1 We thank all the reviewers for their feedback!

2 Our paper formalizes a data acquisition problem when one cannot verify the true labels of the collected data. Under-
3 standing this problem is important as high-quality data are crucial for high-quality output of machine learning systems.
4 The writing of our submission focused more on the basics to ensure clarity for the general, diverse Neurips readers.
5 Most of the technical results were either deferred to the appendix or compressed to fit in the page limit. We thus want to
6 elaborate a bit more on some of our technical contributions.

7 One of our major technical contributions is the explicit sensitivity guarantee of the peer-prediction style mechanisms.
8 Sensitivity is the key property we want for a data acquisition mechanism, which prevents data providers from harmful
9 misreporting. In this work, we provide explicit and verifiable conditions for sensitivity, which is absent in the previous
10 work. In [1], checking whether the mechanism will be sensitive requires knowing whether a system of linear equations
11 (which has exponential size in our problem) has a unique solution. So it is not clear how likely the mechanisms will be
12 sensitive. In our data acquisition setting, we are able to give much stronger guarantees:

13 1. When $\theta$ has a finite support, we give a sufficient condition for sensitivity. The condition only uses the
14    distribution of $\theta$ and a single data point (Corollary 5.1). The basic message is that when there is enough
15    "correlation" between other people's data and $\theta$, the mechanism will be sensitive. Corollary 5.1 quantifies the
16    "correlation" using the k-rank of the distribution matrix. It is arguably not difficult to have enough "correlation":
17    a naive relaxation of Corollary 5.1 says that assuming different $\theta$ lead to different data distributions, the
18    mechanism will be sensitive if the total number of other people's data points $\geq |\Theta| - 1$.[1]

19 2. For an exponential family distribution, we give the explicit necessary and sufficient condition for the mechanism
20    to be sensitive (Theorem 5.3), which is based on a function of the normalization term of the exponential family
21    probability.

22 This kind of sensitivity guarantee is possible because of the special structure of the reports (or the signals): each dataset
23 consists of i.i.d. samples (conditioning on theta), despite the fact that the signal space for our problem is much larger
24 (the number of possible realizations of a dataset is exponentially large).

25 Besides sensitivity, another important property of our mechanisms is budget feasibility.

26 1. For the one-time acquisition, the log PMI payment rule requires PMI to be bounded. We give a polynomial-time
27    method (a trivial method would take exponential time) to find the bounds for finite-size $\Theta$ (Appendix C.2).

28 2. For the multiple-time mechanism, budget feasibility can be guaranteed for any underlying distribution. This is
29    achieved by carefully choosing the convex function $f$ in the $f$-mutual information gain to have a bounded
30    derivative.

31 If the paper is accepted, we will add more explanation about our technical results and factor out some omitted key steps
32 in the proofs. We thank all the comments on the writing/related work/typos etc. We will revise the paper accordingly.

33 **Other minor comments:**

34 **About the proof of Theorem 5.1.** Yes, showing budget-fixed from budget bounded is straightforward based on our
35    normalization. We will make it clear in the final version. Thanks for pointing this out.

36 **About Example 3.2.** It is a claim about knowledge: we may not need to know the entire prior to compute $p(\theta|\mathbf{x}_i, y_i)$.
37    In this example, it is right we need to know $p(\theta)$, but we do not necessarily need $p(\mathbf{x}_i, y_i|\theta)$ because $p(y_i|\mathbf{x}_i, \theta)$
38    will be sufficient.

# References

40 [1] Yuqing Kong and Grant Schoenebeck. Water from two rocks: Maximizing the mutual information. In *Proceedings*
41    *of the 2018 ACM Conference on Economics and Computation*, pages 177–194, 2018.

---

[1] Here we have found a typo in our statement of Corollary 5.1, the ineuqality should be $\sum_{j \neq i}(rank_k(G_j) - 1) N_j + 1 \geq |\Theta|$
rather than $\sum_{j \neq i} rank_k(G_j) (N_j - 1) + 1 \geq |\Theta|$. We will correct this in our final version.