

1 **Reply to Reviewer 2, 3 and 4**

2 Novelty of the analysis: Besides keeping the estimator of  $\nabla f(\mathbf{x}_k, \mathbf{y}_k)$  sufficiently accurate like SARAH/SPIDER for  
3 convex/nonconvex minimization, the minimax problem also requires  $\mathbf{y}_k$  to be close to  $\mathbf{y}^*(\mathbf{x}_k)$ , which leads to a more  
4 challenging analysis. Our analysis is based on a recursive relationship between  $\Delta_k = \mathbb{E}[\|\mathbf{v}_k - \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|_2^2 +$   
5  $\|\mathbf{u}_k - \nabla_{\mathbf{y}} f(\mathbf{x}_k, \mathbf{y}_k)\|_2^2]$  and  $\delta_k = \mathbb{E} \|\mathcal{G}_{\lambda, \mathbf{y}}(\mathbf{x}_k, \mathbf{y}_k)\|_2^2$  (defined in line 364 of Appendix B) as in Lemma 13 (Corollary 2,  
6 line 418, Appendix B), which guarantees both  $\Delta_k$  and  $\delta_k$  to be smaller than  $\mathcal{O}(\kappa^{-2}\varepsilon^2)$  (line 434-435, Appendix B). This  
7 is a nontrivial ingredient of the analysis. To achieve the desired convergence rate, all of the stepsizes, mini-batch sizes,  
8 number of the inner iterations, and the orders of  $\Delta_k$  and  $\delta_k$  must be balanced carefully. In comparison, SARAH/SPIDER  
9 for convex/nonconvex minimization only needs to consider the estimator of gradient, which does not involve the extra  
10 complexity in minimax problems.

11 **Reply to Reviewer 2 and 4**

12 Algorithm is complicated: As mentioned in the setting of experiments (line 589, Appendix F), we can select  $q = m =$   
13  $\lceil n/S_2 \rceil$  heuristically and the empirical result show it performs well in practice. We agree that it is worth to see weather  
14 there exists a simpler variant of SREDA which also holds the theoretical guarantee.

15 **Reply to Reviewer 3 and 5**

16 Optimal dependency on  $\varepsilon$ : The lower bound means any stochastic first-order algorithm requires at least  $\mathcal{O}(\varepsilon^{-3})$  calls of  
17 stochastic first-order oracle ( $\nabla_{\mathbf{x}} F(\mathbf{x}, \mathbf{y}; \xi), \nabla_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}; \xi)$ ) to find  $\varepsilon$ -stationary point of  $\Phi(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y})$   
18 (line 93, Definition 1). As we mentioned in line 180-185, we can consider the special case of minimax problem whose  
19 objective function has the form  $f(\mathbf{x}, \mathbf{y}) = g(\mathbf{x}) + h(\mathbf{y})$  where  $g$  is possibly nonconvex and  $h$  is strongly-concave,  
20 which leads to minimizing on  $\mathbf{x}$  and maximizing on  $\mathbf{y}$  are independently. Consequently, finding  $\mathcal{O}(\varepsilon)$ -stationary point  
21 of the corresponding  $\Phi(\mathbf{x})$  can be reduced to finding  $\mathcal{O}(\varepsilon)$ -stationary point of nonconvex function  $g(\mathbf{x})$ , which is  
22 based on the stochastic first order-oracle  $\nabla_{\mathbf{x}} F(\mathbf{x}, \mathbf{y}; \xi) = \nabla g(\mathbf{x}; \xi)$  (this equality holds for any  $\mathbf{y}$  and we also have  
23  $\nabla_{\mathbf{x}} F(\mathbf{x}, \mathbf{y}^*(\mathbf{x}); \xi) = \nabla g(\mathbf{x}; \xi)$ ). Hence, the analysis of stochastic nonconvex miminization problem [5] based on  
24  $\nabla g(\mathbf{x}; \xi)$  can directly lead to the  $\mathcal{O}(\varepsilon^{-3})$  lower bound for our minimax problem.

25 **Reply to Reviewer 2**

26 Relations to compositional optimization: We thank the reviewer for pointing out this valuable reference  
27 (arXiv:1908.11468). We are happy to cite this paper and compare it with SREDA. Both this work and SREDA  
28 use variance reduction to address nonconvex multi-level (two-level) optimization problem, however, their settings are  
29 quite different. We can reformulate our minimax problem (2) as compositional problem:

$$\min_{\mathbf{x}, \mathbf{y}} f_2(f_1(\mathbf{x}, \mathbf{y})) + \Psi(\mathbf{y}), \quad (\text{a})$$

30 where  $f_1(\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \arg \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}))$ ,  $f_2(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}, \mathbf{y})$  and  $\Psi(\mathbf{y})$  is the indicator function of  $\mathcal{Y}$ . The two-level  
31 nested-SPIDER (arXiv:1908.11468) requires to access the stochastic gradient of  $f_1$  and  $f_2$  to solve the compositional  
32 problem. For the two-level formulation (a) of our minimax problem, it is not natural to access the (stochastic) gradient  
33 of  $f_1$  as an oracle since we can not provide the closed form of  $\arg \max_{\mathbf{y} \in \mathbb{R}} f(\mathbf{x}, \mathbf{y})$  and its (stochastic) gradient in  
34 general. A more reasonable way is to solve it by accessing the (stochastic) gradient of  $f$  like SREDA.

35 **Reply to Reviewer 3**

36 1. Problem formulation in experiments: Due to the space limitation, we gave the problem formulation in Appendix F.  
37 We are happy to follow the reviewer’s suggestion and include it in the main text if the paper is accepted.

38 2. Minor remarks: Thanks for the suggestion. We will simplify the presentation of Theorem 2 and cite the papers about  
39 non-stochastic algorithms for minimax problems.

40 **Reply to Reviewer 4**

41 1 & 2. Please see “Reply to Reviewer 2, 3 and 4” and “Reply to Reviewer 2 and 4” above.

42 3. We thank the reviewer for pointing out this valuable reference, which is complementary to our work. We will cite it  
43 and compare to SREDA. First, the convergence result of epoch-GDA is based on the measure of nearly  $\varepsilon$ -stationary  
44 point because it also considers an additional constraint on  $\mathbf{x}$ , while SREDA is based on  $\varepsilon$ -stationary point. Second, the  
45 stochastic first-order oracle complexity of SREDA depends on  $\mathcal{O}(\kappa^3\varepsilon^{-3})$ , while epoch-GDA depends on  $\tilde{\mathcal{O}}(\kappa^2\varepsilon^{-4})$ .

46 **Reply to Reviewer 5**

47 Presentation of Algorithm 4: We thank the reviewer for this suggestion, and will add the arguments in the expression of  
48 the function ConcaveMaximizer.