

1 We thank all the reviewers for their feedback and comments. They will be incorporated into the revised version.

2 We would like to highlight several points regarding our contribution. In the paper we consider the robust linear
3 regression problem which has drawn significant interest in the last decades. Even in the simpler setting where only the
4 responses are obviously corrupted there have been numerous works proposing novel algorithms that recover the gold
5 model θ^* . Though this setting may seem simple, it has considerable applications. Nonetheless, currently none of the
6 proposed algorithms are suitable for modern large-scale problems. This observation is the incentive behind our work:
7 we propose the first highly scalable and efficient algorithm for the robust linear regression problem with oblivious
8 response corruptions. We show that running SGD on the ℓ_1 loss outperforms all current algorithms, theoretically
9 and empirically. Not only is it currently the unique suitable large-scale algorithm, but it also presents several major
10 advantages even in the offline setting: a) there are no hyperparameters to tune, the knowledge of σ or η are not necessary
11 b) the $O(1/n)$ dominant convergence terms are independent of the conditioning κ of the data c) the algorithm adapts to
12 the difficulty of the adversary since the rate depends on the effective corruption proportion $\tilde{\eta}$. We also want to underline
13 that the algorithm is extremely simple to implement: it represents only two lines of code.

14 Concerning the Gaussian assumption on the data, we agree that it is restrictive and that it does not always represent real
15 world situations. However we highlight that this assumption (or similar specific structural assumptions on the design
16 matrix and the noise) is made in the previous papers that deal with the framework that we consider. The key for our
17 analysis is that the continuous density of (x, ε) smooths the objective f , this enables a proof of type Polyak-Juditsky to
18 obtain the final rate. The more restrictive Gaussian assumption is made in order to obtain a closed form for f which
19 enables simpler (but not simple !) computations to obtain the final convergence rate. We believe that the proof could
20 be extended to any sub-Gaussian continuous data, however this would require considerably heavier and cumbersome
21 computations (see the figure for an example where the covariates are sampled from a zero-mean uniform distribution
22 with covariance H and the noise from a uniform distribution with variance σ^2).

23 **To Reviewer 1.** As pointed out, the dependency on η (or rather $\tilde{\eta}$) we obtain
24 might not be optimal and is still an interesting open question. It is however the
25 best dependency which has been obtained since the effective outlier proportion
26 $\tilde{\eta}$ is strictly smaller than η . We agree that real data experiments would have been
27 a plus, they will be added to the revised version.

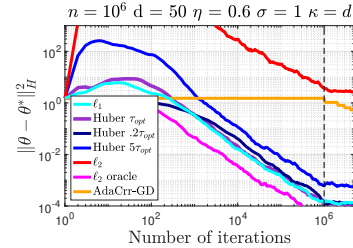
28 **To Reviewer 2.** The centered data assumption is a classical assumption which
29 is often made in machine learning, note that it is for instance made in previous
30 papers that deal with the framework we consider. This assumption is reasonable
31 since it is of common practice to pre-process and center the features. We highlight
32 the fact that even for the Gaussian linear regression problem the minimax rate of
33 estimation is $O(1/n)$, in the sense that no estimator that uses the same information
34 can improve upon this rate. Therefore, in the setting we consider, the rate we
35 obtain cannot be improved.

36 **To Reviewer 3.** The linear model is one of the simplest model we could have considered and linear regression is
37 certainly amongst the oldest and most fundamental statistical methods. However: (a) it is still intensively used in
38 practice and studied in theory, (b) more complex models can be seen as an extension of it and (c) there are still many
39 interesting open questions. Indeed even for this simple model there were still no efficient algorithms which could deal
40 with large datasets, hence the motivation of our work. The extension of the algorithm to broader and more general
41 models is undoubtedly interesting, however it is outside the scope of our present work.

42 **To Reviewer 4.** As seen in Theorem 4, the terms depending on the variance σ go to 0 as $\sigma \rightarrow 0$. Hence the extreme
43 case where there is no 'nice' noise is not pathological and the algorithm still performs well. On the other hand, the
44 performance of the algorithm is indeed inevitably dependent of $\tilde{\eta}$ (and not on η as for other robust algorithms): we
45 cannot expect any algorithm to behave well as $\tilde{\eta}$ goes to 1 since this corresponds to having all the data corrupted. The
46 dependency $(1 - \tilde{\eta})^{-2}$ we obtain is to our knowledge the best which has been obtained.

47 We thank the reviewer for the interesting missing reference which will be added to the revised version. The literature on
48 robust regression is vast and we have tried to be as exhaustive as possible. The mentioned paper proposes an interesting
49 EM-algorithm, however they do not give any convergence guarantee. Furthermore every M-step of their algorithm is
50 computationally heavy since it requires solving a weighted LS problem on the whole dataset, making it inappropriate
51 for large-scale problems.

52 As pointed out, the Tukey loss is often used to deal with corruptions, however its non-convexity makes it much harder
53 to analyse and optimise. Furthermore, as for the Huber loss, the Tukey loss requires tuning an extra parameter and we
54 believe it will not lead to better experimental results for the linear model with corrupted response. However the Tukey
55 loss (as well as other redescending M-estimators) is well-known for being more robust to corruption in the features x .
56 We let further investigations as future work.



Online robust regression on non-Gaussian data. Same adversarial setting as in Figure 1 in the paper.