

1 We thank all the reviewers for their valuable feedback and appreciating our contributions. First we would like to
2 emphasize the technical novelty of our upper bound and lower bound as Reviewer #1, Reviewer # 3 and Reviewer #4
3 commented on the technical novelty of our theoretical results.

4 **Technical novelty of the upper bound.** In the exploration phase, Jin et al. [2020] set reward to be 1 for significant
5 states and 0 for other states. Note their technique cannot be used in Linear MDPs because there are possibly infinitely
6 many states and thus one needs to take the structure of linear functions into account. In this paper, we use UCB bonus
7 *as the reward signal* in the exploration phase. To our knowledge, this idea is new in the literature. We also would like to
8 thank Reviewer # 4 for a detailed description of our main algorithmic ideas.

9 **Technical novelty of the lower bound.** We have discussed the differences between our lower bound and that in [Du
10 et al. 2020] in Line 276 - 279. We acknowledge that in our hard instance, we use a similar feature extractor as that in
11 [Du et al. 2020]. However, all other aspects of the hard instance construction are significantly different from that in [Du
12 et al. 2020]. For example, for the hard instance in [Du et al. 2020], only a single state-action pair has non-zero reward
13 value, which is not case in our hard instance. Note that such distinction is crucial, since in our hard instance the optimal
14 Q -function is *exactly* linear, whereas the the optimal Q -function is only *approximately* linear in the hard instance in [Du
15 et al. 2020]. Moreover, we focus on the reward-free setting while Du et al. [2020] focused on the standard RL setting.

16 Below we address specific concerns from each reviewer.

17 ——— **To Reviewer #1** ———

18 **Lack of rigor.** We have introduced necessary background on MDP in Section 2.1, including the state space, the action
19 space, the transition operator, the reward distribution, the Q -function, etc. We have also provided necessary definitions
20 related to linear function approximation in Section 2.2. Our descriptions mostly follow existing works. We will expand
21 this part to make the paper clearer.

22 **Q^* is linear on the suboptimal action.** In our construction, when defining the reward functions, we first define the
23 optimal Q -function (Q^*) as a specific linear function (see Line 292), and then define the reward values according to
24 the Bellman equations (see Line 296). Therefore, the optimal Q function must be linear for both optimal actions and
25 suboptimal actions in our hard instances.

26 **Relation to prior work.** We will discuss the suggested paper in the next version. Thanks for the suggestion.

27 ——— **To Reviewer #2** ———

28 **Extension to more general settings.** Even in the standard RL setting, going beyond linear MDPs is hard. See the open
29 problems in [Du et al. 2020]. Therefore, we believe it is highly non-trivial to obtain more general results.

30 ——— **To Reviewer #3** ———

31 **More emphasize on the lower bound.** Thanks for the suggestion. We will emphasize more on the lower bound and
32 the implied conceptual messages in the final version.

33 **Why do you need optimism in the planning phase.** Optimism in the planning phase is used when we prove Lemma
34 3.3. It also guarantees the correctness of the first inequality in Line 247-248.

35 ——— **To Reviewer #4** ———

36 We would like to thank the reviewer for the detailed description of our key ideas in our algorithm. The understanding is
37 correct.

38 **Experiments.** Thanks for the suggestion. We will consider adding empirical results in the next version.

39 **Related work.** Thanks for the references. We will add more discussion in the next version.

40 **Agent just gets samples from the reward function.** If we only have samples, we can change Line 6 in Algorithm 2 to
41 $w_h \leftarrow (\Lambda_h)^{-1} \sum_{\tau=1}^K \phi(s_h^\tau, a_h^\tau)(V_{h+1}(s_{h+1}^\tau) + r_h^\tau(s_h^\tau, a_h^\tau))$ where $r_h^\tau(s_h^\tau, a_h^\tau)$ is the *sampled* reward value, and remove
42 $r_h(\cdot, \cdot)$ from Line 7. Our theoretical results still hold after this modification, and we will add a discussion on this.

43 **Linearity approximately holds.** This is an interesting question and we will list it as a future direction.

44 **Line 182-183.** This is correct.

45 **The effect of increasing/decreasing c_β .** c_β needs to be larger than a universal constant in order to guarantee optimism.
46 Once c_β is larger than that constant, the sample complexity decreases as c_β decreases.