
Faster Perturbed Stochastic Gradient Methods for Finding Local Minima

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Escaping from saddle points and finding local minimum is a central problem in
2 nonconvex optimization. Perturbed gradient methods are perhaps the simplest
3 approach for this problem. However, to find $(\epsilon, \sqrt{\epsilon})$ -approximate local minima, the
4 existing best stochastic gradient complexity for this type of algorithms is $\tilde{O}(\epsilon^{-3.5})$,
5 which is not optimal. In this paper, we propose Pullback, a faster perturbed
6 stochastic gradient framework for finding local minima. We show that Pullback
7 with stochastic gradient estimators such as SARAH/SPIDER and STORM can
8 find (ϵ, ϵ_H) -approximate local minima within $\tilde{O}(\epsilon^{-3} + \epsilon_H^{-6})$ stochastic gradient
9 evaluations (or $\tilde{O}(\epsilon^{-3})$ when $\epsilon_H = \sqrt{\epsilon}$). The core idea of our framework is a
10 step-size “pullback” scheme to control the average movement of the iterates, which
11 leads to faster convergence to the local minima.

12 1 Introduction

13 In this paper, we focus on the following optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) := \mathbb{E}_{\xi}[f(\mathbf{x}; \xi)], \quad (1.1)$$

14 where $f(\mathbf{x}; \xi) : \mathbb{R}^d \rightarrow \mathbb{R}$ is a stochastic function indexed by some random vector ξ , and it is
15 differentiable and possibly nonconvex. We consider the case where only the stochastic gradients
16 $\nabla f(\mathbf{x}; \xi)$ are accessible. (1.1) can unify a variety of stochastic optimization problems, such as finite-
17 sum optimization and online optimization. Since in general, finding global minima of a nonconvex
18 function could be an NP-hard problem [12], one often seeks to finding an (ϵ, ϵ_H) -approximate local
19 minimum \mathbf{x} , i.e., $\|\nabla F(\mathbf{x})\|_2 \leq \epsilon$ and $\lambda_{\min}(\nabla^2 F(\mathbf{x})) \geq -\epsilon_H$, where $\nabla F(\mathbf{x})$ is the gradient of F
20 and $\lambda_{\min}(\nabla^2 F(\mathbf{x}))$ is the smallest eigenvalue of the Hessian of F at \mathbf{x} . In many machine learning
21 applications such as matrix sensing and completion [5, 11], it suffices to find local minima due to the
22 fact that all local minima are global minima.

23 For the case where f is a deterministic function, it has been shown that vanilla gradient descent fails
24 to find local minima efficiently since the iterates will get stuck at saddle points for exponential time
25 [8]. To address this issue, the simplest idea is to add random noises as a perturbation to the stuck
26 iterates. Jin et al. [13] showed that the simple perturbation step is enough for gradient descent to
27 escape saddle points and find $(\epsilon, \sqrt{\epsilon})$ -approximate local minima within $\tilde{O}(1/\epsilon^2)$ gradient evaluations,
28 which matches the number of gradient evaluations for gradient descent to find ϵ -stationary points
29 [19]. Such matching results suggest that perturbed gradient methods can find local minima efficiently,
30 at least for deterministic optimization. When it comes to stochastic optimization, a natural question
31 arises:

32 *Can perturbed stochastic gradient methods find local minima efficiently?*

33 To answer this question, we first look into existing results of perturbed stochastic gradient methods
34 for finding local minima. Ge et al. [10] showed that perturbed Stochastic gradient descent can find

35 $(\epsilon, \sqrt{\epsilon})$ -approximate local minima within $\tilde{O}(\text{poly}(\epsilon^{-1}))$ stochastic gradient evaluations. Daneshmand
36 et al. [7] showed that under a specific CNC condition, stochastic gradient descent is able to find
37 $(\epsilon, \sqrt{\epsilon})$ -approximate local minima within $\tilde{O}(1/\epsilon^5)$ stochastic gradient evaluations. Later on, Li [17]
38 showed that simple stochastic recursive gradient descent (SSRGD) can find $(\epsilon, \sqrt{\epsilon})$ -approximate
39 local minima within $\tilde{O}(\epsilon^{-3.5})$ stochastic gradient evaluations, which is the state-of-the-arts to date.
40 However, none of these results by perturbed stochastic gradient methods matches the optimal result
41 $\tilde{O}(\epsilon^{-3})$ for finding ϵ -stationary points, achieved by SPIDER [9], SNVRG [30] and STORM [6] (See
42 also Arjevani et al. [4] for the lower bound results). Therefore, whether perturbed stochastic gradient
43 methods can find local minima as efficiently as finding stationary points still remains unknown.

44 In this work, we give an affirmative answer to the above question. We propose a general framework
45 named Pullback, which works together with existing popular stochastic gradient estimators such
46 as SARAH/SPIDER and STORM to find approximate local minima efficiently. We summarize our
47 contributions as follows:

- 48 • We prove that Pullback finds (ϵ, ϵ_H) -approximate local minima within $\tilde{O}(\epsilon^{-3} + \epsilon_H^{-6})$ stochastic
49 gradient evaluations. Specifically, in the classic setting where $\epsilon_H = \sqrt{\epsilon}$, our Pullback together
50 with the SARAH/SPIDER estimator enjoys an $\tilde{O}(\epsilon^{-3})$ stochastic gradient complexity, which
51 outperforms previous best known complexity result $\tilde{O}(\epsilon^{-3.5})$ achieved by Li [17]. Our result
52 also matches the best possible complexity result $\tilde{O}(\epsilon^{-3})$ achieved by negative curvature search
53 based algorithms [9, 31], which suggests that *simple* methods such as perturbed stochastic gradient
54 methods can find local minima as efficiently as the more complicated ones.
- 55 • Besides, we also show that Pullback with a recent proposed STORM estimator is also able to find
56 (ϵ, ϵ_H) -approximate local minima within $\tilde{O}(\epsilon^{-3} + \epsilon_H^{-6})$ stochastic gradient evaluations.
- 57 • At the core of our Pullback is a novel step-size "pullback" scheme to control the average movement
58 of the iterates, which may be of independent interest to other related nonconvex optimization
59 algorithm design.

60 To compare with, we summarized related results of stochastic first-order methods for finding local
61 minima in Table 1.

Table 1: Comparison of of different optimization algorithm for find approximate local minima of non
convex online problems.

Algorithm	Gradient complexity	Classic Setting	Neon2
Neon2+Natasha2 [1]	$\tilde{O}(\epsilon^{-3.25} + \epsilon^{-3}\epsilon_H^{-1} + \epsilon_H^{-5})$	$\tilde{O}(\epsilon^{-3.5})$	needed
Neon2+SCSG [3]	$\tilde{O}(\epsilon^{-10/3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon_H^{-5})$	$\tilde{O}(\epsilon^{-3.5})$	needed
SNVRG ⁺ +Neon2 [31]	$\tilde{O}(\epsilon^{-3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon_H^{-5})$	$\tilde{O}(\epsilon^{-3.5})$	needed
SPIDER-SFO ⁺ (+Neon2)[9]	$\tilde{O}(\epsilon^{-3} + \epsilon^{-2}\epsilon_H^{-2} + \epsilon_H^{-5})$	$\tilde{O}(\epsilon^{-3})$	needed
Perturbed SGD [10]	$\text{Poly}(d, \epsilon^{-1}, \epsilon_H^{-1})$	$\tilde{O}(\text{Poly}(\epsilon^{-1}))$	No
CNC-SGD [7]	$\tilde{O}(\epsilon^{-4} + \epsilon_H^{-10})$	$\tilde{O}(\epsilon^{-5})$	No
SSRGD [17]	$\tilde{O}(\epsilon^{-3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon^{-1}\epsilon_H^{-4})$	$\tilde{O}(\epsilon^{-3.5})$	No
Pullback (This paper)	$\tilde{O}(\epsilon^{-3} + \epsilon_H^{-6})$	$\tilde{O}(\epsilon^{-3})$	No

62 **Notations** We use lower case letters to denote scalars, lower and upper case bold letters to denote
63 vectors and matrices. We use $\|\cdot\|$ to indicate Euclidean norm. We use $\mathbb{B}_x(r)$ to denote a Euclidean
64 ball center at \mathbf{x} with radius r . We also use the standard O and Ω notations. We use $\lambda_{\min}(\mathbf{M})$ to denote
65 the minimum eigenvalue of matrix \mathbf{M} . We say $a_n = O(b_n)$ if and only if $\exists C > 0, N > 0, \forall n >$
66 $N, a_n \leq Cb_n; a_n = \Omega(b_n)$ if $a_n \geq Cb_n$. The notation \tilde{O} is used to hide logarithmic factors.

67 2 Related Work

68 In this section, we review some important related works.

69 **Variance reduction methods for finding stationary points.** Our algorithm takes stochastic gradient
70 estimators as its subroutine. In specific, Johnson and Zhang [14], Xiao and Zhang [28] proposed
71 Stochastic Variance Reduced Gradient (SVRG) for convex optimization in the finite-sum setting.
72 Reddi et al. [25], Allen-Zhu and Hazan [2] analyzed SVRG for nonconvex optimization. Lei et al. [16]

73 proposed a new variance reduction algorithm, dubbed stochastically controlled stochastic gradient
74 (SCSG) algorithm, which finds a ϵ -stationary point within $O(\epsilon^{-10/3})$ stochastic gradient evaluations.
75 Nguyen et al. [21] proposed a SARAH algorithm which uses a recursive gradient estimator for convex
76 optimization, and it was extended to nonconvex optimization in [22]. Fang et al. [9] proposed a
77 SPIDER algorithm with a recursive gradient estimator and proved an $O(\epsilon^{-3})$ stochastic gradient
78 evaluations to find a ϵ -stationary point, which matches a corresponding lower bound. Concurrently,
79 Zhou et al. [30] proposed an SNVRG algorithm with a nested gradient estimator and proved an $\tilde{O}(\epsilon^{-3})$
80 stochastic gradient evaluations to find a ϵ -stationary point. Wang et al. [27] proposed a Spiderboost
81 algorithm with a constant step size, achieves the same $O(\epsilon^{-3})$ gradient complexity. Pham et al.
82 [23] extended SARAH [22] to proximal optimization and proved $O(\epsilon^{-3})$ gradient complexity for
83 finding stationary points. Recently, Cutkosky and Orabona [6] proposed a recursive momentum-based
84 algorithm called STORM and proved an $\tilde{O}(\epsilon^{-3})$ gradient complexity to find ϵ -stationary points.
85 Tran-Dinh et al. [26] proposed a SARAH-SGD algorithm which hybrids both SGD and SARAH
86 algorithm with an $\tilde{O}(\epsilon^{-3})$ gradient complexity when ϵ is small. Li et al. [18] proposed a PAGE
87 algorithm with probabilistic gradient estimator which also attains an $\tilde{O}(\epsilon^{-3})$ gradient complexity.
88 In our work, we employ SARAH/SPIDER and STORM as the gradient estimator in our Pullback
89 framework since they are most representative and simple to use.

90 **Utilizing negative curvature descent to escape from saddle points.** To escape saddle points, a
91 widely used approach is to first compute the direction of the negative curvature of the saddle point and
92 move away along that direction. In stochastic optimization, to find (ϵ, ϵ_H) -approximate local minima,
93 [1] proposed a Natasha algorithm using Hessian-vector product to compute the negative curvature
94 direction with the total computation cost of $\tilde{O}(\epsilon^{-3.25} + \epsilon^{-3}\epsilon_H^{-1} + \epsilon_H^{-5})$. Later, Xu et al. [29] proposed
95 a Neon method which computes the negative curvature direction with perturbed stochastic gradients,
96 whose total computational cost is $\tilde{O}(\epsilon^{-10/3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon_H^{-6})$. [3] proposed a Neon2 negative
97 curvature computation subroutine with $\tilde{O}(\epsilon^{-10/3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon_H^{-5})$ stochastic gradient evaluations.
98 Fang et al. [9] then showed that SPIDER equipped with Neon2 can find (ϵ, ϵ_H) -approximate local
99 minima within $\tilde{O}(\epsilon^{-3} + \epsilon^{-2}\epsilon_H^{-2} + \epsilon_H^{-5})$ stochastic gradient evaluations, while independently Zhou
100 et al. [31] proved that SNVRG equipped with Neon2 can find (ϵ, ϵ_H) -approximate local minima
101 within $\tilde{O}(\epsilon^{-3} + \epsilon^{-2}\epsilon_H^{-3} + \epsilon_H^{-5})$ stochastic gradient evaluations. In contrast to this line of works, our
102 algorithm is simpler since it does not need to use the negative curvature search routine.

103 3 Preliminaries

104 In this section, we present assumptions and definitions that will be used throughout our analysis.

105 We first introduce the standard smoothness and Hessian Lipschitz assumptions.

106 **Assumption 3.1.** For all ξ , $f(\cdot; \xi)$ is L -smooth and its Hessian is ρ -Lipschitz continuous w.r.t. \mathbf{x} ,
107 i.e., for any $\mathbf{x}_1, \mathbf{x}_2$, we have that

$$\|\nabla f(\mathbf{x}_1; \xi) - \nabla f(\mathbf{x}_2; \xi)\|_2 \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|_2, \quad \|\nabla^2 f(\mathbf{x}_1; \xi) - \nabla^2 f(\mathbf{x}_2; \xi)\|_2 \leq \rho\|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

108 This assumption directly implies that the expected objective function $F(\mathbf{x})$ is also L -smooth and its
109 Hessian is ρ -Lipschitz continuous. This assumption is standard for finding approximate local minima
110 in all the results presented in Table 1.

111 **Assumption 3.2.** The squared difference between the stochastic gradient and full gradient is bounded
112 by $\sigma^2 < \infty$, i.e., for any $\mathbf{x}, \xi \in \mathbb{R}^d$, $\|\nabla f(\mathbf{x}; \xi) - \nabla F(\mathbf{x})\|_2^2 \leq \sigma^2$.

113 Assumption 3.2 is standard in online/stochastic optimization for finding second-order stationary
114 points [9, 17], and immediately implies that the variance of the stochastic gradient is bounded by σ^2 .
115 It can be relaxed to be $\|\nabla f(\mathbf{x}; \xi) - \nabla F(\mathbf{x})\|_2$ has a σ -Sub-Gaussian tail.

116 Let $\mathbf{x}_0 \in \mathbb{R}^d$ be the starting point of the algorithm. We assume the gap between the initial function
117 value and the optimal value is bounded.

118 **Assumption 3.3.** We have $\Delta = F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x}) < +\infty$.

119 Next, we give the formal definition of approximate local minima (a.k.a., second-order stationary
120 points).

121 **Definition 3.4.** We call $\mathbf{x} \in \mathbb{R}^d$ an (ϵ, ϵ_H) -approximate local minimum, if

$$\|\nabla F(\mathbf{x})\|_2 \leq \epsilon, \lambda_{\min}(\nabla^2 F(\mathbf{x})) \geq -\epsilon_H.$$

122 The definition of (ϵ, ϵ_H) -approximate local minima is a generalization of the classical $(\epsilon, \sqrt{\epsilon})$ -
 123 approximate local minima studied by Nesterov and Polyak [20], Jin et al. [13].

124 4 The Pullback Framework

125 In this section, we present our main algorithm Pullback. We begin with reviewing the mechanism
 126 of perturbed gradient descent in deterministic optimization, and then we discuss the main difficulty of
 127 extending it to the stochastic optimization case. Finally, we show how we overcome such a difficulty
 128 by presenting our Pullback framework.

129 **How does perturbed gradient descent escape from saddle points?** We review the perturbed
 130 gradient descent [13] (PGD for short) with its proof roadmap, which shows how PGD finds $(\epsilon, \sqrt{\epsilon})$ -
 131 approximate local minima efficiently. In general, the whole process of perturbed gradient descent
 132 can be decomposed into several epochs, and each epoch consists of two non-overlapping phases: the
 133 *gradient descent phase* (GD phase for short) and the *Escape from saddle point phase* (Escape phase
 134 for short). In each epoch, PGD starts with the GD phase by default. In the GD phase, PGD performs
 135 vanilla gradient descent to update its iterate, until at some iterate $\tilde{\mathbf{x}}$, the norm of the gradient
 136 $\|\nabla F(\tilde{\mathbf{x}})\|_2$ is less than the target accuracy $\tilde{O}(\epsilon)$. Then PGD switches to the Escape phase. In the
 137 Escape phase, PGD first adds a uniform random noise (or Gaussian noise) to the current iterate $\tilde{\mathbf{x}}$,
 138 then it runs $\ell_{\text{thres}} = \tilde{O}(\epsilon^{-1/2})$ steps of vanilla gradient descent. PGD then compares the function
 139 value gap between the current iterate and the beginning iterate of Escape phase $\tilde{\mathbf{x}}$. If the gap is less
 140 than a threshold $\mathcal{F} = \tilde{O}(\epsilon^{1.5})$, then PGD outputs $\tilde{\mathbf{x}}$ as the targeted local minimum. Otherwise, PGD
 141 starts a new epoch and performs gradient descent again.

142 To see why PGD can find $(\epsilon, \sqrt{\epsilon})$ -approximate local minima within $\tilde{O}(\epsilon^{-2})$ gradient evaluations, we
 143 do the following calculation. First, when PGD is in the GD phase, the function value decreases $\tilde{O}(\epsilon^2)$
 144 per step (following the standard gradient descent analysis). When PGD is in the Escape phase, the
 145 function value decreases $\mathcal{F}/\ell_{\text{thres}} = \tilde{O}(\epsilon^2)$ per step on average. Therefore, the total number of steps
 146 will be bounded by $\tilde{O}(\epsilon^{-2})$, which is in the same order as GD for finding ϵ -stationary points.

147 **Limitation of existing methods.** However, extending the two-phase PGD algorithm from determinis-
 148 tic optimization to stochastic optimization with a competitive gradient complexity is very challenging.
 149 We take the SSRGD algorithm proposed by Li [17] as an example, which uses SARAH/SPIDER [9] as
 150 its gradient estimator. Unlike deterministic optimization where we can access the exact function value
 151 $F(\mathbf{x})$ and gradient $\nabla F(\mathbf{x})$ defined in (1.1), in the stochastic optimization case we can only access
 152 the stochastic function $f(\mathbf{x}; \xi)$ and the stochastic gradient $\nabla f(\mathbf{x}; \xi)$. Therefore, in order to estimate
 153 the gradient norm $\|\nabla F(\mathbf{x})\|_2$ (which is required at the beginning of Escape phase), a naive approach
 154 (adapted by Li [17]) is to sample a big batch of stochastic gradient $\nabla f(\mathbf{x}; \xi_1), \dots, \nabla f(\mathbf{x}; \xi_B)$ and
 155 uses their mean to approximate $\nabla F(\mathbf{x})$. Standard concentration analysis suggests that in order to
 156 achieve an ϵ -accuracy, the batch size B should be in the order $\tilde{O}(\epsilon^{-2})$. Thus, each Escape phase
 157 leads to a $\mathcal{F} = \tilde{O}(\epsilon^{1.5})$ function value decrease with at least $\tilde{O}(\epsilon^{-2})$ number of stochastic gradient
 158 evaluations, which contributes $\tilde{O}(1/\epsilon^{1.5} \cdot \epsilon^{-2}) = \tilde{O}(\epsilon^{-3.5})$ gradient complexity in the end. This is
 159 already worse than the $O(\epsilon^{-3})$ gradient complexity of SPIDER for finding ϵ -stationary points.

160 **Our approach.** Here we propose our Pullback framework in Algorithm 1, which overcomes the
 161 aforementioned limitation. In detail, Pullback inherits the two-phase structure of PGD and SSRGD,
 162 and it takes either SARAH/SPIDER or STORM [6] as its gradient estimator. The two gradient
 163 estimators are summarized as subroutines GradEst-SPIDER and GradEst-STORM in Algorithms 2
 164 and 3 respectively, and we use \mathbf{d}_t to denote their estimated gradient at iterate \mathbf{x}_t . The key improvement
 165 of Pullback is that, it directly takes the output of the gradient estimator GradEst to estimate the
 166 true gradient $\nabla F(\mathbf{x})$, which avoids sampling a big batch of stochastic gradients as in Li [17] and
 167 thus saves the total gradient complexity. A similar strategy has also been adapted in [9], but for the
 168 negative curvature search subroutine. However, such a strategy leads to a new problem to be solved.

169 Since we use \mathbf{d}_t to directly estimate $\nabla F(\mathbf{x}_t)$, in order to make such an estimation valid, we need to
 170 guarantee that the error between \mathbf{d}_t and $\nabla F(\mathbf{x}_t)$ is small enough, e.g., up to $O(\epsilon)$ accuracy. Notice

Algorithm 1 Pullback

Input: Initial point \mathbf{x}_1 , step size η and η_H , perturbation radius r , threshold parameter ℓ_{thres} , average movement \bar{D} .

```
1:  $\mathbf{d}_1 \leftarrow \text{GradEst}(0, \mathbf{0}, \mathbf{0}, \mathbf{x}_1)$ ,  $s \leftarrow 0$ ,  $t \leftarrow 1$ , FIND  $\leftarrow$  false
2: while FIND = false do
3:    $s \leftarrow s + 1$ ,  $t_s \leftarrow t$ , FIND  $\leftarrow$  true
4:   while  $\|\mathbf{d}_t\|_2 > \epsilon$  do
5:      $\eta_t \leftarrow \eta / \|\mathbf{d}_t\|_2$ , {"PullBack"}
6:      $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_t \mathbf{d}_t$ ,  $\mathbf{d}_{t+1} \leftarrow \text{GradEst}(t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1})$ ,  $t \leftarrow t + 1$ 
7:   end while
8:    $m_s \leftarrow t$ ,  $\boldsymbol{\xi} \sim \text{Uniform } B_0(r)$ ,  $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \boldsymbol{\xi}$ ,  $\mathbf{d}_{t+1} \leftarrow \text{GradEst}(t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1})$ ,  $t \leftarrow t + 1$ 
9:   for  $k = 0, \dots, \ell_{\text{thres}} - 1$  do
10:     $\eta_t \leftarrow \eta_H$ ,  $D \leftarrow \sum_{i=m_s}^t \eta_i^2 \|\mathbf{d}_i\|_2^2$ 
11:    if  $D > (t - m_s + 1)\bar{D}$  then
12:      Set  $\eta_t$  such that  $\sum_{i=m_s}^t \eta_i^2 \|\mathbf{d}_i\|_2^2 = (t - m_s + 1)\bar{D}$  {"PullBack"}
13:       $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_t \mathbf{d}_t$ ,  $\mathbf{d}_{t+1} \leftarrow \text{GradEst}(t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1})$ ,  $t \leftarrow t + 1$ , FIND  $\leftarrow$  false, break
14:    end if
15:     $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_t \mathbf{d}_t$ ,  $\mathbf{d}_{t+1} \leftarrow \text{GradEst}(t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1})$ ,  $t \leftarrow t + 1$ 
16:  end for
17: end while
Output:  $\mathbf{x}_{m_s}$ 
```

Algorithm 2 GradEst-SPIDER($t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1}, b, q, B$)

Input: Big batch size B , mini-batch size b , loop length q

```
1: if  $t \bmod q = 0$  then
2:   Generate  $\boldsymbol{\xi}_{t+1}^1, \dots, \boldsymbol{\xi}_{t+1}^B$ . Set  $\mathbf{d}_{t+1} \leftarrow \sum_{i=1}^B \nabla f(\mathbf{x}_{t+1}; \boldsymbol{\xi}_{t+1}^i) / B$ 
3: else
4:   Generate  $\boldsymbol{\xi}_{t+1}^1, \dots, \boldsymbol{\xi}_{t+1}^b$ . Set  $\mathbf{d}_{t+1} \leftarrow \mathbf{d}_t + \sum_{i=1}^b [\nabla f(\mathbf{x}_{t+1}; \boldsymbol{\xi}_{t+1}^i) - \nabla f(\mathbf{x}_t; \boldsymbol{\xi}_{t+1}^i)] / b$ 
5: end if
Output:  $\mathbf{d}_{t+1}$ 
```

171 that the recursive structure of SARAH/SPIDER and STORM suggests the following error bound:

$$\forall t, \|\mathbf{d}_t - \nabla F(\mathbf{x}_t)\|_2^2 = \tilde{O}\left(\sum_{i=s_t}^{t-1} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2\right), \quad (4.1)$$

172 where s_t is some reference index only related to t . Therefore, in order to make the error
173 $\|\mathbf{d}_t - \nabla F(\mathbf{x}_t)\|_2$ small, it suffices to make the movement of the iterates $\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2$ small either
174 individually or on average. We achieve this goal by our proposed step-size ‘‘Pullback’’ scheme. In de-
175 tail, in the GD phase, when the norm of the estimated gradient $\|\mathbf{d}_t\|_2$ is large, we pull the step-size η_t
176 back to a smaller value via normalization, which forces the movement $\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2 = \eta_t \|\mathbf{d}_t\|_2 = \eta$
177 to be small. Such an approach is also adapted by Fang et al. [9] as a normalized gradient de-
178 scent for finding first-order stationary points. In the Escape phase, which starts at m_s -th step, we
179 record the accumulative squared movement starting from \mathbf{x}_{m_s+1} (after the perturbation step) as
180 $D := \sum_{i=m_s+1}^t \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 = \sum_{i=m_s+1}^t \eta_i^2 \|\mathbf{d}_i\|_2^2$. When the average movement $D/(t - m_s + 1)$
181 is large, we pull the *last* step size η_t back to a smaller value, which forces the average movement
182 $D/(t - m_s + 1)$ to be small. Fortunately, such a simple step-size calibration scheme allows us to
183 well-control the error between \mathbf{d}_t and $\nabla F(\mathbf{x}_t)$, and to reduce the gradient complexity.

184 5 Main Results

185 In this section, we present the main theoretical results. We first present the convergence guarantee of
186 Pullback-SPIDER, which uses GradEst-SPIDER to estimate the gradient \mathbf{d}_t in Algorithm 1.

187 **Theorem 5.1.** Under Assumptions 3.1, 3.2 and 3.3, choose batch size $B = \tilde{O}(\sigma^2 \epsilon^{-2} +$
188 $\sigma^2 \rho^2 \epsilon_H^{-4})$, $b = q = \sqrt{B}$, set step size $\eta = \sigma / (2\sqrt{BL})$, $\eta_H = \tilde{O}(L^{-1})$, perturbation ra-
189 dius $r \leq \min\{\sigma / (2\sqrt{BL}), \log(4/\delta)\eta_H \sigma^2 / (2B\epsilon), \sqrt{2\log(4/\delta)\eta_H \sigma^2 / (BL)}\}$, threshold $\ell_{\text{thres}} =$

Algorithm 3 GradEst-STORM($t, \mathbf{d}_t, \mathbf{x}_t, \mathbf{x}_{t+1}, a, b, B$)

Input: Initial batch size B , mini batch size b and weight parameter a .

- 1: **if** $t = 0$ **then**
- 2: Generate $\xi_{t+1}^1, \dots, \xi_{t+1}^B$. Set $\mathbf{d}_{t+1} \leftarrow \sum_{i=1}^B \nabla f(\mathbf{x}_{t+1}; \xi_{t+1}^i)/B$
- 3: **else**
- 4: Generate $\xi_{t+1}^1, \dots, \xi_{t+1}^b$
- 5: Set $\mathbf{d}_{t+1} \leftarrow (1-a)[\mathbf{d}_t - \sum_{i=1}^b \nabla f(\mathbf{x}_t; \xi_{t+1}^i)/b] + \sum_{i=1}^b \nabla f(\mathbf{x}_{t+1}; \xi_{t+1}^i)/b$
- 6: **end if**

Output: \mathbf{d}_{t+1}

190 $\tilde{O}(1/(\eta_H \epsilon_H))$ and $\bar{D} = \sigma^2/(4BL^2)$. Then with high probability, Pullback-SPIDER can find
191 (ϵ, ϵ_H) -approximate local minima within $\tilde{O}(\sigma L \Delta \epsilon^{-3} + \sigma \rho^3 L \Delta \epsilon_H^{-6})$ stochastic gradient evaluations.

192 **Remark 5.2.** In the classical setting $\epsilon = \sqrt{\epsilon_H}$, our result gives $\tilde{O}(\epsilon^{-3})$ gradient complexity, which
193 outperforms the best existing result $\tilde{O}(\epsilon^{-3.5})$ for perturbed stochastic gradient methods achieved
194 by SSRGD [17]. For sufficiently small ϵ , Arjevani et al. [4] proved the lower bound of gradient
195 complexity $\Omega(\epsilon^{-3} + \epsilon_H^{-5})$ for any first-order stochastic methods to find (ϵ, ϵ_H) -approximate local
196 minima. Our results matches the lower bound $\tilde{O}(\epsilon^{-3})$ when $\epsilon_H \leq \epsilon^{3/5}$. For the general case, there is
197 still a gap in the dependency of ϵ_H between our result and the lower bound, and we leave to close it
198 as future work.

199 Next, we present the convergence guarantee of Pullback-STORM, which uses GradEst-STORM to
200 estimate the gradient \mathbf{d}_t in Algorithm 1.

201 **Theorem 5.3.** Under Assumptions 3.1, 3.2 and 3.3, choose the mini batch size $b = \tilde{O}(\sigma \epsilon^{-1} +$
202 $\sigma \rho \epsilon_H^{-2})$, and initial batch size $B = b^2$, set step size $\eta = \sigma/(2bL)$, $\eta_H = \tilde{O}(L^{-1})$,
203 weight $a = 56^2 \log(4/\delta)/b$, threshold $\ell_{\text{thres}} = \tilde{O}(1/(\eta_H \epsilon_H))$, perturbation radius $r \leq$
204 $\min\{\sigma/(2bL), \log(4/\delta)^2 \eta_H \sigma^2/(\epsilon b^2), \sqrt{2 \log(4/\delta)^2 \eta_H \sigma^2/(b^2 L)}\}$, and $\bar{D} = \sigma^2/(4b^2 L^2)$. Then
205 with high probability, Pullback-STORM can find (ϵ, ϵ_H) -approximate local minima within
206 $\tilde{O}(\sigma L \Delta \epsilon^{-3} + \sigma \rho^3 L \Delta \epsilon_H^{-6})$ stochastic gradient evaluations.

207 **Remark 5.4.** Different from Pullback-SPIDER, the estimation error $\|\mathbf{d}_t - \nabla F(\mathbf{x}_t)\|_2$ of
208 Pullback-STORM is controlled by the weight parameter a . This allows us to come up with a
209 simpler single-loop algorithm instead of a double-loop algorithm.

210 6 Proof Outline of the Main Results

211 Due to the page limit, we only outline the proof of Theorem 5.1 and leave the proof of Theorem 5.3
212 to the appendix.

213 Let ϵ_t denote the difference between true gradient $\nabla F(\mathbf{x}_t)$ and the estimated gradient \mathbf{d}_t , which is
214 $\epsilon_t := \mathbf{d}_t - \nabla F(\mathbf{x}_t)$. The following lemma suggests that the estimation error $\|\epsilon_t\|_2$ can be bounded.

215 **Lemma 6.1.** Under Assumptions 3.1 and 3.2, set $b = q = \sqrt{B}$, $\eta \leq \sigma/(2\sqrt{B}L)$, $r \leq \sigma/(2\sqrt{B}L)$
216 and $\bar{D} \leq \sigma^2/(4BL^2)$, then with probability at least $1 - \delta$, for all t we have

$$\|\epsilon_t\|_2 \leq \sqrt{8 \log(4/\delta)} \sigma / \sqrt{B}.$$

217 Specifically, by the choice of B in Theorem 5.1 we have that $\|\epsilon_t\|_2 \leq \epsilon/2$.

218 *Proof of Lemma 6.1.* By GradEst-SPIDER presented in Algorithm 2 we have

$$\epsilon_{t+1} = \frac{1}{B} \sum_{i=1}^B [\nabla f(\mathbf{x}_{t+1}; \xi_{t+1}^i) - \nabla F(\mathbf{x}_{t+1})], \quad t \bmod q = 0,$$

$$\epsilon_{t+1} = \epsilon_t + \frac{1}{b} \sum_{i=1}^b [\nabla f(\mathbf{x}_{t+1}; \xi_{t+1}^i) - \nabla f(\mathbf{x}_t; \xi_{t+1}^i) - \nabla F(\mathbf{x}_{t+1}) + \nabla F(\mathbf{x}_t)], \quad t \bmod q \neq 0.$$

219 By the L -smoothness in Assumption 3.1 we have

$$\|\nabla f(\mathbf{x}_{t+1}; \xi_{t+1}^i) - \nabla f(\mathbf{x}_t; \xi_{t+1}^i) - \nabla F(\mathbf{x}_{t+1}) + \nabla F(\mathbf{x}_t)\|_2 \leq 2L \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2.$$

220 Then by Assumption 3.2 and Azuma–Hoeffding inequality (See Lemma D.1 for details), with
 221 probability at least $1 - \delta$, we have

$$\forall t > 0, \|\epsilon_{t+1}\|_2^2 \leq 4 \log(4/\delta) \left(\frac{\sigma^2}{B} + \frac{4L^2}{b} \sum_{i=\lfloor t/q \rfloor q}^t \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \right). \quad (6.1)$$

222 Notice that GradEst-SPIDER is parallel with Pullback. Thus we need to further bound (6.1) by
 223 considering iterates in three different cases: (1) for step i in the GD phase, we have $\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \leq \eta^2$
 224 due to the "Pullback" scheme; (2) for $i = m_s$ for some s in the Escape phase, we have $\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \leq$
 225 r^2 ; and (3) for the other steps in Escape phase, we have on average, $\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \leq \bar{D}$. Therefore
 226 we have

$$\|\epsilon_{t+1}\|_2^2 \leq 4 \log(4/\delta) \left(\frac{\sigma^2}{B} + \frac{4L^2}{b} \cdot q \cdot \max\{\eta^2, r^2, \bar{D}\} \right) \leq \frac{8 \log(4/\delta) \sigma^2}{B}.$$

227

□

228 Lemma 6.1 guarantees that with high probability $\|\nabla F(\mathbf{x}_t)\|_2 \leq \|\mathbf{d}_t\|_2 + \epsilon$, which ensures
 229 $\|\nabla F(\mathbf{x}_{m_s})\|_2 \leq 2\epsilon$ when the algorithm terminates. Next lemma bounds the function value de-
 230 crease in the GD phase, which is also valid for Pullback-STORM.

231 **Lemma 6.2.** Suppose the event in Lemma 6.1 holds, $\eta \leq \epsilon/(2L)$, then for any s , we have

$$F(\mathbf{x}_{t_s}) - F(\mathbf{x}_{m_s}) \geq \frac{(m_s - t_s)\eta\epsilon}{8}.$$

232 The choice of η in Theorem 5.1 further implies that the loss decreases by at least $\sigma\epsilon/(16\sqrt{B}L)$ per
 233 step on average.

234 *Proof of Lemma 6.2.* For any $t_s \leq t < m_s$, we can show the following property (See Lemma D.2),

$$F(\mathbf{x}_{t+1}) \leq F(\mathbf{x}_t) - \frac{\eta_t}{2} \|\mathbf{d}_t\|_2^2 + \frac{\eta_t}{2} \|\epsilon_t\|_2^2 + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2. \quad (6.2)$$

235 Plugging the update rule $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{d}_t$ into (6.2) gives,

$$\begin{aligned} F(\mathbf{x}_{t+1}) &= F(\mathbf{x}_t) - \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2 \left(\frac{1}{2\eta_t} - \frac{L}{2} \right) + \frac{\eta_t \|\epsilon_t\|_2^2}{2} \\ &\leq F(\mathbf{x}_t) - \eta^2 \left(\frac{1}{2\eta_t} - \frac{L}{2} \right) + \frac{\eta_t \epsilon^2}{8}, \\ &\leq F(\mathbf{x}_t) - \frac{\eta\epsilon}{8} \end{aligned}$$

236 where the first inequality holds due to the fact that $\eta_t = \eta/\|\mathbf{d}_t\|_2$ and $\|\epsilon_t\|_2 \leq \epsilon/2$, and the second
 237 inequality is by $\eta_t = \eta/\|\mathbf{d}_t\|_2 \leq \eta/\epsilon \leq 1/(2L)$. □

238 Following Lemma shows that if \mathbf{x}_{m_s} is a saddle point, then with high probability, the algorithm will
 239 break during the Escape phase and set FIND←false. Thus, whenever \mathbf{x}_{m_s} is not a local minima, the
 240 algorithm cannot terminate.

241 **Lemma 6.3.** Under Assumptions 3.1 and 3.2, set perturbation radius $r \leq L\eta_H\epsilon_H/(C\rho)$, step
 242 size $\eta_H \leq \min\{1/(16L \log(\eta_H\epsilon_H\sqrt{d}LC^{-1}\rho^{-1}\delta^{-1}r^{-1})), 1/(8CL \log \ell_{\text{thres}})\} = \tilde{O}(L^{-1})$, $\ell_{\text{thres}} =$
 243 $2 \log(\eta_H\epsilon_H\sqrt{d}LC^{-1}\rho^{-1}\delta^{-1}r^{-1})/(\eta_H\epsilon_H) = \tilde{O}(\eta_H^{-1}\epsilon_H^{-1})$, and $\bar{D} < C^2L^2\eta_H^2\epsilon_H^2/(\rho^2\ell_{\text{thres}}^2)$, where
 244 $C = O(\log(d\ell_{\text{thres}}/\delta)) = \tilde{O}(1)$. We also set $b = q = \sqrt{B} \geq 16 \log(4/\delta)/(\eta_H^2\epsilon_H^2)$. Then for
 245 any s , when $\lambda_{\min}(\nabla^2 F(\mathbf{x}_{m_s})) \leq -\epsilon_H$, with probability at least $1 - 2\delta$ algorithm breaks in the
 246 Escape phase.

247 *Proof of Lemma 6.3.* Let $\{\mathbf{x}_t\}, \{\mathbf{x}'_t\}$ be two coupled sequences by running Pullback-SPIDER from
 248 $\mathbf{x}_{m_s+1}, \mathbf{x}'_{m_s+1}$ with $\mathbf{x}_{m_s+1} - \mathbf{x}'_{m_s+1} = r_0 \mathbf{e}_1$, where $\mathbf{x}_{m_s+1}, \mathbf{x}'_{m_s+1} \in \mathbb{B}_{\mathbf{x}_{m_s}}(r)$. Here $r_0 = \delta r/\sqrt{d}$
 249 and \mathbf{e}_1 denotes the smallest eigenvector direction of Hessian $\nabla^2 F(\mathbf{x}_{m_s})$.

250 When $\lambda_{\min}(\nabla^2 F(\mathbf{x}_{m_s})) \leq -\epsilon_H$, under the parameter choice in Lemma 6.3, we can show that
 251 at least one of two sequence will escape the saddle point (See Lemma D.3). To be specific, with
 252 probability at least $1 - \delta$,

$$\max_{m_s < t < m_s + \ell_{\text{thres}}} \{\|\mathbf{x}_t - \mathbf{x}_{m_s+1}\|_2, \|\mathbf{x}'_t - \mathbf{x}'_{m_s+1}\|_2\} \geq \frac{L\eta_H\epsilon_H}{C\rho}. \quad (6.3)$$

253 (6.3) suggests that for any two points $\mathbf{x}_{m_s+1}, \mathbf{x}'_{m_s+1}$ satisfying $\mathbf{x}_{m_s+1} - \mathbf{x}'_{m_s+1} = r_0\mathbf{e}_1$, at least one
 254 of them will generate a sequence of iterates which finally move more than $L\eta_H\epsilon_H/(C\rho)$. Thus, let
 255 $\mathcal{S} \subseteq \mathbb{B}_{m_s}(r)$ be the set of \mathbf{x}_{m_s+1} which will not generate a sequence of iterates moving more than
 256 $\frac{L\eta_H\epsilon_H}{C\rho}$, then in the direction \mathbf{e}_1 , the "thickness" of \mathcal{S} is smaller than r_0 . Simple integration shows that
 257 the ratio between the volume of \mathcal{S} and $\mathbb{B}_{m_s}(r)$ is bounded by δ . Therefore, since \mathbf{x}_{m_s+1} is generated
 258 from \mathbf{x}_{m_s} by adding a uniform random noise in ball $\mathbb{B}_{m_s}(r)$, we conclude that the probability for
 259 \mathbf{x}_{m_s+1} locating in \mathcal{S} is less than δ . Applying union bound, we get with probability at least $1 - 2\delta$,

$$\exists m_s < t < m_s + \ell_{\text{thres}}, \|\mathbf{x}_t - \mathbf{x}_{m_s+1}\|_2 \geq \frac{L\eta_H\epsilon_H}{C\rho}. \quad (6.4)$$

260 Denote \mathcal{E} as the event that the algorithm does not break in the Escape phase. Then under \mathcal{E} , for any
 261 $m_s < t < m_s + \ell_{\text{thres}}$, we have

$$\|\mathbf{x}_t - \mathbf{x}_{m_s+1}\|_2 \leq \sum_{i=m_s+1}^{t-1} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2 \leq \sqrt{(t - m_s) \sum_{i=m_s}^{t-1} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2} \leq (t - m_s)\sqrt{\bar{D}},$$

262 where the first inequality is due to the triangle inequality and the second inequality is due to Cauchy-
 263 Schwarz inequality. Thus, by the choice of ℓ_{thres} and \bar{D} , we have

$$\|\mathbf{x}_t - \mathbf{x}_{m_s+1}\|_2 \leq (t - m_s)\sqrt{\bar{D}} \leq \ell_{\text{thres}}\sqrt{\bar{D}} < C \cdot \frac{L\eta_H\epsilon_H}{\rho}.$$

264 Then by (6.4), we know that $\mathbb{P}(\mathcal{E}) \leq 2\delta$. Therefore when $\lambda_{\min}(\nabla^2 F(\mathbf{x}_{m_s})) \leq -\epsilon_H$, with probability
 265 at least $1 - 2\delta$, Pullback breaks in the Escape phase. \square

266 Next lemma bounds the decreasing value of the function during the Escape phase if the algorithm
 267 breaks in the Escape phase (i.e. FIND is false).

268 **Lemma 6.4** (localization). Suppose the result of Lemma 6.1 holds, and
 269 set the step size $\eta_H \leq 1/(L\sqrt{128\log(4/\delta)})$, perturbation radius $r \leq$
 270 $\min\{\log(4/\delta)\eta_H\sigma^2/(2B\epsilon), \sqrt{2\log(4/\delta)\eta_H\sigma^2/(BL)}\}$, and $\bar{D} = \sigma^2/(4BL^2)$. Suppose the
 271 algorithm breaks in the Escape phase starting at \mathbf{x}_{m_s} , then we have

$$F(\mathbf{x}_{m_s}) - F(\mathbf{x}_{t_{s+1}}) \geq (t_{s+1} - m_s) \frac{\log(4/\delta)\eta_H\sigma^2}{B}.$$

272 *Proof of Lemma 6.4.* For any $m_s < i < t_{s+1}$, we can show the following property (See Lemma D.2),

$$F(\mathbf{x}_{i+1}) \leq F(\mathbf{x}_i) - \frac{\eta_i}{2}\|\mathbf{d}_i\|_2^2 + \frac{\eta_i}{2}\|\epsilon_i\|_2^2 + \frac{L}{2}\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2. \quad (6.5)$$

273 Plugging the update rule $\mathbf{x}_{i+1} = \mathbf{x}_i - \eta_i\mathbf{d}_i$ into (6.5) gives,

$$\begin{aligned} F(\mathbf{x}_{i+1}) &\leq F(\mathbf{x}_i) + \frac{\eta_i}{2}\|\epsilon_i\|_2^2 - \left(\frac{1}{2\eta_i} - \frac{L}{2}\right)\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \\ &\leq F(\mathbf{x}_i) + \frac{\eta_H}{2} \frac{8\log(4/\delta)\sigma^2}{B} - \frac{1}{4\eta_H}\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2 \end{aligned} \quad (6.6)$$

274 where the second inequality holds due to Lemma 6.1 and $\eta_i \leq \eta_H \leq 1/(2L)$ for any $m_s < i <$
 275 t_{s+1} . Telescoping (6.6) from $i = m_s + 1$ to $t_{s+1} - 1$, we have

$$F(\mathbf{x}_{t_{s+1}}) \leq F(\mathbf{x}_{m_s+1}) + 4\eta_H \log(4/\delta)(t_{s+1} - m_s - 1) \frac{\sigma^2}{B} - \frac{1}{4\eta_H} \sum_{i=m_s+1}^{t_{s+1}-1} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2.$$

276 Finally, we have

$$\begin{aligned}
F(\mathbf{x}_{m_s+1}) - F(\mathbf{x}_{t_{s+1}}) &\geq \sum_{i=m_s+1}^{t_{s+1}-1} \frac{\|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2}{4\eta_H} - 4\log(4/\delta)(t_{s+1} - m_s - 1)\eta_H \frac{\sigma^2}{B} \\
&= (t_{s+1} - m_s - 1) \left(\frac{\bar{D}}{4\eta_H} - \frac{4\log(4/\delta)\eta_H\sigma^2}{B} \right) \\
&= (t_{s+1} - m_s - 1) \left(\frac{\sigma^2}{16\eta_H B L^2} - \frac{4\log(4/\delta)\eta_H\sigma^2}{B} \right) \\
&\geq (t_{s+1} - m_s - 1) \frac{4\log(4/\delta)\eta_H\sigma^2}{B}, \tag{6.7}
\end{aligned}$$

277 where the last inequality is by the choice of $\eta_H \leq 1/(L\sqrt{128\log(4/\delta)})$. For $i = m_s$, we have (See
278 Lemma D.2)

$$F(\mathbf{x}_{m_s+1}) \leq F(\mathbf{x}_{m_s}) + (\|\mathbf{d}_{m_s}\|_2 + \|\epsilon_{m_s}\|_2 + Lr/2)r. \tag{6.8}$$

279 Plugging $\|\mathbf{d}_{m_s}\|_2 \leq \epsilon$ and $\|\epsilon_{m_s}\|_2 \leq \epsilon/2$ into (6.8) gives,

$$F(\mathbf{x}_{m_s+1}) \leq F(\mathbf{x}_{m_s}) + (4\epsilon + Lr/2)r \leq F(\mathbf{x}_{m_s}) + \frac{2\log(4/\delta)\eta_H\sigma^2}{B}, \tag{6.9}$$

280 where the last inequality is by the choice $r \leq \min \{ \log(4/\delta)\eta_H\sigma^2/(2B\epsilon), \sqrt{2\log(4/\delta)\eta_H\sigma^2/(BL)} \}$.
281 Combining (6.7) and (6.9) and applying $t_{s+1} - m_s \geq 2$ gives,

$$F(\mathbf{x}_{m_s}) - F(\mathbf{x}_{t_{s+1}}) \geq [4(t_{s+1} - m_s - 1) - 2] \frac{\log(4/\delta)\eta_H\sigma^2}{B} \geq (t_{s+1} - m_s) \frac{\log(4/\delta)\eta_H\sigma^2}{B}.$$

282

□

283 Now, we can provide the proof of Theorem 5.1 .

284 *Proof of Theorem 5.1.* The analysis can be divided into two phases, i.e., GD phase and Escape phase.
285 The function value will decrease at different rates in different phases.

286 **GD phase:** In this phase, $\|\mathbf{d}_t\|_2 \geq \epsilon$ and $\|\epsilon_t\|_2 \leq \epsilon/2$ due to Lemma 6.1. Thus the gradients of the
287 function are large $\|\nabla F(\mathbf{x})\|_2 \geq \epsilon/2$. Lemma 6.2 further shows that the loss decreases by at least
288 $\sigma\epsilon/(16\sqrt{BL})$ on average.

289 **Escape phase:** In this phase, the starting point \mathbf{x}_{m_s} satisfies $\|\nabla F(\mathbf{x}_{m_s})\|_2 \leq \|\mathbf{d}_{m_s}\|_2 + \|\epsilon_t\|_2 \leq 2\epsilon$.
290 If \mathbf{x}_{m_s} is a saddle point with $\lambda_{\min}(\nabla^2 F(\mathbf{x}_{m_s})) \leq -\epsilon_H$, then by Lemma 6.3, with high probability
291 Pullback-SPIDER will break Escape phase, set FIND←False and begin a new GD phase. Further
292 by Lemma 6.4, the loss will decrease by at least $\log(4/\delta)\eta_H\sigma^2/B$ per step on average.

293 **Sample Complexity:** Note that the total amount for function value can decrease is at most $\Delta =$
294 $F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x}) < +\infty$. So the algorithm must end and find an (ϵ, ϵ_H) -approximate local
295 minimum within $\tilde{O}(\sqrt{BL}\Delta\sigma^{-1}\epsilon^{-1} + BL\Delta\sigma^{-2})$ iterations. Notice that on average we sample
296 $\max\{b, B/q\} = \sqrt{B}$ examples per iteration, so the total sample complexity is $\tilde{O}(BL\Delta\sigma^{-1}\epsilon^{-1} +$
297 $B^{3/2}L\Delta\sigma^{-2})$. Plugging in the choice of $B = \tilde{O}(\sigma^2\epsilon^{-2} + \sigma^2\rho^2\epsilon_H^{-4})$ in Theorem 5.1, we have the
298 total gradient complexity

$$\tilde{O}\left(\frac{\sigma L \Delta}{\epsilon^3} + \frac{\sigma \rho^2 L \Delta}{\epsilon \epsilon_H^4} + \frac{\sigma \rho^3 L \Delta}{\epsilon_H^6}\right) = \tilde{O}\left(\frac{\sigma L \Delta}{\epsilon^3} + \frac{\sigma \rho^3 L \Delta}{\epsilon_H^6}\right),$$

299 where the equation is due to the Young's inequality.

300

□

301 7 Conclusions

302 In this paper, we propose a perturbed stochastic gradient framework named Pullback for finding
303 local minima. Pullback can find (ϵ, ϵ_H) -approximate local minima within $\tilde{O}(\epsilon^{-3} + \epsilon_H^{-6})$ stochastic
304 gradient evaluations, which matches the best possible complexity results in the classical $\epsilon_H = \sqrt{\epsilon}$
305 setting. Our results show that simple perturbed gradient methods can be as efficient as more
306 sophisticated algorithms for finding local minima.

307 **References**

- 308 [1] Allen-Zhu, Z. (2018). Natasha 2: Faster non-convex optimization than sgd. In *Advances in*
309 *neural information processing systems*, pages 2675–2686.
- 310 [2] Allen-Zhu, Z. and Hazan, E. (2016). Variance reduction for faster non-convex optimization. In
311 *International Conference on Machine Learning*, pages 699–707.
- 312 [3] Allen-Zhu, Z. and Li, Y. (2018). Neon2: Finding local minima via first-order oracles. In *Advances*
313 *in Neural Information Processing Systems*, pages 3716–3726.
- 314 [4] Arjevani, Y., Carmon, Y., Duchi, J. C., Foster, D. J., Sekhari, A., and Sridharan, K. (2020).
315 Second-order information in non-convex stochastic optimization: Power and limitations. In
316 *Conference on Learning Theory*, pages 242–299. PMLR.
- 317 [5] Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2016). Global optimality of local search for low
318 rank matrix recovery. *arXiv preprint arXiv:1605.07221*.
- 319 [6] Cutkosky, A. and Orabona, F. (2019). Momentum-based variance reduction in non-convex sgd.
320 In *Advances in Neural Information Processing Systems*, pages 15210–15219.
- 321 [7] Daneshmand, H., Kohler, J., Lucchi, A., and Hofmann, T. (2018). Escaping saddles with
322 stochastic gradients. *arXiv preprint arXiv:1803.05999*.
- 323 [8] Du, S. S., Jin, C., Lee, J. D., Jordan, M. I., Póczos, B., and Singh, A. (2017). Gradient descent
324 can take exponential time to escape saddle points. *arXiv preprint arXiv:1705.10412*.
- 325 [9] Fang, C., Li, C. J., Lin, Z., and Zhang, T. (2018). Spider: Near-optimal non-convex optimization
326 via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing*
327 *Systems*, pages 689–699.
- 328 [10] Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015). Escaping from saddle points—online stochastic
329 gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842.
- 330 [11] Ge, R., Jin, C., and Zheng, Y. (2017). No spurious local minima in nonconvex low rank
331 problems: A unified geometric analysis. In *International Conference on Machine Learning*, pages
332 1233–1242. PMLR.
- 333 [12] Hillar, C. J. and Lim, L.-H. (2013). Most tensor problems are np-hard. *Journal of the ACM*
334 (*JACM*), 60(6):1–39.
- 335 [13] Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. (2017). How to escape saddle
336 points efficiently. In *International Conference on Machine Learning*, pages 1724–1732. PMLR.
- 337 [14] Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive
338 variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323.
- 339 [15] Kallenberg, O. and Sztencel, R. (1991). Some dimension-free features of vector-valued martin-
340 gales. *Probability Theory and Related Fields*, 88(2):215–247.
- 341 [16] Lei, L., Ju, C., Chen, J., and Jordan, M. I. (2017). Non-convex finite-sum optimization via scsg
342 methods. In *Advances in Neural Information Processing Systems*, pages 2348–2358.
- 343 [17] Li, Z. (2019). Ssrgd: Simple stochastic recursive gradient descent for escaping saddle points. In
344 *Advances in Neural Information Processing Systems*, pages 1521–1531.
- 345 [18] Li, Z., Bao, H., Zhang, X., and Richtárik, P. (2020). Page: A simple and optimal probabilistic
346 gradient estimator for nonconvex optimization. *arXiv preprint arXiv:2008.10898*.
- 347 [19] Nesterov, Y. (1998). Introductory lectures on convex programming volume i: Basic course.
348 *Lecture notes*, 3(4):5.
- 349 [20] Nesterov, Y. and Polyak, B. T. (2006). Cubic regularization of newton method and its global
350 performance. *Mathematical Programming*, 108(1):177–205.

- 351 [21] Nguyen, L. M., Liu, J., Scheinberg, K., and Takáč, M. (2017a). Sarah: A novel method for
352 machine learning problems using stochastic recursive gradient. In *International Conference on*
353 *Machine Learning*, pages 2613–2621. PMLR.
- 354 [22] Nguyen, L. M., Liu, J., Scheinberg, K., and Takáč, M. (2017b). Stochastic recursive gradient
355 algorithm for nonconvex optimization. *arXiv preprint arXiv:1705.07261*.
- 356 [23] Pham, N. H., Nguyen, L. M., Phan, D. T., and Tran-Dinh, Q. (2020). Proxsarah: An efficient
357 algorithmic framework for stochastic composite nonconvex optimization. *Journal of Machine*
358 *Learning Research*, 21(110):1–48.
- 359 [24] Pinelis, I. (1994). Optimum bounds for the distributions of martingales in banach spaces. *The*
360 *Annals of Probability*, pages 1679–1706.
- 361 [25] Reddi, S. J., Hefny, A., Sra, S., Póczos, B., and Smola, A. (2016). Stochastic variance reduction
362 for nonconvex optimization. In *International Conference on Machine Learning*, pages 314–323.
- 363 [26] Tran-Dinh, Q., Pham, N. H., Phan, D. T., and Nguyen, L. M. (2019). Hybrid stochastic gradient
364 descent algorithms for stochastic nonconvex optimization. *arXiv preprint arXiv:1905.05920*.
- 365 [27] Wang, Z., Ji, K., Zhou, Y., Liang, Y., and Tarokh, V. (2019). Spiderboost and momentum: Faster
366 variance reduction algorithms. *Advances in Neural Information Processing Systems*, 32:2406–
367 2416.
- 368 [28] Xiao, L. and Zhang, T. (2014). A proximal stochastic gradient method with progressive variance
369 reduction. *SIAM Journal on Optimization*, 24(4):2057–2075.
- 370 [29] Xu, Y., Jin, R., and Yang, T. (2017). First-order stochastic algorithms for escaping from saddle
371 points in almost linear time. *arXiv preprint arXiv:1711.01944*.
- 372 [30] Zhou, D., Xu, P., and Gu, Q. (2018). Stochastic nested variance reduction for nonconvex
373 optimization. *Advances in Neural Information Processing Systems*, 31:3921–3932.
- 374 [31] Zhou, D., Xu, P., and Gu, Q. (2020). Stochastic nested variance reduction for nonconvex
375 optimization. *Journal of machine learning research*.

376 **Checklist**

- 377 1. For all authors...
- 378 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
379 contributions and scope? [Yes]
- 380 (b) Did you describe the limitations of your work? [Yes] We note in Remark 5.2 that there
381 is a gap in the dependency of ϵ_H between our result and the lower bound. We will
382 explore this in future work.
- 383 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Our
384 work studies the theoretical aspect of optimization algorithms, thus it does not have
385 any negative social impact.
- 386 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
387 them? [Yes]
- 388 2. If you are including theoretical results...
- 389 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 390 (b) Did you include complete proofs of all theoretical results? [Yes]
- 391 3. If you ran experiments...
- 392 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
393 mental results (either in the supplemental material or as a URL)? [N/A]
- 394 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
395 were chosen)? [N/A]
- 396 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
397 ments multiple times)? [N/A]
- 398 (d) Did you include the total amount of compute and the type of resources used (e.g., type
399 of GPUs, internal cluster, or cloud provider)? [N/A]
- 400 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 401 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 402 (b) Did you mention the license of the assets? [N/A]
- 403 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 404
- 405 (d) Did you discuss whether and how consent was obtained from people whose data you're
406 using/curating? [N/A]
- 407 (e) Did you discuss whether the data you are using/curating contains personally identifiable
408 information or offensive content? [N/A]
- 409 5. If you used crowdsourcing or conducted research with human subjects...
- 410 (a) Did you include the full text of instructions given to participants and screenshots, if
411 applicable? [N/A]
- 412 (b) Did you describe any potential participant risks, with links to Institutional Review
413 Board (IRB) approvals, if applicable? [N/A]
- 414 (c) Did you include the estimated hourly wage paid to participants and the total amount
415 spent on participant compensation? [N/A]