

1 We thank the reviewers for the positive feedback. We will answer the questions in the following.

2 **Reviewer1.** *Q1: Robustness of the proposed methods against Byzantine workers.*

3 The robustness of Sign and Median SGD both come from the fact that when performing variable updates, the mean of
4 gradients is replaced by the median of gradients, which is less sensitive to extreme values (i.e., those that are possibly
5 provided by Byzantine workers). As noticed by the reviewer, the perturbation mechanism in Section 4 is gradually
6 converting the estimated statistic from median to mean by adding more noise. Thus, the more noise is added, the less
7 robust the estimated statistic will be. This implies that the robustness of Noisy Sign and Median SGD Byzantine workers
8 depends on the the amount of artificial noise (b in Algorithm 3 and 4). On heterogeneous data with Byzantine workers,
9 the performance of the median-based algorithms can be affected by mainly two factors. One is the gap between median
10 and mean, and the other one is the misleading information provided by Byzantine workers. When the possible effect of
11 Byzantine workers is relatively small (e.g. a small number of Byzantine workers) compared with the median-mean gap,
12 some noise is still preferred to reduce the median-mean gap even though this could amplify the effect of Byzantine
13 workers (there could be a trade-off). For non-heterogeneous data (e.g. iid setting), the median could be very close to the
14 mean and the noise scale b in Noisy Sign and Median SGD should be set small or even 0, the effect of such a small
15 noise might be negligible for both convergence and robustness. We will add a more detailed discussion to the paper.

16 **Reviewer 2.** *Q1: What claims can be made in terms of privacy.*

17 This is a very interesting question. The noise added to the gradient might provide some level of differential privacy
18 while improving the utility of the algorithms in practice. If one knows the upper bound on L_2 norm of all possible
19 gradients, one can use it to calculate the differential privacy cost ϵ using standard privacy accountant (e.g. Abadi et al.
20 [2016]). However, in many cases the L_2 norm bound on gradients is unknown or hard to compute, and gradient clipping
21 is required. It is unclear how the utility of the algorithms will be affected if gradient clipping is applied. We will discuss
22 this question in the paper.

23 **Reviewer 4.** *Q1: The increased variance in the treatment.*

24 Indeed, as mentioned by the reviewer, the perturbation mechanism applied in Noisy Sign and Median SGD can
25 increase variance of the original median estimator while making median closer to mean. This mechanism can be
26 understood as trading bias with variance. One needs to properly deal with the added variance in the optimization
27 algorithms. For example, the optimal learning rate of Noisy Sign and Median SGD can depend on the added variance as
28 implied by Theorem 5 and Theorem 6. In practice, such increased variance can change the best learning rate found by
29 hyperparameter search strategies. We will add more discussion on the effect of increased variance.

30 **References**

31 Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learn-
32 ing with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications*
33 *Security*, pages 308–318, 2016.