1   We thank the reviewers for appreciating our work, and for their insightful and constructive comments.

2   

3   **1. Regarding the innovation of this paper.**

4   **A1:** For the comparison with [15], Ge et al. [15] uses the SVRG estimator $v_t = \frac{1}{b}\sum_{i\in I_b}\left(\nabla f_i(x_t) - \nabla f_i(\tilde{x})\right) + \nabla f(\tilde{x})$

5 which reuses a fixed snapshot full gradient $\nabla f(\tilde{x})$, and achieves a suboptimal convergence $\tilde{O}(n^{2/3}/\epsilon^2)$. Here we use the

6 stochastic recursive gradient estimator $v_t = \frac{1}{b}\sum_{i\in I_b}\left(\nabla f_i(x_t) - \nabla f_i(x_{t-1})\right) + v_{t-1}$ to obtain the improved optimal

7 convergence $\tilde{O}(n^{1/2}/\epsilon^2)$ for escaping saddle points which matches the lower bound $\Omega(n^{1/2}/\epsilon^2)$. Moreover, we also

8 have discussed with the authors of [15], we cannot directly add a perturbation step to SPIDER [11] for escaping saddle

9 points. The original SPIDER [11] uses an additional negative-curvature search subroutine (Neon2 [4]) to escape saddle

10 points. Thus, in this paper, we use a few simple but effective modifications to obtain a simple algorithm for achieving

11 the new convergence results. There are three differences between our SSRGD algorithm and SPIDER. 1) SPIDER [11]

12 uses the normalized gradient update $x_{t+1} = x_t - \eta(v_t/\|v_t\|)$, while our SSRGD uses the direct update $x_{t+1} = x_t - \eta v_t$,

13 where $v_t$ is the gradient direction estimator. 2) SPIDER uses a very small step size $\eta = O(\epsilon/L)$, while our SSRGD

14 uses the more standard step size $\eta = O(1/L)$. 3) SPIDER uses an additional negative-curvature search subroutine (e.g.,

15 Neon2 [4]) to escape saddle points, while our SSRGD directly escapes saddle points by adding a uniform perturbation

16 sometimes. These key differences are crucial for achieving the new convergence result $\tilde{O}(n^{1/2}/\epsilon^2)$ and also show that

17 our SSRGD is more attractive in practice. Intuitively, in the large gradient situations (e.g., at the beginning phase),

18 SPIDER goes very slowly due to the normalized gradient update. Also note that SPIDER can still go slowly even

19 in the small gradient situations due to its tiny step size. As a result, SPIDER may get worse convergence results in

20 some situations (see the following response A2). Moreover, our SSRGD which uses direct update (non-normalized),

21 standard step size and simple perturbation step (no negative-curvature search subroutine) is more attractive and easy to

22 implement in practice. Besides, we believe our simple/natural analysis will be useful for future work.

23   **2. Regarding the convergence results.**

24   **A2:** We would like to point out that our convergence rate is not worse than SPIDER [11]. Moreover, our SSRGD is

25 better than SPIDER if $\delta \leq 1/\sqrt{n}$ (i.e., high second-order accuracy situation) due to the $1/\delta^5$ term in SPIDER (see our

26 Table 1). Concretely, for achieving an $(\epsilon, \delta)$-second-order stationary point ($\|\nabla f(x)\| \leq \epsilon$ and $\lambda_{\min}(\nabla^2 f(x)) \geq -\delta$),

27 the last three algorithms in our Table 1 can obtain the best rate $\tilde{O}(n^{1/2}/\epsilon^2)$ in a large range of $\epsilon$ and $\delta$. SNVRG+Neon2

28 [32] and our SSRGD achieve the rate $\tilde{O}(n^{1/2}/\epsilon^2)$ if $1/\epsilon \geq 1/\delta^2$ and $n \leq 1/\epsilon$. SPIDER+Neon2 [11] achieves the rate

29 $\tilde{O}(n^{1/2}/\epsilon^2)$ if $1/\epsilon \geq 1/\delta^2$ and $n \geq 1/\epsilon$. Similarly, for the online case (Table 2), SNVRG+Neon2 [32], SPIDER+Neon2

30 [11] and our SSRGD can also obtain the best rate $\tilde{O}(1/\epsilon^3)$ in a large range of $\epsilon$ and $\delta$ for achieving an $(\epsilon, \delta)$-second-

31 order stationary point, i.e., SNVRG+Neon2 [32] and our SSRGD achieve the best rate $\tilde{O}(1/\epsilon^3)$ if $1/\delta^3 \leq 1/\epsilon$, and

32 SPIDER+Neon2 [11] achieves the best rate $\tilde{O}(1/\epsilon^3)$ if $1/\delta^2 \leq 1/\epsilon$.

33   

34   **Future directions:** The achieved convergence rate $O(n^{1/2}/\epsilon^2)$ matches the lower bound $\Omega(n^{1/2}/\epsilon^2)$ in finite-sum

35 case. However, whether $\Omega(1/\epsilon^3)$ is the lower bound for online case is unknown. This should be an important future

36 problem. The combination of the stabilized trick proposed in Ge et al. [15] is also indeed an interesting extension.

37   **Experiments:** We will try to add some experiments although previous related papers also did not do the experiments.

38   

39   **Regarding the difficulties to avoid using the negative-curvature search subroutine:** Previous algorithms use the

40 negative-curvature search subroutine (e.g. Neon2) to find the approximate smallest eigenvector of the Hessian (i.e.,

41 decreasing direction) near saddle points. However, if we just use a uniform perturbation step, it is not very easy to

42 obtain the smallest eigenvector direction since we only roughly have $O(r/\sqrt{d})$ amount in the smallest eigenvector

43 direction in a uniform perturbation ball with radius $r$ in $\mathbb{R}^d$ space. In our analysis, by simultaneously bounding the

44 coupled distance and coupled variance, we can show that the amount in the smallest eigenvector direction will increase

45 exponentially and thus avoid the negative-curvature search subroutine.

46   

47   **Subroutine:** The subroutine Neon [31] extracts the negative curvature based on a perturbation step, power method and

48 Nesterov's accelerated gradient method. Neon2 [4] is based on a perturbation step, Oja's algorithm and Chebyshev

49 polynomial. Thus they are at least more complicated in implementation than just a perturbation step as our SSRGD

50 used. Please see A2 in Response to Reviewer #1 for the convergence comparison among the last three algorithms.

51   **Optimal among simple algorithms:** Note that Du et al. [10] showed that it is necessary to add the perturbation step

52 for gradient descent to escape saddle points efficiently. Our SSRGD indeed just adds a perturbation step on the simple

53 recursive gradient descent to achieve the optimal convergence $\tilde{O}(n^{1/2}/\epsilon^2)$. **Improvement/Comparison:** Intuitively, it

54 can improve upon SVRG [21] since it uses more aggressive recursive gradient estimator instead of the conservative

55 SVRG estimator which uses a fixed full gradient snapshot. Please see A1 in Response to Reviewer #1 for more details

56 and a comparison with SPIDER [11], and see Response to Reviewer #3 for the discussion of saddle point iterations.