

1 Thanks for the reviews! Response to individual reviewers are below.

2 **R1.** You are right that the contribution is theoretical, but our formulation of the representation learning problem and the  
3 computational oracles assumed are heavily motivated by practice. The FLAMBE algorithm is not difficult to implement  
4 with modern deep learning libraries. In addition, modern RL algorithms can be employed, in the planning step (Eq.  
5 (2) in Algorithm 2, see e.g. the recent preprint <https://arxiv.org/abs/2007.08459>), as well as in the maximum  
6 likelihood step (e.g., using VAEs as in recent model-based RL works such as <https://arxiv.org/abs/2005.05960>).  
7 To summarize, our theory directly motivates new strategies for representation learning and exploration that we believe  
8 will be empirically effective, and we are excited about experimenting with these approaches in future work.

9 **R2.** We believe the algorithm is actually quite precisely written and intuitive. Perhaps the simpler algorithm (and  
10 simpler analysis) for the simplex features setting provided in the appendix is useful for added intuition. The algorithm  
11 at a high-level can be seen as doing MLE to estimate features, and using an algorithm for low-rank MDPs with known  
12 features to compute a good policy. Inducting these two at each step  $h = 0, \dots, H$  with some careful attention to details  
13 gives the FLAMBE algorithm. For specifics:

- 14 • We believe you may have missed the role of the exploratory policy  $\rho_h$  which is computed in iteration  $h - 1$  of the  
15 algorithm (via a call to Algorithm 2) and is used to collect the data for the MLE step in iteration  $h$ . Concretely, our  
16 model-based planning step (call to Alg 2) returns an exploratory policy  $\rho_h^{\text{pre}}$ , which is executed for  $h - 1$  steps,  
17 followed by *two* uniform actions at steps  $h$  and  $h + 1$  to collect the samples for MLE estimation at level  $h + 1$ . We  
18 need two random actions as planning stays one step behind model learning, as briefly discussed in lines 291–295.  
19 We *do not* assume access to a generative model; we are in the online exploration setting and collect data by the  
20 learned exploratory policies  $\rho_h$ . The use of such exploratory policies to cover states well is common to many  
21 provable methods (e.g., Du et al, 2019b, Misra et al., 2019).
- 22 • Loop over episodes missing – The only interaction with the environment happens in the line “Collect n triples ...”  
23 in Algorithm 1. To collect one triple, we execute a full episode, so the total number of episodes is  $nH$ .
- 24 • Objective for elliptical potential – This quadratic objective originated in the linear bandits literature and is quite  
25 standard in linear RL problems with known features, e.g., it appears in the analysis of Jin et al.
- 26 • The sampling oracle is used in Algorithm 2, to optimize Eq (2) in the learned model. This can be done with  
27 dynamic programming-type algorithms, such as Optimistic LSVI, which is quite straightforward and described in  
28 Appendix C, Lemma 5. We *do not* say that Algorithm 2 corresponds to optimistic LSVI.

29 Thanks for the references. We agree the count-based exploration work is relevant. Note that the other works mentioned  
30 do not consider representation learning in the context of exploration, which is our focus, and so they are less relevant.  
31 We will add more discussion to address some of the confusions above,

32 **R3.** Assumption 1 in block MDPs – You are right, the realizability assumptions in our work and block MDP results  
33 are not equivalent, but there is some subtlety here. We consider “model-based realizability” while a weaker notion may  
34 be simply that  $\phi^* \in \Phi$ , with no assumption on  $\mu^*$ . The realizability assumption for block MDP results is between  
35 these two: Du et al.’s and Misra et al.’s assumption is equivalent to realizability of  $\phi^*$  and *the support* of  $\mu^*$  (since the  
36 next state is also decodable). We do not see a natural analog to this intermediate assumption for the general low rank  
37 setting, but we agree that considering the weaker assumption (only  $\phi^* \in \Phi$ ) is a nice question for future work. Note  
38 also that when we discuss expressiveness, we are focusing only on the dynamics assumptions, and not on realizability  
39 requirements. We can expand on this in the final version.

40 We can add some discussions around practical issues, comparison to other works, and on limitations of low rank MDPs.

41 **R5.** Thanks for the pointers to bisimulation! We can discuss more in the final version. Briefly, our learned representa-  
42 tions are related to a state abstraction notion, called Kinematic Inseparability (from Misra et al. 2019), which is finer  
43 than bisimulation, but remains meaningful in the reward-free setting. Note that in the absence of rewards, every MDP  
44 admits a trivial bisimulation that aggregates all states together. Indeed, it is not possible to learn a bisimulation with  
45 polynomial sample complexity while exploring in a sparse reward problem, as formally proved in Modi et al. 2020  
46 (Proposition B.1). Thus our abstraction is less coarse, but learnable in the exploration setting.

47 Regarding reward-free exploration: You are right that system identification may be overkill for easy RL problems (e.g.,  
48 dense rewards, local exploration suffices, etc.). However, we are interested in *provable* sample efficiency, which means  
49 we must consider *hard* RL problems. In this case, the distinction between reward-free and reward-sensitive learning is  
50 less significant, since even in the reward-sensitive setting the agent must explore the entire environment so it can certify  
51 that there is no “hidden rewards” anywhere.

52 You are also hinting at a separation between model-based and value/policy-based algorithms. This seems plausible, and  
53 we agree that developing provable model-free algorithms for low rank MDPs is an exciting direction for future work.