

1 We are grateful to the reviewers for their thoughtful comments which have improved the work. Below are responses to
 2 each and every point of the reviewers (where we refer to the same references in the reviewer 2’s comments):

3 **Reviewer 1:** We thank the reviewer for the positive comments. **1a.** The dependence of the constant C in Theorem 2
 4 on the parameters (such as α) is complicated due to the delicate nature of our chaining method, so we did not write
 5 down the constant explicitly. We will carefully address which parameters that the constant C depends on. **1b.** We
 6 appreciate the reviewer’s excellent suggestions on the writing structure, we will rewrite lines 29-30 and move lines
 7 52-55 (the motivation of introducing distributional robustness) to follow the formal introduction of our model. **1c.**
 8 We thank the reviewer for pointing out the typos in our paper. The \inf in line 45 refers to the infimum of a set, so
 9 we don’t have to add a subscript for it. We will follow the reviewer’s advice to correct the other typos. **2a.** Although
 10 the Weibull-type condition is already a significant generalization of related assumptions which are standard in convex
 11 regression formulations, we agree that weakening to finite moments assumption is a very interesting direction to explore
 12 in future work. We also thank the reviewer for bringing up the ideas of robust statistics. In our problem, we care about
 13 the shape of the estimator, it is not clear how estimators such as Huber-type M-estimators can be modified to preserve
 14 the convex shape, but this is also an interesting direction to explore. **2b.** Please see our response to **4** of Reviewer 2.

15 **Reviewer 2:** We thank the reviewer for the insightful comments. **1.** We appreciate the reviewer for bringing [2] to us,
 16 and we agree that it is important to illustrate the similarities and differences between [2, Theorem 1] and our result. The
 17 main difference is, the loss function in [2, Theorem 1] is required to be a composition of a Lipschitz function $l : \mathbb{R} \rightarrow \mathbb{R}$
 18 and a linear function of x, y , due to the non-linear nature of our convex regression model, this result does not apply to our
 19 case directly. Besides, our approach also allows l be function maps \mathbb{R}^2 to \mathbb{R} . However, we agree that the proof of the two
 20 results share a similar idea: we apply the dual representation of the standard Wasserstein-based distributionally robust
 21 optimization problem, and then lower bound the dual parameter λ using the structure of the loss function l . Moreover,
 22 we will follow the reviewer’s advice to cite other results similar to Lemma 1. **2.** The new interpretation of equation (2)
 23 that relates our problem to covariate shifting is very interesting, we will add this point to our submission. **3.** We thank
 24 the reviewer for giving us the chance to illustrate our contribution beyond the results in [4]. First of all, their results only
 25 focus on problems with **finite** dimensional decision space (or feasible set \mathbb{X} , in the words of [4]), while in our case the
 26 decision space \mathcal{F}_n is clearly **infinite** dimensional. Secondly, to the best of our understanding, the concentration result
 27 ([4, Theorem 3.4], or [1, 6]) allows us to choose a proper δ_n such that the Wasserstein ball centered at P_n (the empirical
 28 measure) with radius δ_n covers the true underlying probability measure P (see [4, Theorem 3.5]) with high probability,
 29 and as a consequence we can appropriately choose δ_n to ensure the **consistency** of the estimator of the optimal decision
 30 variable (see [4, Theorem 3.6(2)]). However, it is non-trivial to identify the convergence rate of the estimator. For
 31 example, for finite dimensional decisions, such convergence rates should match the canonical rate $O(n^{-1/2})$. The choice
 32 of δ_n suggested in [4] does not provide the canonical rate, so it doesn’t seem that the results in [4] are directly applicable
 33 to recover convergence rates for estimators; not even in the finite dimensional case which is the environment of [4], let
 34 alone the infinite dimensional case, which is the setting of our paper. **4.** We thank the reviewer for raising this issue, and
 35 we are agree with the reviewer’s opinion from the optimization point of view. However, there is a technical statistical
 36 reason behind our formulation. On the one hand, the penalty term $\delta_n \|\nabla f\|_\infty$ captures the impact of the uncertainty
 37 set $\{P : D(P, P_n) \leq \delta_n\}$. On the other hand, the $\log n$ constraint introduced in \mathcal{F}_n is related to a compactification
 38 argument applied to the decision space. This issue is particularly important in the current setting of infinite dimensional
 39 decisions. The current formulation provides one possible tradeoff of these two effects that guarantee the estimator \hat{f}_n
 40 converges to f_* with order $\tilde{O}(n^{-1/d})$. Furthermore, the $\log n$ constraint already relaxes the typical assumption that
 41 $\|\nabla f\|_\infty < C$ (in which C needs to be known apriori). Finally, we will provide a formal proof of the optimality of
 42 piecewise affine functions. **5.** We appreciate the reviewer’s suggestions of adding experiments on real world data. We
 43 consider a public dataset from United States Environmental Protection Agency, which was suggested by [R. Mazumder,
 44 A. Choudhury, G. Iyengar, B. Sen, A Computational Framework for Multivariate Convex Regression and its Variants].
 45 The dataset consists of 600 air market data of California in the first quarter
 46 of 2019. The response was the amount of heat input with the covariates
 47 corresponding to the amounts of emissions of SO2, NOx, CO2 (in tons) and the
 48 NOX rate. Empirical evidence suggests that relationship between the response
 49 and the log transformation of each individual covariate can be modeled well by
 50 a convex fit, so we do the log transformation on covariates and then standardize
 51 the data. Since we never know f^* in real data, we can not evaluate our
 52 method in the same way as the submitted paper. Instead, we randomly split
 53 the dataset into a training set with 400 data and a test set with 200 data, and
 54 we implement three different approaches: DRCR (our estimator), LSE (standard convex regression estimator) and LR
 55 (linear regression). We repeat the experiment 10 times and then compare the average training l_1 loss and average test l_1
 56 error. We summarize the results in the table on the right, it is clear that our method outperforms both LSE and LR.

Method	Training loss	Test error
DRCR	0.1238	0.1294
LSE	0.1485	0.1516
LR	0.1691	0.1692

57 **Reviewer 3:** We thank the reviewer for the positive comments.