

1 We thank all reviewers for their valuable time and feedback. We believe that we have addressed the main concerns
2 about the assumptions, related work, and evaluation.

3 CONTENT

4 **R2**, learning the features. The linear MDP assumption $\Pi\Phi = \Phi M$ and $\mathbf{r} = \Phi w$ implies that the transition matrix Π is
5 low-rank and that rewards are linear in the features, so these assumptions can be used to guide feature learning. For
6 example, for discrete actions we can set $\phi(s, a) = \phi(s) \otimes \delta(a)$, and make $\phi(s|\theta)$ the output of a neural network. The
7 objective would be to minimize $L(\theta, K, w) := E_{(s,a,r,s') \sim \mathcal{D}} [c_1 \|\phi(s'|\theta) - K\phi(s, a|\theta)\|^2 + c_2 (r - \phi(s, a|\theta)^\top w)^2]$.
8 Note that multiple recent works (offline and online) simply assume a linear MDP with known features in analysis.

9 **R2**, violation of assumptions. In terms of exploration, we do not require the logging policy to cover all states, but just to
10 span the features (feature covariance matrix should be full-rank, assumption A3). If this is not the case, the error would
11 depend on the missing subspace. In terms of linearity, for non-linear MDPs our method would incur an approximation
12 error. Separating the linearity of dynamics and rewards (rather than implicitly assuming both with a linear Q-function)
13 at least allows us to drop the linear-reward assumption if we learn the full distribution (simulated in Figure 1, middle).

14 **R3, R4**, necessity of entropy. Our finite-sample guarantees for J_π actually hold for any distribution satisfying the
15 constraint (see Remark 1) as long as the MDP is linear, and we can also use the simple estimate in Equation (10).
16 Empirically, we find entropy to be a good regularizer when learning the full distribution $d_\pi(s, a) = \mu_\pi(s) \otimes \pi(a|s)$.
17 Learning d_π is useful in the case of non-linear rewards, or $E_d[r(s, a)]$ cannot be computed in closed form. Appendix D
18 gives a justification of the maximum-entropy objective for MDPs with sufficiently random dynamics. More generally,
19 maximizing entropy is equivalent to minimizing KL-divergence to the uniform distribution. We can use the same
20 KL-divergence formulation to impose different distribution priors when available.

21 **R3**, comparison to bounds in Duan and Wang (2020). These bounds are expressed in terms of χ^2 divergence, and for
22 linear f , they are a function of the spectrum of the feature covariance matrix Σ . This also the case with our bound - it
23 scales inversely with the smallest eigenvalue of Σ . The results are admittedly different otherwise (we will clarify this in
24 the paper). Note also that the bounds of Duan and Wang (2020) scale with the horizon / effective horizon, and are thus
25 infinite in our setting. Furthermore, their discounted infinite-horizon bound scales as $N^{-1/2}$ where N is the number of
26 trajectories, whereas our bound scales as $T^{-1/2}$ where T is the number of transitions (possibly in a single trajectory).

27 EXPERIMENTS

28 **R1, R3**, we agree the the methods are not always well-separated in the sense of 95% confidence intervals, and BRM and
29 FQI can perform well. However, note that some of the evaluated environments violate our assumptions (Taxi, Acrobot),
30 and on these our method performs at least as well as the baselines. The only evaluated environment satisfying linearity
31 and ergodicity is LQ control (Figure 2 left), where our approach is clearly better. Concretely, the model-based method
32 is clearly not sensitive to ε_1 (controlling level of suboptimality of the target policy), while BRM and FQI deteriorate
33 with larger ε_1 , and the curves are well-separated for $\varepsilon_1 \geq 0.2$.

34 **R2**, for Acrobot, it is an interesting observation that MaxEnt and Model often underestimate the expected reward
35 (actually so does FQI). However, for smaller true J_π , the same methods overestimate. This may be an artifact of
36 covariance regularization for linear regression in the model / Q-function.

37 **R3**, we chose not to explicitly demonstrate divergence of FQI and biasedness of BRM since these issues are well
38 known (see references in Section 4).

39 **STYLE. R2**, thank you very much for the suggestions. We agree about Section 3.3 and in retrospect should have saved
40 a discussion of policy improvement for future work. We will remove it and move some of the Appendix into the paper.

41 RELATED WORK

42 **R2**, thank you for the suggestions. We will add references to maximum-entropy approaches in RL and IRL. In RL,
43 entropy is typically a regularizer on the policy, rather than the state-action distribution. IRL approaches are more similar
44 to ours, maximizing the entropy of distribution over paths s.t. feature expectations match demonstration data.

45 **R4**, thank you, we will improve the discussion. We actually do refer to the “Breaking the Curse of the Horizon” work
46 (lines 24,65, 204-208). Their formulation also relies on the relaxed feature expectation constraint. They minimize the
47 constraint violation while maximizing over features, and have a consistency guarantee. We assume oracle features and a
48 linear MDP, satisfy an approximate constraint (given by the model estimate), and show finite-sample guarantees.

49 **R3, R4**. We will add a reference to GenDICE. In the paper, we refer to “RL via Fenchel-Rockafellar Duality” by
50 Nachum and Dai (2020), which provides a unified view of the DICE papers, including GenDICE. Section 7 there
51 discusses the undiscounted setting via different DICE methods. We compare to the basic Lagrangian formulation in our
52 paper (lines 195-205), and will add a note on GenDICE (regularized Lagrangian).