

1 We thank all reviewers for appreciating our technical contributions: a novel algorithm and a sound theoretical analysis.  
2 Please find our response to each reviewer below.

3 **To Reviewer #1:**

4 *-Definition 3.1:* We will state in Definition 3.1 explicitly that we assume  $\gamma > 0$  to avoid confusion. Thanks for pointing  
5 out. We will also state that we assume  $\gamma > 0$  again before Theorem 6.1 to emphasize this assumption.

6 *-Assuming  $\gamma > 0$  is restrictive:* Assuming  $\gamma$  strictly positive is not a restrictive assumption for the *finite action setting*  
7 considered in this paper. First, in the contextual linear bandit literature, this assumption is widely discussed. See, e.g.,  
8 [1,2]. The notion, context, in the bandit literature is essentially  $\phi(s)$  in our paper and the number of contexts can also  
9 be infinite. Second, there are many natural environments in RL which satisfy this assumption. For example, in many  
10 environments, states can be classified as good states and bad states. In these environments, an agent can obtain a reward  
11 only if it is in a good state. There are also two kinds of actions: good actions and bad actions. If the agent is in a good  
12 state and chooses a good action, the agent will transit to a good state. If the agent chooses a bad action, the agent will  
13 transit to a bad state. If the agent is in a bad state, whatever action the agent chooses, the agent will transit to a bad state.  
14 Note that for this kind of environments, there is a strictly positive gap between good actions and bad actions when the  
15 agent is in good states and there is no difference between good actions and bad actions when the agent is in bad states.  
16 In this case,  $\gamma$  is strictly positive, since by Definition 3.1, we take the minimum over all state-action pairs with *strictly*  
17 positive gap. These environments are natural generalizations of the combination lock environment [3]. Some Atari  
18 games, e.g. Freeway, have a similar flavor as these environments.

19 *-What if  $\gamma$  is unknown:* While our algorithm requires  $\gamma$  as an input for deciding hyper-parameters, it is easy to see from  
20 the proof that as long as we choose some  $\gamma' \leq \gamma$  and choose hyper-parameters according to  $\gamma'$ , the sample complexity  
21 bound still holds. Therefore, we can just do a grid search over  $[0, 1]$  to find such  $\gamma'$ . We will add more discussion.

22 *-Section 5:* We will put an overview of our algorithm at the beginning of Section 5. The design of our algorithm is  
23 actually fairly natural as we explain below. Since we know  $Q^*$  belongs to a function class, the natural idea is to use  
24 regression to recover it (Line 9 of Algorithm 2), for which we need unbiased samples. To obtain unbiased samples for  
25 level  $h$ , we need to make sure that we execute the optimal policy for level  $h' = h + 1, \dots, H$  (this accounts for our  
26 definition in Eq. 1). Therefore, in Line 1-2 of Algorithm 1, we learn  $Q^*$  from the last level to the first level (learning  
27 at the last level is equivalent to the stochastic bandit problem, where we have unbiased samples). The tricky part is  
28 that the  $Q$ -function we learned, say at level  $h' = h + 1, \dots, H$ , is only guaranteed to be optimal with respect to the  
29 roll-in policy we used to gather samples. When we learn  $Q^*$  at level  $h$ , the roll-in policy for level  $h' = h + 1, \dots, H$   
30 can change. To address this problem, we need to check that the learned  $Q$ -function is still optimal for the new roll-in  
31 policy (Algorithm 3 and the DSEC oracle). We have explained the intuition of the DSEC oracle in Section 1.

32 *-Abstract / introduction:* We follow one of the common writing styles that first states our precise result and gives a  
33 description of our main techniques in the abstract, and then puts an informal main theorem and a technical overview in  
34 the introduction. We are happy to adjust our writing style.

35 *-Regularizer in Oracle 4.1 / minor concerns / typos:* We will fix them accordingly. Thanks for pointing out.

36 **To Reviewer #2:**

37 *- $Q$ -function is not exactly linear:* This is a great question. Our algorithm can still be applied if the best linear predictor  
38 has approximation error smaller than  $\frac{\gamma}{\text{poly}(H,d)}$ . The proof is essentially the same. We will add more discussion.

39 *-Most state-action pairs have a large gap:* This is an interesting question. We do believe the sample complexity can be  
40 improved with proper quantification on “most state-action pairs”. We will discuss this in the final version.

41 *-Need to collect on-the-go reward:* We agree this incurs higher computation complexity. We will list improving the  
42 computation complexity as an open problem.

43 *-Proof of Theorem 6.1:* In Line 401 and 402 we use the induction hypothesis of later levels. We will add more details.

44 *-Notation  $\gamma$ :* Thanks for pointing out. We will change the notation to  $\Delta$ .

45 **To Reviewer #3:**

46 *-More general setting / generalization of Lemma 4:* Thanks for the question! We have discussed going beyond linear  
47 function approximation in Section 7. A natural extension of Lemma 4 is to study the Fisher information matrix, though  
48 this generalization is not trivial. Once we know the number of calls to the DSEC oracle is bounded (via a generalization  
49 of Lemma 4), we can then reuse the current proof framework.

50 *-Optimality / connections with stochastic bandits:* We will discuss more about connections with stochastic bandits.  
51 Thanks for the suggestion. While Line 6-8 is not optimal for the stochastic bandit case, it not obvious how to improve  
52 that part since Algorithm 3 also depends on the data collection policy. We will add more discussion.

53 [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits.

54 [2] Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit feedback.

55 [3] Sham M Kakade. On the sample complexity of reinforcement learning.