

1 Many thanks for your valuable comments and suggestions. We will address the minor comments (such as to fix the  
2 typos, add the discussion about the parameters, and broaden the related works section) in the next version of the paper.  
3 We now respond to the major comments are as follows.

4 **To Reviewer 2:** We believe that the proposed techniques could be helpful in other RL settings such as continuous  
5 state spaces. Take RL with the linear model as an example. In this case, to learn  $Q^*$ , it suffices to learn the  $d$ -  
6 dimensional vector  $w^*$  (which we name as pseudo value function) because  $Q_h^*(s, a) = \phi(s, a)^\top w_h^*$ , where  $\phi(s, a)$  is  
7 the  $d$ -dimensional feature vector of  $(s, a)$ . To reduce the estimation variance via reference-advantage decomposition,  
8 we first aim at learning a rough estimation of  $w^*$  as the reference pseudo value function, using time steps that are  
9 only polynomially depending on  $A, H, d, 1/\epsilon, \ln(1/p), \ln T$ , where  $p$  is the failure probability, and  $\epsilon$  is the accuracy  
10 parameter. More formally, we believe that the key is to prove an analogue of Lemma 5 for the linear model. We will  
11 add this discussion in the next version of the paper.

12 We will also discuss the work on policy certificates (Dann et al., 2019) in related work section.

13

14 **To Reviewer 3:** At line 397, we mentioned that  $[c_1, c_2, c_3] = [2, 2, 5]$  is sufficient. The safe choice in the current  
15 version of the paper for  $c_4$  is  $c_4 = 57603$ , which seems quite large because we bound  $(1 + \frac{1}{H})$  by 2 in many places  
16 for convenience. When  $H$  is sufficiently large (say greater or equal to 200), we can bring the choice of  $c_4$  to around  
17 200. We will add this discussion at the next version of the paper.

18

19 **To Reviewer 4:** Thank you for your suggestion. We will include the table for all notations in the next version of the  
20 paper.