

1 We would like to thank all the reviewers for the detailed and insightful comments.

2 **Reviewer 1**

3 *Regarding applying our technique to other methods:* The clipping trick and much of the subsequent analysis can be
4 used to analyze other methods, including UCB-VI, UBEV, and other model-based methods. In fact, Simchowitz and
5 Jamieson (2019) use the clipping trick to analyze model-based algorithms in the tabular setting, not the model-free one
6 we analyze here. We do not include this in our paper because the model-based algorithms incur an S^2 dependence in
7 the “lower order” K -independent term. This degrades the rate substantially in the nonparametric case, since effectively
8 we set S in terms of the number of episodes K and the zooming dimension. So this term is not actually “lower order”
9 in the nonparametric case, which leads to a worse final bound.

10 *Regarding future work:* Based on the above remarks, an intriguing question is whether model-based methods can
11 achieve similar “zooming” guarantees via adaptive partitioning. A very recent paper of Sinclair et al. (<https://arxiv.org/abs/2007.00717>)
12 make some progress in this direction, but does not obtain zooming-dimension
13 dependence. We also mention two other questions: Can our techniques and adaptive guarantees extend to the infinite
14 horizon discounted setting? Can adaptive discretization be combined with function approximation, analogous to the
15 policy zooming approach of Krishnamurthy et al. (2019) for contextual bandits? This would allow us to use function
16 approximation to deal with large state spaces and adaptive discretization to deal with large action spaces, which may be
17 quite effective in practice. We can mention these directions in the final version.

18 **Reviewer 2.** Thanks for the kind words!

19 **Reviewer 3.** Thanks for catching the typos. We will fix them in the final version.

20 *Regarding linear methods:* This is a great question, but quite subtle! Remember, our method applies when Q^* is
21 Lipschitz continuous. This is of course implied by Q^* being linear in some known features (assuming boundedness), but
22 linearity is a much stronger assumption. A helpful analogy to think about is regression: if one knew the true regression
23 function were linear, one should not use a nonparametric method. But nonparametric methods are applicable much
24 more generally. It’s the standard estimation error/approximation error tradeoff.

25 However, the RL setting is more subtle. If Q^* is Lipschitz, then work from Lipschitz bandits suggests that our method
26 is optimal. If Q^* is linear (and that is all we assume), then our method achieves the $K^{\frac{d+1}{d+2}}$ regret rate, but actually we
27 do not know of any better guarantees for this setting. In particular, the recent $\text{poly}(d)\sqrt{K}$ results for linear function
28 approximation require much stronger assumptions, such as linear MDP or low inherent Bellman error. Q^* linear is a
29 much weaker assumption than e.g., linear MDP, and it is not clear whether $\text{poly}(d)\sqrt{K}$ rates are achievable here at all.

30 **Reviewer 4** Thanks for catching the typos. We will fix them in the final version.

31 *Regarding zooming dimension vs covering dimension:* We encourage the reviewer to examine the experiments of
32 Sinclair et al., (2019), which shows (1) that adaptive partitioning performs much better than uniform discretization and
33 (2) provides concrete examples and visualizations of this benefit. The zooming dimension gives a tighter bound for
34 Sinclair et al.’s adaptive partitioning scheme and captures the improvements observed in the experiments. Note that we
35 analyze exactly their algorithm but their bound does not capture these improvements.

36 In addition, (1) the example we give in Figure 1 of our paper is actually quite general, as any problem where the
37 near-optimal actions concentrate onto a low dimensional manifold has smaller zooming dimension than covering
38 dimension; (2) Theorem 1 is even more refined than the zooming dimension bound; (3) zooming dimension and
39 related quantities like “near optimality dimension” are widely studied in the bandit literature (20+ papers in the top ML
40 conferences) and it is natural to extend these notions to RL.

41 *Regarding the clipping trick:* Please see Lemma 11 and 13 in the supplement for a mathematical statement. In words,
42 instead of trying to bound the number of balls (which would be quite large), we notice that the regret incurred is
43 determined by the sum of *clipped* surpluses for this ball, where we clip at level $\text{gap}/(H + 1)$. This is because we only
44 play a bad ball when there is a large error at a later time step, so we can credit this mistake to the later error and clip the
45 current surplus. Using this decomposition, clipping will quickly take effect for “bad” balls, which allows us to bound
46 the regret incurred by this ball.