

1 We thank all the reviewers for their thoughtful reviews and for pointing out some typos in the manuscript.

2 **Responses to Reviewer-1’s comments:**

3 “...the restriction on  $\delta$  seems crude. I cannot see where these restrictions come from...”

4 First, we note that our technical results remain the same without this restriction on  $\delta$ . However, this restriction is made  
5 to comply with a standard requirement on  $\delta$  in the literature of differential privacy to ensure the robustness of the  
6 definition. In particular, it is standard to require  $\delta$  to be much smaller than  $1/n$  to rule out trivial mechanisms, e.g., the  
7 one that selects one individual uniformly at random and publishes her record in the clear (note that such mechanism  
8 satisfies  $(0, 1/n)$ -differential privacy but is blatantly non-private). This requirement on  $\delta$  is discussed in several early  
9 references on differential privacy including the textbook by Dwork & Roth (towards the end of Sec. 2.3.3) and the  
10 survey by Salil Vadhan, and several others. In fact, in some of these references  $\delta$  is even assumed to be a negligible  
11 function of  $n$  (i.e., smaller than the inverse of any polynomial in  $n$ ).

12 “There seems to be an assumption  $\sigma \geq 1$  in Lemma 3 of [ACG+16], however I do not see this appearing in the article.”

13 First, we note that the lower bound on  $\sigma$  in [ACG+16, Lemma 3] is because the bound on the norm of  $f$  in that lemma  
14 (which represents the Lipschitz constant) is assumed to be 1. So, when the Lipschitz constant is  $L$ , then by simple  
15 re-normalization it is easy to see that the analogous lower bound on  $\sigma$  is  $L$ . Second, as pointed out in the supplementary  
16 document, due to different normalization of the noise in the gradient update step (noise in our case is normalized by the  
17 batch size  $m$ ), our setting of  $\sigma$  is smaller than that of [ACG+16] by a factor of  $m$ . Taking these different normalizations  
18 into account, we note that the analogous lower bound on  $\sigma$  is  $L/m$ , and it is indeed satisfied in our case. To see this,  
19 note that given the setting of  $m$  in Step 2 of Algorithm 1 and the fact that  $\epsilon \leq 1$ , the setting of  $\sigma$  in Step 1 implies that  
20  $\sigma \geq 2L/m$ .

21 “...is it not entirely correct to say (as in Definition 1): ‘for any pair of datasets  $S$  and  $S'$  differ in exactly one data  
22 point...’ Then the results of [ACG+16] would not hold. Perhaps you could comment on the neighbouring relation of  $S$   
23 and  $S'$ .”

24 The definition of differential privacy with respect to the adjacency notion involving addition/removal of one element  
25 in the dataset is equivalent (up to a factor of 2 in the privacy parameters) to the definition w.r.t. the adjacency notion  
26 involving replacement of one element. This follows from a simple triangle-inequality style argument since replacing  
27 one element  $z$  by another element  $z'$  can be carried out via a removal step (of  $z$ ) followed by an addition step (of  $z'$ ).  
28 Hence, the same techniques and results of [ACG+16] still apply in our case (after renormalizing the privacy parameters  
29 by a factor of 2).

30 **Response to Reviewer-2’s comments:**

31 “I still think this work lacks the experimental part. As I know, most of the recent work on the central  $(\epsilon, \delta)$  DP-ERM has  
32 experimental study such as [1-6]. ... should provide some experimental study in order to say the improvement.”

33 Our work is the first one to focus on differentially private stochastic convex optimization (as opposed to previous works  
34 on private empirical risk minimization). Our primary contribution is theoretical and we believe that the fundamental  
35 nature of the question we resolve makes the work interesting for the community. At the same time, our analysis concerns  
36 standard algorithms such as mini-batch DP SGD and objective perturbation for which one can easily find experimental  
37 results in the literature. The only new algorithmic aspect is the use of prox step (to get Moreau-Yosida smoothing). This  
38 step is necessary in the worst case but will not make much difference for simple linear models used in most experiments.  
39 Thus we do not think that experiments are likely to give additional insights into the question we investigate.