

---

# A Generic Approach for Escaping Saddle points

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A central challenge to using first-order methods for optimizing nonconvex problems  
2 is the presence of saddle points. First-order methods often get stuck at saddle points,  
3 greatly deteriorating their performance. Typically, to escape from saddles one has  
4 to use second-order methods. However, most works on second-order methods rely  
5 extensively on expensive Hessian-based computations, making them impractical in  
6 large-scale settings. To tackle this challenge, we introduce a generic framework that  
7 minimizes Hessian based computations while at the same time provably converging  
8 to second-order critical points. Our framework carefully alternates between a first-  
9 order and a second-order subroutine, using the latter only close to saddle points,  
10 and yields convergence results competitive to the state-of-the-art. Empirical results  
11 suggest that our strategy also enjoys good practical performance.

## 12 1 Introduction

13 We study nonconvex *finite-sum* problems of the form

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

14 where neither  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  nor the individual functions  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  ( $i \in [n]$ ) are necessarily convex.  
15 We operate in a general nonconvex setting except for few smoothness assumptions like Lipschitz  
16 continuity of the gradient and Hessian. Optimization problems of this form arise naturally in machine  
17 learning and statistics as empirical risk minimization (ERM) and M-estimation respectively.

18 In the large-scale settings, algorithms based on first-order information of functions  $f_i$  are typically  
19 favored as they are relatively inexpensive and scale seamlessly. An algorithm widely used in practice  
20 is stochastic gradient descent (SGD), which has the iterative update:

$$x_{t+1} = x_t - \eta_t \nabla f_{i_t}(x_t), \quad (2)$$

21 where  $i_t \in [n]$  is a randomly chosen index and  $\eta_t$  is a learning rate. Under suitable selection of the  
22 learning rate, we can show that SGD converges to a point  $x$  that, in expectation, satisfies the stationarity  
23 condition  $\|\nabla f(x)\| \leq \epsilon$  in  $O(1/\epsilon^4)$  iterations [14]. This result has two critical weaknesses: (i) It does  
24 not ensure convergence to local optima or second-order critical points; (ii) The rate of convergence of  
25 the SGD algorithm is slow.

26 For general nonconvex problems, one has to settle for a more modest goal than sub-optimality, as  
27 finding the global minimizer of finite-sum nonconvex problem will be in general intractably hard.  
28 Unfortunately, SGD does not even ensure second-order critical conditions such as local optimality  
29 since it can get stuck at saddle points. This issue has recently received considerable attention in the  
30 ML community, especially in the context of deep learning [8–10]. These works argue that saddle  
31 points are highly prevalent in most optimization paths, and are the primary obstacle for training large  
32 deep networks. To tackle this issue and achieve a second-order critical point for which  $\|\nabla f\| \leq \epsilon$   
33 and  $\nabla^2 f \succeq -\sqrt{\epsilon}\mathbb{I}$ , we need algorithms that either use the Hessian explicitly or exploit its structure.

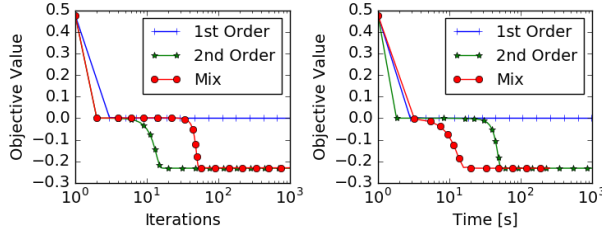


Figure 1: First order methods like GD can potentially get stuck at saddle points. Second-order methods can escape it in very few iterations (as observed in the left plot) but at the cost of expensive Hessian based iterations (see time plot to the right). The proposed framework, which is a novel mix of the two strategies, can escape saddle points *faster* in time by carefully trading off computation and iteration complexity.

34 A key work that explicitly uses Hessians to obtain faster convergence rates is the cubic regularization  
 35 (CR) method [28]. In particular, Nesterov and Polyak [28] showed that CR requires  $O(1/\epsilon^{3/2})$   
 36 iterations to achieve the second-order critical conditions. However, each iteration of CR is expensive  
 37 as it requires computing the Hessian and solving multiple linear systems, each of which has complexity  
 38  $O(d^\omega)$  ( $\omega$  is the matrix multiplication constant), thus, undermining the benefit of its faster convergence.  
 39 Recently, Agarwal et al. [3] designed an algorithm to solve the CR more efficiently, however, it still  
 40 exhibits slower convergence in practice compared to first-order methods. Both of these approaches  
 41 use Hessian based optimization in each iteration, which make them slow in practice.

42 A second line of work focuses on using Hessian information (or its structure) whenever the method  
 43 gets stuck at stationary points that are not second-order critical. To our knowledge, the first work  
 44 in this line is [13], which shows that for a class of functions that satisfy a special property called  
 45 “strict-saddle” property, a noisy variant of SGD can converge to a point close to a local minimum. For  
 46 this class of functions, points close to saddle points have a Hessian with a large negative eigenvalue,  
 47 which proves instrumental in escaping saddle points using an isotropic noise. While such a noise-  
 48 based method is appealing as it only uses first-order information, it has a very bad dependence on the  
 49 dimension  $d$ , and furthermore, the result only holds when the strict-saddle property is satisfied [13].  
 50 More recently, Carmon et al. [6] presented a new faster algorithm that alternates between first-order  
 51 and second-order subroutines. However, their algorithm is designed for the simple case of  $n = 1$   
 52 in (1) and hence, can be expensive in practice.

53 Inspired by this line of work, we develop a general framework for finding second-order critical points.  
 54 The key idea of our framework is to use first-order information for the most part of the optimization  
 55 process and invoke Hessian information only when stuck at stationary points that are not second-order  
 56 critical. We summarize the key idea and main contributions of this paper below.

57 **Main Contributions:** We develop an algorithmic framework for converging to second-order crit-  
 58 ical points and provide convergence analysis for it. Our framework carefully alternates between  
 59 two subroutines that use gradient and Hessian information, respectively, and ensures second-order  
 60 criticality. Furthermore, we present two instantiations of our framework and provide convergence  
 61 rates for them. In particular, we show that a simple instance of our framework, based on SVRG,  
 62 achieves convergence rates competitive with the current state-of-the-art methods; thus highlighting  
 63 the simplicity and applicability of our framework. Finally, we demonstrate the empirical performance  
 64 of a few algorithms encapsulated by our framework and show their superior performance.

65 **Related Work.** There is a vast literature on algorithms for solving optimization problems of the  
 66 form (1). A classical approach for solving such optimization problems is SGD, which dates back  
 67 at least to the seminal work of [35]. Since then, SGD has been a subject of extensive research,  
 68 especially in the convex setting [5, 20, 24, 30]. Recently, new faster methods, called variance  
 69 reduced (VR) methods, have been proposed for convex finite-sum problems. VR methods attain faster  
 70 convergence by reducing the variance in the stochastic updates of SGD, see e.g., [11, 12, 17, 19,  
 71 36, 37]. Accelerated variants of these methods achieve the lower bounds proved in [2, 21], thereby  
 72 settling the question of their optimality. Furthermore, [31] developed an asynchronous framework for  
 73 VR methods and demonstrated their benefits in parallel environments.

74 Most of the aforementioned prior works study stochastic methods in convex or very specialized  
 75 nonconvex settings that admit theoretical guarantees on sub-optimality. For the general nonconvex  
 76 setting, it is only recently that non-asymptotic convergence rate analysis for SGD and its variants was  
 77 obtained in [14], who showed that SGD ensures  $\|\nabla f\| \leq \epsilon$  (in expectation) in  $O(1/\epsilon^4)$  iterations.  
 78 A similar rate for parallel and distributed SGD was shown in [23]. For these problems, Reddi et al.  
 79 [32, 33, 34] proved faster convergence rates that ensure the same optimality criteria in  $O(n + n^{2/3}/\epsilon^2)$ ,  
 80 which is an order  $n^{1/3}$  faster than GD. While these methods ensure convergence to *stationary* points at

81 a faster rate, the question of convergence to local minima (or in general to second-order critical points)  
 82 has not been addressed. To our knowledge, convergence rates to second-order critical points (defined  
 83 in Definition 1) for general nonconvex functions was first studied by [28]. However, each iteration of  
 84 the algorithm in [28] is prohibitively expensive since it requires eigenvalue decompositions, and hence,  
 85 is unsuitable for large-scale high-dimensional problems. More recently, Agarwal et al. [3], Carmon  
 86 et al. [6] presented algorithms for finding second-order critical points by tackling some practical  
 87 issues that arise in [28]. However, these algorithms are either only applicable to a restricted setting  
 88 or heavily use Hessian based computations, making them unappealing from a practical standpoint.  
 89 *Noisy* variants of first-order methods have also been shown to escape saddle points (see [13, 16, 22]),  
 90 however, these methods have strong dependence on either  $n$  or  $d$ , both of which are undesirable.

## 91 2 Background & Problem Setup

92 We assume that each of the functions  $f_i$  in (1) is  $L$ -smooth, i.e.,  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$   
 93 for all  $i \in [n]$ . Furthermore, we assume that the Hessian of  $f$  in (1) is Lipschitz, i.e., we have

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq M\|x - y\|, \quad (3)$$

94 for all  $x, y \in \mathbb{R}^d$ . Such a condition is typically necessary to ensure convergence of algorithms to the  
 95 second-order critical points [28]. In addition to the above smoothness conditions, we also assume  
 96 that the function  $f$  is bounded below, i.e.,  $f(x) \geq B$  for all  $x \in \mathbb{R}^d$ .

97 In order to measure stationarity of an iterate  $x$ , similar to [14, 27, 28], we use the condition  
 98  $\|\nabla f(x)\| \leq \epsilon$ . In this paper, we are interested in convergence to second-order critical points. Thus, in  
 99 addition to stationarity, we also require the solution to satisfy the Hessian condition  $\nabla^2 f(x) \succeq -\gamma\mathbb{I}$   
 100 [28]. For iterative algorithms, we require both  $\epsilon, \gamma \rightarrow 0$  as the number of iterations  $T \rightarrow \infty$ . When  
 101 all saddle points are non-degenerate, such a condition implies convergence to a local optimum.

102 **Definition 1.** An algorithm  $\mathcal{A}$  is said to obtain a point  $x$  that is a  $(\epsilon, \gamma)$ -second order critical point if  
 103  $\mathbb{E}[\|\nabla f(x)\|] \leq \epsilon$  and  $\nabla^2 f(x) \succeq -\gamma\mathbb{I}$ , where the expectation is over any randomness in  $\mathcal{A}$ .

104 We must exercise caution while interpreting results pertaining to  $(\epsilon, \gamma)$ -second order critical points.  
 105 Such points need not be close to any local minima either in objective function value, or in the domain  
 106 of (1). For our algorithms, we use only an Incremental First-order Oracle (IFO) [2] and an Incremental  
 107 Second-order Oracle (ISO), defined below.

108 **Definition 2.** An IFO takes an index  $i \in [n]$  and a point  $x \in \mathbb{R}^d$ , and returns the pair  $(f_i(x), \nabla f_i(x))$ .  
 109 An ISO takes an index  $i \in [n]$ , point  $x \in \mathbb{R}^d$  and vector  $v \in \mathbb{R}^d$  and returns the vector  $\nabla^2 f_i(x)v$ .

110 IFO and ISO calls are typically cheap, with ISO call being relatively more expensive. In many  
 111 practical settings that arise in machine learning, the time complexity of these oracle calls is linear in  
 112  $d$  [4, 29]. For clarity and clean comparison, the dependence of time complexity on Lipschitz constant  
 113  $L, M$ , initial point and any polylog factors in the results is hidden.

## 114 3 Generic Framework

115 In this section, we propose a generic framework for escaping saddle points while solving non-  
 116 convex problems of form (1). One of the primary difficulties in reaching a second-order critical  
 117 point is the presence of saddle points. To evade such points, one needs to use properties  
 118 of both gradients and Hessians. To this end, our framework is based on two core subroutines:  
 119 GRADIENT-FOCUSED-OPTIMIZER and HESSIAN-FOCUSED-OPTIMIZER.

120 The idea is to use these two subroutines, each focused on different aspects of the optimiza-  
 121 tion procedure. GRADIENT-FOCUSED-OPTIMIZER focuses on using gradient information for de-  
 122 creasing the function. On its own, the GRADIENT-FOCUSED-OPTIMIZER might not converge to  
 123 a local minimizer since it can get stuck at a saddle point. Hence, we require the subroutine  
 124 HESSIAN-FOCUSED-OPTIMIZER to help avoid saddle points. A natural idea is to interleave these  
 125 subroutines to obtain a second-order critical point. But it is not even clear if such a procedure even  
 126 converges. We propose a carefully designed procedure that effectively balances these two subroutines,  
 127 which not only provides meaningful theoretical guarantees, but remarkably also translates into strong  
 128 empirical gains in practice.

129 Algorithm 1 provides pseudocode of our framework. Observe that the algorithm is still  
 130 abstract, since it does not specify the subroutines GRADIENT-FOCUSED-OPTIMIZER and  
 131 HESSIAN-FOCUSED-OPTIMIZER. These subroutines determine the crucial update mechanism of the  
 132 algorithm. We will present specific instance of these subroutines in the next section, but we assume  
 133 the following properties to hold for these subroutines.

---

**Algorithm 1** Generic Framework
 

---

```

1: Input - Initial point:  $x^0$ , total iterations  $T$ , error threshold parameters  $\epsilon, \gamma$  and probability  $p$ 
2: for  $t = 1$  to  $T$  do
3:    $(y^t, z^t) = \text{GRADIENT-FOCUSED-OPTIMIZER}(x^{t-1}, \epsilon)$  (refer to G.1 and G.2)
4:   Choose  $u^t$  as  $y^t$  with probability  $p$  and  $z^t$  with probability  $1 - p$ 
5:    $(x^{t+1}, \tau^{t+1}) = \text{HESSIAN-FOCUSED-OPTIMIZER}(u^t, \epsilon, \gamma)$  (refer to H.1 and H.2)
6:   if  $\tau^{t+1} = \emptyset$  then
7:     Output set  $\{x^{t+1}\}$ 
8:   end if
9: end for
10: Output set  $\{y^1, \dots, y^T\}$ 

```

---

134 • **GRADIENT-FOCUSED-OPTIMIZER**: Suppose  $(y, z) = \text{GRADIENT-FOCUSED-OPTIMIZER}(x, n, \epsilon)$ ,  
 135 then there exists positive function  $g : \mathbb{N} \times \mathbb{R}^+ \rightarrow \mathbb{R}^+$ , such that

136 **G.1**  $\mathbb{E}[f(y)] \leq f(x)$ ,

137 **G.2**  $\mathbb{E}[\|\nabla f(y)\|^2] \leq \frac{1}{g(n, \epsilon)} \mathbb{E}[f(x) - f(z)]$ .

138 Here the outputs  $y, z \in \mathbb{R}^d$ . The expectation in the conditions above is over any randomness  
 139 that is a part of the subroutine. The function  $g$  will be critical for the overall rate of Algorithm 1.  
 140 Typically, **GRADIENT-FOCUSED-OPTIMIZER** is a first-order method, since the primary aim of  
 141 this subroutine is to focus on gradient based optimization.

142 • **HESSIAN-FOCUSED-OPTIMIZER**: Suppose  $(y, \tau) = \text{HESSIAN-FOCUSED-OPTIMIZER}(x, n, \epsilon, \gamma)$   
 143 where  $y \in \mathbb{R}^d$  and  $\tau = \{\emptyset, \diamond\}$ . If  $\tau = \emptyset$ , then  $y$  is a  $(\epsilon, \gamma)$ -second order critical point with  
 144 probability at least  $1 - q$ . Otherwise if  $\tau = \diamond$ , then  $y$  satisfies the following condition:

145 **H.1**  $\mathbb{E}[f(y)] \leq f(x)$ ,

146 **H.2**  $\mathbb{E}[f(y)] \leq f(x) - h(n, \epsilon, \gamma)$  when  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$  for some function  $h : \mathbb{N} \times \mathbb{R}^+ \times$   
 147  $\mathbb{R}^+ \rightarrow \mathbb{R}^+$ .

148 Here the expectation is over any randomness in subroutine **HESSIAN-FOCUSED-OPTIMIZER**.  
 149 The two conditions ensure that the objective function value, in expectation, never increases and  
 150 furthermore, decreases with a certain rate when  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$ . In general, this subroutine  
 151 utilizes the Hessian or its properties for minimizing the objective function. Typically, this is the  
 152 most expensive part of the Algorithm 1 and hence, needs to be invoked judiciously.

153 The key aspect of these subroutines is that they, in expectation, never increase the objective function  
 154 value. The functions  $g$  and  $h$  will determine the convergence rate of Algorithm 1. In order to provide  
 155 a concrete implementation, we need to specify the aforementioned subroutines. Before we delve into  
 156 those details, we will provide a generic convergence analysis for Algorithm 1.

### 157 Convergence Analysis

158 **Theorem 1.** Let  $\Delta = f(x^0) - B$  and  $\theta = \min((1 - p)\epsilon^2 g(n, \epsilon), ph(n, \epsilon, \gamma))$ . Also, let set  $\Gamma$   
 159 be the output of Algorithm 1 with **GRADIENT-FOCUSED-OPTIMIZER** satisfying **G.1** and **G.2** and  
 160 **HESSIAN-FOCUSED-OPTIMIZER** satisfying **H.1** and **H.2**. Furthermore,  $T$  be such that  $T > \Delta/\theta$ .

161 Suppose the multiset  $S = \{i_1, \dots, i_k\}$  are  $k$  indices selected independently and uniformly randomly  
 162 from  $\{1, \dots, |\Gamma|\}$ . Then the following holds for the indices in  $S$ :

163 1.  $y^t$ , where  $t \in \{i_1, \dots, i_k\}$ , is a  $(\epsilon, \gamma)$ -critical point with probability at least  $1 - \max(\Delta/(T\theta), q)$ .

164 2. If  $k = O(\log(1/\zeta)/\min(\log(\Delta/(T\theta)), \log(1/q)))$ , with at least probability  $1 - \zeta$ , at least one  
 165 iterate  $y^t$  where  $t \in \{i_1, \dots, i_k\}$  is a  $(\epsilon, \gamma)$ -critical point.

166 The proof of the result is presented in Appendix A. The key point regarding the above result  
 167 is that the overall convergence rate depends on the magnitude of both functions  $g$  and  $h$ . The-  
 168 orem 1 shows that the slowest amongst the subroutines **GRADIENT-FOCUSED-OPTIMIZER** and  
 169 **HESSIAN-FOCUSED-OPTIMIZER** governs the overall rate of Algorithm 1. Thus, it is important to  
 170 ensure that both these procedures have good convergence. Also, note that the optimal setting for  $p$   
 171 based on the result above satisfies  $1/p = 1/\epsilon^2 g(n, \epsilon) + 1/h(n, \epsilon, \gamma)$ . We defer further discussion of  
 172 convergence to next section, where we present more specific convergence and rate analysis.

---

**Algorithm 2** SVRG( $x^0, \epsilon$ )

---

1: **Input:**  $x_m^0 = x^0 \in \mathbb{R}^d$ , epoch length  $m$ , step sizes  $\{\eta_i > 0\}_{i=0}^{m-1}$ , iterations  $T_g, S = \lceil T_g/m \rceil$   
2: **for**  $s = 0$  **to**  $S - 1$  **do**  
3:    $\tilde{x}^s = x_0^{s+1} = x_m^s$   
4:    $g^{s+1} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x}^s)$   
5:   **for**  $t = 0$  **to**  $m - 1$  **do**  
6:     Uniformly randomly pick  $i_t$  from  $\{1, \dots, n\}$   
7:      $v_t^{s+1} = \nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s) + g^{s+1}$   
8:      $x_{t+1}^{s+1} = x_t^{s+1} - \eta_t v_t^{s+1}$   
9:   **end for**  
10: **end for**  
11: **Output:**  $(y, z)$  where  $y$  is Iterate  $x_a$  chosen uniformly random from  $\{x_t^{s+1}\}_{t=0}^{m-1}\}_{s=0}^{S-1}$  and  $z = x_m^S$ .

---

## 173 4 Concrete Instantiations

174 We now present specific instantiations of our framework in this section. Before we state our key  
175 results, we discuss an important subroutine that is used as GRADIENT-FOCUSED-OPTIMIZER for  
176 rest of this paper: SVRG. We give a brief description of the algorithm in this section and show that it  
177 meets the conditions required for a GRADIENT-FOCUSED-OPTIMIZER. SVRG [17, 32] is a stochastic  
178 algorithm recently shown to be very effective for reducing variance in finite-sum problems. We seek  
179 to understand its benefits for nonconvex optimization, with a particular focus on the issue of escaping  
180 saddle points. Algorithm 2 presents SVRG’s pseudocode.

181 Observe that Algorithm 2 is an epoch-based algorithm. At the start of each epoch  $s$ , a full gradient  
182 is calculated at the point  $\tilde{x}^s$ , requiring  $n$  calls to the IFO. Within its inner loop SVRG performs  $m$   
183 stochastic updates. Suppose  $m$  is chosen to be  $O(n)$  (typically used in practice), then the total IFO  
184 calls per epoch is  $\Theta(n)$ . Strong convergence rates have been proved Algorithm 2 in the context  
185 of convex and nonconvex optimization [17, 32]. The following result shows that SVRG meets the  
186 requirements of a GRADIENT-FOCUSED-OPTIMIZER.

187 **Lemma 1.** *Suppose  $\eta_t = \eta = 1/4Ln^{2/3}$ ,  $m = n$  and  $T_g = T_\epsilon$ , which depends on  $\epsilon$ , then Algorithm 2*  
188 *is a GRADIENT-FOCUSED-OPTIMIZER with  $g(n, \epsilon) = T_\epsilon/40Ln^{2/3}$ .*

189 In rest of this section, we discuss approaches using SVRG as a GRADIENT-FOCUSED-OPTIMIZER.  
190 In particular, we propose and provide convergence analysis for two different methods with different  
191 HESSIAN-FOCUSED-OPTIMIZER but which use SVRG as a GRADIENT-FOCUSED-OPTIMIZER.

### 192 4.1 Hessian descent

193 The first approach is based on directly using the eigenvector corresponding to the smallest eigenvalue  
194 as a HESSIAN-FOCUSED-OPTIMIZER. More specifically, when the smallest eigenvalue of the Hessian  
195 is negative and reasonably large in magnitude, the Hessian information can be used to ensure  
196 descent in the objective function value. The pseudo-code for the algorithm is given in Algorithm 3.

197 The key idea is to utilize the minimum eigenvalue information in order to make a descent step.  
198 If  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$  then the idea is to use this information to take a descent step. Note the  
199 subroutine is designed in a fashion such that the objective function value never increases. Thus, it  
200 naturally satisfies the requirement H.1 of HESSIAN-FOCUSED-OPTIMIZER. The following result  
201 shows that HESSIANDESCENT is a HESSIAN-FOCUSED-OPTIMIZER.

202 **Lemma 2.** *HESSIANDESCENT is a HESSIAN-FOCUSED-OPTIMIZER with  $h(n, \epsilon, \gamma) = \frac{\rho}{24M^2} \gamma^3$ .*

203 The proof of the result is presented in Appendix C. With SVRG as GRADIENT-FOCUSED-OPTIMIZER  
204 and HESSIANDESCENT as HESSIAN-FOCUSED-OPTIMIZER, we show the following key result:

205 **Theorem 2.** *Suppose SVRG with  $m = n$ ,  $\eta_t = \eta = 1/4Ln^{2/3}$  for all  $t \in \{1, \dots, m\}$  and  $T_g =$   
206  $40Ln^{2/3}/\epsilon^{1/2}$  is used as GRADIENT-FOCUSED-OPTIMIZER and HESSIANDESCENT is used as  
207 HESSIAN-FOCUSED-OPTIMIZER with  $q = 0$ , then Algorithm 1 finds a  $(\epsilon, \sqrt{\epsilon})$ -second order critical  
208 point in  $T = O(\Delta/\min(p, 1-p)\epsilon^{3/2})$  with probability at least 0.9.*

209 The result directly follows from using Lemma 1 and 2 in Theorem 1. The result shows that the iteration  
210 complexity of Algorithm 1 in this case is  $O(\Delta/\epsilon^{3/2} \min(p, 1-p))$ . Thus, the overall IFO complexity

---

**Algorithm 3** HESSIANDESCENT  $(x, \epsilon, \gamma)$ 

---

- 1: Find  $v$  such that  $\|v\| = 1$ , and with probability at least  $\rho$  the following inequality holds:  $\langle v, \nabla^2 f(x)v \rangle \leq \lambda_{\min}(\nabla^2 f(x)) + \frac{\gamma}{2}$ .
  - 2: Set  $\alpha = |\langle v, \nabla^2 f(x)v \rangle|/M$ .
  - 3:  $u = x - \alpha \text{sign}(\langle v, \nabla f(x) \rangle)v$ .
  - 4:  $y = \arg \min_{z \in \{u, x\}} f(z)$
  - 5: **Output:**  $(y, \diamond)$ .
- 

211 of SVRG algorithm is  $(n + T_g) \times T = O(n/\epsilon^{3/2} + n^{2/3}/\epsilon^2)$ . Since each IFO call takes  $O(d)$  time, the  
212 overall time complexity of all GRADIENT-FOCUSED-OPTIMIZER steps is  $O(nd/\epsilon^{3/2} + n^{2/3}d/\epsilon^2)$ .  
213 To understand the time complexity of HESSIANDESCENT, we need the following result [3].

214 **Proposition 1.** *The time complexity of finding  $v \in \mathbb{R}^d$  that  $\|v\| = 1$ , and with probability at least  $\rho$*   
215 *the following inequality holds:  $\langle v, \nabla^2 f(x)v \rangle \leq \lambda_{\min}(\nabla^2 f(x)) + \frac{\gamma}{2}$  is  $O(nd + n^{3/4}d\gamma^{1/2})$ .*

216 Note that each iteration of Algorithm 1 in this case has just linear dependence on  $d$ . Since the total  
217 number of HESSIANDESCENT iterations is  $O(\Delta/\min(p, 1-p)\epsilon^{3/2})$  and each iteration has the  
218 complexity of  $O(nd + n^{3/4}d/\epsilon^{1/4})$ , using the above remark, we obtain an overall time complexity  
219 of HESSIANDESCENT is  $O(nd/\epsilon^{3/2} + n^{3/4}d/\epsilon^{7/4})$ . Combining this with the time complexity of  
220 SVRG, we get the following result.

221 **Corollary 1.** *The overall running time of Algorithm 1 to find a  $(\epsilon, \sqrt{\epsilon})$ -second order critical point,*  
222 *with parameter settings used in Theorem 2, is  $O(nd/\epsilon^{3/2} + n^{3/4}d/\epsilon^{7/4} + n^{2/3}d/\epsilon^2)$ .*

223 Note that the dependence on  $\epsilon$  is much better in comparison to that of Noisy SGD used in [13].  
224 Furthermore, our results are competitive with [3, 6] in their respective settings, but with a much  
225 simpler algorithm and analysis. We also note that our algorithm is faster than the one proposed in [16],  
226 which has a time complexity of  $O(nd/\epsilon^2)$ .

## 227 4.2 Cubic Descent

228 In this section, we show that the cubic regularization method in [28] can be used as  
229 HESSIAN-FOCUSED-OPTIMIZER. More specifically, here HESSIAN-FOCUSED-OPTIMIZER approx-  
230 imately solves the following optimization problem:

$$y = \arg \min_z \langle \nabla f(x), z - x \rangle + \frac{1}{2} \langle z - x, \nabla^2 f(x)(z - x) \rangle + \frac{M}{6} \|z - x\|^3, \quad (\text{CUBICDESCENT})$$

231 and returns  $(y, \diamond)$  as output. The following result can be proved for this approach.

232 **Theorem 3.** *Suppose SVRG (same as Theorem 2) is used as GRADIENT-FOCUSED-OPTIMIZER and*  
233 *CUBICDESCENT is used as HESSIAN-FOCUSED-OPTIMIZER with  $q = 0$ , then Algorithm 1 finds a*  
234  *$(\epsilon, \sqrt{\epsilon})$ -second order critical point in  $T = O(\Delta/\min(p, 1-p)\epsilon^{3/2})$  with probability at least 0.9.*

235 In principle, Algorithm 1 with CUBICDESCENT as HESSIAN-FOCUSED-OPTIMIZER can con-  
236 verge without the use of GRADIENT-FOCUSED-OPTIMIZER subroutine at each iteration since  
237 it essentially reduces to the cubic regularization method of [28]. However, in practice, we  
238 would expect GRADIENT-FOCUSED-OPTIMIZER to perform most of the optimization and  
239 HESSIAN-FOCUSED-OPTIMIZER to be used for far fewer iterations. Using the method developed  
240 in [28] for solving CUBICDESCENT, we obtain the following corollary.

241 **Corollary 2.** *The overall running time of Algorithm 1 to find a  $(\epsilon, \sqrt{\epsilon})$ -second order critical point,*  
242 *with parameter settings used in Theorem 3, is  $O(nd^\omega/\epsilon^{3/2} + n^{2/3}d/\epsilon^2)$ .*

243 Here  $\omega$  is the matrix multiplication constant. The dependence on  $\epsilon$  is weaker in comparison to  
244 Corollary 1. However, each iteration of CUBICDESCENT is expensive (as seen from the factor  $d^\omega$   
245 in the corollary above) and thus, in high dimensional settings typically encountered in machine learning,  
246 this approach can be expensive in comparison to HESSIANDESCENT.

## 247 4.3 Practical Considerations

248 The focus of this section was to demonstrate the wide applicability of our framework; wherein using  
249 a simple instantiation of this framework, we could achieve algorithms with fast convergence rates.

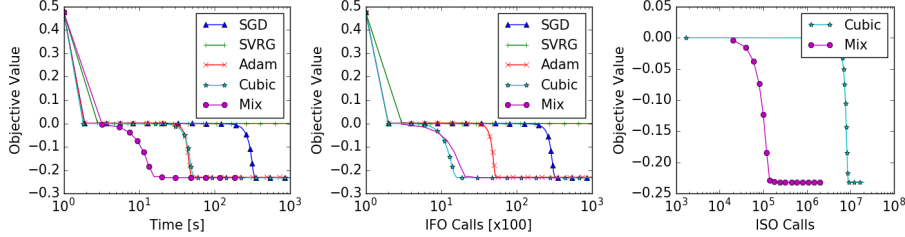


Figure 2: Comparison of various methods on a synthetic problem. Our mix framework successfully escapes saddle point and uses relatively few ISO calls in comparison to CUBICDESCENT.

250 To further achieve good empirical performance, we had to slightly modify these procedures. For  
 251 HESSIAN-FOCUSED-OPTIMIZER, we found stochastic, adaptive and inexact approaches for solving  
 252 HESSIANDESCENT and CUBICDESCENT work well in practice. Due to lack of space, the exact  
 253 description of these modifications is deferred to Appendix F. Furthermore, in the context of deep  
 254 learning, empirical evidence suggests that first-order methods like ADAM [18] exhibit behavior that  
 255 is in congruence with properties G.1 and G.2. While theoretical analysis for a setting where ADAM  
 256 is used as GRADIENT-FOCUSED-OPTIMIZER is still unresolved, we nevertheless demonstrate its  
 257 performance through empirical results in the following section.

## 258 5 Experiments

259 We now present empirical results for our saddle point avoidance technique with an aim to highlight  
 260 three aspects: (i) the framework successfully escapes non-degenerate saddle points, (ii) the framework  
 261 is fast, and (iii) the framework is practical on large-scale problems. All the algorithms are implemented  
 262 on TensorFlow [1]. In case of deep networks, the Hessian-vector product is evaluated using the trick  
 263 presented in [29]. We run our experiments on a commodity machine with Intel<sup>®</sup> Xeon<sup>®</sup> CPU  
 264 E5-2630 v4 CPU, 256GB RAM, and NVidia<sup>®</sup> Titan X (Pascal) GPU.

265 **Synthetic Problem** To demonstrate the fast escape from a saddle point by the proposed method, we  
 266 consider the following simple nonconvex finite-sum problem:

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n x^T A_i x + b_i^T x + \|x\|_{10}^{10} \quad (4)$$

267 Here the parameters are designed such that  $\sum_i b_i = 0$  and  $\sum_i A_i$  matrix has exactly one negative  
 268 eigenvalue of  $-0.001$  and other eigenvalues randomly chosen in the interval  $[1, 2]$ . The total number  
 269 of examples  $n$  is set to be 100,000 and  $d$  is 1000. It is not hard to see that this problem has a non-  
 270 degenerate saddle point at the origin. This allows us to explore the behaviour of different optimization  
 271 algorithms in the vicinity of the saddle point. In this experiment, we compare a mix of SVRG  
 272 and HESSIANDESCENT (as in Theorem 2) with SGD (with constant step size), ADAM, SVRG and  
 273 CUBICDESCENT. The parameter of these algorithms is chosen by grid search so that it gives the  
 274 best performance. The subproblem of CUBICDESCENT was solved with gradient descent [6] until  
 275 the gradient norm of the subproblem is reduced below  $10^{-3}$ . We study the progress of optimization,  
 276 i.e., decrease in function value with wall clock time, IFO calls, and ISO calls. All algorithms were  
 277 initialized with the same starting point very close to origin.

278 The results are presented in Figure 2, which shows that our proposed mix framework was the *fastest*  
 279 to escape the saddle point in terms of wall clock time. We observe that performance of the first order  
 280 methods suffered severely due to the saddle point. Note that SGD eventually escaped the saddle  
 281 point due to inherent noise in the mini-batch gradient. CUBICDESCENT, a second-order method,  
 282 escaped the saddle point faster in terms of iterations using the Hessian information. But operating on  
 283 Hessian information is expensive as a result this method was slow in terms of wall clock time. The  
 284 proposed framework, which is a mix of the two strategies, inherits the best of both worlds by using  
 285 cheap gradient information most of the time and reducing the use of relatively expensive Hessian  
 286 information (ISO calls) by 100x. This resulted in *faster* escape from saddle point in terms of wall  
 287 clock time.

288 **Deep Networks** To investigate the practical performance of the framework for deep learning prob-  
 289 lems, we applied it to two deep autoencoder optimization problems from [15] called “CURVES” and

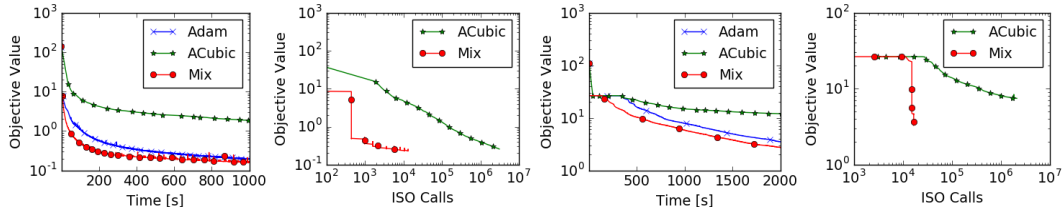


Figure 3: Comparison of various methods on CURVES and MNIST Deep Autoencoder. Our mix approach converges faster than the baseline methods and uses relatively few ISO calls in comparison to APPROXCUBICDESCENT.

290 “MNIST”. Due to their high difficulty, performance on these problems has become a standard bench-  
 291 mark for neural network optimization methods, e.g. [25, 26, 38, 39]. The “CURVES” autoencoder  
 292 consists of an encoder with layers of size (28x28)-400-200-100- 50-25-6 and a symmetric decoder  
 293 totaling in 0.85M parameters. The six units in the code layer were linear and all the other units were  
 294 logistic. The network was trained on 20,000 images and tested on 10,000 new images. The data set  
 295 contains images of curves that were generated from three randomly chosen points in two dimensions.  
 296 The “MNIST” autoencoder consists of an encoder with layers of size (28x28)-1000-500-250-30 and  
 297 a symmetric decoder, totaling in 2.8M parameters. The thirty units in the code layer were linear and  
 298 all the other units were logistic. The network was trained on 60,000 images and tested on 10,000 new  
 299 images. The data set contains images of handwritten digits 0-9. The pixel intensities were normalized  
 300 to lie between 0 and 1.<sup>1</sup>

301 As an instantiation of our framework, we use a mix of ADAM, which is popular in deep learning  
 302 community, and an APPROXCUBICDESCENT for the practical reasons mentioned in Section 4.3. This  
 303 method with ADAM and APPROXCUBICDESCENT. The parameters of these algorithms were chosen  
 304 to produce the best generalization on a held out test set. The regularization parameter  $M$  was chosen  
 305 as the smallest value such that the function value does not fluctuate in the first 10 epochs. We use  
 306 the initialization suggested in [25] and a mini-batch size of 1000 for all the algorithms. We report  
 307 objective function value against wall clock time and ISO calls.

308 The results are presented in Figure 3, which shows that our proposed mix framework was the *fastest*  
 309 to escape the saddle point in terms of wall clock time. ADAM took considerably more time to escape  
 310 the saddle point, especially in the case of MNIST. While APPROXCUBICDESCENT escaped the  
 311 saddle point in relatively fewer iterations, each iteration required considerably large number of ISO  
 312 calls; as a result, the method was extremely slow in terms of wall clock time, despite our efforts to  
 313 improve it via approximations and code optimizations. On the other hand, our proposed framework,  
 314 seamlessly balances these two methods, thereby, resulting in the fast decrease of training loss.

## 315 6 Discussion

316 In this paper, we examined a generic strategy to escape saddle points in nonconvex finite-sum problems  
 317 and presented its convergence analysis. The key intuition is to alternate between a first-order and  
 318 second-order based optimizers; the latter is mainly intended to escape points that are only stationary  
 319 but are not second-order critical points. We presented two different instantiations of our framework  
 320 and provided their detailed convergence analysis. While both our methods explicitly use the Hessian  
 321 information, one can also use noisy first-order methods as HESSIAN-FOCUSED-OPTIMIZER (see for  
 322 e.g. noisy SGD in [13]). In such a scenario, we exploit the negative eigenvalues of the Hessian to  
 323 escape saddle points by using isotropic noise, and do not explicitly use ISO. For these methods, under  
 324 strict-saddle point property [13], we can show convergence to local optima within our framework.

325 We primarily focused on obtaining second-order critical points for nonconvex finite-sums (1). This  
 326 does not necessarily imply low test error or good generalization capabilities. Thus, we should be  
 327 careful when interpreting the results presented in this paper. A detailed discussion or analysis of  
 328 these issues is out of scope of this paper. While a few prior works argue for convergence to local  
 329 optima, the exact connection between generalization and local optima is not well understood, and is  
 330 an interesting open problem. Nevertheless, we believe the techniques presented in this paper can be  
 331 used alongside other optimization tools for faster and better nonconvex optimization.

<sup>1</sup>Data available at: [www.cs.toronto.edu/~jmartens/digs3pts\\_1.mat,mnist\\_all.mat](http://www.cs.toronto.edu/~jmartens/digs3pts_1.mat,mnist_all.mat)



## References

- 332
- 333 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,  
334 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,  
335 Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser,  
336 Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray,  
337 Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul  
338 Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden,  
339 Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale  
340 machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>.  
341 Software available from tensorflow.org.
- 342 [2] Alekh Agarwal and Leon Bottou. A lower bound for the optimization of finite sums.  
343 *arXiv:1410.0723*, 2014.
- 344 [3] Naman Agarwal, Zeyuan Allen Zhu, Brian Bullins, Elad Hazan, and Tengyu Ma. Finding  
345 approximate local minima for nonconvex optimization in linear time. *CoRR*, abs/1611.01146,  
346 2016.
- 347 [4] Naman Agarwal, Brian Bullins, and Elad Hazan. Second order stochastic optimization in linear  
348 time. *CoRR*, abs/1602.03943, 2016.
- 349 [5] Léon Bottou. Stochastic gradient learning in neural networks. *Proceedings of Neuro-Nimes*, 91  
350 (8), 1991.
- 351 [6] Yair Carmon, John C. Duchi, Oliver Hinder, and Aaron Sidford. Accelerated methods for  
352 non-convex optimization. *CoRR*, abs/1611.00756, 2016.
- 353 [7] C. Cartis and K. Scheinberg. Global convergence rate analysis of unconstrained optimization  
354 methods based on probabilistic models. *Mathematical Programming*, pages 1–39, 2017. ISSN  
355 1436-4646. doi: 10.1007/s10107-017-1137-4.
- 356 [8] Anna Choromanska, Mikael Henaff, Michaël Mathieu, Gérard Ben Arous, and Yann LeCun.  
357 The loss surface of multilayer networks. *CoRR*, abs/1412.0233, 2014.
- 358 [9] Yann Dauphin, Harm de Vries, and Yoshua Bengio. Equilibrated adaptive learning rates for non-  
359 convex optimization. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett,  
360 editors, *Advances in Neural Information Processing Systems 28*, pages 1504–1512. Curran  
361 Associates, Inc., 2015.
- 362 [10] Yann N. Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and  
363 Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-  
364 convex optimization. In *Proceedings of the 27th International Conference on Neural Information  
365 Processing Systems, NIPS’14*, pages 2933–2941, 2014.
- 366 [11] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient  
367 method with support for non-strongly convex composite objectives. In *NIPS 27*, pages 1646–  
368 1654, 2014.
- 369 [12] Aaron J Defazio, Tibério S Caetano, and Justin Domke. Finito: A faster, permutable incremental  
370 gradient method for big data problems. *arXiv:1407.2710*, 2014.
- 371 [13] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points - online  
372 stochastic gradient for tensor decomposition. In *Proceedings of The 28th Conference on  
373 Learning Theory, COLT 2015*, pages 797–842, 2015.
- 374 [14] Saeed Ghadimi and Guanghui Lan. Stochastic first- and zeroth-order methods for nonconvex  
375 stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013. doi: 10.1137/  
376 120880811.
- 377 [15] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with  
378 neural networks. *science*, 313(5786):504–507, 2006.
- 379 [16] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M. Kakade, and Michael I. Jordan. How to escape  
380 saddle points efficiently. *CoRR*, abs/1703.00887, 2017.

- 381 [17] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance  
382 reduction. In *NIPS 26*, pages 315–323, 2013.
- 383 [18] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*,  
384 abs/1412.6980, 2014.
- 385 [19] Jakub Konečný, Jie Liu, Peter Richtárik, and Martin Takáč. Mini-Batch Semi-Stochastic  
386 Gradient Descent in the Proximal Setting. *arXiv:1504.04407*, 2015.
- 387 [20] Harold Joseph Kushner and Dean S Clark. *Stochastic approximation methods for constrained  
388 and unconstrained systems*, volume 26. Springer Science & Business Media, 2012.
- 389 [21] Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method.  
390 *arXiv:1507.02000*, 2015.
- 391 [22] Kfir Y. Levy. The power of normalization: Faster evasion of saddle points. *CoRR*,  
392 abs/1611.04831, 2016.
- 393 [23] Xiangru Lian, Yijun Huang, Yuncheng Li, and Ji Liu. Asynchronous Parallel Stochastic Gradient  
394 for Nonconvex Optimization. In *NIPS*, 2015.
- 395 [24] Lennart Ljung. Analysis of recursive stochastic algorithms. *Automatic Control, IEEE Transactions on*,  
396 22(4):551–575, 1977.
- 397 [25] James Martens. Deep learning via hessian-free optimization. In *Proceedings of the 27th  
398 International Conference on Machine Learning (ICML-10)*, pages 735–742, 2010.
- 399 [26] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored ap-  
400 proximate curvature. In *International Conference on Machine Learning*, pages 2408–2417,  
401 2015.
- 402 [27] Yurii Nesterov. *Introductory Lectures On Convex Optimization: A Basic Course*. Springer,  
403 2003.
- 404 [28] Yurii Nesterov and Boris T Polyak. Cubic regularization of newton method and its global  
405 performance. *Mathematical Programming*, 108(1):177–205, 2006.
- 406 [29] Barak A. Pearlmutter. Fast exact multiplication by the hessian. *Neural Computation*, 6(1):  
407 147–160, January 1994. ISSN 0899-7667.
- 408 [30] BT Poljak and Ya Z Tsyypkin. Pseudogradient adaptation and training algorithms. *Automation  
409 and Remote Control*, 34:45–67, 1973.
- 410 [31] Sashank Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex J Smola. On variance  
411 reduction in stochastic gradient descent and its asynchronous variants. In *NIPS 28*, pages  
412 2629–2637, 2015.
- 413 [32] Sashank J. Reddi, Ahmed Hefny, Suvrit Sra, Barnabás Póczos, and Alexander J. Smola. Stochas-  
414 tic variance reduction for nonconvex optimization. In *Proceedings of the 33rd International  
415 Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*,  
416 pages 314–323, 2016.
- 417 [33] Sashank J. Reddi, Suvrit Sra, Barnabás Póczos, and Alexander J. Smola. Fast incremental  
418 method for nonconvex optimization. *CoRR*, abs/1603.06159, 2016.
- 419 [34] Sashank J. Reddi, Suvrit Sra, Barnabás Póczos, and Alexander J. Smola. Fast stochastic methods  
420 for nonsmooth nonconvex optimization. *CoRR*, 2016.
- 421 [35] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical  
422 Statistics*, 22:400–407, 1951.
- 423 [36] Mark W. Schmidt, Nicolas Le Roux, and Francis R. Bach. Minimizing Finite Sums with the  
424 Stochastic Average Gradient. *arXiv:1309.2388*, 2013.
- 425 [37] Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized  
426 loss. *The Journal of Machine Learning Research*, 14(1):567–599, 2013.

- 427 [38] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of  
428 initialization and momentum in deep learning. In *International conference on machine learning*,  
429 pages 1139–1147, 2013.
- 430 [39] Oriol Vinyals and Daniel Povey. Krylov subspace descent for deep learning. In *AISTATS*, pages  
431 1261–1268, 2012.

## Appendix: A Generic Approach for Escaping Saddle points

### 433 A Proof of Theorem 1

434 The case of  $\tau = \emptyset$  can be handled in a straightforward manner, so let us focus on the case where  
 435  $\tau = \diamond$ . We split our analysis into cases, each analyzing the change in objective function value  
 436 depending on second-order criticality of  $y^t$ .

437 We start with the case where the gradient condition of second-order critical point is violated and then  
 438 proceed to the case where the Hessian condition is violated.

439 **Case I:**  $\mathbb{E}[\|\nabla f(y^t)\|] \geq \epsilon$  for some  $t > 0$

440 We first observe the following:  $\mathbb{E}[\|\nabla f(y^t)\|^2] \geq (\mathbb{E}[\|\nabla f(y^t)\|])^2 \geq \epsilon^2$ . This follows from a straight-  
 441 forward application of Jensen's inequality. From this inequality, we have the following:

$$\epsilon^2 \leq \mathbb{E}[\|\nabla f(y^t)\|^2] \leq \frac{1}{g(n, \epsilon)} \mathbb{E}[f(x^{t-1}) - f(z^t)]. \quad (5)$$

This follows from the fact that  $y^t$  is the output of GRADIENT-FOCUSED-OPTIMIZER subroutine, which satisfies the condition that for  $(y, z) = \text{GRADIENT-FOCUSED-OPTIMIZER}(x, n, \epsilon)$ , we have

$$\mathbb{E}[\|\nabla f(y)\|^2] \leq \frac{1}{g(n, \epsilon)} \mathbb{E}[f(x) - f(z)].$$

From Equation (5), we have

$$\mathbb{E}[f(z^t)] \leq \mathbb{E}[f(x^{t-1})] - \epsilon^2 g(n, \epsilon).$$

442 Furthermore, due to the property of non-increasing nature of GRADIENT-FOCUSED-OPTIMIZER, we  
 443 also have  $\mathbb{E}[y^t] \leq \mathbb{E}[f(x^{t-1})]$ .

444 We now focus on the HESSIAN-FOCUSED-OPTIMIZER subroutine. From the property of  
 445 HESSIAN-FOCUSED-OPTIMIZER that the objective function value is non-increasing, we have  
 446  $\mathbb{E}[f(x^t)] \leq \mathbb{E}[f(u^t)]$ . Therefore, combining with the above inequality, we have

$$\begin{aligned} \mathbb{E}[f(x^t)] &\leq \mathbb{E}[f(u^t)] \\ &= p\mathbb{E}[f(y^t)] + (1-p)\mathbb{E}[f(z^t)] \\ &\leq p\mathbb{E}[f(x^{t-1})] + (1-p)(\mathbb{E}[f(x^{t-1})] - \epsilon^2 g(n, \epsilon)) \\ &= \mathbb{E}[f(x^{t-1})] - (1-p)\epsilon^2 g(n, \epsilon). \end{aligned} \quad (6)$$

447 The first equality is due to the definition of  $u^t$  in Algorithm 1. Therefore, when the gradient condition  
 448 is violated, irrespective of whether  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$  or  $\nabla^2 f(y^t) \succeq -\gamma \mathbb{I}$ , the objective function  
 449 value always decreases by at least  $\epsilon^2 g(n, \epsilon)$ .

450 **Case II:**  $\mathbb{E}[\|\nabla f(y^t)\|] < \epsilon$  and  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$  for some  $t > 0$

In this case, we first note that for  $y = \text{HESSIAN-FOCUSED-OPTIMIZER}(x, n, \epsilon, \gamma)$  and  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$ , we have  $\mathbb{E}[f(y)] \leq f(x) - h(n, \epsilon, \gamma)$ . Observe that  $x^t = \text{HESSIAN-FOCUSED-OPTIMIZER}(u^t, n, \epsilon, \gamma)$ . Therefore, if  $u^t = y^t$  and  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$ , then we have

$$\mathbb{E}[f(x^t)|u^t = y^t] \leq f(y^t) - h(n, \epsilon, \gamma) \leq f(x^{t-1}) - h(n, \epsilon, \gamma).$$

451 The second inequality is due to the non-increasing property of GRADIENT-FOCUSED-OPTIMIZER.  
 452 On the other hand, if  $u^t = z^t$ , we have hand, if we have  $\mathbb{E}[f(x^t)|u^t = z^t] \leq f(z^t)$ . This is due to the  
 453 non-increasing property of HESSIAN-FOCUSED-OPTIMIZER. Combining the above two inequalities  
 454 and using the law of total expectation, we get

$$\begin{aligned} \mathbb{E}[f(x^t)] &= p\mathbb{E}[f(x^t)|u^t = y^t] + (1-p)\mathbb{E}[f(x^t)|u^t = z^t] \\ &\leq p(\mathbb{E}[f(y^t)] - h(n, \epsilon, \gamma)) + (1-p)\mathbb{E}[f(z^t)] \\ &\leq p(\mathbb{E}[f(x^{t-1})] - h(n, \epsilon, \gamma)) + (1-p)\mathbb{E}[f(x^{t-1})] \\ &= \mathbb{E}[f(x^{t-1})] - ph(n, \epsilon, \gamma). \end{aligned} \quad (7)$$

455 The second inequality is due to the non-increasing property of GRADIENT-FOCUSED-OPTIMIZER.  
 456 Therefore, when the hessian condition is violated, the objective function value always decreases by at  
 457 least  $ph(n, \epsilon, \gamma)$ .

458 **Case III:**  $\mathbb{E}[\|\nabla f(y^t)\|] < \epsilon$  and  $\nabla^2 f(y^t) \succeq -\gamma\mathbb{I}$  for some  $t > 0$

459 This is the favorable case for the algorithm. The only condition to note is that the objective function  
 460 value will be non-increasing in this case too. This is, again, due to the non-increasing properties of  
 461 subroutines GRADIENT-FOCUSED-OPTIMIZER and HESSIAN-FOCUSED-OPTIMIZER. In general,  
 462 greater the occurrence of this case during the course of the algorithm, higher will the probability that  
 463 the output of our algorithm satisfies the desired property.

The key observation is that Case I & II cannot occur large number of times since each of these cases strictly decreases the objective function value. In particular, from Equation (6) and (7), it is easy to see that each occurrence of Case I & II the following holds:

$$\mathbb{E}[f(x^t)] \leq \mathbb{E}[f(x^{t-1})] - \theta,$$

464 where  $\theta = \min((1-p)\epsilon^2g(n, \epsilon), ph(n, \epsilon, \gamma))$ . Furthermore, the function  $f$  is lower bounded by  
 465  $B$ , thus, Case I & II cannot occur more than  $(f(x^0) - B)/\theta$  times. Therefore, the probability of  
 466 occurrence of Case III is at least  $1 - (f(x^0) - B)/(T\theta)$ , which completes the first part of the proof.

467 The second part of the proof simply follows from first part. As seen above, the probability of Case I  
 468 & II is at most  $(f(x^0) - B)/T\theta$ . Therefore, probability that an element of the set  $S$  falls in Case III  
 469 is at least  $1 - ((f(x^0) - B)/T\theta)^k$ , which gives us the required result for the second part.

## 470 B Proof of Lemma 1

471 *Proof.* The proof follows from the analysis in [32] with some additional reasoning. We need to show  
 472 two properties: **G.1** and **G.2**, both of which are based on objective function value. To this end, we  
 473 start with an update in the  $s^{\text{th}}$  epoch. We have the following:

$$\begin{aligned} \mathbb{E}[f(x_{t+1}^{s+1})] &\leq \mathbb{E}[f(x_t^{s+1}) + \langle \nabla f(x_t^{s+1}), x_{t+1}^{s+1} - x_t^{s+1} \rangle + \frac{L}{2}\|x_{t+1}^{s+1} - x_t^{s+1}\|^2] \\ &\leq \mathbb{E}[f(x_t^{s+1}) - \eta_t \|\nabla f(x_t^{s+1})\|^2 + \frac{L\eta_t^2}{2}\|v_t^{s+1}\|^2]. \end{aligned} \quad (8)$$

The first inequality is due to  $L$ -smoothness of the function  $f$ . The second inequality simply follows from the unbiasedness of SVRG update in Algorithm 2. For the analysis of the algorithm, we need the following Lyapunov function:

$$A_t^{s+1} := \mathbb{E}[f(x_t^{s+1}) + \mu_t \|x_t^{s+1} - \tilde{x}^s\|^2].$$

This function is a combination of objective function and the distance of the current iterate from the latest snapshot  $\tilde{x}^s$ . Note that the term  $\mu_t$  is introduced only for the analysis and is not part of the algorithm (see Algorithm 2). Here  $\{\mu_t\}_{t=0}^m$  is chosen such the following holds:

$$\mu_t = \mu_{t+1}(1 + \eta_t\beta_t + 2\eta_t^2L^2) + \eta_t^2L^3,$$

474 for all  $t \in \{0, \dots, m-1\}$  and  $\mu_m = 0$ . For bounding the Lyapunov function  $A$ , we need the following  
 475 bound on the distance of the current iterate from the latest snapshot:

$$\begin{aligned} \mathbb{E}[\|x_{t+1}^{s+1} - \tilde{x}^s\|^2] &= \mathbb{E}[\|x_{t+1}^{s+1} - x_t^{s+1} + x_t^{s+1} - \tilde{x}^s\|^2] \\ &= \mathbb{E}[\|x_{t+1}^{s+1} - x_t^{s+1}\|^2 + \|x_t^{s+1} - \tilde{x}^s\|^2 + 2\langle x_{t+1}^{s+1} - x_t^{s+1}, x_t^{s+1} - \tilde{x}^s \rangle] \\ &= \mathbb{E}[\eta_t^2\|v_t^{s+1}\|^2 + \|x_t^{s+1} - \tilde{x}^s\|^2 - 2\eta_t\mathbb{E}[\langle \nabla f(x_t^{s+1}), x_t^{s+1} - \tilde{x}^s \rangle]] \\ &\leq \mathbb{E}[\eta_t^2\|v_t^{s+1}\|^2 + \|x_t^{s+1} - \tilde{x}^s\|^2] + 2\eta_t\mathbb{E}\left[\frac{1}{2\beta_t}\|\nabla f(x_t^{s+1})\|^2 + \frac{1}{2}\beta_t\|x_t^{s+1} - \tilde{x}^s\|^2\right]. \end{aligned} \quad (9)$$

476 The second equality is due to the unbiasedness of the update of SVRG. The last inequality follows  
 477 from a simple application of Cauchy-Schwarz and Young's inequality. Substituting Equation (8) and

478 Equation (9) into the Lypunov function  $A_{t+1}^{s+1}$ , we obtain the following:

$$\begin{aligned}
A_{t+1}^{s+1} &\leq \mathbb{E}[f(x_t^{s+1}) - \eta_t \|\nabla f(x_t^{s+1})\|^2 + \frac{L\eta_t^2}{2} \|v_t^{s+1}\|^2] \\
&\quad + \mathbb{E}[\mu_{t+1}\eta_t^2 \|v_t^{s+1}\|^2 + \mu_{t+1} \|x_t^{s+1} - \tilde{x}^s\|^2] \\
&\quad + 2\mu_{t+1}\eta_t \mathbb{E}\left[\frac{1}{2\beta_t} \|\nabla f(x_t^{s+1})\|^2 + \frac{1}{2}\beta_t \|x_t^{s+1} - \tilde{x}^s\|^2\right] \\
&\leq \mathbb{E}[f(x_t^{s+1}) - \left(\eta_t - \frac{\mu_{t+1}\eta_t}{\beta_t}\right) \|\nabla f(x_t^{s+1})\|^2] \\
&\quad + \left(\frac{L\eta_t^2}{2} + \mu_{t+1}\eta_t^2\right) \mathbb{E}[\|v_t^{s+1}\|^2] + (\mu_{t+1} + \mu_{t+1}\eta_t\beta_t) \mathbb{E}[\|x_t^{s+1} - \tilde{x}^s\|^2]. \tag{10}
\end{aligned}$$

479 To further bound this quantity, we use Lemma 3 to bound  $\mathbb{E}[\|v_t^{s+1}\|^2]$ , so that upon substituting it in  
480 Equation (10), we see that

$$\begin{aligned}
A_{t+1}^{s+1} &\leq \mathbb{E}[f(x_t^{s+1})] - \left(\eta_t - \frac{\mu_{t+1}\eta_t}{\beta_t} - \eta_t^2 L - 2\mu_{t+1}\eta_t^2\right) \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] \\
&\quad + [\mu_{t+1}(1 + \eta_t\beta_t + 2\eta_t^2 L^2) + \eta_t^2 L^3] \mathbb{E}[\|x_t^{s+1} - \tilde{x}^s\|^2] \\
&\leq A_t^{s+1} - \left(\eta_t - \frac{\mu_{t+1}\eta_t}{\beta_t} - \eta_t^2 L - 2\mu_{t+1}\eta_t^2\right) \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2].
\end{aligned}$$

481 The second inequality follows from the definition of  $\mu_t$  and  $A_t^{s+1}$ . Since  $\eta_t = \eta = 1/(4Ln^{2/3})$  for  
482  $j > 0$  and  $t \in \{0, \dots, j-1\}$ ,

$$A_j^{s+1} \leq A_0^{s+1} - v_n \sum_{t=0}^{j-1} \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2], \tag{11}$$

where

$$v_n = \left(\eta_t - \frac{\mu_{t+1}\eta_t}{\beta_t} - \eta_t^2 L - 2\mu_{t+1}\eta_t^2\right).$$

We will prove that for the given parameter setting  $v_n > 0$  (see the proof below). With  $v_n > 0$ , it is easy to see that  $A_j^{s+1} \leq A_0^{s+1}$ . Furthermore, note that  $A_0^{s+1} = \mathbb{E}[f(x_0^{s+1}) + \mu_0 \|x_0^{s+1} - \tilde{x}^s\|^2] = \mathbb{E}[f(x_0^{s+1})]$  since  $x_0^{s+1} = \tilde{x}^s$  (see Algorithm 2). Also, we have

$$\mathbb{E}[f(x_j^{s+1}) + \mu_j \|x_j^{s+1} - \tilde{x}^s\|^2] \leq \mathbb{E}[f(x_0^{s+1})]$$

483 and thus, we obtain  $\mathbb{E}[f(x_j^{s+1})] \leq \mathbb{E}[f(x_0^{s+1})]$  for all  $j \in \{0, \dots, m\}$ . Furthermore, using simple  
484 induction and the fact that  $x_0^{s+1} = x_m^s$  for all epoch  $s \in \{0, \dots, S-1\}$ , it easy to see that  
485  $\mathbb{E}[f(x_j^{s+1})] \leq f(x^0)$ . Therefore, with the definition of  $y$  specified in the output of Algorithm 2, we  
486 see that the condition **G.1** of GRADIENT-FOCUSED-OPTIMIZER is satisfied for SVRG algorithm.

487 We now prove that  $v_n > 0$  and also **G.2** of GRADIENT-FOCUSED-OPTIMIZER is satisfied for SVRG  
488 algorithm. By using telescoping the sum with  $j = m$  in Equation (11), we obtain

$$\sum_{t=0}^{m-1} \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] \leq \frac{A_0^{s+1} - A_m^{s+1}}{v_n}.$$

489 This inequality in turn implies that

$$\sum_{t=0}^{m-1} \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] \leq \frac{\mathbb{E}[f(\tilde{x}^s) - f(\tilde{x}^{s+1})]}{v_n}, \tag{12}$$

490 where we used that  $A_m^{s+1} = \mathbb{E}[f(x_m^{s+1})] = \mathbb{E}[f(\tilde{x}^{s+1})]$  (since  $\mu_m = 0$ ), and that  $A_0^{s+1} = \mathbb{E}[f(\tilde{x}^s)]$   
491 (since  $x_0^{s+1} = \tilde{x}^s$ ). Now sum over all epochs to obtain

$$\frac{1}{T_g} \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] \leq \frac{\mathbb{E}[f(x^0) - f(x_m^S)]}{T_g v_n}. \tag{13}$$

492 Here we used the the fact that  $\tilde{x}^0 = x^0$ . To obtain a handle on  $v_n$  and complete our analysis, we will  
493 require an upper bound on  $\mu_0$ . We observe that  $\mu_0 = \frac{L}{16n^{4/3}} \frac{(1+\theta)^{m-1}}{\theta}$  where  $\theta = 2\eta^2 L^2 + \eta\beta$ . This  
494 is obtained using the relation  $\mu_t = \mu_{t+1}(1 + \eta\beta + 2\eta^2 L^2) + \eta^2 L^3$  and the fact that  $\mu_m = 0$ . Using  
495 the specified values of  $\beta$  and  $\eta$  we have

$$\theta = 2\eta^2 L^2 + \eta\beta = \frac{1}{8n^{4/3}} + \frac{1}{4n} \leq \frac{3}{4n}.$$

496 Using the above bound on  $\theta$ , we get

$$\begin{aligned}\mu_0 &= \frac{L}{16n^{4/3}} \frac{(1+\theta)^m - 1}{\theta} = \frac{L((1+\theta)^m - 1)}{2(1+2n^{1/3})} \\ &\leq \frac{L((1+\frac{3}{4n})^{\lfloor 4n/3 \rfloor} - 1)}{2(1+2n^{1/3})} \leq n^{-1/3}(L(e-1)/4),\end{aligned}\quad (14)$$

497 wherein the second inequality follows upon noting that  $(1+\frac{1}{l})^l$  is increasing for  $l > 0$  and  $\lim_{l \rightarrow \infty} (1+\frac{1}{l})^l = e$  (here  $e$  is the Euler's number). Now we can lower bound  $v_n$ , as

$$v_n = \min_t (\eta - \frac{\mu_{t+1}\eta}{\beta} - \eta^2 L - 2\mu_{t+1}\eta^2) \geq (\eta - \frac{\mu_0\eta}{\beta} - \eta^2 L - 2\mu_0\eta^2) \geq \frac{1}{40Ln^{2/3}}.$$

499 The first inequality holds since  $\mu_t$  decreases with  $t$ . The second inequality holds since (a)  $\mu_0/\beta$  can be  
500 upper bounded by  $(e-1)/4$  (follows from Equation (14)), (b)  $\eta^2 L \leq \eta/4$  and (c)  $2\mu_0\eta^2 \leq (e-1)\eta/8$   
501 (follows from Equation (14)). Substituting the above lower bound in Equation (13), we obtain the  
502 following:

$$\frac{1}{T_g} \sum_{s=0}^{S-1} \sum_{t=0}^{m-1} \mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] \leq \frac{40Ln^{2/3}\mathbb{E}[f(x^0) - f(x_m^S)]}{T_g}.\quad (15)$$

503 From the definition of  $(y, z)$  in output of Algorithm 2 i.e.,  $y$  is Iterate  $x_a$  chosen uniformly random  
504 from  $\{x_t^{s+1}\}_{t=0}^{m-1}\}_{s=0}^{S-1}$  and  $z = x_m^S$ , it is clear that Algorithm 2 satisfies the **G.2** requirement  
505 of GRADIENT-FOCUSED-OPTIMIZER with  $g(n, \epsilon) = T_g/40Ln^{2/3}$ . Since both **G.1** and **G.2** are  
506 satisfied for Algorithm 2, we conclude that SVRG is a GRADIENT-FOCUSED-OPTIMIZER.  $\square$

## 507 C Proof of Lemma 2

*Proof.* The first important observation is that the function value never increases because  $y = \arg \min_{z \in \{u, x\}} f(z)$  i.e.,  $f(y) \leq f(x)$ , thus satisfying **H.1** of HESSIAN-FOCUSED-OPTIMIZER. We now analyze the scenario where  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$ . Consider the event where we obtain  $v$  such that

$$\langle v, \nabla^2 f(x)v \rangle \leq \lambda_{\min}(\nabla^2 f(x)) + \frac{\gamma}{2}.$$

508 This event (denoted by  $\mathcal{E}$ ) happens with at least probability  $\rho$ . Note that, since  $\lambda_{\min}(\nabla^2 f(x)) \leq -\gamma$ ,  
509 we have  $\langle v, \nabla^2 f(x)v \rangle \leq -\frac{\gamma}{2}$ . In this case, we have the following relationship:

$$\begin{aligned}f(y) &\leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x) + \frac{M}{6}\|y - x\|^3 \\ &= f(x) - \alpha|\langle \nabla f(x), v \rangle| + \frac{\alpha^2}{2}v^T \nabla^2 f(x)v + \frac{M\alpha^3}{6}\|v\|^3 \\ &\leq f(x) + \frac{\alpha^2}{2}v^T \nabla^2 f(x)v + \frac{M\alpha^3}{6} \\ &\leq f(x) - \frac{1}{2M^2}|v^T \nabla^2 f(x)v|^3 + \frac{1}{6M^2}|v^T \nabla^2 f(x)v|^3 \\ &= f(x) - \frac{1}{3M^2}|v^T \nabla^2 f(x)v|^3 \leq f(x) - \frac{1}{24M^2}\gamma^3.\end{aligned}\quad (16)$$

The first inequality follows from the  $M$ -lipschitz continuity of the Hessian  $\nabla^2 f(x)$ . The first equality follows from the update rule of HESSIANDESCENT. The second inequality is obtained by dropping the negative term and using the fact that  $\|v\| = 1$ . The second equality is obtained by substituting  $\alpha = \frac{|v^T \nabla^2 f(x)v|}{M}$ . The last inequality is due to the fact that  $\langle v, \nabla^2 f(x)v \rangle \leq -\frac{\gamma}{2}$ . In the other scenario where

$$\langle v, \nabla^2 f(x)v \rangle \leq \lambda_{\min}(\nabla^2 f(x)) + \frac{\gamma}{2},$$

510 we can at least ensure that  $f(y) \leq f(x)$  since  $y = \arg \min_{z \in \{u, x\}} f(z)$ . Therefore, we have

$$\begin{aligned}\mathbb{E}[f(y)] &= \rho\mathbb{E}[f(y)|\mathcal{E}] + (1-\rho)\mathbb{E}[f(y)|\mathcal{E}^c] \\ &\leq \rho\mathbb{E}[f(y)|\mathcal{E}] + (1-\rho)f(x) \\ &\leq \rho[f(x) - \frac{\rho}{24M^2}\gamma^3] + (1-\rho)f(x) \\ &= f(x) - \frac{\rho}{24M^2}\gamma^3.\end{aligned}\quad (17)$$

511 The last inequality is due to Equation (16). Hence, HESSIAN-FOCUSED-OPTIMIZER satisfies **H.2** of  
 512 HESSIAN-FOCUSED-OPTIMIZER with  $h(n, \epsilon, \gamma) = \frac{\rho}{24M^2}\gamma^3$ , thus concluding the proof.  $\square$

### 513 **D Proof of Theorem 3**

514 First note that cubic method is a descent method (refer to Theorem 1 of [28]); thus, **H.1** is trivially sat-  
 515 isfied. Furthermore, cubic descent is a HESSIAN-FOCUSED-OPTIMIZER with  $h(n, \epsilon, \gamma) = \frac{2\gamma^3}{81M^3}\gamma^3$ .  
 516 This, again, follows from Theorem 1 of [28]. The result easily follows from the aforementioned  
 517 observations.

### 518 **E Other Lemmas**

519 The following bound on the variance of SVRG is useful for our proof [32].

520 **Lemma 3.** [32] Let  $v_t^{s+1}$  be computed by Algorithm 2. Then,

$$\mathbb{E}[\|v_t^{s+1}\|^2] \leq 2\mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] + 2L^2\mathbb{E}[\|x_t^{s+1} - \tilde{x}^s\|^2].$$

521 *Proof.* We use the definition of  $v_t^{s+1}$  to get

$$\begin{aligned} \mathbb{E}[\|v_t^{s+1}\|^2] &= \mathbb{E}[\|(\nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s)) + \nabla f(\tilde{x}^s)\|^2] \\ &= \mathbb{E}[\|(\nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s)) + \nabla f(\tilde{x}^s) - \nabla f(x_t^{s+1}) + \nabla f(x_t^{s+1})\|^2] \\ &\leq 2\mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] + 2\mathbb{E}\left[\|\nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s) - \mathbb{E}[\nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s)]\|^2\right] \end{aligned}$$

522 The inequality follows from the simple fact that  $(a + b)^2 \leq a^2 + b^2$ . From the above inequality, we  
 523 get the following:

$$\begin{aligned} \mathbb{E}[\|v_t^{s+1}\|^2] &\leq 2\mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] + 2\mathbb{E}\|\nabla f_{i_t}(x_t^{s+1}) - \nabla f_{i_t}(\tilde{x}^s)\|^2 \\ &\leq 2\mathbb{E}[\|\nabla f(x_t^{s+1})\|^2] + 2L^2\mathbb{E}[\|x_t^{s+1} - \tilde{x}^s\|^2] \end{aligned}$$

524 The first inequality follows by noting that for a random variable  $\zeta$ ,  $\mathbb{E}[\|\zeta - \mathbb{E}[\zeta]\|^2] \leq \mathbb{E}[\|\zeta\|^2]$ . The  
 525 last inequality follows from  $L$ -smoothness of  $f_{i_t}$ .  $\square$

### 526 **F Approximate Cubic Regularization**

527 Cubic regularization method of [19] is designed to operate on full batch, i.e., it does not exploit the  
 528 finite-sum structure of the problem and requires the computation of the gradient and the Hessian on  
 529 the entire dataset to make an update. However, such full-batch methods do not scale gracefully with  
 530 the size of data and become prohibitively expensive on large datasets. To overcome this challenge,  
 531 we devised an approximate cubic regularization method described below:

532 1. Pick a mini-batch  $\mathcal{B}$  and obtain the gradient and the hessian based on  $\mathcal{B}$ , i.e.,

$$g = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla f_i(x) \quad H = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla^2 f_i(x) \quad (18)$$

533 2. Solve the sub-problem

$$v^* = \arg \min_v \langle g, v \rangle + \frac{1}{2} \langle v, H v \rangle + \frac{M}{6} \|v\|^3 \quad (19)$$

534 3. Update:  $x \leftarrow x + v^*$

535 We found that this mini-batch training strategy, which requires the computation of the gradient and  
 536 the Hessian on a small subset of the dataset, to work well on a few datasets (CURVES, MNIST,  
 537 CIFAR10). A similar method has been analysed in [7].

538 Furthermore, in many deep-networks, adaptive per-parameter learning rate helps immensely [18].  
 539 One possible explanation for this is that the scale of the gradients in each layer of the network



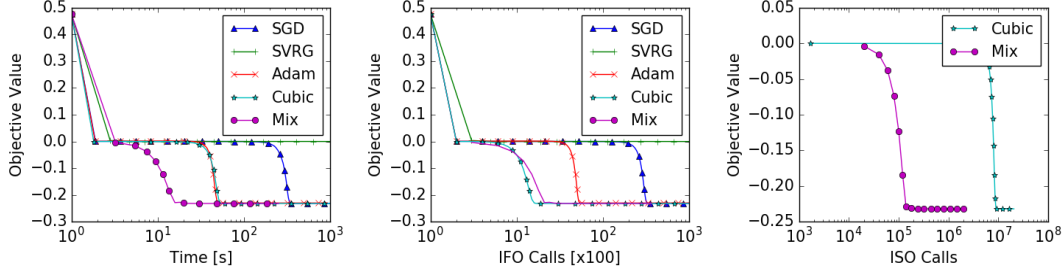


Figure 4: Comparison of various methods on a synthetic problem. Our mix framework successfully escapes saddle point.

540 often differ by several orders of magnitude. A well-suited optimization method should take this into  
 541 account. This is the reason for popularity of methods like ADAM or RMSPROP in the deep learning  
 542 community. On similar lines, to account for different per-parameter behaviour in cubic regularization,  
 543 we modify the sub-problem by adding a diagonal matrix  $M_d$  in addition to the scalar regularization  
 544 coefficient  $M$ , i.e.,

$$\min_v \langle g, v \rangle + \frac{1}{2} \langle v, H v \rangle + \frac{1}{6} M \|M_d v\|^3. \quad (20)$$

545 Also we devised an adaptive rule to obtain the diagonal matrix as  $M_d = \text{diag}((s + 10^{-12})^{1/9})$ , where  
 546  $s$  is maintained as a moving average of third order polynomial of the mini-batch gradient  $g$ , in a  
 547 fashion similar to RMSPROP and ADAM:

$$s \leftarrow \beta s + (1 - \beta)(|g|^3 + 2g^2), \quad (21)$$

548 where  $|g|^3$  and  $g^2$  are vectors such that  $[|g|^3]_i = |g_i|^3$  and  $[g^2]_i = g_i^2$  respectively for all  $i \in [n]$ . The  
 549 experiments reported on CURVES and MNIST in this paper utilizes both the above modifications  
 550 to the cubic regularization, with  $\beta$  set to 0.9. We refer to this modified procedure as ACubic in our  
 551 results.

## 552 G Experiment Details

553 In this section we provide further experimental details and results to aid reproducibility.

### 554 G.1 Synthetic Problem

555 The parameter selection for all the methods were carried as follows:

- 556 1. SGD: The scalar step-size was determined by a grid search.
- 557 2. ADAM: We performed a grid search over  $\alpha$  and  $\varepsilon$  parameters of ADAM tied together, i.e.,  $\alpha = \varepsilon$ .
- 558 3. SVRG: The scalar step-size was determined by a grid search.
- 559 4. CUBICDESCENT: The regularization parameter  $M$  was chosen by grid search. The sub-problem  
 560 was solved with gradient descent [6] with the step-size of solver to be  $10^{-2}$  and run till the gradient  
 561 norm of the sub-problem is reduced below  $10^{-3}$ .

562 **Further Observations** The results are presented in Figure 4. The other first order methods like  
 563 ADAM with higher noise could escape relatively faster whereas SVRG with reduced noise stayed  
 564 stuck at the saddle point.

### 565 G.2 Deep Networks

566 **Methods** The parameter selection for all the methods were carried as follows::

- 567 1. ADAM: We performed a grid search over  $\alpha$  and  $\varepsilon$  parameters of ADAM so as to produce the best  
 568 generalization on a held out test set. We found it to be  $\alpha = 10^{-3}$ ,  $\varepsilon = 10^{-3}$  for CURVES and  
 569  $\alpha = 10^{-2}$ ,  $\varepsilon = 10^{-1}$  for MNIST.
- 570 2. APPROXCUBICDESCENT: The regularization parameter  $M$  was chosen as the largest value such  
 571 function value does not jump in first 10 epochs. We found it to be  $M = 10^3$  for both CURVES  
 572 and MNIST. The sub-problem was solved with gradient descent [6] with the step-size of solver to  
 573 be  $10^{-3}$  and run till the gradient norm of the sub-problem is reduced below 0.1.

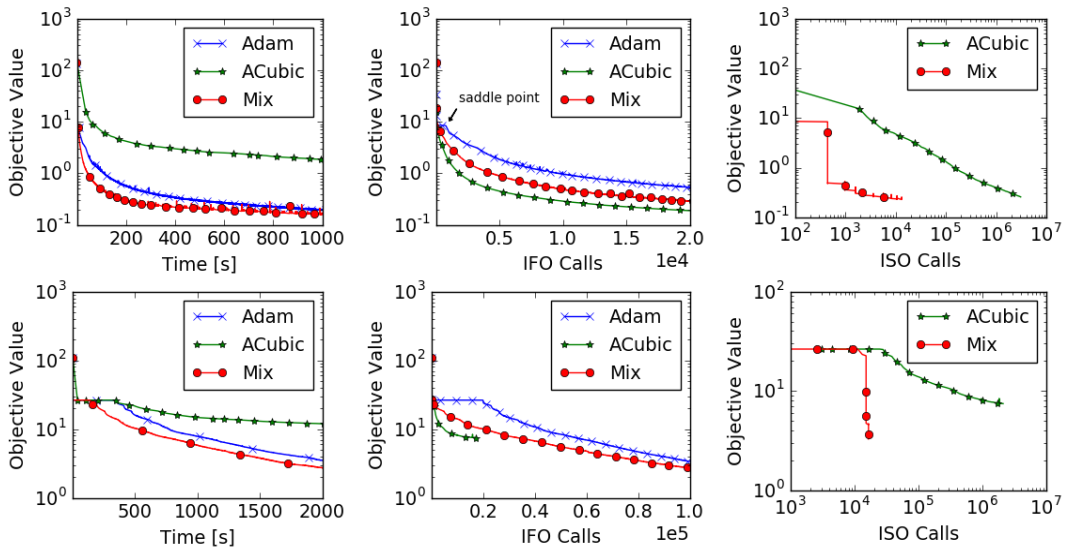


Figure 5: Comparison of various methods on a Deep Autoencoder on CURVES (top) and MNIST (bottom). Our mix approach converges faster than the baseline methods and uses relatively few ISO calls in comparison to APPROXCUBICDESCENT