

1 Firstly, we thank the reviewers for their valuable comments. We will address the comments of each reviewer in what
2 follows.

3 **R1, R4:** *The results are very specific to the particular model:* Indeed it is the case that our theoretical results assume
4 that data providers are constrained in ℓ_2 -norm, and that both the learner and the data providers are interested in solving
5 linear least squares problems. However, it is fairly easy to see that our theoretical results generalise to kernel ridge
6 regression, and thus, our theoretical results hold for a far wider array of learning models than linear predictors. In
7 addition, we believe that our work is an important first step in relaxing the overly pessimistic assumptions of adversarial
8 machine learning. Whilst it is not reasonable in practice to assume that data is sampled i.i.d. from the distribution of
9 interest, neither is it reasonable to assume that data is provided with the sole intention of hindering learning. By taking
10 the incentives of data providers directly into account during the optimisation process, we hope to better reflect the
11 reality of sampled data in practice. As previously stated, we believe our work forms a first step in achieving this goal.

12 **R1:** *It is not clear from the paper how their model differs from previously considered ones:* In this paper, we study a
13 specific subclass of SPGs in which both the learner and data providers are interested in solving least squares linear
14 regression problems with their own data labels. Whilst Brückner and Scheffer (2011) give an algorithm which converges
15 to local optima for general SPGs, we give an algorithm for this specific subclass of SPGs which converges to **globally**
16 **optimal solutions**. Note that SPGs are bilevel optimisation problems, which are, in general, NP-hard. With this in mind,
17 we believe our results are novel and significant, as we have provided a practically efficient algorithm for a large subset
18 of bilevel optimisation problems. Algorithms for specific subclasses of SPGs have been considered before, but to the
19 best of our knowledge, our paper is the first to consider SPGs for linear least squares regression. Closely related to
20 our work is the problem of robust regression in which the challenge is to choose a model which performs well in the
21 presence of worst-case noise. The subclass of SPGs we consider allow us to model data providers with more nuanced
22 motivations, which we believe are more likely to arise in practice. We will highlight these differences from previous
23 work in future versions of the paper.

24 **R2:** *Why would the learner ever evaluate on both the manipulated and unmanipulated data in practice?:* We believe that
25 our work has applications whenever a learner has access to a reliable, but costly, verification process. In practice, many
26 machine learning models make decisions which affect those who provide data. Thus, data providers may manipulate
27 the data they submit in order to obtain a preferred outcome. In cases in which the learner has access to a verification
28 process, the learner can recover the originally sampled data as well as the goals of the data providers. Unfortunately,
29 this verification process may be too expensive to use on every single data point. However, the learner may be able to
30 verify a sample which can be used to learn the motivations of data providers and select a model which anticipates the
31 manipulations that data providers are likely to make. We believe that our theoretical model captures this dynamic. Note
32 that using verified data points to improve learning has been applied extensively in adversarial machine learning contexts
33 (for examples, see Charikar *et al.* (2017) and Raghavendra and Yau (2020)). One practical example of such a setting
34 is insurance fraud. An insurance company may gather information from a customer to better evaluate potential risk.
35 However, a smart customer, who knows the information they provide could be used to decide the cost of their insurance,
36 may lie when submitting their information. Whilst insurance companies can often verify information regarding their
37 customers, verification is often expensive, requiring a significant deployment of staff and/or resources.

38 **R2:** *I am guessing it is the "interior point" line - this seems suspiciously bad. Is this a strong baseline / was the baseline*
39 *correctly executed with proper hyper-parameter choice?:* The reviewer is correct in their assumption that the interior
40 point line in figure 1 corresponds to the method of Brückner and Scheffer. We shall update the legend of figure 1
41 to be more clear. This approach involves reformulating SPGs into a single-level optimisation problem which can be
42 solved via conventional optimisation techniques. In their work, Brückner and Scheffer highlight a number of SPGs for
43 which this problem reformulation is simple and can be easily solved. For the family of SPGs we consider, the problem
44 reformulation is nontrivial and nonconvex. The error tolerances used for both our algorithm and the interior point
45 method we use to solve this problem reformulation are the same. We believe that the poor performance of the interior
46 point method for higher values of γ reflects ill-conditioning issues present in the nonconvex objective, but we cannot
47 confirm this.

48 **R2:** *Is there any way to contextualize the MSE in Figure 1?:* A data label in the Medical Personal Costs dataset
49 corresponds to the medical charges for a given individual. Since these charges range from \$1000 to \$63,000, we
50 numerically scale the data labels by dividing them by 100 before passing them to each algorithm. For values of $\gamma > 0.4$,
51 we observe that, for modest data providers, our algorithm is at least more accurate by \$45 on average. For severe data
52 providers, we observe an even greater difference. For example, when $\gamma > 0.5$, the predictions made by our algorithm
53 are at least \$120 more accurate on average. Considering the modest data provider only sought to reduce their charges
54 by \$100, and the severe data provider only sought to reduce their charges by \$300, we believe these differences are
55 fairly significant. We shall provide more contextualisation for our empirical experiments in future versions of the paper.