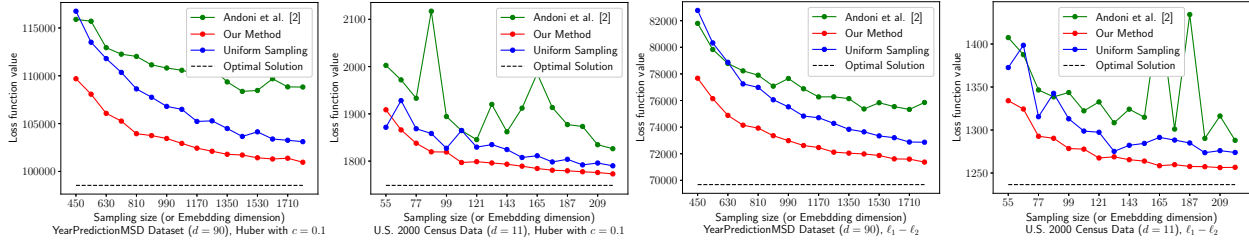


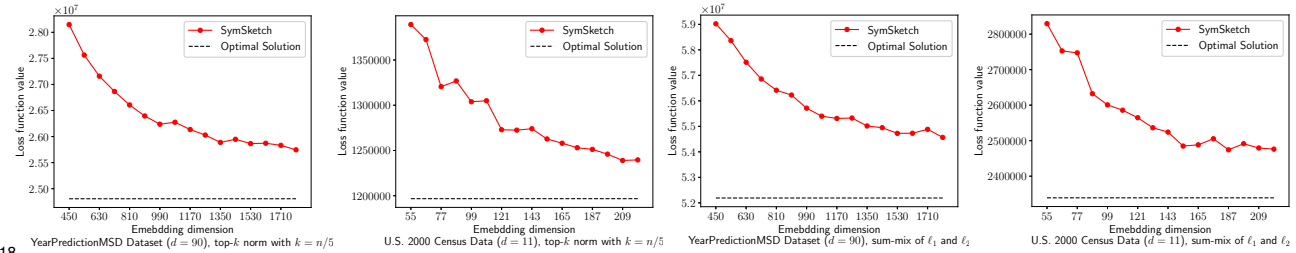
1 We thank all reviewers for their comments, and will incorporate suggestions in the final version. Although the goal of  
 2 this paper is theoretical, we perform experiments to resolve reviewers’ concern about practicality of our methods.

3 **Experiment Setup.** We compare the proposed algorithms with baseline algorithms on the U.S. 2000 Census Data  
 4 containing  $n = 5 \times 10^6$  rows and  $d = 11$  columns and UCI YearPredictionMSD dataset which has  $n = 515,345$  rows  
 5 and  $d = 90$  columns. All algorithms are implemented in Python 3.7. To solve the optimization problems induced by  
 6 the regression problems and their sketched versions, we invoke the `minimize` function in `scipy.optimize`. Each  
 7 experiment is repeated for *25 times*, and the mean of the loss function value is reported. In all experiments, we vary the  
 8 sampling size or embedding dimension from  $5d$  to  $20d$ , and observe their effects on the quality of approximation.

9 **Experiments on Orlicz norm.** We compare our algorithm in Section 2 with uniform sampling and the embedding  
 10 in [2]. We also calculate the optimal solution to verify the approximation ratio. We try Orlicz norms induced by  
 11 two different  $G$  functions: Huber with  $c = 0.1$  and “ $\ell_1 - \ell_2$ ”. See Table 1 in our submission for definitions. Our  
 12 experimental results given below clearly demonstrate the practicality of our algorithm. In both datasets, our algorithm  
 outperforms both baseline algorithms by a significant margin, and achieves the best accuracy in almost all settings.



13 **Experiments on symmetric norm.** We compare our algorithm in Section 3 (SymSketch) with the optimal solution to  
 14 verify the approximation ratio. We try two different symmetric norms: top- $k$  norm with  $k = n/5$  and sum-mix of  $\ell_1$  and  
 15  $\ell_2$  norm ( $\|x\|_1 + \|x\|_2$ ). See Line 58-60 in our submission for definitions of these norms. As shown below, SymSketch  
 16 achieves reasonable approximation ratios with moderate embedding dimension. In particular, the algorithm achieves an  
 17 approximation ratio of 1.25 when the embedding dimension is only  $5d$ .



18 **(Reviewer #1) Assumption 1.** Our sampling algorithm in fact works for  $\ell_p$  norms when  $p > 2$ . In general, suppose  
 19 the function  $G : \mathbb{R} \rightarrow \mathbb{R}$  satisfies that for all  $0 < x < y$ ,  $G(y)/G(x) \leq C_G(y/x)^p$ , for the Orlicz norm induced by  
 20  $G$ , given a well-conditioned basis with condition number  $\kappa_G$ , our sampling algorithm returns a matrix with roughly  
 21  $O((\sqrt{d}\kappa_G)^p \cdot d/\epsilon^2)$  rows such that Theorem 1 holds. However it is not clear how to calculate well-conditioned bases  
 22 in input-sparsity time when  $p > 2$ . Our current method fails since it requires an oblivious subspace with  $\text{poly}(d)$   
 23 distortion, and it is known that such embedding does not exist when  $p > 2$  [9]. Since we focus on input-sparsity time  
 24 algorithms in this paper, we did not include the  $p > 2$  case. We will add more discussion on this in the final version.

25 **(Reviewer #2) Results on symmetric norm.** We disagree that this is an incremental improvement. First, the previous  
 26 embedding with  $d \log n$  distortion only works for Orlicz norms, and in this paper we give the first subspace embedding  
 27 for general symmetric norms. Second, the construction in [7] is only for streaming algorithms. To construct a subspace  
 28 embedding, we need to show that (i) norms of all vectors in a subspace are preserved and (ii) there is a simple estimator  
 29 in the sketch space. Neither of them can be satisfied by the construction in [7].

30 **Comparison with [11].** First, our definitions for Orlicz norm leverage score and well-conditioned basis, as given in  
 31 Definition 2 and 3, are different from all previous works and are closely related to the Orlicz norm under consideration.  
 32 The algorithm in [11], on the other hand, simply uses  $\ell_p$  leverage scores. Under our definition, we can prove that the  
 33 sum of leverage scores is bounded by  $O(C_G d \kappa_G^2)$  (Lemma 4), whose proof requires a novel probabilistic argument. In  
 34 contrast, the upper bound on sum of leverage scores in [11] is  $O(\sqrt{nd})$  (Lemma 38 in [11]). Thus, the algorithm in  
 35 [11] runs in an iterative manner since in each round the algorithm can merely reduce the dimension from  $n$  to  $O(\sqrt{nd})$ ,  
 36 while our algorithm is one-shot. We will of course add a more detailed comparison with [11] in the final version.

37 **(Reviewer #3)** The uniqueness of  $\alpha$  follows from the assumption that  $G$  is strictly increasing, in which case the two  
 38 definitions are equivalent. The assumption that  $G$  is strictly increasing was also implicitly made in Andoni et al. [2].  
 39 It is indeed an interesting problem to generalize our techniques to other problems, e.g., classification problems and  
 40 non-linear regression problems. We leave this as a future work.