

1 We thank all reviewers for their thoughtful comments.

2 **Experimental setup (R1, R3):** The reviewers brought up an excellent suggestion to deploy our algorithm in production
3 on real users, for example using a user-interactive system. We agree this is an appropriate next step for development.
4 However, applying contextual bandits to continuous-action domains is a relatively new approach; the large-scale
5 contextual bandit datasets for the examples in the introduction have not been collected yet.

6 Simulating contextual bandit learning using regression/classification datasets is a standard protocol for academic
7 evaluation of bandit algorithms [see e.g. A. Bietti et al. A Contextual Bandit Bake-off]. Note that we have chosen
8 the regression datasets with the criterion of having million scale examples with unique regression values, so it is very
9 relevant to a real contextual bandit setting.

10 **Approximate ERM guarantee of Train_tree (R2, R4):** At a high level, our analysis of the Train_tree algorithm
11 needs to account for compounding errors accumulated in each of the nodes, which results in a $O(K\sqrt{1/n})$ rate slower
12 than the $O(\sqrt{K/n})$ rate achieved by ERM. We will clarify this in the final version.

13 **Realizability assumption (R3, R4):** Indeed we assume realizability in Theorems 6 and 7. This is a typical situation in
14 ML theory when the theory does not necessarily capture *all* cases when things work, but instead gives justification and
15 assurance that they do work in practice.

16 Realizability has been widely used in many successful contextual bandit algorithms (e.g. LinUCB, OFUL).

17 Regarding checking realizability: This argument applies to essentially all supervised learning settings, where standard
18 surrogate loss minimization can fail without realizability or related unverifiable assumptions.

19 We now answer specific questions raised by individual reviewers:

20 **R1:** We agree that in our online contextual bandit learning setting, tuning hyperparameters h and K can be problematic.
21 Note that all other hyperparameters such as learning rate, the greedy parameter, and the penalty term in line 5
22 of Algorithm 3 (CATS_Off) are fixed in our experiments and for all of the datasets. For practical applications, we
23 recommend using CATS_Off to perform off-policy optimization to find the best (h, K) combination, whose effectiveness
24 is demonstrated in Figure 1 (right) in the paper. Note that in CATS_Off, we only use input logged data to select the best
25 h, K hyperparameters.

26 Regarding the choice of h and K in bias-variance trade-off: Tuning K has the usual bias-complexity trade-off in
27 standard statistical learning. If we are competing with the best (non-smoothed) policy in a class, then tuning h also
28 incurs a bias-variance tradeoff: (1) the smoothed loss more closely approximates the true loss for smaller h ; (2) a
29 smaller h makes achieving a low h -smoothed regret bound harder.

30 **R2:** The reviewer's understanding regarding the dTree baseline and the filter tree algorithm is correct. We will clarify
31 this in the final version. We will elaborate more on line 208-213 in the final version.

32 **R3:** The reviewer mentioned our accuracy is similar to the dLinear baseline. This is true, but our algorithm beats
33 dLinear in terms of time complexity. As R4 has acknowledged, our algorithm wins over dLinear computationally
34 and wins over dTree statistically.

35 The reviewer has suggested three papers to compare our algorithm with. Note that [1] was published *after* the submission
36 deadline so we could not have possibly compared to it. [2] and [3] are also not directly relevant to our work since (1)
37 they impose structure on contexts, rather than actions, and moreover the structure is *linear* rather than Lipschitz; (2)
38 they work with discrete action spaces. We thank the reviewer for bringing these papers to our attention and will cite
39 them in our final draft.

40 We will polish the writing on the sections that the reviewer highlighted in the final version.

41 **R4:** The idea of smoothing was initially proposed in [35], but their proposed algorithm is computationally intractable.
42 Our contribution is proposing a computationally efficient algorithm for continuous action spaces in the smoothing
43 framework. Our online algorithm CATS, its off-policy optimization version CATS_Off, and their statistical guarantees
44 as well as our careful implementation design with $O(\log K)$ are all new and constitute the contributions of this work.

45 See Figure 1 of [35] for an example of a discontinuous loss function, where the best single action is impossible to find
46 within a finite number of rounds of contextual bandit learning.

47 Regarding the reason we need to guess width but the location can be automatically adjusted: If the learner sets a high
48 discretization level K then it can take actions in $[0,1]$ with fine granularity.