

---

# Unified Convergence Theory of Stochastic and Variance-Reduced Cubic Newton Methods

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study stochastic Cubic Newton methods for solving general possibly non-  
2 convex minimization problems. We propose a new framework, which we call  
3 the *helper framework*, that provides a unified view of the stochastic and variance-  
4 reduced second-order algorithms equipped with global complexity guarantees. It  
5 can also be applied to learning with auxiliary information. Our helper framework  
6 offers the algorithm designer high flexibility for constructing and analysis of the  
7 stochastic Cubic Newton methods, allowing arbitrary size batches, and the use  
8 of noisy and possibly biased estimates of the gradients and Hessians, incorporat-  
9 ing both the variance reduction and the lazy Hessian updates. We recover the  
10 best-known complexities for the stochastic and variance-reduced Cubic Newton,  
11 under weak assumptions on the noise and avoiding artificial logarithms. A direct  
12 consequence of our theory is the new lazy stochastic second-order method, which  
13 significantly improves the arithmetic complexity for large dimension problems. We  
14 also establish complexity bounds for the classes of gradient-dominated objectives,  
15 that include convex and strongly convex problems. For Auxiliary Learning, we  
16 show that using a helper (auxiliary function) can outperform training alone if a  
17 given similarity measure is small.

## 18 1 Introduction

19 In many fields of machine learning, it is common to optimize a function  $f(\mathbf{x})$  that can be expressed  
20 as a finite sum:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) \right\}, \quad (1)$$

21 or, more generally, as an expectation over some given probability distribution:  $f(\mathbf{x}) = \mathbb{E}_{\zeta} [f(\mathbf{x}, \zeta)]$ .  
22 When  $f$  is non-convex, this problem is especially difficult, since finding a global minimum is NP-hard  
23 in general [15]. Hence, the reasonable goal is to look for approximate solutions. The most prominent  
24 family of algorithms for solving large-scale problems of the form (1) are the *first-order methods*, such  
25 as the Stochastic Gradient Descent (SGD) [26, 17]. They employ only stochastic gradient information  
26 about the objective  $f(\mathbf{x})$  and guarantee the convergence to a stationary point, which is a point with a  
27 small gradient norm.

28 Nevertheless, when the objective function is non-convex, a stationary point may be a saddle point or  
29 even a local maximum, which is not desirable. Another common issue is that first-order methods  
30 typically have a slow convergence rate, particularly when the problem is *ill-conditioned*. Therefore,  
31 they may not be suitable when high precision for the solution is required.

32 To address these challenges, we can take into account *second-order information* (the Hessian matrix)  
33 and apply Newton’s method (see, e.g. [20]). Among the many versions of this algorithm, the Cubic  
34 Newton method [21] is one of the most theoretically established. With the Cubic Newton method, we

35 can guarantee *global convergence* to an approximate *second-order* stationary point (in contrast, the  
36 pure Newton method without regularization can even diverge when it starts far from a neighborhood  
37 of the solution). For a comprehensive historical overview of the different variants of Newton’s  
38 method, see [25]. Additionally, the rate of convergence of the Cubic Newton is *provably better* than  
39 those for the first-order methods.

40 Therefore, theoretical guarantees of the Cubic Newton method seem to be very appealing for practical  
41 applications. However, the basic version of the Cubic Newton requires the exact gradient and Hessian  
42 information in each step, which can be very expensive to compute in the large scale setting. To  
43 overcome this issue, several techniques have been proposed:

- 44 • One popular approach is to use inexact *stochastic gradient and Hessian estimates* with sub-  
45 sampling [34, 18, 33, 22, 13, 6, 1]. This technique avoids using the full oracle information,  
46 but typically it has a slower convergence rate compared to the exact Cubic Newton.
- 47 • *Variance reduction* techniques [36, 30] combine the advantages of stochastic and exact  
48 methods, achieving an improved rates by recomputing the full gradient and Hessian  
49 information at some iterations.
- 50 • *Lazy Hessian* updates [27, 10] utilize a simple idea of reusing an old Hessian for several  
51 iterations of a second-order scheme. Indeed, since the cost of computing one Hessian is  
52 usually much more expensive than one gradient, it can improve the arithmetic complexity of  
53 our methods.
- 54 • In addition, exploiting the special structure of the function  $f$  (if known) can also be helpful.  
55 For instance, some studies [21, 19] consider *gradient-dominated objectives*, a subclass  
56 of non-convex functions that have improved convergence rates and can even be shown to  
57 converge to the global minimum. Examples of such objectives include convex and star-  
58 convex functions, uniformly convex functions, and functions satisfying the PL condition  
59 [24] as a special case.

60 In this work, we revise the current state-of-the-art convergence theory for the stochastic Cubic  
61 Newton method and propose a unified and improved complexity guarantees for different versions of  
62 the method, which combine all the advanced techniques listed above.

63 Our developments are based on the new *helper framework* for the second-order optimization, that we  
64 present in Section 3. For the first-order optimization, a similar in-spirit techniques called *learning*  
65 *with auxiliary information* was developed recently in [8, 31]. Thus, our results can also be seen as a  
66 generalization of the Auxiliary Learning paradigm to the second-order optimization. However, note  
67 that in our second-order case, we have more freedom for choosing the "helper functions" (namely, we  
68 use one for the gradients and one for the Hessians). That brings more flexibility into our methods and  
69 it allows, for example, to use the lazy Hessian updates.

70 Our new helper framework provides us with a unified view of the stochastic and variance-reduced  
71 methods and can be used by an algorithm designed to construct new methods. Thus, we show how to  
72 recover already known versions of the stochastic Cubic Newton with the best convergence rates, as  
73 well as present the new *Lazy Stochastic Second-Order Method*, which significantly improves the total  
74 arithmetic complexity for large-dimension problems.

## 75 **Contributions.**

- 76 • We introduce the *helper framework* which we argue encompasses multiple methods in a  
77 unified way. Such methods include stochastic methods, variance reduction, Lazy methods,  
78 core sets, and semi-supervised learning.
- 79 • This framework covers previous versions of the variance-reduced stochastic Cubic Newton  
80 methods with known rates. Moreover, it provides us with new algorithms that employ *Lazy*  
81 *Hessian* updates and significantly improves the arithmetic complexity (for high dimensions),  
82 by using the same Hessian snapshot for several steps of the method.
- 83 • In the case of Auxiliary learning we provably show a benefit from using auxiliary tasks  
84 as helpers in our framework. In particular, we can replace the smoothness constant by a  
85 similarity constant which might be smaller.
- 86 • Moreover, our analysis works both for the general class of non-convex functions, as well as  
87 for the class of gradient-dominated problems, that includes convex and uniformly convex  
88 functions. Hence, in particular, we are the first to establish the convergence rates of the

89 stochastic Cubic Newton algorithms with variance reduction for the gradient-dominated  
 90 case.

## 91 2 Notation and Assumptions

92 For simplicity, we consider the finite-sum optimization problem (1), while it can be also possible  
 93 to generalize our results to arbitrary expectations. We assume that our objective  $f$  is bounded  
 94 from below and denote  $f^* := \inf_{\mathbf{x}} f(\mathbf{x})$ , and use the following notation:  $F_0 := f(\mathbf{x}_0) - f^*$ , for  
 95 some initial  $\mathbf{x}_0 \in \mathbb{R}^d$ . We denote by  $\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$ ,  $\mathbf{x} \in \mathbb{R}^d$ , the standard Euclidean norm for  
 96 vectors, and the spectral norm for symmetric matrices,  $\|\mathbf{H}\| := \max\{\lambda_{\max}(\mathbf{H}), -\lambda_{\min}(\mathbf{H})\}$ , where  
 97  $\mathbf{H} = \mathbf{H}^\top \in \mathbb{R}^{d \times d}$ . We will also use  $x \wedge y$  to denote  $\min(x, y)$ .

98 Throughout this work, we make the following smoothness assumption on the objective  $f$  :

**Assumption 1 (Lipschitz Hessian)** *The Hessian of  $f$  is Lipschitz continuous, for some  $L > 0$ :*

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$

99 Our goal is to explore the potential of using the Cubically regularized Newton methods to solve  
 100 problem (1). At each iteration, being at a point  $\mathbf{x} \in \mathbb{R}^d$ , we compute the next point  $\mathbf{x}^+$  by solving  
 101 the subproblem of the form

$$\mathbf{x}^+ \in \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left\{ \Omega_{M, \mathbf{g}, \mathbf{H}}(\mathbf{y}, \mathbf{x}) := \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{H}(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{M}{6} \|\mathbf{y} - \mathbf{x}\|^3 \right\}. \quad (2)$$

102 Here,  $\mathbf{g}$  and  $\mathbf{H}$  are estimates of the gradient  $\nabla f(\mathbf{x})$  and the Hessian  $\nabla^2 f(\mathbf{x})$ , respectively. Note that  
 103 solving (2) can be done efficiently even for non-convex problems (see [9, 21, 5]). Generally, the cost  
 104 of computing  $\mathbf{x}^+$  is  $\mathcal{O}(d^3)$  arithmetic operations, which are needed for evaluating an appropriate  
 105 factorization of  $\mathbf{H}$ . Hence, it is of a similar order as the cost of the classical Newton's step.

We will be interested to find a second-order stationary point to (1). We call  $(\varepsilon, c)$ -approximate  
 second-order local minimum a point  $\mathbf{x}$  that satisfies:

$$\|\nabla f(\mathbf{x})\| \leq \varepsilon \quad \text{and} \quad \lambda_{\min}(\nabla^2 f(\mathbf{x})) \geq -c\sqrt{\varepsilon},$$

where  $\varepsilon, c > 0$  are given tolerance parameters. Let us define the following accuracy measure (see  
 [21]):

$$\mu_c(\mathbf{x}) := \max\left(\|\nabla f(\mathbf{x})\|^{3/2}, \frac{-\lambda_{\min}(\nabla^2 f(\mathbf{x}))^3}{c^{3/2}}\right), \quad \mathbf{x} \in \mathbb{R}^d, c > 0.$$

106 Note that this definition implies that if  $\mu_c(\mathbf{x}) \leq \varepsilon^{3/2}$  then  $\mathbf{x}$  is an  $(\varepsilon, c)$ -approximate local minimum.

107 **Computing gradients and Hessians.** It is clear that computing the Hessian matrix can be  
 108 much more expensive than computing the gradient vector. We denote the corresponding arith-  
 109 metic complexities by *HessCost* and *GradCost*. We will make and follow the convention that  
 110  $HessCost = d \times GradCost$ , where  $d$  is the dimension of the problem. For example, this is known to  
 111 hold for neural networks using the backpropagation algorithm [16]. However, if the Hessian has a  
 112 sparse structure, the cost of computing the Hessian can be cheaper [23]. Then, we can replace  $d$  with  
 113 the *effective dimension*  $d_{\text{eff}} := \frac{HessCost}{GradCost} \leq d$ .

## 114 3 Second-Order Optimization with Helper Functions

115 In this section, we extend the helper framework previously introduced in [8] for first-order optimiza-  
 116 tion methods to second-order optimization.

117 **General principle.** The general idea is the following: imagine that, besides the objective function  $f$   
 118 we have access to a help function  $h$  that we think is similar in some sense (that will be defined later)  
 119 to  $f$  and thus it should help to minimize it.

120 Note that many optimization algorithms can be framed in the following sequential way. For a current  
 121 state  $\mathbf{x}$ , we compute the next state  $\mathbf{x}^+$  as:

$$\mathbf{x}^+ \in \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left\{ \hat{f}_{\mathbf{x}}(\mathbf{y}) + Mr_{\mathbf{x}}(\mathbf{y}) \right\},$$

122 where  $\hat{f}_{\mathbf{x}}(\cdot)$  is an approximation of  $f$  around current point  $\mathbf{x}$ , and  $r_{\mathbf{x}}(\mathbf{y})$  is a regularizer that encodes  
 123 how accurate the approximation is, and  $M > 0$  is a regularization parameter. In this work, we  
 124 are interested in cubically regularized second-order models of the form (2) and we use  $r_{\mathbf{x}}(\mathbf{y}) :=$   
 125  $\frac{1}{6} \|\mathbf{y} - \mathbf{x}\|^3$ .

126 Now let us look at how we can use a helper  $h$  to construct the approximation  $\hat{f}$ . We notice that we  
 127 can write

$$f(\mathbf{y}) := \underbrace{h(\mathbf{y})}_{\text{cheap}} + \underbrace{f(\mathbf{y}) - h(\mathbf{y})}_{\text{expensive}}$$

128 We discuss the actual practical choices of the helper function  $h$  below. We assume now that we can  
 129 afford the second-order approximation for the cheap part  $h$  around the current point  $\mathbf{x}$ . However,  
 130 approximating the part  $f - h$  can be expensive (as for example when the number of elements  $n$  in  
 131 finite sum (1) is huge), or even impossible (due to lack of data). Thus, we would prefer to approximate  
 132 the expensive part less frequently. For this reason, let us introduce an extra *snapshot point*  $\tilde{\mathbf{x}}$  that  
 133 is updated less often than  $\mathbf{x}$ . Then, we use it to approximate  $f - h$ . Another question that we still  
 134 need to ask is *what order should we use for the approximation of  $f - h$ ?* We will see that order 0  
 135 (approximating by a constant) leads as to the basic stochastic methods, while for orders 1 and 2 we  
 136 equip our methods with the variance reduction.

137 Combining the two approximations for  $h$  and  $f - h$  we get the following model of our objective  $f$ :

$$\hat{f}_{\mathbf{x}, \tilde{\mathbf{x}}}(\mathbf{y}) = C(\mathbf{x}, \tilde{\mathbf{x}}) + \mathcal{G}(h, \mathbf{x}, \tilde{\mathbf{x}}), \mathbf{y} - \mathbf{x} + \frac{1}{2} \langle \mathcal{H}(h, \mathbf{x}, \tilde{\mathbf{x}})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \quad (3)$$

138 where  $C(\mathbf{x}, \tilde{\mathbf{x}})$  is a constant,  $\mathcal{G}(h, \mathbf{x}, \tilde{\mathbf{x}})$  is a linear term, and  $\mathcal{H}(h, \mathbf{x}, \tilde{\mathbf{x}})$  is a matrix. Note that if  
 139  $\tilde{\mathbf{x}} \equiv \mathbf{x}$ , then the best second-order model of the form (3) is the Taylor polynomial of degree two for  
 140  $f$  around  $\mathbf{x}$ , and that would give us the exact Newton-type method. However, when the points  $\mathbf{x}$  and  
 141  $\tilde{\mathbf{x}}$  are different, we obtain much more freedom in constructing our models.

142 For using this model in our cubically regularized method (2), we only need to define the gradient  
 143  $\mathbf{g} = \mathcal{G}(h, \mathbf{x}, \tilde{\mathbf{x}})$  and the Hessian estimates  $\mathbf{H} = \mathcal{H}(h, \mathbf{x}, \tilde{\mathbf{x}})$ , and we can also treat them differently  
 144 (using two different helpers  $h_1$  and  $h_2$ , correspondingly). Thus we come to the following general  
 145 second-order (meta)algorithm. We perform  $S$  rounds, the length of each round is  $m \geq 1$ , which is  
 our key parameter:

---

**Algorithm 1** Cubic Newton with helper functions

---

**Input:**  $\mathbf{x}_0 \in \mathbb{R}^d$ ,  $S$ ,  $m \geq 1$ ,  $M > 0$ .

```

1: for  $t = 0, \dots, Sm - 1$  do
2:   if  $t \bmod m = 0$  then
3:     Update  $\tilde{\mathbf{x}}_t$  (using previous states  $\mathbf{x}_{i < t}$ )
4:   else
5:      $\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1}$ 
6:   Form helper functions  $h_1, h_2$ 
7:   Compute the gradient  $\mathbf{g}_t = \mathcal{G}(h_1, \mathbf{x}_t, \tilde{\mathbf{x}}_t)$ , and the Hessian  $\mathbf{H}_t = \mathcal{H}(h_2, \mathbf{x}_t, \tilde{\mathbf{x}}_t)$ 
8:   Compute the cubic step  $\mathbf{x}_{t+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^d} \Omega_{M, \mathbf{g}_t, \mathbf{H}_t}(\mathbf{y}, \mathbf{x}_t)$ 
return  $\mathbf{x}_{out}$  using the history  $(\mathbf{x}_i)_{0 \leq i \leq Sm}$ 

```

---

146  
 147 In Algorithm 1 we update the snapshot  $\tilde{\mathbf{x}}$  regularly every  $m$  iterations. The two possible options are

$$\tilde{\mathbf{x}}_t = \mathbf{x}_{t \bmod m} \quad (\text{use the last iterate}) \quad (4)$$

148 or

$$\tilde{\mathbf{x}}_t = \arg \min_{i \in \{t-m+1, \dots, t\}} f(\mathbf{x}_i) \quad (\text{use the best iterate}) \quad (5)$$

149 Clearly, option (5) is available only in case we can efficiently estimate the function values. However,  
 150 we will see that it serves us with better global convergence guarantees, for the gradient-dominated  
 151 functions.

152 It remains only to specify how we choose the helpers  $h_1$  and  $h_2$ . We need to assume that they are  
 153 somehow similar to  $f$ . Let us present several efficient choices that lead to implementable second-order  
 154 schemes.

155 **3.1 Basic Stochastic Methods**

If the objective function  $f$  is very "expensive" (for example of the form (1) with  $n \rightarrow \infty$ ), one option is to ignore the part  $f - h$  i.e. to approximate it by a zeroth-order approximation:  $f(\mathbf{y}) - h(\mathbf{y}) \approx f(\tilde{\mathbf{x}}) - h(\tilde{\mathbf{x}})$ . Since it is just a constant, we do not need to update  $\tilde{\mathbf{x}}$ . In this case, we have:

$$\mathcal{G}(h_1, \mathbf{x}, \tilde{\mathbf{x}}) := \nabla h_1(\mathbf{x}), \quad \mathcal{H}(h_2, \mathbf{x}, \tilde{\mathbf{x}}) := \nabla^2 h_2(\mathbf{x}).$$

156 To treat this choice of the helpers and motivated by the form of the errors in Lemma 5, we assume the  
157 following similarity assumptions:

**Assumption 2 (Bounded similarity)** *Let for some  $\delta_1, \delta_2 \geq 0$ , it holds*

$$\mathbb{E}_{h_1}[\|\mathcal{G}(h_1, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\|^{3/2}] \leq \delta_1^{3/2}, \quad \mathbb{E}_{h_2}[\|\mathcal{H}(h_2, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla^2 f(\mathbf{x})\|^3] \leq \delta_2^3, \quad \forall \mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{R}^d.$$

158 Under this assumption, we prove the following theorem:

**Theorem 1** *Under Assumptions 1 and 2, and  $M \geq L$ , for an output of Algorithm 1  $\mathbf{x}_{out}$  chosen uniformly at random from  $(\mathbf{x}_i)_{0 \leq i \leq S_m}$ , we have:*

$$\mathbb{E}[\mu_M(\mathbf{x}_{out})] = \mathcal{O}\left(\frac{\sqrt{M}F_0}{S_m} + \frac{\delta_2^3}{M^{3/2}} + \delta_1^{3/2}\right).$$

159 We see that according to this result, we can get  $\mathbb{E}[\mu_M(\mathbf{x}_{out})] \leq \varepsilon^{3/2}$  only for  $\varepsilon > \delta_1$ . In other words,  
160 we can converge only to a certain *neighbourhood around a stationary point*, that is determined by the  
161 error  $\delta_1$  of the stochastic gradients.

162 However, as we will show next, this seemingly pessimistic dependence leads to the same rate of  
163 classical subsampled Cubic Newton methods discovered in [18, 33, 34].

164 Let us discuss now the specific case of stochastic optimization, where  $f$  has the specific form (1),  
165 with  $n$  potentially being very large. In this case, it is customary to sample batches at random and  
166 assume the noise to be bounded in expectation. Precisely speaking, if we assume the standard  
167 assumption that for one index sampled uniformly at random, we have  $\mathbb{E}_i\|\nabla f(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2 \leq \sigma_g^2$   
168 and  $\mathbb{E}_i\|\nabla^2 f(\mathbf{x}) - \nabla^2 f_i(\mathbf{x})\|^3 \leq \sigma_h^3$ , then it is possible to show that for

$$h_1 = \frac{1}{b_g} \sum_{i \in \mathcal{B}_g} f_i \quad \text{and} \quad h_2 = \frac{1}{b_h} \sum_{i \in \mathcal{B}_h} f_i, \quad (6)$$

169 for batches  $\mathcal{B}_g, \mathcal{B}_h \subseteq [n]$  sampled uniformly at random and of sizes  $b_g$  and  $b_h$  respectively, Assump-  
170 tion 2 is satisfied with [28]:  $\delta_1 = \frac{\sigma_g}{\sqrt{b_g}}$  and  $\delta_2 = \tilde{\mathcal{O}}\left(\frac{\sigma_h}{\sqrt{b_h}}\right)$ . Note that we can use the same random  
171 subsets of indices  $\mathcal{B}_g, \mathcal{B}_h$  for all iterations.

**Corollary 1** *In Algorithm 1, let us choose  $M = L$  and  $m = 1$ , with basic helpers (6). Then, according to Theorem 1, for any  $\varepsilon > 0$ , to reach an  $(\varepsilon, L)$ -approximate second-order local minimum, we need at most  $S = \frac{\sqrt{L}F_0}{\varepsilon^{3/2}}$  iterations with  $b_g = \left(\frac{\sigma_g}{\varepsilon}\right)^2$  and  $b_h = \frac{\sigma_h^2}{\varepsilon}$ . Therefore, the total arithmetic complexity of the method becomes*

$$\mathcal{O}\left(\frac{\sigma_g^2}{\varepsilon^{7/2}} + \frac{\sigma_h^2}{\varepsilon^{5/2}} d_{\text{eff}}\right) \times \text{GradCost}.$$

172 It improves upon the complexity  $\mathcal{O}\left(\frac{1}{\varepsilon^4}\right) \times \text{GradCost}$  of the first-order SGD for non-convex optimiza-  
173 tion [12], unless  $d_{\text{eff}} > \frac{1}{\varepsilon^{3/2}}$  (high cost of computing the Hessians).

174 **3.2 Let the Objective Guide Us**

175 If the objective  $f$  is such that we can afford to access its gradients and Hessians from time to time  
176 (functions of the form (1) with  $n < \infty$  and "reasonable"), then we can do better than the previous  
177 chapter. In this case, we can afford to use a better approximation of the term  $f(\mathbf{y}) - h(\mathbf{y})$ . From a  
178 theoretical point of view, we can treat the case when  $f$  is only differentiable once, and thus we can  
179 only use a first-order approximation of  $f - h$ , in this case, we will only be using the hessian of the  
180 helper  $h$  but only gradients of  $f$ . However, in our case, if we assume we have access to gradients then

181 we can also have access to the Hessians of  $f$  as well (from time to time). For this reason, we consider  
 182 a second-order approximation of the term  $f - h$ , if we follow the procedure that we described above  
 183 we find:

$$\mathcal{G}(h_1, \mathbf{x}, \tilde{\mathbf{x}}) := \nabla h_1(\mathbf{x}) - \nabla h_1(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}}) + (\nabla^2 f(\tilde{\mathbf{x}}) - \nabla^2 h_1(\tilde{\mathbf{x}}))(\mathbf{x} - \tilde{\mathbf{x}}) \quad (7)$$

$$\mathcal{H}(h_2, \mathbf{x}, \tilde{\mathbf{x}}) := \nabla^2 h_2(\mathbf{x}) - \nabla^2 h_2(\tilde{\mathbf{x}}) + \nabla^2 f(\tilde{\mathbf{x}}) \quad (8)$$

184 We see that there is an explicit dependence on the snapshot  $\tilde{\mathbf{x}}$  and thus we need to address the question  
 185 of how this snapshot point should be updated in Algorithm 1. In general, we can update it with a  
 186 certain probability  $p \sim \frac{1}{m}$ , and we can use more advanced combinations of past iterates (like the  
 187 average). However, for our purposes, we simply choose option 4 (i.e. the last iterate), thus it is only  
 188 updated once every  $m$  iterations.

189 We also need to address the question of the measure of similarity in this case. Since we are using a  
 190 second-order approximation of  $f - h$ , it is very logical to compare them using the difference between  
 191 their third derivatives or equivalently, the Hessian Lipschitz constant of their difference. Precisely we  
 192 make the following similarity assumption :

**Assumption 3 (Lipschitz similarity)** *Let for some  $\delta_1, \delta_2 \geq 0$ , it holds,  $\forall \mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{R}^d$ :*

$$\mathbb{E}_{h_1} [\|\mathcal{G}(h_1, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\|^{3/2}] \leq \delta_1^{3/2} \|\mathbf{x} - \tilde{\mathbf{x}}\|^3,$$

$$\mathbb{E}_{h_2} [\|\mathcal{H}(h_2, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla^2 f(\mathbf{x})\|^3] \leq \delta_2^3 \|\mathbf{x} - \tilde{\mathbf{x}}\|^3.$$

193 In particular, if  $f - h_1$  and  $f - h_2$  have  $\delta_1$  and  $\delta_2$  Lipschitz Hessians respectively then  $h_1$  and  $h_2$   
 194 satisfy Assumption 3.

195 Under this assumption, we show that the errors resulting from the use of the snapshot can be  
 196 successfully balanced by choosing  $M$  satisfying:

$$4\left(\frac{\delta_1}{M}\right)^{3/2} + 73\left(\frac{\delta_2}{M}\right)^3 \leq \frac{1}{24m^3}. \quad (9)$$

197 And we have the following theorem.

**Theorem 2** *For  $f, h_1, h_2$  verifying Assumptions 1,3. For a regularization parameter  $M$  chosen such that  $M \geq L$  and (9) is satisfied. For an output of Algorithm 1  $\mathbf{x}_{out}$  chosen uniformly at random from  $(\mathbf{x}_i)_{0 \leq i \leq Sm:=T}$ , we have:*

$$\mathbb{E}[\mu_M(\mathbf{x}_{out})] = \mathcal{O}\left(\frac{\sqrt{M}F_0}{Sm}\right),$$

198 In particular, we can choose  $M = \max(L, 32\delta_1 m^2, 16\delta_2 m)$  which gives

$$\mathbb{E}[\mu_M(\mathbf{x}_{out})] = \mathcal{O}\left(\frac{\sqrt{L}F_0}{Sm} + \frac{\sqrt{\delta_2}F_0}{S\sqrt{m}} + \frac{\sqrt{\delta_1}F_0}{S}\right). \quad (10)$$

199 Based on the choices of the helpers  $h_1$  and  $h_2$  we can have many algorithms. We discuss these in  
 200 the following sections. We start by discussing variance reduction and Lazy Hessians which rely on  
 201 sampling batches randomly, then move to core-sets which try to find, more intelligently, representative  
 202 weighted batches of data, after this, we discuss semi-supervised learning and how unlabeled data can  
 203 be used to engineer the helpers. More generally, auxiliary learning tries to leverage auxiliary tasks in  
 204 training a given main task, the auxiliary tasks can be treated as helpers.

### 205 3.3 Variance Reduction and Lazy Hessians

206 The following lemma demonstrates that we can create helper functions  $h$  with lower similarity to the  
 207 main function  $f$  of the form (1) by employing sampling and averaging.

**Lemma 1** Let  $f = \frac{1}{n} \sum_{i=1}^n f_i$  such that all  $f_i$  are twice differentiable and have  $L$ -Lipschitz Hessians. Let  $\mathcal{B} \subset \{1, \dots, n\}$  be of size  $b$  and sampled with replacement uniformly at random, and define  $h_{\mathcal{B}} = \frac{1}{b} \sum_{i \in \mathcal{B}} f_i$ , then  $h_{\mathcal{B}}$  satisfies Assumption 3 with  $\delta_1 = \frac{L}{\sqrt{b}}$  and  $\delta_2 = \mathcal{O}\left(\frac{\sqrt{\log(d)L}}{\sqrt{b}}\right)$ .

208 **Choice of the parameter  $m$  in Algorithm 1.** Minimizing the total arithmetic cost, we choose  
 209  $m = \arg \min_m \#Grad(m, \varepsilon) + d \#Hess(m, \varepsilon)$ , where  $\#Grad(m, \varepsilon)$  and  $\#Hess(m, \varepsilon)$  denote  
 210 the number of gradients and Hessians required to find an  $\varepsilon$  stationary point.

211 Now we are ready to discuss several special cases that are direct consequences from Theorem 2.

212 First, note that choosing  $h_1 = h_2 = f$  gives the classical Cubic Newton method [21], whereas  
 213 choosing  $h_1 = f$  and  $h_2 = 0$ , gives the Lazy Cubic Newton [10]. In both cases, we recuperate the  
 214 known rates of convergence.

215

**General variance reduction.** If we sample batches  $\mathcal{B}_g$  and  $\mathcal{B}_h$  of sizes  $b_g$  and  $b_h$  consecutively at random and choose

$$h_1 = \frac{1}{b_g} \sum_{i \in \mathcal{B}_g} f_i \quad \text{and} \quad h_2 = \frac{1}{b_h} \sum_{i \in \mathcal{B}_h} f_i,$$

216 and use these helpers along with the estimates (7), (8), we obtain the *Variance Reduced Cubic*  
 217 *Newton* algorithm [36, 30]. According to Lemma 1, this choice corresponds to  $\delta_1 = \frac{L}{\sqrt{b_g}}$  and

218  $\delta_2 = \tilde{\mathcal{O}}\left(\frac{L}{\sqrt{b_h}}\right)$ . For  $b_g \sim m^4 \wedge n$ ,  $b_h \sim m^2 \wedge n$  and  $M = L$ , we have the non-convex convergence  
 219 rate  $\mathcal{O}\left(\frac{\sqrt{LF_0}}{Sm}\right)$ , which is the same as that of the cubic Newton algorithm but with a smaller cost per

220 iteration. Minimizing the total arithmetic cost, we can choose  $m = \arg \min_m \frac{dn + d(m^3 \wedge nm) + (m^5 \wedge nm)}{m}$ .

221 Let us denote by  $g^{VR}(n, d)$  the corresponding optimal value. Then we reach an  $(\varepsilon, L)$ -approximate  
 222 second-order local minimum in at most  $\mathcal{O}\left(\frac{g^{VR}(n, d)}{\varepsilon^{3/2}}\right) \times GradCost$  arithmetic operations.

**Variance reduction with Lazy Hessians.** We can also use lazy updates for Hessians combined with variance-reduced gradients. This corresponds to choosing

$$h_1 = \frac{1}{b_g} \sum_{i \in \mathcal{B}_g} f_i \quad \text{and} \quad h_2 = 0,$$

223 which implies (according to Lemma 1) that  $\delta_1 = \frac{L}{\sqrt{b_g}}$  and  $\delta_2 = L$ . In this case, we need  $b_g \sim m^2$

224 to obtain a convergence rate of  $\mathcal{O}\left(\frac{\sqrt{LF_0}}{S\sqrt{m}}\right)$ , which matches the convergence rate of the Lazy Cubic

225 Newton method while using stochastic gradients. We choose this time  $m = \arg \min_m \frac{nd + (m^3 \wedge mn)}{\sqrt{m}}$ ,

226 as before. Let us denote  $g^{Lazy}(n, d)$  the corresponding minimum. Then we guarantee to reach an  
 227  $(\varepsilon, mL)$ -approximate second-order local minimum in at most  $\mathcal{O}\left(\frac{g^{Lazy}(n, d)}{\varepsilon^{3/2}}\right) \times GradCost$  operations.

228 **To be lazy or not to be?** We show that  $g^{Lazy}(n, d) \sim (nd)^{5/6} \wedge n\sqrt{d}$  and  $g^{VR}(n, d) \sim (nd)^{4/5} \wedge$   
 229  $(n^{2/3}d + n)$ . In particular, for  $d \geq n^{2/3}$  we have  $g^{Lazy}(n, d) \leq g^{VR}(n, d)$  and thus for  $d \geq n^{2/3}$   
 230 *it is better to use Lazy Hessians* than variance-reduced Hessians from a gradient equivalent cost  
 231 perspective. We note also that for the Lazy approach, we can keep a factorization of the Hessian (this  
 232 factorization induces most of the cost of solving the cubic subproblem) and thus it is as if we only  
 233 need to solve the subproblem once every  $m$  iterations, so the Lazy approach has a big advantage  
 234 compared to the general approach, and the advantage becomes even bigger for the case of large  
 235 dimensions.

236 Note that according to the theory, we could use the same random batches  $\mathcal{B}_g, \mathcal{B}_h \subseteq [n]$  generated  
 237 once for all iterations. However, using the resampled batches can lead to a more stable convergence.

### 238 3.4 Other Applications

239 The result in (10) is general enough that it can include many other applications that are only limited  
 240 by our imagination. To cite a few such applications there are:

241 **Core sets.** [3] The idea of core sets is simple: can we summarize a potentially big data set using  
 242 only a few (weighted) important examples? Many reasons such as redundancy make the answer yes.  
 243 Devising approaches to find such core sets is outside of the scope of this work, but in general, we  
 244 can see from (10) that if we have batches  $\mathcal{B}_g, \mathcal{B}_h$  such that they are  $(\delta_1, 1)$  and  $(\delta_2, 2)$  similar to  $f$   
 245 respectively, then we can keep reusing the same batch  $\mathcal{B}_g$  for at least  $\sqrt{\frac{L}{\delta_1}}$  times, and  $\mathcal{B}_h$  for  $\frac{L}{\delta_2}$  all  
 246 the while guaranteeing an improved rate. So then if we can design such small batches with small  $\delta_1$   
 247 and  $\delta_2$  then we can keep reusing them, and enjoy the improved rate without needing large batches.

248 **Auxiliary learning.** [4, 2, 32] study how a given task  $f$  can be trained in the presence of auxiliary  
 249 (related) tasks. Our approach can be indeed used for auxiliary learning by treating the auxiliaries as  
 250 helpers. If we compare (10) to the rate that we obtained without the use of the helpers:  $\mathcal{O}(\frac{\sqrt{LF_0}}{S})$ , we  
 251 see that we have a better rate using the helpers/auxiliary tasks when  $\frac{1}{m} + \frac{\sqrt{\delta_2}}{\sqrt{mL}} + \frac{\sqrt{\delta_1}}{\sqrt{L}} \leq 1$ .

252 **Semi-supervised learning.**[35] Semi-supervised learning is a machine learning approach that com-  
 253 bines the use of both labeled data and unlabeled data during training. In general, we can use the  
 254 unlabeled data to construct the helpers, we can start for example by using random labels for the  
 255 helpers and improving the labels with training. There are at least two special cases where our theory  
 256 implies improvement by only assigning random labels to the unlabeled data. In fact, for both regular-  
 257 ized least squares and logistic regression, we notice that the Hessian is independent of the labels (only  
 258 depends on inputs) and thus if the unlabeled data comes from the same distribution as the labeled  
 259 data, then we can use it to construct helpers which, at least theoretically, have  $\delta_1 = \delta_2 = 0$ . Because  
 260 the Hessian is independent of the labels, we can technically endow the unlabeled data with random  
 261 labels. Theorem 2 would imply in this case  $\mathbb{E}[\mu_L(\mathbf{x}_{out})] = \mathcal{O}(\frac{\sqrt{LF_0}}{Sm})$ , where  $S$  is the number of  
 262 times we use labeled data and  $S(m-1)$  is the number of unlabeled data.

## 263 4 Gradient-Dominated Functions

264 We consider now the class of gradient-dominated functions defined below.

**Assumption 4**  $(\tau, \alpha)$ -*gradient dominated.* A function  $f$  is called gradient dominated on set if  
 it holds, for some  $\alpha \geq 1$  and  $\tau > 0$ :

$$f(\mathbf{x}) - f^* \leq \tau \|\nabla f(\mathbf{x})\|^\alpha, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (11)$$

265 Examples of functions satisfying this assumption are convex functions ( $\alpha = 1$ ) and strongly convex  
 266 functions ( $\alpha = 2$ ), see Appendix D.1. For such functions, we can guarantee convergence (in  
 267 expectation) to a *global minimum*, i.e. we can find a point  $\mathbf{x}$  such that  $f(\mathbf{x}) - f^* \leq \varepsilon$ .

268 The Gradient-dominance property is interesting because many non-convex functions have been shown  
 269 to satisfy it [29, 14, 19]. Furthermore, besides convergence to a global minimum, we get accelerated  
 270 rates.

271 We note that for  $\alpha > 3/2$  (and only for this case), we needed to assume the following (stronger)  
 272 inequality:

$$\mathbb{E}f(\mathbf{x}_t) - f^* \leq \tau \mathbb{E}[\|\nabla f(\mathbf{x}_t)\|]^\alpha, \quad (12)$$

273 where the expectation is taken with respect to the iterates  $(\mathbf{x}_t)$  of our algorithms. This is a stronger  
 274 assumption than (11). To avoid using this stronger assumption, we can assume that the iterates belong  
 275 to some compact set  $Q \subset \mathbb{R}^d$  and that the gradient norm is uniformly bounded:  $\forall \mathbf{x} \in Q : \|\nabla f(\mathbf{x})\| \leq$   
 276  $G$ . Then, a  $(\tau, \alpha)$ -gradient dominated on set  $Q$  function is also a  $(\tau G^{\alpha-3/2}, 3/2)$ -gradient dominated  
 277 on this set for any  $\alpha > 3/2$ .

278 In Theorem 3 we extend the results of Theorem 1 to gradient-dominated functions.

**Theorem 3** Under Assumptions 1,2,4, for  $M \geq L$  and  $T := Sm$  we have:

- For  $1 \leq \alpha \leq 3/2$ :  $\mathbb{E}[f(\mathbf{x}_T)] - f^* = \mathcal{O}\left(\left(\frac{\alpha\sqrt{M}\tau^{3/(2\alpha)}}{(3-2\alpha)T}\right)^{\frac{2\alpha}{3-2\alpha}} + \tau\frac{\delta_2^{2\alpha}}{M^\alpha} + \tau\delta_1^\alpha\right)$ .

- For  $3/2 < \alpha \leq 2$ , let  $h_0 = \mathcal{O}\left(\frac{F_0}{(\sqrt{M}\tau^{\frac{3}{2\alpha}})^{\frac{2\alpha}{3-2\alpha}}}\right)$ , then for  $T \geq t_0 = \mathcal{O}(h_0^{\frac{3-2\alpha}{2\alpha}} \log(h_0))$  we have:

$$E[f(\mathbf{x}_T)] - f^* = \mathcal{O}\left(\left(\sqrt{M}\tau^{\frac{3}{2\alpha}}\right)^{\frac{2\alpha}{3-2\alpha}} \left(\frac{1}{2}\right)^{\left(\frac{2\alpha}{3}\right)^{T-t_0}} + \tau\frac{\delta_2^{2\alpha}}{M^\alpha} + \tau\delta_1^\alpha\right).$$

279 Theorem 3 shows (up to the noise level) for  $1 \leq \alpha < 3/2$  a sublinear rate, for  $\alpha = 3/2$  a linear rate  
 280 (obtained by taking the limit  $\alpha \rightarrow 3/2$ ) and a superlinear rate for  $\alpha > 3/2$ .

281 We do the same thing for Theorem 2 which we extend in Theorem 4. In this case, we need to set the  
 282 snapshot line 3 in Algorithm 1) as in 5 i.e. the snapshot corresponds to the state with the smallest  
 283 value of  $f$  during the last  $m$  iterations.

**Theorem 4** Under Assumptions 1,3,4, for  $M = \max(L, 34\delta_1 m^2, 11\delta_2 m)$ , we have:

- For  $1 \leq \alpha \leq 3/2$ :  $\mathbb{E}[f(\mathbf{x}_{Sm})] - f^* = \mathcal{O}\left(\left(\frac{\alpha\sqrt{M}\tau^{3/(2\alpha)}}{(3-2\alpha)Sm}\right)^{\frac{2\alpha}{3-2\alpha}}\right)$ .

- For  $3/2 < \alpha \leq 2$ , let  $h_0 = \mathcal{O}\left(\frac{F_0}{\left(\frac{\sqrt{M}}{m}\tau^{\frac{3}{2\alpha}}\right)^{\frac{2\alpha}{3-2\alpha}}}\right)$ , then for  $S \geq s_0 = \mathcal{O}(h_0^{\frac{3-2\alpha}{2\alpha}} \log(h_0))$  we have:

$$\mathbb{E}[f(\mathbf{x}_{Sm})] - f^* = \left(\left(\frac{\sqrt{M}}{m}\tau^{\frac{3}{2\alpha}}\right)^{\frac{2\alpha}{3-2\alpha}} \left(\frac{1}{2}\right)^{\left(\frac{2\alpha}{3}\right)^{S-s_0}}\right)$$

284 Again, the same behavior is observed as for Theorem 3 but this time without noise (variance reduction  
 285 is working). To the best of our knowledge, this is the first time such analysis is made. As a direct  
 286 consequence of our results, we obtain new global complexities for the variance-reduced and lazy  
 287 variance-reduced Cubic Newton methods on the class of gradient-dominated functions.

288 To compare the statements of Theorems 3 and 4, for convex functions (i.e.  $\alpha = 1$ ), Theorem 3  
 289 guarantees convergence to a  $\varepsilon$ -global minimum in at most  $\mathcal{O}\left(\frac{1}{\varepsilon^{5/2}} + \frac{d}{\varepsilon^{3/2}}\right)$  *GradCost*, whereas  
 290 Theorem 4 only needs  $\mathcal{O}\left(\frac{g(n,d)}{\sqrt{\varepsilon}}\right)$  *GradCost*, where  $g(n,d)$  is either  $g^{Lazy}(n,d) = (nd)^{5/6} \wedge n\sqrt{d}$   
 291 or  $g^{VR}(n,d) = (nd)^{4/5} \wedge (n^{2/3}d + n)$ . See the Appendix D.3 for more details.

## 292 5 Limitations and possible extensions

293 **Estimating similarity between the helpers and the main function.** While we show in this work  
 294 that we can have an improvement over training alone, this supposes that we know the similarity  
 295 constants  $\delta_1, \delta_2$ , hence it will be interesting to have approaches that can adapt to such constants.

296 **Engineering helper functions.** Building helper task with small similarities is also an interesting idea.  
 297 Besides the examples in supervised learning and core-sets that we provide, it is not evident how to do  
 298 it in a generalized way.

299 **Using the helper to regularize the cubic subproblem.** We note that while we proposed to approxi-  
 300 mate the ‘‘cheap’’ part as well in Section 3, one other theoretically viable approach is to keep it intact  
 301 and approximately solve a ‘‘proximal type’’ problem involving  $h$ , this will lead to replacing  $L$  by  $\delta$ ,  
 302 but the subproblem is even more difficult to solve. However our theory suggests that we don’t need  
 303 to solve this subproblem exactly, we only need  $m \geq \frac{L}{\delta}$ . We do not treat this case here.

## 304 6 Conclusion

305 In this work, we proposed a general theory for using auxiliary information in the context of the  
 306 cubically regularized Newton’s method. Our theory encapsulates the classical stochastic methods as  
 307 well as variance reduction and Lazy methods. For auxiliary learning, we showed a provable benefit  
 308 compared to training alone. Besides studying the convergence for general non-convex functions  
 309 for which we show convergence to approximate local minima, we also study gradient-dominated  
 310 functions, for which convergence is accelerated and is to approximate global minima.

## References

- [1] A. Agafonov, D. Kamzolov, P. Dvurechensky, A. Gasnikov, and M. Takáč. Inexact tensor methods and their application to stochastic convex optimization. *arXiv preprint arXiv:2012.15636*, 2020.
- [2] N. Aviv, A. Idan, M. Haggai, C. Gal, and F. Ethan. Auxiliary learning by implicit differentiation. *ICLR 2021*.
- [3] O. Bachem, M. Lucic, and K. Andreas. Practical coresets constructions for machine learning. *arXiv:1703.06476 [stat.ML]*<https://arxiv.org/abs/1703.06476>, 2017.
- [4] S. Baifeng, H. Judy, S. Kate, D. Trevor, and X. Huijuan. Auxiliary task reweighting for minimum-data learning. *34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada*.
- [5] C. Cartis, N. I. Gould, and P. L. Toint. Adaptive cubic regularisation methods for unconstrained optimization. Part I: motivation, convergence and numerical results. *Mathematical Programming*, 127(2):245–295, 2011.
- [6] C. Cartis and K. Scheinberg. Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Mathematical Programming*, 169:337–375, 2018.
- [7] C.-C. Chang and C.-J. Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- [8] E. M. Chayti and S. P. Karimireddy. Optimization with access to auxiliary information. *arXiv:2206.00395 [cs.LG]*, 2022.
- [9] A. R. Conn, N. I. Gould, and P. L. Toint. *Trust region methods*. SIAM, 2000.
- [10] N. Doikov, E. M. Chayti, and M. Jaggi. Second-order optimization with lazy hessians. *arXiv:2212.00781 [math.OC]*, 2022.
- [11] N. Doikov and Y. Nesterov. Minimizing uniformly convex functions by cubic regularization of newton method. *Journal of Optimization Theory and Applications* 189:317–339, 2021.
- [12] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [13] S. Ghadimi, H. Liu, and T. Zhang. Second-order methods with cubic regularization under inexact information. *arXiv preprint arXiv:1710.05782*, 2017.
- [14] M. Hardt and T. Ma. Identity matters in deep learning. *arXiv preprint arXiv:1611.04231*, 2016.
- [15] C. J. Hillar and L.-H. Lim. Most tensor problems are np-hard. *Journal of the ACM (JACM)* 60 45., 2013.
- [16] H. J. Kelley. Gradient theory of optimal flight paths. *Ars Journal*, 30(10):947–954, 1960.
- [17] J. Kiefer and J. Wolfowitz. Stochastic estimation of the maximum of a regression function. *Ann. Math. Statist. Volume 23, Number 3*, 462-466, 1952.
- [18] J. M. Kohler and A. Lucchi. Sub-sampled cubic regularization for non-convex optimization. *arXiv preprint arXiv:1705.05933*, 2017.
- [19] S. Masiha, S. Salehkaleybar, N. He, N. Kiyavash, and P. Thiran. Stochastic second-order methods improve best-known sample complexity of sgd for gradient-dominated functions. In *NeurIPS 2022 - Advances in Neural Information Processing Systems*, volume 35, pages 10862–10875, 2022.
- [20] Y. Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- [21] Y. Nesterov and B. Polyak. Cubic regularization of Newton’s method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.

- 355 [22] T. Nilesh, S. Mitchell, J. Chi, R. Jeffrey, and J. Michael I. Stochastic cubic regularization for  
356 fast nonconvex optimization. *Part of Advances in Neural Information Processing Systems 31*,  
357 2018.
- 358 [23] J. Nocedal and S. Wright. *Numerical optimization*. Springer Science & Business Media, 2006.
- 359 [24] B. T. Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel'noi Matematiki*  
360 *i Matematicheskoi Fiziki*, 3(4):643–653,, 1963.
- 361 [25] B. T. Polyak. Newton’s method and its use in optimization. *European Journal of Operational*  
362 *Research*, 181(3):1086–1096, 2007.
- 363 [26] H. Robbins and S. Monro. A stochastic approximation method the annals of mathematical  
364 statistics. *Vol. 22, No. 3. pp. 400-407*, 1951.
- 365 [27] V. Shamanskii. A modification of Newton’s method. *Ukrainian Mathematical Journal*,  
366 19(1):118–122, 1967.
- 367 [28] J. A. Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends®*  
368 *in Machine Learning*, 8(1-2):1–230, 2015.
- 369 [29] L. uanzhi and Y. Yang. Convergence analysis of two-layer neural networks with relu activation.  
370 *Advances in neural information processing systems*, 30, 2017.
- 371 [30] Z. Wang, Z. Yi, L. Yingbin, and L. Guanghui. Stochastic variance-reduced cubic regularization  
372 for nonconvex optimization. *AISTATS*, 2019.
- 373 [31] B. Woodworth, K. Mishchenko, and F. Bach. Two losses are better than one: Faster optimization  
374 using a cheaper proxy. *arXiv preprint arXiv:2302.03542*, 2023.
- 375 [32] L. Xingyu, S. B. Harjatin, K. George, and H. David. Adaptive auxiliary task weighting for  
376 reinforcement learning. *33rd Conference on Neural Information Processing Systems (NeurIPS*  
377 *2019), Vancouver, Canada*.
- 378 [33] P. Xu, F. Roosta-Khorasani, and M. W. Mahoney. Newton-type methods for non-convex  
379 optimization under inexact Hessian information. *arXiv preprint arXiv:1708.07164* ., 2017.
- 380 [34] P. Xu, J. Yang, F. Roosta-Khorasani, and M. W. Mahoney. Sub-sampled newton methods with  
381 non-uniform sampling. 2016.
- 382 [35] X. Yang, Z. Song, I. King, and Z. Xu. A survey on deep semi-supervised learning. Technical  
383 report, 2021.
- 384 [36] D. Zhou, P. Xu, and Q. Gu. Stochastic variance-reduced cubic regularization methods. *Journal*  
385 *of Machine Learning Research 20 1-47*, 2019.

## 386 0 Experiments

### 387 0.1 To be lazy or not

To verify our findings from Subsection 3.3, we consider a logistic regression problem on the ‘‘a9a’’ data set [7]. Specifically

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-b_i \mathbf{x}^\top \mathbf{a}_i)) + \frac{\lambda}{2} \|\mathbf{x}\|^2,$$

388 where  $\{(\mathbf{a}_i, b_i)\}_{i=1}^N$  are samples from our data, and  $\lambda \geq 0$  is a regularization parameter.

389 We consider the variance-reduced cubic Newton method from [36] (referred to as ‘‘full VR’’), its lazy  
390 version where we do not update the snapshot Hessian (‘‘Lazy VR’’), the stochastic Cubic Newton  
391 method (‘‘SCN’’), the Cubic Newton algorithm (‘‘CN’’), Gradient Descent with line search (‘‘GD’’)  
392 and Stochastic Gradient Descent (‘‘SGD’’). We report the results in terms of time and gradient arithmetic  
393 computations needed to arrive at a given level of convergence.

394 Figure 1 shows how the lazy version saves both time and arithmetic computations without sacrificing  
395 the convergence precision.

Logistic regression: a9a,  $d = 123$ ,  $n = 32561$ , L2-regularization

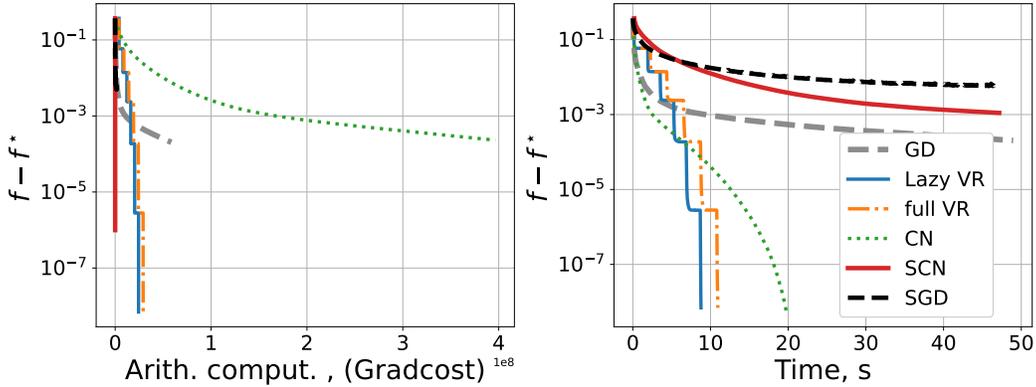


Figure 1: Comparison of the convergence of different algorithms. We see that ‘‘Lazy VR’’ has the same convergence speed as its full version ‘‘full VR’’ and the cubic Newton method ‘‘CN’’ while needing less time and fewer arithmetic computations.

396 The *Gradcost* is computed using the convention that computing one hessian is  $d$  times as expensive  
397 as computing one gradient.

### 398 0.2 Auxiliary Learning

399 Our goal is to show that the helper framework is very general and that it goes beyond the variance  
400 reduction and lazy Hessian computations. For the previously considered problem of training the  
401 logistic regression (using the same ‘‘a9a’’ data set), we suppose that we also have access to unlabeled  
402 data (in this sense this becomes semi-supervised learning). Specifically, we have a labeled dataset  
403  $\mathcal{D}_l = \{(\mathbf{a}_i, b_i)\}_{i=1}^{N_l}$  and an unlabeled data set  $\mathcal{D}_u = \{\mathbf{a}_i\}_{i=N_l+1}^{N_l+N_u}$ , we suppose that both data sets are  
404 sampled from the same distribution  $\mathcal{P}_{(\mathbf{a},b)}$ .

Our goal is to minimize

$$f(\mathbf{x}) = \mathbb{E}_{(\mathbf{a},b) \sim \mathcal{P}_{(\mathbf{a},b)}} [\log(1 + \exp(-b\mathbf{x}^\top \mathbf{a}))].$$

405 A simple computation shows that the Hessian of  $f$  only depends on  $\mathcal{P}_{\mathbf{a}}$ , and, for this reason, we can  
406 use unlabeled data to construct a good approximation of the true Hessian (if we can sample from  $\mathcal{P}_{\mathbf{a}}$ ,  
407 we construct the exact Hessian and thus have a helper  $h$  with  $\delta_1 = \delta_2 = 0$ ).

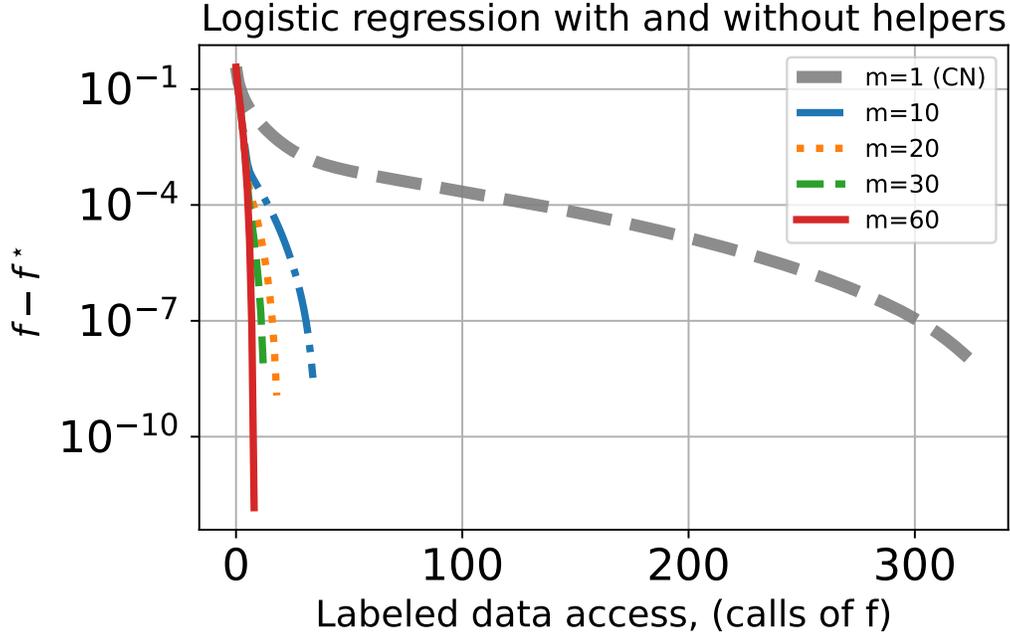


Figure 2: Cubic Newton method with and without using the helper function  $h$ . For  $m = 1$  this is simply the classic Cubic Newton method. To give an intuitive meaning to the plot,  $\frac{1}{m}$  is the percentage of labeled data used during training. We can clearly see that using our approach we benefit a lot from the helper function  $h$ .

Let

$$h(\mathbf{x}) = \mathbb{E}_{\mathbf{a} \sim \mathcal{P}_a, b \sim \text{Random}\{\pm 1\}} [\log(1 + \exp(-b\mathbf{x}^\top \mathbf{a}))],$$

408 where  $\text{Random}\{\pm 1\}$  is any distribution on labels. In our experiments, we use uniform distribution.

409 Figure 2 shows that, indeed we can benefit a lot from using this helper function.

### 410 0.3 Additional experiments

411 We go back to comparing the algorithms in 0.1. We consider now non-convex problems.

First we consider logistic regression with a non-convex regularizer  $\text{Reg}(\mathbf{x}) = \sum_{i=1}^d \frac{\mathbf{x}_i^2}{1+\mathbf{x}_i^2}$ . Precisely speaking we minimize

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-b_i \mathbf{x}^\top \mathbf{a}_i)) + \lambda \text{Reg}(\mathbf{x}).$$

412 Figure 3 shows the results in this case. Again we see that “lazy VR” reduces both time and gradient  
413 equivalent computations without sacrificing the convergence speed.

Logistic regression: a9a, d = 123, n = 32561, Non-convex

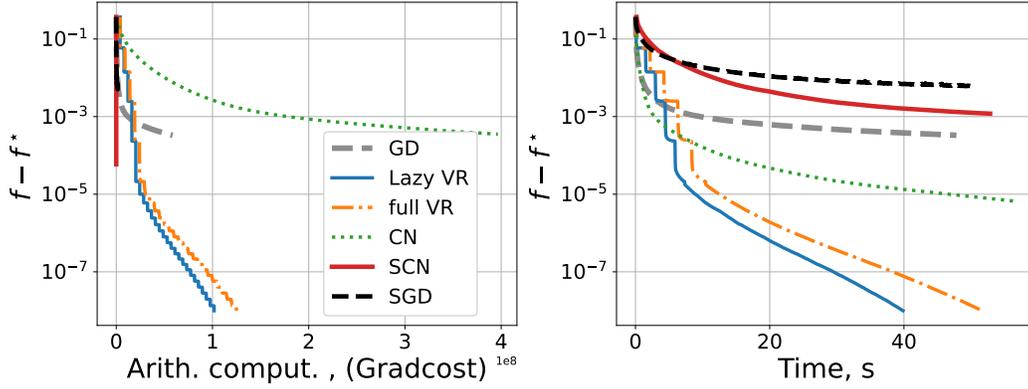


Figure 3: Comparison of the convergence of different algorithms. We see that using our approach we benefit a lot from the helper function  $h$ .

Second, we consider a simple diagonal neural network with L2 loss with data generated from a normal distribution. specifically, we want to minimize

$$f(\mathbf{x} := (\mathbf{u}, \mathbf{v})) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{a}_i^\top \mathbf{u} \circ \mathbf{v} - b_i\|^2 + \frac{\lambda}{2} \|\mathbf{x}\|^2,$$

414 where  $\circ$  is the element-wise vector product.

L2 loss, Diagonal Neural Network: d=100, n=10000

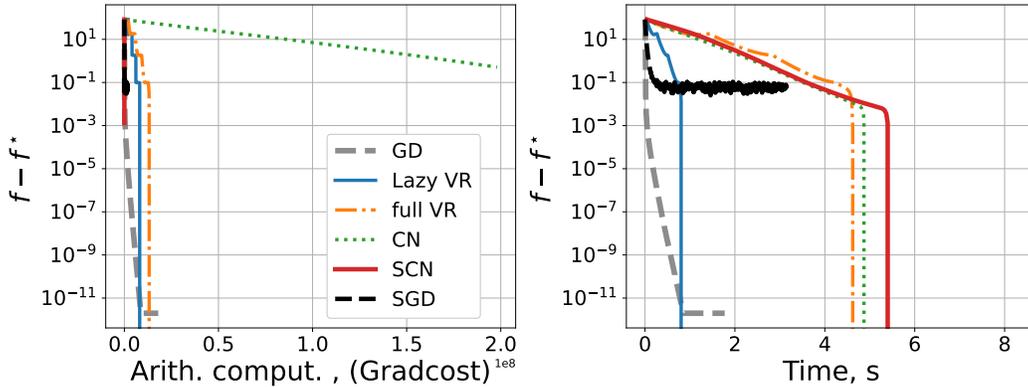


Figure 4: Comparison of the convergence of the different algorithms. Except for gradient descent (“GD”) which performs very well in this case, again the same conclusions as in Figure 2 with respect to “Lazy VR” can be said.

415 Figure 4 shows again that compared to other second-order methods, “Lazy VR” has considerable  
 416 time and computation savings. It also has a close performance to gradient descent with line search  
 417 which performs very well in this case.

#### 418 0.4 Reproducibility

419 We will make our code available with all the details necessary for re-  
 420 producing our results in [https://anonymous.4open.science/r/  
 421 Unified-Convergence-Theory-of-Cubic-Newton-s-method-E4C0/README](https://anonymous.4open.science/r/Unified-Convergence-Theory-of-Cubic-Newton-s-method-E4C0/README).  
 422 md.

423 **A Theoretical Preliminaries**

We consider the general problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

424 Where  $f$  is twice differentiable with  $L$ -Lipschitz Hessian i.e.:

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (13)$$

425 As a direct consequence of (13) (see [21, 20]) we have for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ :

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})\| \leq \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2, \quad (14)$$

426

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{6}\|\mathbf{y} - \mathbf{x}\|^3. \quad (15)$$

427 For  $\mathbf{x}$  and  $\mathbf{x}^+$  defined as in Equation (2) i.e.

$$\mathbf{x}^+ \in \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left\{ \Omega_{M, \mathbf{g}, \mathbf{H}}(\mathbf{y}, \mathbf{x}) := \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{H}(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{M}{6} \|\mathbf{y} - \mathbf{x}\|^3 \right\}. \quad (16)$$

428 The optimality condition of (16) ensures that :

$$\langle \mathbf{g}, \mathbf{x}^+ - \mathbf{x} \rangle + \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{M}{2} r^3 = 0, \quad (17)$$

429 where we denoted  $r = \|\mathbf{x}^+ - \mathbf{x}\|$ .

430 It is also known that the solution to (16) verifies :

$$\mathbf{H} + \frac{M}{2} r \mathbb{I} \succeq 0 \quad (18)$$

431 We start by proving the following Theorem

**Theorem 5** For any  $\mathbf{x} \in \mathbb{R}^d$ , let  $\mathbf{x}^+$  be defined by (2). Then, for  $M \geq L$  we have:

$$f(\mathbf{x}) - f(\mathbf{x}^+) \geq \frac{1}{1008\sqrt{M}} \mu_M(\mathbf{x}^+) + \frac{M\|\mathbf{x} - \mathbf{x}^+\|^3}{72} - \frac{4\|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2}}{\sqrt{M}} - \frac{73\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3}{M^2}.$$

Using (15) with  $\mathbf{y} = \mathbf{x}^+$  and  $\mathbf{x} = \mathbf{x}$  and for  $M \geq L$  we have:

$$\begin{aligned} f(\mathbf{x}^+) &\stackrel{(15)}{\leq} f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{6} r^3 \\ &\stackrel{(17)+(18)}{\leq} f(\mathbf{x}) - \frac{6M-4L}{24} r^3 + \langle \nabla f(\mathbf{x}) - \mathbf{g}, \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle (\nabla^2 f(\mathbf{x}) - \mathbf{H})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \\ &\stackrel{M \geq L}{\leq} f(\mathbf{x}) - \frac{M}{12} r^3 + \langle \nabla f(\mathbf{x}) - \mathbf{g}, \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle (\nabla^2 f(\mathbf{x}) - \mathbf{H})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \end{aligned}$$

432 Using Young's inequality  $xy \leq \frac{x^p}{p} + \frac{y^q}{q} \forall x, y \in \mathbb{R} \forall p, q > 0$  s.t  $\frac{1}{p} + \frac{1}{q} = 1$  we have:

$$\langle \nabla f(\mathbf{x}) - \mathbf{g}, \mathbf{x}^+ - \mathbf{x} \rangle \leq \frac{M}{36} r^3 + \frac{2\sqrt{12}}{3\sqrt{M}} \|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2},$$

433 and

$$\frac{1}{2} \langle (\nabla^2 f(\mathbf{x}) - \mathbf{H})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \leq \frac{M}{36} r^3 + \frac{72}{M^2} \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3.$$

434 Mixing all these ingredients, we get

**Lemma 2** For any  $M \geq L$ , it holds

$$f(\mathbf{x}) - f(\mathbf{x}^+) \geq \frac{M}{36} r^3 - \frac{3}{\sqrt{M}} \|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2} - \frac{72}{M^2} \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3. \quad (19)$$

Using (14) we have:

$$\|\nabla f(\mathbf{x}^+) - \mathbf{g} - \mathbf{H}(\mathbf{x}^+ - \mathbf{x}) + \mathbf{g} - \nabla f(\mathbf{x}) + (\mathbf{H} - \nabla^2 f(\mathbf{x}))(\mathbf{x}^+ - \mathbf{x})\| \leq \frac{L}{2}r^2,$$

435 applying the triangular inequality we get for  $M \geq L$  :

$$\begin{aligned} \|\nabla f(\mathbf{x}^+)\| &\leq \frac{L}{2}r^2 + \|\mathbf{g} + \mathbf{H}(\mathbf{x}^+ - \mathbf{x})\| + \|\nabla f(\mathbf{x}) - \mathbf{g}\| + \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|r \\ &\leq \frac{L + 2M}{2}r^2 + \|\nabla f(\mathbf{x}) - \mathbf{g}\| + \frac{1}{2M}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^2 \\ &\leq \frac{3M}{2}r^2 + \|\nabla f(\mathbf{x}) - \mathbf{g}\| + \frac{1}{2M}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^2 \end{aligned}$$

436 By the convexity of  $x \mapsto x^{3/2}$  we have for any  $(a_i) \geq 0$  :  $(\sum_i a_i x_i)^{3/2} \leq (\sum_i a_i)^{1/2} \sum_i a_i x_i^{3/2}$ ,  
437 applying this to the above inequality we get

**Lemma 3** For any  $M \geq L$ , it holds

$$\frac{1}{\sqrt{M}}\|\nabla f(\mathbf{x}^+)\|^{3/2} \leq 3Mr^3 + \frac{2}{\sqrt{M}}\|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2} + \frac{1}{M^2}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3 \quad (20)$$

We can also bound the smallest eigenvalue of the Hessian. Using the smoothness of the Hessian we have:

$$\begin{aligned} \nabla^2 f(\mathbf{x}^+) &\succeq \nabla^2 f(\mathbf{x}) - L\|\mathbf{x}^+ - \mathbf{x}\|\mathbb{I} \\ &\succeq \mathbf{H} + \nabla^2 f(\mathbf{x}) - \mathbf{H} - Lr\mathbb{I} \\ &\succeq \mathbf{H} - \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|\mathbb{I} - Lr\mathbb{I} \\ &\stackrel{(18)}{\succeq} -\frac{Mr}{2}\mathbb{I} - \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|\mathbb{I} - Lr\mathbb{I} \end{aligned}$$

Which means for  $M \geq L$  we have:

$$-\lambda_{\min}(\nabla^2 f(\mathbf{x}^+)) \leq \frac{3Mr}{2} + \|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|$$

438 Then the convexity of  $x \mapsto x^3$  leads to the following lemma :

**Lemma 4** For any  $M \geq L$ , it holds

$$\frac{-\lambda_{\min}(\nabla^2 f(\mathbf{x}^+))^3}{M^2} \leq 14Mr^3 + \frac{4}{M^2}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3 \quad (21)$$

439 Now the quantity  $\mu_M(\mathbf{x}) = \max(\|\nabla f(\mathbf{x})\|^{3/2}, \frac{-\lambda_{\min}(\nabla^2 f(\mathbf{x}^+))^3}{M^{3/2}})$  which we can be bounded using  
440 Lemmas 3 and 4 :

$$\frac{1}{\sqrt{M}}\mu(\mathbf{x}^+) \leq 14Mr^3 + \frac{2}{\sqrt{M}}\|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2} + \frac{4}{M^2}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3. \quad (22)$$

441 Combining Lemma 2 and (22) we get the inequality given in Theorem 5:

$$f(\mathbf{x}) - f(\mathbf{x}^+) \geq \frac{1}{1008\sqrt{M}}\mu_M(\mathbf{x}^+) + \frac{M}{72}r^3 - \frac{4}{\sqrt{M}}\|\nabla f(\mathbf{x}) - \mathbf{g}\|^{3/2} - \frac{73}{M^2}\|\nabla^2 f(\mathbf{x}) - \mathbf{H}\|^3.$$

## 442 B More on Section 3.1

### 443 B.1 Similarity using sampling

444 One common approach for constructing gradient and Hessian estimates is sub-sampling. The idea  
445 behind sub-sampling is simple: for an objective of the form in (1), we randomly sample two batches  
446  $\mathcal{B}_g$  and  $\mathcal{B}_h$  of sizes  $b_g$  and  $b_h$  consecutively from the distribution  $\mathcal{D}$  and define:

$$\mathbf{g}_{t, \mathcal{B}_g} = \frac{1}{b_g} \sum_{i \in \mathcal{B}_g} \nabla f_i(\mathbf{x}_t) \quad \text{and} \quad \mathbf{H}_{t, \mathcal{B}_h} = \frac{1}{b_h} \sum_{i \in \mathcal{B}_h} \nabla^2 f_i(\mathbf{x}_t) \quad (23)$$

447 In this particular scenario, the ‘‘elementary’’ estimates  $\nabla f(\mathbf{x}_t, \zeta)$  and  $\nabla^2 f(\mathbf{x}_t, \zeta)$  are unbiased, and  
448 we will assume that they satisfy  $\mathbb{E}_i \|\nabla f(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2 \leq \sigma_g^2$  and  $\mathbb{E}_i \|\nabla^2 f(\mathbf{x}) - \nabla^2 f_i(\mathbf{x})\|^3 \leq \sigma_h^3$ .

**Lemma 5** For the estimators defined in (23) we have:

$$\mathbb{E}\|\nabla f(\mathbf{x}_t) - \mathbf{g}_{t, \mathcal{B}_g}\|^2 \leq \frac{\sigma_g^2}{b_g} \quad \text{and} \quad \mathbb{E}\|\nabla^2 f(\mathbf{x}_t) - \mathbf{H}_{t, \mathcal{B}_h}\|^3 \leq \mathcal{O}(\log(d)^{3/2} \frac{\sigma_h^3}{b_h^{3/2}}),$$

where  $\mathcal{O}$  hides constant multiplicative factors.

449 Lemma 5 demonstrates how the utilization of batching can decrease the noise. To simplify things, we  
450 can keep in mind this straightforward rule:

$$\text{If we employ a batch of size } b_a, \text{ then we need to modify } \sigma_a \text{ by } \frac{\sigma_a}{\sqrt{b_a}} \text{ for } a \in \{g, h\}.$$

452 Lemma 5 is based on the following two Lemmas :

**Lemma 6 (Lyapunov's inequality)** For any random variable  $X$  and any  $0 < s < t$  we have

$$\mathbb{E}[|X|^s]^{1/s} \leq \mathbb{E}[|X|^t]^{1/t}.$$

453 and

**Lemma 7** Suppose that  $q \geq 2$ ,  $p \geq 2$ , and fix  $r \geq \max(q, 2 \log(p))$ . Consider i.i.d. random self-adjoint matrices  $Y_1, \dots, Y_N$  with dimension  $p \times p$ ,  $\mathbb{E}[Y_i] = 0$ . It holds that:

$$\left[ \mathbb{E} \left[ \left\| \sum_{i=1}^N Y_i \right\|_2^q \right] \right]^{1/q} \leq 2\sqrt{er} \left\| \left( \sum_{i=1}^N \mathbb{E}[Y_i^2] \right)^{1/2} \right\|_2 + 4er \mathbb{E}[\max_i \|Y_i\|_2^q]^{1/q}.$$

454 Lemma 7 is taken from [36].

455 Now if we have  $X_1, \dots, X_b \in \mathbb{R}^d$ ,  $b$  i.i.d vector-valued random variables such that  $\mathbb{E}[X_i] = \mu$  and  
456  $\mathbb{E}[\|X_i - \mu\|^2] \leq \sigma^2$  then by applying Lemma 6 we get :

$$\mathbb{E} \left[ \left\| \frac{1}{b} \sum_i X_i - \mu \right\|^{3/2} \right] \leq \mathbb{E} \left[ \left\| \frac{1}{b} \sum_i X_i - \mu \right\|^2 \right]^{3/4} \leq \frac{\sigma^{3/2}}{b^{3/4}}.$$

457 When we have  $b$  i.i.d matrix-valued random variables  $Y_1, \dots, Y_b \in \mathbb{R}^{d \times d}$  such that  $\mathbb{E}[Y_i] = \mu$ ,  
458  $\mathbb{E}[\|Y_i - \mu\|^2] \leq \sigma_2^2$  and  $\mathbb{E}[\|Y_i - \mu\|^3] \leq \sigma_3^3$  (by Jensen's inequality  $\sigma_2 \leq \sigma_3$ ), then by applying  
459 Lemma 7 we get:

$$\begin{aligned} \mathbb{E} \left[ \left\| \frac{1}{b} \sum_i Y_i - \mu \right\|^3 \right] &\leq \left( 2\sqrt{\frac{2e \log(d)}{b}} \sigma_2 + \frac{8e \log(d)}{b} \sigma_3 \right)^3 \\ &= \mathcal{O} \left( \frac{\sigma_3^3}{b^{3/2}} \right) \end{aligned}$$

460 These last two inequalities are identical to the statement of Lemma 5.

## 461 B.2 Proof of Theorem 1

462 We use here  $\delta_1 = \sigma_g$  and  $\delta_2 = \sigma_h$ .

463 Combining both Theorem 5, Assumption 2 and Lemma 6 we get:

$$\begin{aligned} \mathbb{E}f(\mathbf{x}_t) - \mathbb{E}f(\mathbf{x}_{t+1}) &\geq \frac{1}{1008\sqrt{M}} \mathbb{E}\mu_M(\mathbf{x}_{t+1}) - \frac{4}{\sqrt{M}} \mathbb{E}\|\nabla f(\mathbf{x}_t) - \mathbf{g}_t\|^{3/2} - \frac{73}{M^2} \mathbb{E}\|\nabla^2 f(\mathbf{x}_t) - \mathbf{H}_t\|^3 \\ &\stackrel{\text{Lemma 6}}{\geq} \frac{1}{1008\sqrt{M}} \mathbb{E}\mu_M(\mathbf{x}_{t+1}) - \frac{4}{\sqrt{M}} \sigma_g^{3/2} - \frac{73}{M^2} \sigma_h^3 \end{aligned}$$

464 By summing the above inequality from  $t = 0$  to  $t = T - 1$  and rearranging we get:

$$\frac{1}{1008T} \sum_{t=1}^T \mathbb{E}\mu_M(\mathbf{x}_t) \leq \sqrt{M} \frac{\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_T)}{T} + 4\sigma_g^{3/2} + \frac{73}{M^{3/2}} \sigma_h^3$$

465 All is left is to use the fact that  $\mathbb{E}f(\mathbf{x}_0) - \mathbb{E}f(\mathbf{x}_T) \leq \mathbb{E}f(\mathbf{x}_0) - f^* = F_0$ , and by definition of  $\mathbf{x}_{out}$

466 :  $\mathbb{E}\mu_M(\mathbf{x}_{out}) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}\mu_M(\mathbf{x}_t)$ , thus :

$$\frac{1}{1008} \mathbb{E}\mu_M(\mathbf{x}_{out}) \leq \frac{\sqrt{M} F_0}{T} + \frac{73}{M^{3/2}} \sigma_h^3 + 4\sigma_g^{3/2}$$

467 **C Helper**

468 **C.1 Proof of Lemma 1**

We have

$$f = \frac{1}{n} \sum_{i=1}^n f_i$$

469 and we suppose that all the  $f_i$ 's have  $L$ -Lipschitz Hessians, so  $f$  too has an  $L$ -Lipschitz Hessian.  
470 Thus  $f_i - f$  has  $2L$ -Lipschitz Hessian.

Applying (1) and 14 to  $f_i - f$  we get

$$\|\mathcal{G}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\| \leq L\|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$

and

$$\|\mathcal{H}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla^2 f(\mathbf{x})\| \leq 2L\|\mathbf{x} - \tilde{\mathbf{x}}\|$$

471 We note also the if  $i$  is chosen at random then  $\mathbb{E}_i \mathcal{G}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) = \nabla f(\mathbf{x})$  and  $\mathbb{E}_i \mathcal{H}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) = \nabla^2 f(\mathbf{x})$ .

By using the properties of variance we have

$$\mathbb{E}_{\mathcal{B}} \|\mathcal{G}(f_{\mathcal{B}}, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\|^2 \leq \frac{\mathbb{E}_i \|\mathcal{G}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\|^2}{b} \leq \frac{L^2}{b} \|\mathbf{x} - \tilde{\mathbf{x}}\|^4$$

Now all that is left is to apply Lemmas 6 and 7 we get

$$\mathbb{E}_{\mathcal{B}} \|\mathcal{G}(f_{\mathcal{B}}, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla f(\mathbf{x})\|^3 \leq \frac{L^3}{b^{3/2}} \|\mathbf{x} - \tilde{\mathbf{x}}\|^3$$

472 And

$$\begin{aligned} \mathbb{E}_{\mathcal{B}} \|\mathcal{H}(f_{\mathcal{B}}, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla^2 f(\mathbf{x})\|^3 &\leq \left(2\sqrt{\frac{2e \log(d)}{b}} + \frac{8e \log(d)}{b}\right)^3 \mathbb{E}_i \|\mathcal{H}(f_i, \mathbf{x}, \tilde{\mathbf{x}}) - \nabla^2 f(\mathbf{x})\|^3 \\ &\leq \left(2\sqrt{\frac{2e \log(d)}{b}} + \frac{8e \log(d)}{b}\right)^3 L^3 \|\mathbf{x} - \tilde{\mathbf{x}}\|^3 \end{aligned}$$

473 **C.2 Proof of Theorem 2**

We use Theorem 5 and denote denoting  $r_{i+1} = \|\mathbf{x}_{i+1} - \mathbf{x}_i\|$  then by the definition of the similarity of  $h_1, h_2$  to  $f$  we have:

$$\mathbb{E}f(\mathbf{x}_{sm+i}) - \mathbb{E}f(\mathbf{x}_{sm+i+1}) \geq \frac{1}{216\sqrt{M}} \mathbb{E}\mu_M(\mathbf{x}_{sm+i+1}) + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right]$$

474 We sum from  $i = 0$  to  $i = m - 1$

$$\mathbb{E}f(\mathbf{x}_{sm}) - \mathbb{E}f(\mathbf{x}_{(s+1)m}) \geq \sum_{i=0}^{m-1} \frac{1}{216\sqrt{M}} \mathbb{E}\mu_M(\mathbf{x}_{sm+i+1}) + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right]$$

475 We note that  $\|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\| \leq \sum_{j=1}^{i-1} r_{sm+j}$ , this means

$$\mathbb{E}f(\mathbf{x}_{sm}) - \mathbb{E}f(\mathbf{x}_{(s+1)m}) \geq \sum_{i=0}^{m-1} \frac{1}{216\sqrt{M}} \mathbb{E}\mu_M(\mathbf{x}_{sm+i+1}) + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \left(\sum_{j=1}^{i-1} r_{sm+j}\right)^3\right]$$

476 We apply now the following inequality (from [10])  $\sum_{k=1}^{m-1} \left(\sum_{i=1}^k r_i\right)^3 \leq \frac{m^3}{3} \sum_{k=1}^{m-1} r_k^3$  true for positive

477 numbers  $\{r_k\}_{k \geq 1}$  and any  $m \geq 1$ . This inequality means :

$$\sum_{i=0}^{m-1} \left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \left(\sum_{j=1}^{i-1} r_{sm+j}\right)^3\right] \geq \left(\frac{M}{72} - \frac{m^3}{3} \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right)\right) \sum_{i=0}^{m-1} r_{sm+i+1}^3$$

478 The above quantity is thus positive if  $\frac{M}{72} - \frac{m^3}{3} \left( \frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2} \right) \geq 0$ .

479 Equivalently, for  $M$  satisfying

$$\boxed{4\left(\frac{\delta_1}{M}\right)^{3/2} + 73\left(\frac{\delta_2}{M}\right)^3 \leq \frac{1}{24m^3}} \quad (24)$$

We have:

$$\mathbb{E}f(x_{sm}) - \mathbb{E}f(x_{(s+1)m}) \geq \frac{m}{216\sqrt{M}} \frac{1}{m} \sum_{i=0}^{m-1} \mathbb{E}\mu(x_{sm+i+1}).$$

We sum from  $s = 0$  to  $s = S - 1$  which gives :

$$\frac{1}{216Sm} \sum_{s=0}^{S-1} \sum_{i=0}^{m-1} \mathbb{E}\mu(x_{sm+i+1}) \leq \frac{\sqrt{M}(f(x_0) - f^*)}{Sm}$$

And thus by definition of  $\mathbf{x}_{out}$  we have:

$$\mathbb{E}\mu(x_{out}) \leq 216 \frac{\sqrt{M}(f(x_0) - f^*)}{Sm}$$

## 480 D Gradient dominated functions

### 481 D.1 Examples of gradient-dominated functions

482 Let us provide several main examples of functions satisfying (11):

**Example 1** Let  $f$  be convex on a bounded convex set  $Q$  of diameter  $D$ , and let solution  $\mathbf{x}^*$  to (1) belong to  $Q$ . Then, we have:

$$f(\mathbf{x}) - f^* \leq \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \leq D \|\nabla f(\mathbf{x})\|, \quad \forall \mathbf{x} \in Q.$$

Therefore,  $f$  is  $(D, 1)$ -gradient dominated.

**Example 2** Let  $f$  be uniformly convex of degree  $s \geq 2$  with some constant  $\sigma > 0$ :

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\sigma}{s} \|\mathbf{y} - \mathbf{x}\|^s, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Then,  $f$  is  $(\frac{s-1}{s}(\frac{1}{\sigma})^{\frac{1}{s-1}}, \frac{s}{s-1})$ -gradient dominated (see, e.g. [11]).

483 In particular, uniformly convex functions of degree  $s = 2$  are known as *strongly convex*, and we see  
484 that they satisfy condition (11) with  $\tau = \frac{1}{2\sigma}$  and  $\alpha = 2$ . However, the function class (11) is much  
485 wider and it includes also some problems with *non-convex objectives* ([21]).

### 486 D.2 Special cases of Theorem 3

487 For **Convex functions**. Theorem 3 implies that for  $M = \max\left(L, \frac{\sigma_h T}{2D}\right)$  we have the rate

$$\mathbb{E}[f(\mathbf{x}_{out})] - f^* = \mathcal{O}\left(\frac{LD^3}{T^2} + \frac{\sigma_h D^2}{T} + \sigma_g D\right). \quad (25)$$

488 Equation (25) has been obtained by [1] but under the much stronger assumption of almost surely  
489 bounded noise. Using the gradient and Hessian estimates in (23), for  $\varepsilon > 0$  and  $M = L$ , to reach an  
490  $\varepsilon$ -global minimum, we need at most  $T = \mathcal{O}\left(\sqrt{\frac{LD^3}{\varepsilon}}\right)$ ,  $b_h = \mathcal{O}\left(\frac{\sigma_h^2 D}{L\varepsilon}\right)$  and  $b_g = \mathcal{O}\left(\frac{\sigma_g^2 D^2}{\varepsilon^2}\right)$ . In other

491 words, we need at most  $\mathcal{O}\left(\frac{\sigma_g^2 L^{1/2} D^{3/2}}{\varepsilon^{5/2}} + d \frac{\sigma_h^2 D^{5/2}}{L^{1/2} \varepsilon^{3/2}}\right)$  *GradCost*.

492 **s-uniformly convex functions**. For this class of functions, using the estimates in (23) and for  $\varepsilon > 0$   
493 and  $M = L$ , we reach an  $\varepsilon$ -global minimum in at most  $T = \mathcal{O}\left(\frac{\sqrt{L}}{s^2} \log\left(\frac{E_0}{\varepsilon}\right)\right)$ ,  $b_h = \mathcal{O}\left(\frac{\sigma_h^2}{s^{4/3} L \varepsilon^{2/3}}\right)$

494 and  $b_g = \mathcal{O}\left(\frac{\sigma_g^2}{s^{8/3} \varepsilon^{4/3}}\right)$  or equivalently  $\tilde{\mathcal{O}}\left(\frac{\sigma_g^2 \sqrt{L}}{s^{14/3} \varepsilon^{4/3}} + d \frac{\sigma_h^2}{s^{10/3} \sqrt{L} \varepsilon^{2/3}}\right)$  *GradCost*.

495  **$\mu$ -strongly convex functions**. For this class of functions, for  $M = L$ , for any  $\varepsilon > 0$  to get  
496  $\mathbb{E}[f(\mathbf{x}_{out})] - f^* < \varepsilon$  we need at most  $T = \mathcal{O}\left(t_0 + \log \log\left(\frac{\mu^3}{L^2 \varepsilon}\right)\right)$ ,  $b_h = \mathcal{O}\left(\frac{\sigma_h^2}{L \sqrt{\mu \varepsilon}}\right)$  and  $b_g = \mathcal{O}\left(\frac{\sigma_g^2}{\mu \varepsilon}\right)$

497 or  $\tilde{\mathcal{O}}\left(\frac{t_0 \sigma_g^2}{\mu \varepsilon} + d \frac{t_0 \sigma_h^2}{L \sqrt{\mu \varepsilon}}\right)$  *GradCost*.

498 **D.3 A special case of Theorem 4**

499 Since we have here many special cases depending on the value of  $\alpha$  and the choices of the helpers, we  
 500 will only consider the case of **convex functions** i.e.  $\alpha = 1, \tau = D$ , but it is easy to apply Theorem 4  
 501 to other cases (like uniformly convex and strongly convex functions).

Theorem 4 implies for convex functions the following:

$$\mathbb{E}[f(\mathbf{x}_{out})] - f^* = \mathcal{O}\left(\frac{\delta_1 D^3}{S^2} + \frac{\delta_2 D^3}{S^2 m} + \frac{LD^3}{S^2 m^2}\right)$$

502 We have the following special cases based on the choice of the helper functions:

- 503 • Cubic Newton [21] corresponds to  $\delta_1 = \delta_2 = 0$ , we get indeed its known rate in the  
 504 convex case. Under the **SOGEO** oracle, we reach an  $\varepsilon$ -global minimum in at most  $\mathcal{O}\left(\frac{nd}{\sqrt{\varepsilon}}\right)$   
 505 *GradCost*.
- 506 • Convex Lazy Cubic Newton (which was not considered in [10]) corresponds to  $\delta_1 = 0, \delta_2 =$   
 507  $L$ , which gives  $\mathbb{E}[f(\mathbf{x}_{out})] - f^* = \mathcal{O}\left(\frac{LD^3}{S^2 m}\right)$ . Under the **SOGEO** oracle, by choosing  
 508  $m = d$  we reach an  $\varepsilon$ -global minimum in at most  $\mathcal{O}\left(\frac{n\sqrt{d}}{\sqrt{\varepsilon}}\right)$  *GradCost*.
- 509 • Convex variance reduced cubic Newton (also not considered by [36, 30]), corresponds  
 510 to each time sampling  $\mathcal{B}_g, \mathcal{B}_h$  of sizes  $b_g, b_h$  consecutively at random and setting  $h_1 =$   
 511  $\frac{1}{b_g} \sum_{i \in \mathcal{B}_g} f_i, h_2 = \frac{1}{b_h} \sum_{i \in \mathcal{B}_h} f_i$ . According to Lemma 1 we have  $\delta_1 = \frac{L}{\sqrt{b_g}}$  and  $\delta_2 =$   
 512  $\tilde{\mathcal{O}}\left(\frac{L}{\sqrt{b_h}}\right)$  so for  $b_g \sim m^4, b_h \sim m^2$  and  $M = L$  we get  $\mathbb{E}[f(\mathbf{x}_{out})] - f^* = \mathcal{O}\left(\frac{LD^3}{S^2 m^2}\right)$ .  
 513 Again, Under the **SOGEO** oracle, by choosing  $m = (nd)^{1/5} 1_{d \leq n^{2/3}} + n^{1/3} 1_{d \geq n^{2/3}}$ , we  
 514 reach an  $\varepsilon$ -global-minimum in at most  $\mathcal{O}\left(\frac{\min((nd)^{4/5}, n^{2/3} d + n)}{\sqrt{\varepsilon}}\right)$  *GradCost*.
- 515 • Variance reduced cubic newton with Lazy Hessians, in this case, by using sampling, we  
 516 have  $\delta_1 = \frac{L}{\sqrt{b_g}}$  and  $\delta_2 = L$ . If we take  $m = (nd)^{1/3} 1_{d \leq \sqrt{n}} + d 1_{d \geq \sqrt{n}}$  then we reach an  
 517  $\varepsilon$ -global-minimum in at most  $\mathcal{O}\left(\frac{\min((nd)^{5/6}, n\sqrt{d})}{\sqrt{\varepsilon}}\right)$  *GradCost*. Again this improves the  
 518 complexity of Lazy Cubic Newton.

519 **D.4 Proof of Theorem 3**

From the previous proof we have

$$\mathbb{E}f(\mathbf{x}_t) - \mathbb{E}f(\mathbf{x}_{t+1}) \geq \frac{1}{1008\sqrt{M}} \mathbb{E}\|\nabla f(\mathbf{x}_{t+1})\|^{3/2} - \frac{4}{\sqrt{M}} \sigma_g^{3/2} - \frac{73}{M^2} \sigma_h^3$$

By the definition of  $(\tau, \alpha)$ -gradient dominated functions we have

$$f(\mathbf{x}_t) - f^* \leq \tau \|\nabla f(\mathbf{x}_t)\|^\alpha$$

So

$$\mathbb{E}\|\nabla f(\mathbf{x}_{t+1})\|^{3/2} \geq \mathbb{E}\left(\frac{f(\mathbf{x}_{t+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}}.$$

If  $\alpha \leq 3/2$ , then by Jensen's inequality we have

$$\mathbb{E}\|\nabla f(\mathbf{x}_{t+1})\|^{3/2} \geq \left(\frac{\mathbb{E}f(\mathbf{x}_{t+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}}.$$

520 For  $\alpha > 3/2$  we need to assume that  $\mathbb{E}f(\mathbf{x}_t) - f^* \leq \tau \mathbb{E}[\|\nabla f(\mathbf{x}_t)\|]^\alpha$  which will give us also

521  $\mathbb{E}\|\nabla f(\mathbf{x}_{t+1})\|^{3/2} \geq \left(\frac{\mathbb{E}f(\mathbf{x}_{t+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}}.$

522 We will consider the sequence  $F_t = \mathbb{E}f(\mathbf{x}_t) - f^*$  and denote  $\gamma = \frac{3}{2\alpha}$ ,  $C = \frac{1}{1008\sqrt{M}\tau^\gamma}$  and

523  $a = \frac{4}{\sqrt{M}} \sigma_g^{3/2} + \frac{25}{M^2} \sigma_h^3$ . Then the sequence  $(F_t)$  satisfies :

$$F_t - F_{t+1} \geq CF_{t+1}^\gamma - a \quad (26)$$

- Case  $\gamma = 1$  then

$$F_{t+1} \geq \frac{F_t + a}{C + 1}$$

So by recurrence, we have:

$$F_t \leq (1 + C)^{-t} F_0 + \sum_{i=0}^{t-1} (1 + C)^{-i} \frac{a}{1 + C} \leq (1 + C)^{-t} F_0 + \frac{a}{C}$$

We note that  $(1 + C)^{-t} \leq \exp(\frac{-Ct}{1+C})$  So

$$F_t \leq \exp(\frac{-Ct}{1+C}) F_0 + \frac{a}{C}$$

- Case  $\gamma \in (1, 2]$  then let  $\tilde{F}_t = \frac{F_t}{C^{1/(1-\gamma)}}$  and  $\tilde{a} = \frac{a}{C^{1/(1-\gamma)}}$  for which we have

$$\tilde{F}_t - \tilde{F}_{t+1} \geq \tilde{F}_{t+1}^\gamma - \tilde{a}$$

524 Now let  $x = \tilde{a}^{1/\gamma}$ , and  $\delta_t = \tilde{F}_t - x$  so :

$$\delta_t - \delta_{t+1} \geq (\delta_{t+1} + x)^\gamma - x^\gamma \geq \delta_{t+1}^\gamma$$

525 Where we used in the last inequality the fact that  $(x + y)^\gamma \geq x^\gamma + y^\gamma$  for  $\gamma \geq 1$  and  $x, y \geq 0$ .

All in all  $\delta_t = \frac{F_t - (\frac{a}{C})^{1/\gamma}}{C^{1/(1-\gamma)}}$  and

$$\delta_t - \delta_{t+1} \geq \delta_{t+1}^\gamma$$

If  $\delta_{t+1} \geq \delta_t/2$  then

$$\frac{1}{(\gamma - 1)\delta_{t+1}^{\gamma-1}} - \frac{1}{(\gamma - 1)\delta_t^{\gamma-1}} \geq \frac{\delta_t^{\gamma-1} - \delta_{t+1}^{\gamma-1}}{(\gamma - 1)\delta_t^{\gamma-1}\delta_{t+1}^{\gamma-1}}$$

By concavity of  $x \mapsto x^{\gamma-1}$  (since  $\gamma \leq 2$ ) we get :

$$\frac{1}{(\gamma - 1)\delta_{t+1}^{\gamma-1}} - \frac{1}{(\gamma - 1)\delta_t^{\gamma-1}} \geq \frac{\delta_t - \delta_{t+1}}{(\gamma - 1)\delta_t\delta_{t+1}^{\gamma-1}} \geq \frac{\delta_{t+1}}{\delta_t} \geq 1/2$$

If  $\delta_{t+1} \leq \delta_t/2$  then

$$\frac{1}{(\gamma - 1)\delta_{t+1}^{\gamma-1}} - \frac{1}{(\gamma - 1)\delta_t^{\gamma-1}} \geq \frac{1}{(\gamma - 1)\delta_t^{\gamma-1}}(2^{\gamma-1} - 1) \geq \frac{2^{\gamma-1} - 1}{(\gamma - 1)\delta_0^{\gamma-1}}$$

526 Using the fact that  $(\delta_t)$  is decreasing.

In all cases we have:

$$\frac{1}{(\gamma - 1)\delta_{t+1}^{\gamma-1}} - \frac{1}{(\gamma - 1)\delta_t^{\gamma-1}} \geq \max(1/2, \frac{2^{\gamma-1} - 1}{(\gamma - 1)\delta_0^{\gamma-1}}) := D$$

By summing from  $t = 0$  to  $t = T - 1$  we get :

$$\frac{1}{(\gamma - 1)\delta_T^{\gamma-1}} \geq DT$$

In other words

$$\delta_T \leq \left( \frac{1}{(\gamma - 1)DT} \right)^{\frac{1}{\gamma-1}}$$

527 - Case  $\gamma < 1$ : then we have

$$F_{t+1} \leq \left( \frac{F_t - F_{t+1} + a}{C} \right)^{1/\gamma}$$

By convexity of  $x \mapsto x^{1/\gamma}$  we get

$$F_{t+1} \leq 2^{1/\gamma-1} \left( \frac{F_t - F_{t+1}}{C} \right)^{1/\gamma} + 2^{1/\gamma-1} \left( \frac{a}{C} \right)^{1/\gamma}$$

Let  $\delta_t = \frac{F_t - 2^{1/\gamma-1} \left( \frac{a}{C} \right)^{1/\gamma}}{2^{1/\gamma} C^{1/(1-\gamma)}}$  then we have

$$\delta_{t+1} \leq (\delta_t - \delta_{t+1})^{1/\gamma}$$

528 The sequence  $(\delta_t)$  is decreasing thus  $\delta_{t+1} \leq \delta_t^{1/\gamma}$  which guarantees a superlinear rate the moment  
529  $\delta_t < 1$ .

530 In fact, we can show that at the beginning  $(\delta_t)$  will decrease at least at a linear rate, and thus it will be  
531 at some point  $< 1$ .

532 We have  $\frac{\delta_t}{\delta_{t+1}} \geq 1 + \frac{\delta_t - \delta_{t+1}}{\delta_{t+1}} \geq 1 + \frac{1}{\delta_{t+1}^{1-\gamma}} \geq 1 + \frac{1}{\delta_0^{1-\gamma}}$

533 Which means  $\delta_{t+1} \leq \left(1 + \frac{1}{\delta_0^{1-\gamma}}\right)^{-1} \delta_t = \left(1 - \frac{1}{1 + \delta_0^{1-\gamma}}\right) \delta_t \leq \exp\left(-\frac{1}{1 + \delta_0^{1-\gamma}}\right) \delta_t$ .

534 To have  $\delta_t \leq 1/2$  we need  $t \geq t_0 = (1 + \delta_0^{1-\gamma}) \log(2\delta_0)$  so that for  $t \geq t_0$  we enjoy a superlinear  
535 rate and we have  $\delta_t \leq \left(\frac{1}{2}\right)^{\left(\frac{1}{\gamma}\right)^{t-t_0}}$

536 This finishes the proof.

#### 537 **D.5 Proof of Theorem 4**

538 In Theorem 4 we made the choice of updating the snapshot in the following way  $\tilde{\mathbf{x}}_{s+1} =$   
539  $\mathbf{x}_{\arg \min_{i \in \{0, \dots, m-1\}} f(\mathbf{x}_{sm+i})}$  which means that  $f(\tilde{\mathbf{x}}_{s+1}) \leq f(\mathbf{x}_{sm+i})$  for all  $i \in \{0, \dots, m-1\}$ .

For  $s \in \{0, \dots, S-1\}$  and  $i \in \{0, \dots, m-1\}$  We have the following inequality

$$\mathbb{E}f(\mathbf{x}_{sm+i}) - \mathbb{E}f(\mathbf{x}_{sm+i+1}) \geq \frac{1}{216\sqrt{M}} \mathbb{E}\|\nabla f(\mathbf{x}_{sm+i+1})\|^{3/2} + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right]$$

540 By definition of gradient-dominated functions we have  $\mathbb{E}\|\nabla f(\mathbf{x}_{t+1})\|^{3/2} \geq \left(\frac{\mathbb{E}f(\mathbf{x}_{t+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}}$ .

541 So

$$\begin{aligned} \mathbb{E}f(\mathbf{x}_{sm+i}) - \mathbb{E}f(\mathbf{x}_{sm+i+1}) &\geq \frac{1}{216\sqrt{M}} \left(\frac{\mathbb{E}f(\mathbf{x}_{sm+i+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}} \\ &\quad + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right] \\ &\geq \frac{1}{216\sqrt{M}} \left(\frac{\mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}} + \mathbb{E}\left[\frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right] \end{aligned}$$

Summing the above inequality from  $i = 0$  to  $i = m-1$  and remarking that  $\tilde{\mathbf{x}}_s = \mathbf{x}_{sm}$  we get

$$\mathbb{E}f(\tilde{\mathbf{x}}_s) - \mathbb{E}f(\mathbf{x}_{(s+1)m}) \geq \frac{m}{216\sqrt{M}} \left(\frac{\mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}} + \mathbb{E}\left[\sum_{i=0}^{m-1} \frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right]$$

By definition of  $\tilde{\mathbf{x}}_{s+1}$  we have  $f(\tilde{\mathbf{x}}_{s+1}) \leq f(\tilde{\mathbf{x}}_{sm+i})$  for all  $i \in \{0, \dots, m-1\}$  which leads to

$$\mathbb{E}f(\tilde{\mathbf{x}}_s) - \mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) \geq \frac{m}{216\sqrt{M}} \left(\frac{\mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) - f^*}{\tau}\right)^{\frac{3}{2\alpha}} + \mathbb{E}\left[\sum_{i=0}^{m-1} \frac{M}{72} r_{sm+i+1}^3 - \left(\frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2}\right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3\right]$$

For  $M$  satisfying (24) we have  $\sum_{i=0}^{m-1} \frac{M}{72} r_{sm+i+1}^3 - \left( \frac{4\delta_1^{3/2}}{\sqrt{M}} + \frac{73\delta_2^3}{M^2} \right) \|\mathbf{x}_{sm+i} - \mathbf{x}_{sm}\|^3 \geq 0$  thus we have

$$\mathbb{E}f(\tilde{\mathbf{x}}_s) - \mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) \geq \frac{m}{216\sqrt{M}} \left( \frac{\mathbb{E}f(\tilde{\mathbf{x}}_{s+1}) - f^*}{\tau} \right)^{\frac{3}{2\alpha}}.$$

Let's define  $F_s = \mathbb{E}f(\tilde{\mathbf{x}}_s) - f^*$  then

$$F_s - F_{s+1} \geq CF_{s+1}^\gamma$$

542 which is a special case of the sequence 26 with  $a = 0$ , thus we can apply our findings from before  
 543 and replace  $a$  by 0.