alpha + \delta_t)}$, $C_{\text{IFC}} = 2L + \frac{2L\alpha^2 p}{(L + (1-p)(\alpha + \delta_t))\delta_t}$. Under these choices, we
 228 have $\frac{\alpha}{\delta_t} = O\left(n + \sqrt{n(L/\delta_t + 1)}\right)$, $C_{\text{IDC}} = O\left((L + \delta_t)^2\right)$, and $C_{\text{IFC}} = O(L)$.*

⁶Note that we maintain the full gradient $\nabla f^{\delta_t}(\tilde{x}_k)$ and $\nabla f(\tilde{x}_k) = \nabla f^{\delta_t}(\tilde{x}_k) - \delta_t(\tilde{x}_k - x_0)$.

⁷It is due to the special potential function of BS-SVRG (see (27)), which does not contain these two terms.

229 Under the choices of τ_x and τ_z , the α above is the optimal choice in our analysis. Then, we can
 230 characterize the progress of the inner loop in the following proposition:

231 **Proposition 5.2** (The inner loop of Algorithm 4). *Using the parameters specified in Proposition 5.1,
 232 after running the inner loop (Step 4-12) of Algorithm 4 for k iterations, we can conclude that*

233 (i) under IDC, i.e., $\|x_0 - x^*\| \leq R_0$ for some $x^* \in \mathcal{X}^*$,

$$\mathbb{E} [\|\nabla f(\tilde{x}_k)\|] \leq \left(\delta_t + \left(1 + \frac{\delta_t}{\alpha}\right)^{-k} \sqrt{C_{\text{IDC}}} \right) R_0,$$

234 (ii) under IFC, i.e., $f(x_0) - f(x^*) \leq \Delta_0$,

$$\mathbb{E} [\|\nabla f(\tilde{x}_k)\|] \leq \left(\sqrt{2\delta_t} + \left(1 + \frac{\delta_t}{\alpha}\right)^{-k} \sqrt{C_{\text{IFC}}} \right) \sqrt{\Delta_0}.$$

235 The above results motivate the design of Step 10 and 11. For example, in the IDC case, when the
 236 inner loop breaks at Step 10, using (i) above, we obtain $\mathbb{E} [\|\nabla f(\tilde{x}_k)\|] \leq 2\delta_t R_0$. Then, by discussing
 237 the relative size of δ_t and a certain constant, we can estimate the complexity of Algorithm 4. The
 238 same methodology is used for the IFC case.

239 **Theorem 5.1** (IDC case). *Denote $\delta_{\text{IDC}}^* = \frac{\epsilon q}{2R_0}$ for some $q \in (0, 1)$ and let the outer iteration $t = \ell$
 240 be the first time⁸ $\delta_\ell \leq \delta_{\text{IDC}}^*$. The following assertions hold:*

- 241 (i) At outer iteration ℓ , Algorithm 4 terminates with probability at least $1 - q$.⁹
 242 (ii) The total expected oracle complexity of the $\ell + 1$ outer loops is

$$O \left(\left(n \log \frac{LR_0}{\epsilon q} + \sqrt{\frac{nLR_0}{\epsilon q}} \right) \log \frac{LR_0}{\epsilon q} \right).$$

243 **Theorem 5.2** (IFC case). *Denote $\delta_{\text{IFC}}^* = \frac{\epsilon^2 q^2}{8\Delta_0}$ for some $q \in (0, 1)$ and let the outer iteration $t = \ell$ be
 244 the first time $\delta_\ell \leq \delta_{\text{IFC}}^*$. The following assertions hold:*

- 245 (i) At outer iteration ℓ , Algorithm 4 terminates with probability at least $1 - q$.
 246 (ii) The total expected oracle complexity of the $\ell + 1$ outer loops is

$$O \left(\left(n \log \frac{\sqrt{L\Delta_0}}{\epsilon q} + \frac{\sqrt{nL\Delta_0}}{\epsilon q} \right) \log \frac{\sqrt{L\Delta_0}}{\epsilon q} \right).$$

247 Compared with regularized Katyusha in Table 1, the adaptive regularization approach drops the need
 248 to estimate R_0 or Δ_0 at the cost of a mere $\log \frac{1}{\epsilon}$ factor in the non-dominant term (if ϵ is small).

249 6 Discussion

250 In this work, we proposed several simple and practical schemes that complement existing works
 251 (Table 1). Admittedly, the new schemes are currently only limited to the unconstrained Euclidean
 252 setting, because our techniques heavily rely on the interpolation conditions (1) and (2). On the other
 253 hand, methods such as OGM [30], TM [51] and ITEM [55, 10], which also rely on these conditions,
 254 are still not known to have their proximal variants. We list a few future directions as follows.

255 (1) It is not clear how to naturally connect the parameters of M-OGM-G (Algorithm 2) to OGM-G
 256 (Algorithm 1). The parameters of both algorithms seem to be quite restrictive and hardly generalizable
 257 due to the special construction in (4). Does there exist an optimal method for minimizing the gradient
 258 norm that has a proper potential function (at each iteration)?

259 (2) Is this new ‘‘momentum’’ in OGM-G beneficial for training neural nets? Other classic momentum
 260 schemes such as NAG [40] or heavy-ball momentum method [49] are extremely effective for this
 261 task [53], and they were also originally proposed for convex objectives.

262 (3) Can we directly accelerate SARAH (L2S)? By extending OGM-G? It seems that existing stochastic
 263 acceleration techniques fail to accelerate SARAH (or result in poor dependence on n as in [16]).

⁸We assume that ϵ is small such that $\max\{\delta_{\text{IDC}}^*, \delta_{\text{IFC}}^*\} \leq \delta_0 = L$ for simplicity. In this case, $\ell > 0$.

⁹If Algorithm 4 does not terminate at outer iteration ℓ , it terminates at the next outer iteration with probability at least $1 - q/\beta$. That is, it terminates with higher and higher probability. The same goes for the IFC case.

264 **References**

- 265 [1] Z. Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. *Journal of*
 266 *Machine Learning Research*, 18(1):8194–8244, 2017. 2, 6, 26
- 267 [2] Z. Allen-Zhu. How to make the gradients small stochastically: Even faster convex and noncon-
 268 vex sgd. In *Advances in Neural Information Processing Systems*, pages 1157–1167, 2018. 1, 2,
 269 8
- 270 [3] Z. Allen-Zhu and Y. Yuan. Improved SVRG for Non-Strongly-Convex or Sum-of-Non-Convex
 271 Objectives. In *Proceedings of The 33rd International Conference on Machine Learning*, pages
 272 1080–1089, 2016. 7
- 273 [4] Z. Allen-Zhu, Y. Li, R. M. de Oliveira, and A. Wigderson. Much Faster Algorithms for Matrix
 274 Scaling. In C. Umans, editor, *58th IEEE Annual Symposium on Foundations of Computer*
 275 *Science*, pages 890–901, 2017. 1
- 276 [5] A. Auslender and M. Teboulle. Interior gradient and proximal methods for convex and conic
 277 optimization. *SIAM Journal on Optimization*, 16(3):697–725, 2006. 6
- 278 [6] N. Bansal and A. Gupta. Potential-Function Proofs for Gradient Methods. *Theory of Computing*,
 279 15(4):1–32, 2019. 4
- 280 [7] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Lower bounds for finding stationary points
 281 ii: first-order methods. *Mathematical Programming*, 185(1-2), 2021. 1, 2, 3
- 282 [8] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions*
 283 *on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at [http:](http://www.csie.ntu.edu.tw/~cjlin/libsvm)
 284 [//www.csie.ntu.edu.tw/~cjlin/libsvm](http://www.csie.ntu.edu.tw/~cjlin/libsvm). 13, 14
- 285 [9] M. B. Cohen, A. Madry, D. Tsipras, and A. Vladu. Matrix Scaling and Balancing via Box
 286 Constrained Newton’s Method and Interior Point Methods. In *IEEE 58th Annual Symposium on*
 287 *Foundations of Computer Science*, pages 902–913. IEEE, 2017. 1
- 288 [10] A. d’Aspremont, D. Scieur, and A. Taylor. Acceleration methods. *arXiv preprint*
 289 *arXiv:2101.09545*, 2021. 9
- 290 [11] D. Davis and D. Drusvyatskiy. Complexity of finding near-stationary points of convex functions
 291 stochastically. *arXiv preprint arXiv:1802.08556*, 2018. 8
- 292 [12] A. Defazio. On the Curved Geometry of Accelerated Optimization. In *Advances in Neural*
 293 *Information Processing Systems*, volume 32, pages 1764–1773, 2019. 6
- 294 [13] A. Defazio, F. R. Bach, and S. Lacoste-Julien. SAGA: A Fast Incremental Gradient Method With
 295 Support for Non-Strongly Convex Composite Objectives. In *Advances in Neural Information*
 296 *Processing Systems*, pages 1646–1654, 2014. 14
- 297 [14] J. Diakonikolas and C. Guzmán. Complementary Composite Minimization, Small Gradients in
 298 General Norms, and Applications to Regression Problems. *arXiv preprint arXiv:2101.11041*,
 299 2021. 1
- 300 [15] J. Diakonikolas and P. Wang. Potential Function-based Framework for Making the Gradients
 301 Small in Convex and Min-Max Optimization. *arXiv preprint arXiv:2101.12101*, 2021. 1, 4
- 302 [16] D. Driggs, M. J. Ehrhardt, and C.-B. Schönlieb. Accelerating variance-reduced stochastic
 303 gradient methods. *Mathematical Programming*, 2020. doi: 10.1007/s10107-020-01566-2. 6, 9
- 304 [17] Y. Drori. The exact information-based complexity of smooth convex minimization. *Journal of*
 305 *Complexity*, 39:1–16, 2017. 1, 6
- 306 [18] Y. Drori and A. Taylor. On the oracle complexity of smooth strongly convex minimization.
 307 *arXiv preprint arXiv:2101.09740*, 2021. 1
- 308 [19] Y. Drori and M. Teboulle. Performance of first-order methods for smooth convex minimization:
 309 a novel approach. *Mathematical Programming*, 145(1-2):451–482, 2014. 1, 3, 5, 6
- 310 [20] D. Dua and C. Graff. UCI machine learning repository, 2017. URL [http://archive.ics.](http://archive.ics.uci.edu/ml)
 311 [uci.edu/ml](http://archive.ics.uci.edu/ml). 13, 14
- 312 [21] C. Fang, C. J. Li, Z. Lin, and T. Zhang. SPIDER: Near-Optimal Non-Convex Optimization via
 313 Stochastic Path-Integrated Differential Estimator. In *Advances in Neural Information Processing*
 314 *Systems*, pages 687–697, 2018. 1

- 315 [22] D. J. Foster, A. Sekhari, O. Shamir, N. Srebro, K. Sridharan, and B. Woodworth. The Complexity
316 of Making the Gradient Small in Stochastic Convex Optimization. In *Proceedings of the Thirty-
317 Second Conference on Learning Theory*, pages 1319–1345, 2019. 1, 2, 8
- 318 [23] R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping From Saddle Points — Online Stochastic
319 Gradient for Tensor Decomposition. In *Proceedings of The 28th Conference on Learning
320 Theory*, pages 797–842, 2015. 1
- 321 [24] S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex
322 stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on
323 Optimization*, 22(4):1469–1492, 2012. 6
- 324 [25] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic
325 programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013. 1
- 326 [26] S. Ghadimi and G. Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic
327 programming. *Mathematical Programming*, 156(1-2):59–99, 2016. 1
- 328 [27] M. Ito and M. Fukuda. Nearly optimal first-order methods for convex optimization under
329 gradient norm measure: An adaptive regularization approach. *Journal of Optimization Theory
330 and Applications*, 188(3):770–804, 2021. 1, 2, 8
- 331 [28] C. Jin, R. Ge, P. Netrapalli, S. M. Kakade, and M. I. Jordan. How to Escape Saddle Points
332 Efficiently. In *Proceedings of the 34th International Conference on Machine Learning*, pages
333 1724–1732, 2017. 1
- 334 [29] R. Johnson and T. Zhang. Accelerating Stochastic Gradient Descent using Predictive Variance
335 Reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013. 3, 14
- 336 [30] D. Kim and J. A. Fessler. Optimized first-order methods for smooth convex minimization.
337 *Mathematical Programming*, 159(1):81–107, 2016. 1, 6, 9
- 338 [31] D. Kim and J. A. Fessler. Another Look at the Fast Iterative Shrinkage/Thresholding Algorithm
339 (FISTA). *SIAM Journal on Optimization*, 28(1):223–250, 2018. 1, 3
- 340 [32] D. Kim and J. A. Fessler. Generalizing the optimized gradient method for smooth convex
341 minimization. *SIAM Journal on Optimization*, 28(2):1920–1950, 2018. 1, 2, 3
- 342 [33] D. Kim and J. A. Fessler. Optimizing the efficiency of first-order methods for decreasing the
343 gradient of smooth convex functions. *Journal of Optimization Theory and Applications*, 188(1):
344 192–219, 2021. 1, 2, 3, 4, 5
- 345 [34] D. Kovalev, S. Horváth, and P. Richtárik. Don’t jump through hoops and remove those loops:
346 SVRG and Katyusha are better without the outer loop. In *Algorithmic Learning Theory*, pages
347 451–467. PMLR, 2020. 6, 8
- 348 [35] G. Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*,
349 133(1-2):365–397, 2012. 6
- 350 [36] G. Lan, Z. Li, and Y. Zhou. A unified variance-reduced accelerated gradient method for
351 convex optimization. In *Advances in Neural Information Processing Systems*, volume 32, pages
352 10462–10472, 2019. 6, 7
- 353 [37] B. Li, M. Ma, and G. B. Giannakis. On the Convergence of SARAH and Beyond. In *Proceedings
354 of the Twenty Third International Conference on Artificial Intelligence and Statistics*, pages
355 223–233, 2020. 2, 14, 27
- 356 [38] Z. Li. ANITA: An Optimal Loopless Accelerated Variance-Reduced Gradient Method. *arXiv
357 preprint arXiv:2103.11333*, 2021. 6, 7
- 358 [39] Q. Lin and L. Xiao. An Adaptive Accelerated Proximal Gradient Method and its Homotopy
359 Continuation for Sparse Optimization. In *Proceedings of the 31th International Conference on
360 Machine Learning*, pages 73–81, 2014. 1
- 361 [40] Y. Nesterov. A method for solving the convex programming problem with convergence rate
362 $O(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983. 1, 2, 9
- 363 [41] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer
364 Science & Business Media, 2003. 1
- 365 [42] Y. Nesterov. How to make the gradients small. *Optima. Mathematical Optimization Society
366 Newsletter*, (88):10–11, 2012. 1, 2, 5, 8, 22, 27

- 367 [43] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Program-*
368 *ming*, 140(1):125–161, 2013. 1
- 369 [44] Y. Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018. 3, 18
- 370 [45] Y. Nesterov, A. Gasnikov, S. Guminov, and P. Dvurechensky. Primal–dual accelerated gradient
371 methods with small-dimensional relaxation oracle. *Optimization Methods and Software*, pages
372 1–38, 2020. 2
- 373 [46] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. SARAH: A Novel Method for Ma-
374 chine Learning Problems Using Stochastic Recursive Gradient. In *Proceedings of the 34th*
375 *International Conference on Machine Learning*, pages 2613–2621, 2017. 1, 2
- 376 [47] N. H. Pham, L. M. Nguyen, D. T. Phan, and Q. Tran-Dinh. ProxSARAH: An efficient algorithmic
377 framework for stochastic composite nonconvex optimization. *Journal of Machine Learning*
378 *Research*, 21(110):1–48, 2020. 1
- 379 [48] J. Platt. Sequential minimal optimization: A fast algorithm for training support vector machines.
380 1998. 14
- 381 [49] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *Ussr*
382 *computational mathematics and mathematical physics*, 4(5):1–17, 1964. 4, 9
- 383 [50] U. G. Rothblum and H. Schneider. Scalings of matrices which have prespecified row sums and
384 column sums via optimization. *Linear Algebra and its Applications*, 114:737–764, 1989. 1
- 385 [51] B. V. Scoy, R. A. Freeman, and K. M. Lynch. The Fastest Known Globally Convergent First-
386 Order Method for Minimizing Strongly Convex Functions. *IEEE Control Systems Letters*, 2(1):
387 49–54, 2017. 9
- 388 [52] C. Song, Y. Jiang, and Y. Ma. Variance Reduction via Accelerated Dual Averaging for Finite-
389 Sum Optimization. In *Advances in Neural Information Processing Systems*, volume 33, pages
390 833–844, 2020. 6, 7
- 391 [53] I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and
392 momentum in deep learning. In *Proceedings of the 30th International Conference on Machine*
393 *Learning*, pages 1139–1147, 2013. 9
- 394 [54] A. Taylor and F. Bach. Stochastic first-order methods: non-asymptotic and computer-aided
395 analyses via potential functions. In *Conference on Learning Theory*, pages 2934–2992, 2019. 2,
396 3, 6, 7
- 397 [55] A. Taylor and Y. Drori. An optimal gradient method for smooth strongly convex minimization.
398 *arXiv preprint arXiv:2101.09741*, 2021. 1, 6, 9
- 399 [56] A. B. Taylor, J. M. Hendrickx, and F. Glineur. Smooth strongly convex interpolation and exact
400 worst-case performance of first-order methods. *Mathematical Programming*, 161(1-2):307–345,
401 2017. 3
- 402 [57] P. Tseng. On accelerated proximal gradient methods for convex-concave optimization. <https://www.mit.edu/~dimitrib/PTseng/papers/apgm.pdf>, 2008. Accessed May 1, 2020. 6
- 403
404 [58] B. E. Woodworth and N. Srebro. Tight Complexity Bounds for Optimizing Composite Ob-
405 jectives. In *Advances in Neural Information Processing Systems*, pages 3639–3647, 2016. 3,
406 7
- 407 [59] L. Xiao and T. Zhang. A Proximal Stochastic Gradient Method with Progressive Variance
408 Reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014. 3, 14
- 409 [60] D. Zhou, P. Xu, and Q. Gu. Stochastic Nested Variance Reduction for Nonconvex Optimization.
410 *Journal of Machine Learning Research*, 21:103:1–103:63, 2020. 1
- 411 [61] K. Zhou, F. Shang, and J. Cheng. A Simple Stochastic Variance Reduced Algorithm with Fast
412 Convergence Rates. In *Proceedings of the 35th International Conference on Machine Learning*,
413 pages 5980–5989, 2018. 6
- 414 [62] K. Zhou, Q. Ding, F. Shang, J. Cheng, D. Li, and Z.-Q. Luo. Direct Acceleration of SAGA using
415 Sampled Negative Momentum. In *Proceedings of the Twenty Second International Conference*
416 *on Artificial Intelligence and Statistics*, pages 1602–1610, 2019. 6
- 417 [63] K. Zhou, A. M.-C. So, and J. Cheng. Boosting First-Order Methods by Shifting Objective: New
418 Schemes with Faster Worst-Case Rates. In *Advances in Neural Information Processing Systems*,
419 pages 15405–15416, 2020. 3, 8, 22

420 [64] Z. A. Zhu and L. Orecchia. Linear Coupling: An Ultimate Unification of Gradient and Mirror
421 Descent. In *8th Innovations in Theoretical Computer Science Conference*, volume 67 of *LIPICs*,
422 pages 3:1–3:22, 2017. 6

423 Checklist

- 424 1. For all authors...
- 425 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
426 contributions and scope? [Yes]
 - 427 (b) Did you describe the limitations of your work? [Yes] See Section 6.
 - 428 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We are
429 not aware of clear negative societal impacts since we focus on developing generic
430 algorithms for convex optimization.
 - 431 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
432 them? [Yes]
- 433 2. If you are including theoretical results...
- 434 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See the
435 introduction.
 - 436 (b) Did you include complete proofs of all theoretical results? [Yes]
- 437 3. If you ran experiments...
- 438 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
439 mental results (either in the supplemental material or as a URL)? [Yes]
 - 440 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
441 were chosen)? [Yes] See Appendix A.
 - 442 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
443 ments multiple times)? [Yes] See Figure 1.
 - 444 (d) Did you include the total amount of compute and the type of resources used (e.g., type
445 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix A.
- 446 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 447 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix A.
 - 448 (b) Did you mention the license of the assets? [Yes] LIBSVM [8] is under the BSD license.
 - 449 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - 450 (d) Did you discuss whether and how consent was obtained from people whose data you’re
451 using/curating? [Yes] Details can be found in the online dataset repositories [8, 20].
 - 452 (e) Did you discuss whether the data you are using/curating contains personally identifiable
453 information or offensive content? [Yes] Details can be found in the online dataset
454 repositories [8, 20].
- 455 5. If you used crowdsourcing or conducted research with human subjects...
- 456 (a) Did you include the full text of instructions given to participants and screenshots, if
457 applicable? [N/A]
 - 458 (b) Did you describe any potential participant risks, with links to Institutional Review
459 Board (IRB) approvals, if applicable? [N/A]
 - 460 (c) Did you include the estimated hourly wage paid to participants and the total amount
461 spent on participant compensation? [N/A]