
Provable Benefits of Representational Transfer in Reinforcement Learning

Alekh Agarwal, Yuda Song, Wen Sun, Kaiwen Wang, Mengdi Wang, Xuezhou Zhang*

Abstract

We study the problem of representational transfer in RL, where an agent first pretrains in a number of source tasks to discover a shared representation, which is subsequently used to learn a good policy in a target task. We propose a new notion of task relatedness between source and target tasks, and develop a novel approach for representational transfer under this assumption. Concretely, we show that given a generative access to source tasks, we can discover a representation, using which subsequent linear RL techniques quickly converge to a near-optimal policy, with only online access to the target task. The sample complexity is close to knowing the ground truth features in the target task, and comparable to prior representation learning results in the source tasks. We complement our positive results with lower bounds without generative access, and validate our findings with empirical evaluation on rich observation MDPs that require deep exploration.

1 Introduction

Leveraging historical experiences acquired in learning past skills to accelerate the learning of a new skill is a hallmark of intelligent behavior. In this paper, we study this question in the context of reinforcement learning (RL). Specifically, we consider a setting where the learner is exposed to multiple tasks and ask the following question:

Can we accelerate RL by sharing representations across multiple related tasks?

There is rich empirical literature which studies multiple approaches to this question and various paradigms for instantiating it. For instance, in a multi-task learning scenario, the learner has simultaneous access to different tasks and tries to improve the sample complexity by sharing data across them (Caruana, 1997). Other works study a transfer learning setting, where the learner has access to multiple source tasks during a *pre-training* phase, followed by a target task (Pan and Yang, 2009). The goal is to learn features and/or a policy which can be quickly adapted to succeed in the target task. More generally, the paradigms of meta-learning (Finn et al., 2017), lifelong learning (Parisi et al., 2019) and curriculum learning (Bengio et al., 2009) also consider related questions.

On the theoretical side, questions of representation learning have received an increased recent emphasis owing to their practical significance, both in supervised learning and RL settings. In RL, a limited form of transfer learning across multiple downstream reward functions is enabled by several recent reward-free representation learning approaches (Jin et al., 2020a; Zhang et al., 2020; Wang et al., 2020; Du et al., 2019; Misra et al., 2020; Agarwal et al., 2020; Modi et al., 2021). Inspired by recent treatments of representation transfer in supervised (Maurer et al., 2016; Du et al., 2020) and imitation learning (Arora et al., 2020), some works also study more general task collections in bandits (Hu et al., 2021; Yang et al., 2020, 2022) and RL (Hu et al., 2021; Lu et al., 2021). Almost all these works study settings where the representation is *frozen* after pre-training in the source tasks, and a linear policy or optimal value function approximation is trained in the target task using these learned features. This setting, which we call *representational transfer*, is the main focus of our paper.

A crucial question in formalizing representational transfer settings is the notion of similarity between source and target tasks. Prior works in supervised learning make the stringent assumption that the covariates x follow the same underlying distribution in all the tasks, and only the conditional $P(y|x)$ can vary across tasks (Du et al., 2020). This assumption does not nicely generalize to RL settings, where state distributions are typically policy dependent, and prior attempts to extend this assumption to RL (Lu et al., 2021) result in strong assumptions on the learning setup. Other works (Hu et al.,

*Alphabetical Order. alekhagarwal@google.com, yudas@andrew.cmu.edu, {ws455,kw437}@cornell.edu, {mengdiw,xz7392}@princeton.edu.

2021; Yang et al., 2020, 2022) focus on linear representations only, which limits the expressivity of the feature maps, and does not adequately represent the empirical literature in the field.

Our contributions. In this context, our work makes the following contributions:

- We propose a new linear span assumption of task relatedness for representational transfer, where the target task dynamics can be expressed as a (state-dependent) linear span of source task dynamics, in addition to the dynamics being low-rank under a shared representation. We give examples captured by this assumption, and it generalizes all prior settings for representational transfer in RL. We do not make any linearity assumptions on our feature maps.
- When we have *generative access* to source tasks, we provide a novel algorithm REPTRANSFER that successfully pretrains a representation for downstream *online learning in any target task* (i.e., no generative access in target task) satisfying the linear span assumption, when the source tasks satisfy a common latent reachability assumption. The regret bound of learning in the target task is close to that of learning in a linear MDP equipped with the ground truth features, the strongest possible yardstick in our setup. The additional terms in our regret largely arise out of the distributional mismatch between source and target tasks, which is expected. We complement the theory with an empirical validation of REPTRANSFER on the challenging rich observation combination lock benchmarks (Misra et al., 2020), confirming our theoretical findings.
- Without generative access to source tasks, we show the statistical hardness of representational transfer under the linear span assumption, and confirm this hardness in our empirical evaluation. We show that an additional assumption that every *observed state* is reachable in every source task is sufficient for allowing fully online learning in source tasks.

The new task relatedness assumption, reward-free learning result for low-rank MDPs and our analysis of LSVI-UCB under average case misspecification may be of independent interest.

2 Related Work

Multi-task and Transfer Learning in Supervised Learning. The theoretical benefit of representation learning are well studied under conditions such as the i.i.d. task assumption (Maurer et al., 2016) and the diversity assumption (Du et al., 2020; Tripuraneni et al., 2020). Many works below successfully adopt the frameworks and assumptions to sequential decision making problems.

Multi-task and Transfer Learning in Bandit and small-size MDPs. A sequence of recent works study multi-task linear bandit problems with linear representation functions ($\phi(s) = A s$) (Hu et al., 2021; Yang et al., 2020, 2022). The techniques developed in these works crucially rely on the linear representation function structure and can not be applied to nonlinear function classes. Lazaric et al. (2013) study spectral techniques for online sequential transfer learning. Brunskill and Li (2013) consider multi-task RL under a fixed distribution over finitely many diverse MDPs, while Brunskill and Li (2014) considers transfer learning in semi-MDPs by learning options (temporally extended actions). All these works consider small size tabular models while we focus on large-scale MDPs.

Multi-task and Transfer Learning in RL via representation learning. Beyond tabular MDPs, Arora et al. (2020) and D’Eramo et al. (2019) show benefits of representation learning in imitation learning and planning, but do not address exploration. Lu et al. (2021) study transfer learning in low-rank MDPs with general nonlinear representations, but make a generative model assumption on both the source tasks and the target task, along with other distributional and structural assumptions. We do not require generative access to the target task and make much weaker structural assumptions on the source-target relatedness. Recently and independently, Cheng et al. (2022) also studied transfer learning in low-rank MDPs in the online learning setting, identical to the setting we study in Section 5. However, their analysis relies on an additional assumption that bounds the point-wise TV error with the population TV error, which we show is in fact not necessary. We refer readers to Appendix B for a more detailed comparison with the above two closely related works.

Efficient Representation Learning in RL. Even in the single task setting, efficient representation learning is an active area of study that has seen many recent advances with exploration (Agarwal et al., 2020; Modi et al., 2021; Uehara et al., 2021; Zhang et al., 2022) or without (Ren et al., 2021). Other papers study feature selection (e.g. Farahmand and Szepesvári, 2011; Jiang et al., 2015; Pacchiano et al., 2020; Cutkosky et al., 2021; Lee et al., 2021; Zhang et al., 2021), or focus on learning sparse models (Hao et al., 2021a,b).

3 Preliminaries

Notations: We denote total variation distance of P_1 and P_2 by $k_{TV}(P_1, P_2)$. Given a vector a , we denote $\|a\|_2$ as $\|a\|$. c_0, c_1, \dots are universal constants. We use $\mathbb{E}_{\mathcal{P}}$ to mean $\mathbb{E}_{\mathcal{P}}$ for some universal constant \mathcal{P} . Also, $[K] = \{1, \dots, K\}$ and $\lambda_{\min}(A)$ is the smallest eigenvalue of matrix A . Please see Table 1 for a full list of notations and Table 2 for a list of algorithms.

Low-rank Markov Decision Processes: We consider a finite-horizon episodic Markov Decision Process $\mathcal{M} = (\mathcal{H}; \mathcal{S}; A; \{P_h^0, P_h^1, \dots, P_h^H\}; \{r_h^0, r_h^1, \dots, r_h^H\}; d_0)$, specified by the episode length H , state space \mathcal{S} , discrete action space A of size $|A|$, transition models $P_h^j : \mathcal{S} \times A \rightarrow \Delta(\mathcal{S})$, reward functions $r_h^j : \mathcal{S} \times A \rightarrow [0, 1]$, and an initial distribution $d_0 \in \Delta(\mathcal{S})$. Occasionally, we use P_h^j to denote P_h^j and to denote r_h^j . Starting from an initial state $s_0 \sim d_0$, an agent chooses actions at step h upon observing the state s_h , receives a reward $r_h(s_h; a_h)$, and transitions $s_{h+1} \sim P_h^j(s_h; a_h)$. We assume that d_0 and r_h^j are known. Given a policy $\pi : \mathcal{S} \rightarrow \Delta(A)$ and transition kernel P , we use the notation $\mathbb{E}_{\pi; P}[\cdot]$ to denote expectation under the distribution of trajectories which is executed in P .

The value function $V_{\pi; P, h}(s) = \mathbb{E}_{\pi; P}[\sum_{t=h}^H r_t(s_t; a_t) | s_h = s]$ gives the expected reward-to-go of π under the MDP with transition P , rewards r , starting at state s in step h . Similarly, we define the Q function $Q_{\pi; P, h}(s; a) := r_h(s; a) + \mathbb{E}_{s' \sim P_h(j; s; a)} V_{\pi; P, h+1}(s')$. The expected total reward of a policy π under transition P and rewards r is denoted as $V_{\pi; P} := \mathbb{E}_{s_0 \sim d_0} V_{\pi; P, 0}(s_0)$. We define the state-action occupancy distribution $d_{\pi; P, h}(s; a)$ as the probability of visiting $(s; a)$ at time step h under π and P . $d_{\pi; P, h}(s)$ is the marginal state visitation, which is equal to $\sum_{a \in A} d_{\pi; P, h}(s; a)$.

We study low-rank MDPs defined as follows (Jiang et al., 2017; Agarwal et al., 2020). The conditions on the upper bounds of the norm of P_h^j are just for normalization.

Definition 3.1 (Low-rank MDP) A transition model $P_h^j : \mathcal{S} \times A \rightarrow \Delta(\mathcal{S})$ admits a low rank decomposition with rank $\leq N$ if there exist two unknown embedding functions $\phi_h^j : \mathcal{S} \times A \rightarrow \mathbb{R}^d$, $\psi_h^j : \mathcal{S} \rightarrow \mathbb{R}^d$ such that $\mathbb{E}_{s' \sim P_h^j(s; a)} \langle \phi_h^j(s; a), \psi_h^j(s') \rangle = \langle \phi_h^j(s; a), \psi_h^j(s) \rangle$, where $\|\phi_h^j(s; a)\|_2 \leq 1$ for all $(s; a)$ and for any function $g : \mathcal{S} \rightarrow [0, 1]$, $\mathbb{E}_{s \sim d_h^j} g(s) \leq \mathbb{E}_{s \sim d_h^j} g(s)$. An MDP is a low rank MDP if P_h^j admits such a low rank decomposition for all $j \in [0; 1; \dots; H-1]$.

Low-rank MDPs capture the latent variable model (Agarwal et al., 2020) where ϕ_h^j is a distribution over a discrete latent state space and the block-MDP model (Du et al., 2019) where ϕ_h^j is a one-hot encoding vector. Note that the linear MDP model (Yang and Wang, 2020; Jin et al., 2020b) assumes known ψ_h^j , which significantly simplifies the algorithm design.

Transfer Learning: In contrast to the classic single-task learning setting, in this paper, we explore the setting of transfer learning where learning consists of two phases: (1) the training phase where the agent interacts with K source tasks with dynamics \mathcal{P}_k , and (2) the deployment phase where the agent is deployed into the target task and no longer has access to the source tasks. The performance is measured mainly by the regret incurred in the target task upon deployment, while we also desire small sample complexity in the source tasks. We denote $d_{\pi_k; P_k, h}$.

In order for the pre-training phase to help with learning in the target task, we must make assumptions on the connections between tasks. In this work, we make the following fundamental structural assumption on all the tasks at hand, namely that they share the same underlying representation.

Assumption 3.1 (Common representation) We assume that all tasks \mathcal{P}_k are low-rank MDPs with a shared representation $\phi_h^j(s; a)$ but distinct $\psi_{k; h}^j(s)$, that is $\mathbb{E}_{s' \sim P_{k; h}^j(s; a)} \langle \phi_h^j(s; a), \psi_{k; h}^j(s') \rangle = \langle \phi_h^j(s; a), \psi_{k; h}^j(s) \rangle$.

Assumption 3.2 (Realizability) For any source task $k \in [K]$ and any $h \in [H]$, we assume that the agent has access to realizable function class $\mathcal{F}_{k; h}$ and $\psi_{k; h}^j$, such that $\psi_{k; h}^j \in \mathcal{F}_{k; h}$ and $\psi_{k; h}^j \in \mathcal{F}_{k; h}$.

For normalization, we assume that for all k, h , all $\psi_{k; h}^j$ satisfy $\|\psi_{k; h}^j(s)\|_2 \leq 1$, and for all $k \in [K]$ and any function $g : \mathcal{S} \rightarrow [0, 1]$, $\mathbb{E}_{s \sim d_h^j} g(s) \leq \mathbb{E}_{s \sim d_h^j} g(s)$.

Assumption 3.2 is standard realizability condition in function class which is made in almost all prior works on RL with linear and nonlinear function approximation (e.g., Jiang et al. (2017); Sun et al. (2019); Jin et al. (2020b); Agarwal et al. (2020); Du et al. (2021); Uehara et al. (2021)).

Since all the tasks share a common representation, we expect the agent to learn a good representation during its initial interactions with the source tasks, and then learn a good policy for the target task using this representation in the deployment phase. The availability of multiple source tasks is

beneficial as the agent might not be able to learn a good representation for all parts of the observation space of the target task from one source only. In order to realize this intuition though, we need two additional assumptions which we describe next.

Assumptions for representational transfer. In addition to standard assumptions in the source tasks, we make the following structural and relatedness assumptions on the source and target tasks.

Assumption 3.3 (Feature reachability in the source tasks) We assume that $\gamma := \min_{k \in [K-1], h \in [H]} \max_{\min} E_{\mathcal{P}_k} [\gamma_h^?(s_h; a_h) \gamma_h^?(s_h; a_h)]$ is strictly positive.

Assumption 3.3 intuitively requires that no subspace is unreachable in the source tasks, as otherwise it's impossible to guarantee the quality of the learned representation in that subspace which may contain high rewards in the target task. Note that no reachability is required in the target task.

The next assumption quantifies the relatedness of the target task and the source tasks.

Assumption 3.4 (Relatedness: Point-wise Linear span) For any $h \in [H]$ and $s^0 \in S_{h+1}$, there is a vector $\gamma_h(s^0) \in \mathbb{R}^{K-1}$ such that $\gamma_{k;h}(s^0) = \frac{1}{\sum_{k=1}^{K-1} \gamma_{k;h}(s^0)}$ with $\max_{h,k} \gamma_{k;h}(s^0) \leq \frac{1}{\gamma}$ and $\gamma = \max_{h,k} \max_{s^0 \in S} \gamma_{k;h}(s^0)$.

Assumption 3.4 ensures that if $(s; a)$ is reachable from a $(s; a)$ pair in the target task, then it must be reachable from the same $(s; a)$ pair in at least one of the source tasks. This intuitively is also necessary for transfer learning, as s^0 could be a high rewarding state in the target. A special case is the convex combination, i.e., for any $s^0 \in S; h, \gamma_{k;h}(s^0) = p_k$ with $p_k \geq 0; \sum p_k = 1$, which implies $\max_{k,h} \gamma_{k;h}(s^0) \leq 1$ and $\gamma = 1$. On the other hand, if the target task largely focuses on observations quite rare under the source tasks, then γ can grow large, and this is unavoidable in a transfer learning setting.

We remark that unlike prior work (Du et al., 2020; Tripuraneni et al., 2020; Lu et al., 2021), we do not make any assumption on the data generating distribution in either the source or the target tasks. Instead, our approach is end-to-end, i.e., we collect our own data from scratch for representation learning by doing strategic exploration in the source tasks (see Section B for a detailed discussion). In fact, in Theorem C.1 we show that the above assumptions do not permit successful transfer in the supervised learning setting, thus establishing an interesting separation between supervised and reinforcement learning.

We conclude this section with a couple of examples where our assumptions are satisfied.

Example 3.1 (Mixture of source tasks) Perhaps the simplest example of our assumptions is where $\gamma_{k;h}(s^0; a) = \frac{1}{\sum_{k=1}^{K-1} \gamma_{k;h}(s^0; a)}$ with the coefficients independent of $s^0, k \geq 0$ and $\sum_{k=1}^{K-1} \gamma_{k;h}(s^0; a) = 1$. Such mixtures of base models have been considered in several prior works (Modi et al., 2020; Ayoub et al., 2020). While prior works study the case of arbitrary known base models, we instead allow structured, unknown base models with a shared representation. Here

Example 3.2 (Block MDPs with shared latent dynamics) In this example, each MDP \mathcal{P}_k is a Block MDP (Du et al., 2019) with a shared latent space and a shared decoder $\gamma: S \rightarrow Z$. In a block MDP, given state action pair $(s; a)$, the decoder γ maps s to a latent state z , the next latent state is sampled from the latent transition $\mathcal{P}(z; a)$, and the next state is generated from an emission distribution $\alpha(z)$. Recall that $\alpha(z) > 0$ at only one $z \in Z$ for any $s \in S$ for a block MDP. We assume that the latent transition model $\mathcal{P}(z; a)$ is shared across all the tasks, but the emission process differs across the MDPs. For instance, in a typical navigation example used to motivate Block MDPs, the latent dynamics might correspond to navigating in a shared 2-D map, while emission distributions capture different wall colors or lighting conditions across multiple rooms. Then Assumption 3.3 requires that the agent can visit the entire 2-D map, while Assumption 3.4 requires that the color/lighting conditions of the target task resemble that of at least one source task. The coefficients $\gamma_{k;h}(s^0)$ for any s^0 are non-zero on the source tasks which can generate that observation.

4 Transfer Learning with Generative Access to Source Tasks

We first study a setting where we assume generative model access to the source tasks, while having only online access to the target task

Assumption 4.1 (Generative access to the source tasks) We assume that we have access to generative models for the $K-1$ source tasks. Specifically, for any \mathcal{P}_k with $k \in [K-1]$, we can query any $(s_h; a_h)$ pair, and the generative model will return a next state sample $\mathcal{P}_{k;h}^?(s_h; a_h)$.

Algorithm 1 Transfer learning with generative access (REPTTRANSFER)

PRE-TRAINING PHASE

 Input: exploratory policies $\{g_{k=1}^K\}$, size of cross-sampled datasets, failure probability δ .

- 1: for task pairs i, j , i.e. for all $i, j \in [K-1]$ s.t. $i \neq j$ do . cross sampling procedure
- 2: For each $i \in [H-1]$, sample datasets $\mathcal{D}_{i,j;h}$ containing n i.i.d. $(s; a; s^0)$ tuples sampled as:

$$(s; a) \sim d_{i;h} \times P_{j;h}^?; s \sim P_{j;h}^?(js; a); a \sim U(A); s^0 \sim P_{i;h}^?(js; a); \quad (1)$$

- 3: For all $h \in [H-1]$, learn $b_h = \text{Multi-task REPLEARN}(\bigcup_{j \in [K-1]} \mathcal{D}_{k,j;h}, g_{k \in [K-1]}).$ (Algorithm 3)

DEPLOYMENT PHASE

 Additional Input: number of deployment episodes \bar{S}

- 1: Set $\bar{d} = H \bar{d} + dH \frac{P}{\log(dHT)}$.
 - 2: Run LSVI-UCB $f_{b_h, g_{h=0}^H}; r = r_K; T; \bar{S}$ in the target task $\mathcal{P}_K^?$ (Algorithm 6).
-

Having generative model access is not unrealistic, especially in applications where a high-quality simulation environment is available. Generative access also does not trivialize the challenge of efficient exploration in the source tasks, since there are a potentially infinite number of states and the ground truth representation is unknown. Prior works using generative access typically require either a known representation (so that one can perform D-optimal design to construct an exploratory state-action distribution (Agarwal et al., 2019)), or directly assume access to a diverse state-action distribution which provides coverage and from which one can sample. Neither such a diverse sampling distribution is given in our case. We also note that in Section 5 we will show that without any additional assumptions, generative access in source tasks is necessary.

Algorithm overview. Given the above setup, now we present our algorithm REPTTRANSFER, detailed in Algorithm 1. REPTTRANSFER takes a representation learning approach to the transfer learning problem and operates in two phases. During the pre-training phase, REPTTRANSFER performs reward-free exploration in each of the source tasks to learn task-specific policies which satisfies for all $h = 0; 1; \dots; H-1$, that $E_{k; P_k^?} \int_{\mathcal{H}} (s_h; a_h) \int_{\mathcal{H}} (s_h; a_h) > \min_{\mathcal{H}} I$ for some $\min > 0$. We call such policies \min -exploratory and will use this definition in Lemma 4.1 and Theorem 4.1. In Section 4.1 we present one particular algorithm that finds such with a specific \min , but our main analysis below is modular to the choice of the reward-free exploration algorithm.

Given the exploratory policies $\{g_{k \in [K-1]}\}$, REPTTRANSFER collects a joint dataset across the source tasks using these policies and cross-sampling across pairs of tasks and then learns a single representation b . In the deployment phase, REPTTRANSFER runs optimistic least squares value iteration (Algorithm 6) using the learned representation on the target task.

The cross sampling procedure and learning a shared representation. We now describe the cross sampling procedure in detail. Given an exploratory policy for each source task, the next step is to sample fresh data under cross-sampling procedures using the generative model. Consider a particular $(k; j)$ pair of source environments. For each $i \in [H]$, we first sample $s_{h-1}; a_{h-1} \sim d_{k;h-1}$. Then, in the simulator of task k , we reset to $(s_{h-1}; a_{h-1})$ and perform a transition step s_h , i.e., $s_h \sim P_{j;h-1}(s_{h-1}; a_{h-1})$. Then we reset the simulator for task j to states s_h , sample a uniformly random action a_h and then perform another transition step s_{h+1} , i.e., $s_{h+1} \sim P_k^?(s_h; a_h)$. Such a procedure is only possible under the generative model setting. Given datasets, each collected using the cross-sampling procedure, we perform an MLE-based representation learning procedure (Algorithm 3) jointly on all source tasks. An adaptation of the result of Agarwal et al. (2020, Theorem 21) to the multi-task setting shows that Multi-task REPLEARN (Algorithm 3) achieves the following total variation guarantee with probability at least $1 - \delta$:

$$E_{k,h} \int_{\mathcal{H}} |b_h(s; a) - b_{k;h}(\cdot)|^2 \int_{\mathcal{H}} (s; a) \int_{\mathcal{H}} (s; a) > \frac{2}{TV} N := O((\log(j) + K \log(j))N) : \quad (2)$$

²Given any dataset $\{s; a; r; s^0\}$ feature ϕ , and reward r , LSVI learns a Q function backward, i.e., at step t $w_h = \arg \min_w \int_{s,a,s^0} (w^T \phi(s; a) - V_{h+1}(s^0))^2 + \|w\|^2$ and set $V_h(s) = \max_a (r(s; a) + \gamma V_{h+1}(s^0))$; δ s. UCB, short for Upper Confidence Bound, refers to an exploration bonus added to basic LSVI.

Algorithm 2 REWARDFREE

- 1: Input: MDP $P^?$ with online access, num. LSVI-UCB episodes $N_{\text{LSVI-UCB}}$, num. model-learning episodes $N_{\text{REWARDFREE}}$, failure probability δ .
 - 2: Learn model $\hat{P}_h = (b_h; b_h)_{h=0}^{H-1}$ by running REWARDFREE REP-UCB (Algorithm 4) in $P^?$ for $N_{\text{REWARDFREE}}$ episodes.
 - 3: Set $\hat{P} = dH \log(dHN_{\text{LSVI-UCB}})$.
 - 4: Return \hat{P} = LSVI-UCB $f_{b_h} g_{h=0}^{H-1}; r = 0; N_{\text{LSVI-UCB}}; \text{UNIFORMACTIONS} = \text{TRUE}$ by simulating in the learned model \hat{P} (Algorithm 6). Note this step requires no samples from $P^?$
-

where $D_{k;h}$ denotes the generating distribution of inputted data sets, each with size N . For example, in REPTTRANSFER we have $D_{k;h} = [\prod_{j=1}^K D_{kj;h}]$ and $N = n(K-1)$.

Representational transfer. Next we transfer this guarantee to the target task via Assumption 3.4, a key insight of our work. Specifically, we show that under cross sampling and Assumption 3.4, \hat{P} also linearly approximates the true transition in the target task well, under the occupancy measure of policy. Remarkably, this holds before the agent has ever interacted with the target task.

Lemma 4.1 Suppose Assumption 3.4 and that for all source tasks $k \in [K-1]$ we have a \min -exploratory policy π_k . Then, for any $\epsilon \in (0, 1)$, learning features \hat{P} using the cross-sampling procedure of Algorithm 1 satisfies ϵ w.p. $1 - \delta$. Furthermore, for any $h = 0, 1, \dots, H-1$, there exist $s_h : S \rightarrow \mathbb{R}^d$ such that for any function $g : S \rightarrow [0, 1]$, $\|g(s) - \hat{P}_h(s)\|_2 \leq \epsilon$ and

$$\sup_{P^?} \mathbb{E}_{\pi_k} \|b_h(s_h; a_h) - e_h(s)\|_2 \leq \epsilon \quad \text{TV} := \frac{1}{2} \sum_{j=1}^K \frac{A_j}{A_{\max}} K^{-n} \leq \min :$$

Lemma 4.1 implies that \hat{P} is a feature such that \hat{P}_k is an approximately linear MDP in \hat{P} . Learning in approximately linear MDPs has been studied in Jin et al. (2020b), but under a much stronger error bound rather than the average misspecification here. Our result for downstream task learning under such an average model misspecification case presents the strongest result for learning in an approximately linear MDP, and the result might be of independent interest.

Regret bound for REPTTRANSFER. Putting things together, we obtain the following regret bound for online learning in the target task using our learned representation.

Theorem 4.1 (Regret under generative source access) Suppose Assumptions 3.1-3.4 and 4.1, and suppose the input policies π_k are \min -exploratory. Then, for any $\epsilon \in (0, 1)$, w.p. $1 - \delta$, REPTTRANSFER when deployed in the target task has regret at most $H^2 d^{1.5} \frac{1}{\epsilon} \sqrt{T \log(1/\epsilon)}$, with at most Kn generative accesses per source task, with $O\left(\frac{1}{\min} A_{\max}^3 K T \log \frac{1}{\epsilon} + K \log \frac{1}{\epsilon}\right)$.

Remarkably, Theorem 4.1 shows that with the pre-trained features, we achieve the same regret bound on the target task to the setting of linear MDP with known \hat{P} (Jin et al., 2020b), up to the additional factor. The scaling factor only depends on itself and captures the hardness of transfer learning. For special cases such as convex combination, i.e., state-independent and $\epsilon \in (0, K)$, then $\frac{1}{\min} = 1$. In the worst-case, some dependence on the scales seems unavoidable as we can have a state s^0 such that $\pi_k(s^0) = 1$ and $\pi_1(s^0) = 1$ with $\pi_1(s^0) = 1$. This corresponds to a rarely observed state for the source task encountered often in the target, and our estimates of transitions involving this state can be highly unreliable if it is not seen in any other source, roughly scaling the error between target and source tasks as $\frac{1}{\pi_1(s^0)}$. Obtaining formal lower bounds that capture a matching dependence on structural properties is an interesting question for future research.

4.1 Reward-free exploration in the source environments.

So far we have assumed that we have a reward-free exploration black box algorithm that can provide an exploratory policy π_k for each source task. In this section, we give a detailed algorithm that achieves this goal.

Recall that our transfer learning algorithm, for source task k , relies on an exploratory policy π_k such that the empirical covariance with respect to π_k is lower bounded. This ensures good exploration in the underlying ground truth feature space $P_k^?$. REWARDFREE Algorithm 2 achieves this goal

by first invoking the REWARDFREE REP-UCB algorithm (Algorithm 4) to learn an estimated linear MDP model \hat{P}_k for P_k . Subsequently, we perform reward-free exploration (e.g., using UCB with zero reward) within each \hat{P}_k which involves no further environment interactions.

Lemma 4.2 (Reward-free Exploration) Fix any source task $k \in [K - 1]$. Suppose Assumptions 3.2 and 3.3. Then, for any $\epsilon \in (0, 1)$, w.p. $1 - \epsilon$, REWARDFREE (Algorithm 2) with $N_{\text{LSVI-UCB}} = \lceil e A^3 d^6 H^8 \epsilon^{-2} \rceil$ and $N_{\text{REWARDFREE}} = \lceil e A^3 d^4 H^6 \log \frac{j}{\epsilon} \rceil N_{\text{LSVI-UCB}}^2$ returns a ϵ_{\min} -exploratory policy π_k where $\epsilon_{\min} = e^{-1} A^{-3} d^{-5} H^{-7} \epsilon^{-2}$. The sample complexity here is $N_{\text{REWARDFREE}}$ episodes in the source task.

To the best of our knowledge, Lemma 4.2 is the first result that finds a full-rank policy cover in the low-rank MDP setting, and might be of independent interest. Wagenmaker et al. (2022) recently obtained a related guarantee in the linear MDP setting with known features. Incorporating Lemma 4.2 and Theorem 4.1, we get our final sample complexity bound.

Theorem 4.2 (Regret for REPTTRANSFER with REWARDFREE subroutine) Suppose Assumptions 3.1-3.4 and 4.1 hold, and let $\epsilon \in (0, 1)$. Let $\{f_k, g_{k=1}^k\}$ be exploratory policies learned from running REWARDFREE on each source task with $N_{\text{LSVI-UCB}}$ and $N_{\text{REWARDFREE}}$ set as in Lemma 4.2. Then, w.p. $1 - \epsilon$, running REPTTRANSFER with $\{f_k, g_{k=1}^k\}$ has regret in the target task of $\lceil e H^2 d^{1.5} \sqrt{T \log(1/\epsilon)} \rceil$, with at most $\lceil e A^4 \frac{3}{\max} d^5 H^7 K^2 T^{-2} (\log(j/\epsilon) + K \log j) \rceil$ generative accesses per source task.

5 Transfer Learning with Online Access to Source Tasks

In the previous section, we show that efficient transfer learning is possible under very weak structural assumptions, assuming the help of generative access to the source tasks. One natural question is whether transfer learning is possible with only online access to the source tasks. In that case, the cross-sampling procedure no longer works. Somewhat surprisingly, we show that this is impossible without significantly stronger assumption.

Theorem 5.1 (Impossibility Result) Let M_K be the set of K -task multi-set that satisfies (1) all tasks are Block MDPs; (2) all tasks satisfy Assumption 3.3 and Assumption 3.4; (3) the latent dynamics are exactly the same for all source and target tasks. For any pre-training algorithm \mathcal{A} which outputs a feature $\hat{\phi}$ by interacting with the source tasks $\mathcal{S} \in [K - 1]$, there exists $\mathcal{P}_K^? \in \mathcal{P}_{K,2}^? M_K$, such that with probability at least $1 - \epsilon$, \mathcal{A} will output a feature $\hat{\phi}$, such that for any policy taking the functional form of $\pi(s) = f(\hat{\phi}(s); a) g_{a,2A}; f r(s; a) g_{a,2A}$, we have $\mathbb{E} V_K^? - V_K = 1 - \epsilon$.

The above theorem implies that a representation learned only from online access to source tasks does not enable learning in downstream tasks if the downstream task algorithm is restricted to use the representation as the only information of the state-action pairs (e.g., running UCB with $\hat{\phi}$).

We briefly explain the intuition behind the above lower bound. For a Block MDP, for (s, a) , we can model the ground-truth $q_{z,a}$ as a one-hot encoding $q_{z,a}$ corresponding to the latent state-action pair $(z; a)$ with $z = \phi^?(s)$ being the encoded latent state. The key observation here is that any permutation of $\phi^?$ will also be a perfect feature in terms of characterizing the Block MDPs, since it corresponds to simply permuting the indices of the latent states. Therefore, without cross referencing, the agent could potentially learn different permutations in different source tasks, which would collapse in the target task. A precise constructive proof of Theorem 5.1 can be found in Section C.

Part of the reason that the above example fails is that each source task has its own observed subset of raw states, which permits such permutation to happen. In what follows, we show that, indeed, under an additional assumption on the reachability of raw states, a slight variant of the same algorithm (Algorithm 5 in Section F) can achieve the same regret with only online access to the source tasks.

Assumption 5.1 (Reachability in the raw state space) For all source tasks $k \in [K - 1]$, any policy π and $h = 0; 1; \dots; H - 1$, we have $\inf_{s \in \mathcal{S}; a \in \mathcal{A}} d_{k,h}(s; a)_{\text{raw}} \min_{\pi} \mathbb{E}_{\pi; P_k^?} \sum_h (s_h; a_h) \sum_h (s_h; a_h) > \epsilon$.

Assumption 5.1 implies that for each source task, any policy that achieves a full-rank covariance matrix also achieves global coverage across the raw state-action space. In addition, in order to apply importance sampling (IS) to transfer the TV error from source task to target task, we need to assume

	Source	O-REPTTRANSFER	G-REPTTRANSFER	Oracle	Target
Comblock	1	8006 (294.7)	7790 (267.6)	7048 (164.8)	181450 (147600.2)
Comblock PO	1	1	7764 (145.7)	7336 (270.7)	85000 (14469.8)

(a)

(b)

Figure 1:(a): A visualization of the rich observation comblock environment. Latent states (white and black) emit continuous high-dimensional observations. The reward is sparse (only in white states at H). Each white state has 10 actions where one of them is a good action that leads to the next two white states while the other 9 lead to black states (good action for each white state is different). Once the agent is in black states, it stays stuck in black until the end of the episode. Thus, a random exploration strategy has an exponentially small probability of hitting the goal (i.e., $\frac{1}{(10^H)}$). (b) Top: Number of episodes required to solve the target environment under the setting from Sec. 6.1. Bottom: Number of episodes required to solve the target environment under the setting from Sec. 6.2. An algorithm solves the target task if it can achieve the optimal return (i.e., 1) for 5 consecutive iterations with 50 evaluation runs each. We include the mean and standard deviation (in the brackets) for 5 random seeds. \square denotes that an algorithm can not solve the target task within a fixed sample budget.

that the target task distribution has bounded density. This is true, for example, when S is discrete or when S is compact and has bounded measure.

Assumption 5.2 (Bounded density) For all $(s; h; s; a)$, we have $d_{K;h}(s; a) \leq 1$.

Theorem 5.2 (Regret with online access) Suppose Assumptions 3.1-3.4, 5.1, 5.2 hold. For any $\epsilon \in (0, 1)$, w.p. $1 - \epsilon$, Algorithm 5 with appropriately set parameters achieves a regret in the target task at most $\epsilon^{-1} d^{1.5} H^2 p^T \log(1/\epsilon)$, with at most $\text{poly}(A; \max(d; H; K; T); \epsilon^{-1}; \frac{1}{\text{raw}}; \log(\|j\|_{j=1}^n))$ online queries in the source tasks.

Assumption 5.1 is satisfied in a Block MDP, when, for example, the emission function $(s; z)$ satisfies that $\exists s; z; s.t., \exists (s; z) \in \mathcal{C}$. That is, for any source task, any state in the state space can be generated by at least one latent state.

6 Experiments

For the empirical study, we investigate the benefit of transfer learning under the Block MDP setting, on the challenging Rich Observation Combination Lock (comblock) benchmark. We choose this environment because one must recover the correct feature from rich observations and perform strategic exploration or otherwise pay for exponential sample complexity. We include a visual overview in Fig. 1(a). This environment is uniquely challenging since it requires strategic exploration and latent state discovery at the same time, which results the failures of common deep RL methods (Misra et al., 2020), and theoretical RL approaches based on linear function approximation (and reproducing kernel Hilbert space) (Zhang et al., 2022).

In particular, in this section we study the following questions: i) what are the benefits of representational transfer using multiple source tasks, and ii) whether generative access to source tasks is needed. We design two sets of experiments with various source and target environment configurations. We defer the details of the experiments (the design of vanilla comblock, hyperparameters, etc.) in Appendix I and include the essential design information in the following two sections.

Baselines. We denote \mathcal{C} as the smallest sample complexity of ϵ -VI-UCB using learned features from any of the source tasks; G-REPTTRANSFER as REPTTRANSFER with only online access to the source tasks; S-REPTTRANSFER as REPTTRANSFER with generative access to the source tasks; Oracle as learning in the target task with ground truth features; Target as running BRIEE (Zhang et al., 2022) — the SOTA Block MDP algorithm, in the target task with no pretraining.

³Importantly, BRIEE differs in how it learns features, but uses the same UCB method at each step for policy learning with the learned features.

(a) (b)

Figure 2: (a): Visualization of the decoder source (top) and REPTTRANSFER (bottom). (b): Visualization of the decoder O-REPTTRANSFER (top) and GREPTTRANSFER (bottom). For each baseline, the i -th column in the i -th image denotes the averaged decoded states from the 30 observations generated by latent state, for $i \in \{0, 1, 2\}$ and $h \in \{1, 2\}$ [25], from the corresponding target environment. The optimal decoder should recover the latent states up to a permutation. In Fig a (top), note that the learned features in source task fail to solve the target because of the collapse at timestep 5: both observations from state 0 and 1 are mapped to state 0. Note that in the source task where this feature is trained, such collapse can happen when 0 state have identical latent transition (for detailed discussion we refer to Misra et al. (2020); Zhang et al. (2022)). In Fig b (top), REPTTRANSFER with only online access learns an incorrect decoder when the source tasks' observation spaces are disjoint. This is because the learned feature can decode each source task with a different permutation.

6.1 Comblock without partitioned observations

In this section we start off with an easier setting: we use 5 source tasks, each with $horizon = 3$ latent states per level and 10 actions. The latent transition dynamics are different for each source task, but notice that for comblock, the reachability assumption (c.r. Assumption 3.3) is always satisfied. The emission distribution of all the source and target tasks is identical, so Assumption 5.1 holds here. For the construction of the target task, for each timestep choose the latent transition dynamics from one of the sources uniformly at random (thus Assumption 3.4 is satisfied). We record the number of episodes in the target environment that each method takes to solve the target environment in Table 1(b). We first observe that REPTTRANSFER with either online or generative access can solve the target task (since Assumption 5.1 holds). Second, we observe that directly applying the learned feature from any single source task does not suffice to solve the target environment. This is because the representation learned from a single source task may collapse two latent states into a single one during encoding (e.g., if two latent states at the same time step have exactly identical latent transitions). See the visualization of comparisons between the learned decoders from a single source task and REPTTRANSFER in Fig. 2(a). Third, the result shows that REPTTRANSFER saves order of magnitude of target samples compared with training in the target environment from scratch using the SOTA Block MDPs algorithm BRIE. This set of results verifies the empirical benefits of representation learning from multiple tasks, i.e., resolves ambiguity and speeds up downstream task learning.

6.2 Comblock with partitioned observation space (Comblock-PO)

In this section, following the intuition of our lower bound (Theorem 5.1), we construct a setting where the supports of the emission distributions from each task are completely disjoint, while the emission distribution in the target task is a mixture of all source emissions and the latent dynamics are identical across tasks. Hence Assumption 3.4 holds while Assumption 5.1 fails. So we expect that an algorithm without generative access to source tasks will fail based on Theorem 5.1. Specifically, under disjoint emission supports, a representation can decode the latent state correctly for each source, but permute latent state labels across sources, causing decoding errors on the target. We record the number of target episodes for each method to solve the target task in Table 1(b). We observe that indeed the online version fails while the generative version still succeeds. We show the visualizations of O-REPTTRANSFER and GREPTTRANSFER in Fig. 2(b), and we note that the empirical results exactly match our theoretical results. This ablation verifies that source generative model access is needed unless one has additional stronger distributional assumptions.

References

- Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and algorithms. CS Dept., UW Seattle, Seattle, WA, USA, Tech, Rep, 2019.
- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in Neural Information Processing Systems*, volume 33, pages 20095–20107, 2020.
- Sanjeev Arora, Simon Du, Sham Kakade, Yuping Luo, and Nikunj Saunshi. Provable representation learning for imitation learning via bi-level optimization. *International Conference on Machine Learning* pages 367–376. PMLR, 2020.
- Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin Yang. Model-based reinforcement learning with value-targeted regression. *International Conference on Machine Learning*, pages 463–474. PMLR, 2020.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- Emma Brunskill and Lihong Li. Sample complexity of multi-task reinforcement learning. *arXiv preprint arXiv:1309.6821*, 2013.
- Emma Brunskill and Lihong Li. Pac-inspired option discovery in lifelong reinforcement learning. In *International conference on machine learning*, pages 316–324. PMLR, 2014.
- Rich Caruana. Multitask learning. *Machine learning* 28(1):41–75, 1997.
- Yuan Cheng, Songtao Feng, Jing Yang, Hong Zhang, and Yingbin Liang. Provable benefit of multitask representation learning in reinforcement learning. *arXiv preprint arXiv:2206.05902*, 2022.
- Ashok Cutkosky, Christoph Dann, Abhimanyu Das, Claudio Gentile, Aldo Pacchiano, and Manish Purohit. Dynamic balancing for model selection in bandits and reinforcement learning. *International Conference on Machine Learning*, pages 2276–2285. PMLR, 2021.
- Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Sharing knowledge in multi-task deep reinforcement learning. *International Conference on Learning Representations*, 2019.
- Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford. Provably efficient rl with rich observations via latent state decoding. *International Conference on Machine Learning*, pages 1665–1674. PMLR, 2019.
- Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- Simon S Du, Sham M Kakade, Jason D Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization. *International Conference on Machine Learning*, 2021.
- Amir-massoud Farahmand and Csaba Szepesvári. Model selection in reinforcement learning. *Machine learning*, 85(3):299–332, 2011.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- Botao Hao, Yaqi Duan, Tor Lattimore, Csaba Szepesvári, and Mengdi Wang. Sparse feature selection makes batch reinforcement learning more sample efficient. *International Conference on Machine Learning*, pages 4063–4073. PMLR, 2021a.
- Botao Hao, Tor Lattimore, Csaba Szepesvári, and Mengdi Wang. Online sparse reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 316–324. PMLR, 2021b.

- Jiachen Hu, Xiaoyu Chen, Chi Jin, Lihong Li, and Liwei Wang. Near-optimal representation learning for linear bandits and linear rl. *International Conference on Machine Learning* pages 4349–4358. PMLR, 2021.
- Nan Jiang, Alex Kulesza, and Satinder Singh. Abstraction selection in model-based reinforcement learning. *International Conference on Machine Learning* pages 179–188. PMLR, 2015.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. *International Conference on Machine Learning* pages 1704–1713. PMLR, 2017.
- Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. *International Conference on Machine Learning* pages 4870–4879. PMLR, 2020a.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *Conference on Learning Theory* pages 2137–2143. PMLR, 2020b.
- Alessandro Lazaric, Emma Brunskill, et al. Sequential transfer in multi-armed bandit with finite set of models. *Advances in Neural Information Processing Systems* 26, 2013.
- Jonathan Lee, Aldo Pacchiano, Vidya Muthukumar, Weihao Kong, and Emma Brunskill. Online model selection for reinforcement learning with function approximation. *International Conference on Artificial Intelligence and Statistics* pages 3340–3348. PMLR, 2021.
- Rui Lu, Gao Huang, and Simon S Du. On the power of multitask representation learning in linear mdp. *arXiv preprint arXiv:2106.08053*, 2021.
- Thodoris Lykouris, Max Simchowitz, Alex Slivkins, and Wen Sun. Corruption-robust exploration in episodic reinforcement learning. *Conference on Learning Theory* pages 3242–3245. PMLR, 2021.
- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research* 17(81):1–32, 2016.
- Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. *International conference on machine learning* pages 6961–6971. PMLR, 2020.
- Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. *International Conference on Artificial Intelligence and Statistics* pages 2010–2020. PMLR, 2020.
- Aditya Modi, Jinglin Chen, Akshay Krishnamurthy, Nan Jiang, and Alekh Agarwal. Model-free representation learning and exploration in low-rank mdp. *arXiv preprint arXiv:2102.07035*, 2021.
- Aldo Pacchiano, Christoph Dann, Claudio Gentile, and Peter Bartlett. Regret bound balancing and elimination for model selection in bandits and rl. *arXiv preprint arXiv:2012.13045*, 2020.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22(10):1345–1359, 2009.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks* 113:54–71, 2019.
- Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical exploration for representation learning. *arXiv preprint arXiv:2111.11485*, 2021.
- Wen Sun, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Model-based rl in contextual decision processes: Pac bounds and exponential improvements over model-free approaches. *Conference on learning theory* pages 2898–2933. PMLR, 2019.
- Nilesh Tripuraneni, Michael Jordan, and Chi Jin. On the theory of transfer learning: The importance of task diversity. *Advances in Neural Information Processing Systems* 33:7852–7862, 2020.

- Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and of ine rl in low-rank mdps. In International Conference on Learning Representation, 2021.
- Andrew Wagenmaker, Yifang Chen, Max Simchowitz, Simon S Du, and Kevin Jamieson. Reward-free rl is no harder than reward-aware rl in linear markov decision processes. arXiv preprint arXiv:2201.11206, 2022.
- Ruosong Wang, Simon S Du, Lin Yang, and Russ R Salakhutdinov. On reward-free reinforcement learning with linear function approximation. Advances in neural information processing systems 33:17816–17826, 2020.
- Jiaqi Yang, Wei Hu, Jason D Lee, and Simon Shaolei Du. Impact of representation learning in linear bandits. In International Conference on Learning Representation, 2020.
- Jiaqi Yang, Qi Lei, Jason D Lee, and Simon S Du. Nearly minimax algorithms for linear bandits with shared representation. arXiv preprint arXiv:2203.15664, 2022.
- Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In Proceedings of the 37th International Conference on Machine Learning, pages 10746–10756, 2020.
- Andrea Zanette, Ching-An Cheng, and Alekh Agarwal. Cautiously optimistic policy optimization and exploration with linear function approximation. COLT, 2021.
- Weitong Zhang, Jiafan He, Dongruo Zhou, Amy Zhang, and Quanquan Gu. Provably efficient representation learning in low-rank markov decision processes. arXiv preprint arXiv:2106.11935, 2021.
- Xuezhou Zhang, Yuzhe Ma, and Adish Singla. Task-agnostic exploration in reinforcement learning. Advances in Neural Information Processing Systems 33:11734–11743, 2020.
- Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Wen Sun, and Alekh Agarwal. Efficient reinforcement learning in block mdps: A model-free representation learning approach. arXiv preprint arXiv:2202.00063, 2022.

Appendices

A Notations

Table 1: List of Notations

$S; A; A$	State and action spaces, and $ A_j $.
(S)	The set of distributions supported by S .
$\lambda_{\min}(A)$	Smallest eigenvalue of matrix A .
e_j	One-hot encoding of j , i.e. 0 at each index except the one corresponding to j .
$\ x\ _y$	Length of vector implied from context.
$\min_x f(x)$	$\min_x f(x); y$.
H	Episode length of MDPs, a.k.a. time horizon. We index steps as $0; 1; \dots; H - 1$.
K	There are a total of K tasks. The $r_{1; \dots; K-1}$ are source tasks, and the K -th task is the target task.
d	dimension of the low-rank MDP, i.e. dimension of \mathcal{F} .
$P_{k;h}^?$	Ground truth transition at time h for task k .
r_h	Reward function of the target task (i.e. task).
$E_{;P}[\cdot]$	Expectation under the distribution of trajectories when executed in P . We sometimes omit P when the MDP is clear from context.
$d_{P;h}$	Occupancy distribution of P under transition P at time h .
$d_{k;h}$	Occupancy distribution for the k -th task, i.e. $d_{P_k^?;h}$.
$f_h^?(s; a)$	Embedding function $f(s; a)$ at time h .
$\mathcal{F}_h^?$	Realizable function class for $r_h^?$.
$\ j\ _j$	Defined as $\max_h \ j\ _h$.
$f_{k;h}^?(s^0)$	Emission embedding function $f(s^0)$ at time h for environment k .
$\mathcal{F}_{k;h}^?$	Realizable function class for $r_{k;h}^?$.
$\ j\ _j$	Defined as $\max_{k;h} \ j\ _{k;h}$.
ρ_j	Feature reachability in the source task (Assumption 3.3).
ρ_{\max}	Constants defined in point-wise linear span assumption (Assumption 3.4).
ρ_{raw}	Raw states reachability parameter (Assumption 5.1).

Table 2: List of Algorithms

REPTRANSFER (Algorithm 1)	Transfer learning with generative access.
REWARDFREE (Algorithm 2)	Reward Free exploration in low-rank MDP.
Multi-task REPLEARN (Algorithm 3)	Multi-task Maximum Likelihood (MLE) Representation Learning.
REWARDFREE REP-UCB (Algorithm 4)	Modified REP-UCB for reward-free model learning.
REPTRANSFER for online access (Algorithm 5)	Simplified version of REPTRANSFER if we have reachability under raw states (see Theorem 5.2).
LSVI-UCB (Algorithm 6)	Optimistic Least Squares Value Iteration (with bonus).

B Comparisons to closely related works

Recall that in addition to the commonly made shared representation assumption ([Assumption 3.1](#)), we made two additional structural assumptions: reachability ([Assumption 3.3](#)) and linear span on \mathcal{S}^0 ([Assumption 3.4](#)). The reachability assumption is commonly made in prior works even in single agent RL, e.g. ([Modi et al., 2021](#); [Misra et al., 2020](#)). It ensures that there is no redundant dimensions in the ground-truth representation, which is a reasonable requirement. The linear span assumption is closely related to the diversity assumptions made in prior works of transfer learning in both supervised learning ([Tripuraneni et al., 2020](#)) and reinforcement learning ([Lu et al., 2021](#)).

[Lu et al. \(2021\)](#): In the prior work of [Lu et al. \(2021\)](#), which also studies transfer learning in low-rank MDPs with nonlinear function approximations, they need to make the following assumptions:

1. shared representation (identical to our [Assumption 3.1](#)).
2. task diversity (similar to our [Assumption 3.4](#)).
3. generative model access to both the source and the target tasks. In contrast, we only require generative model access to the source tasks and allow online learning in the target task.
4. a somewhat strong coverage assumption saying that the data covariance matrix (under the generative data distribution) between arbitrary pairs of features^{0,2} must be full rank. In contrast, our analysis only requires coverage in the true features^{0,2} in the source tasks.
5. the existence of an ideal distribution q on which the learned representation can extrapolate. We do not require an assumption of a similar nature. Instead, we show that the data collected from our strategic reward-free exploration phase suffices for successful transfer.
6. the uniqueness for each in the sense of linear-transform equivalence. Two representation functions ϕ and ϕ^0 can yield similar estimation result if and only if they differ by just an invertible linear transformation. In contrast, we do not make any additional structural assumptions on the function class \mathcal{H} beyond realizability.

In summary, our work present a theoretical framework that permits successful representation transfer based on significantly weaker assumptions. We believe that this is a solid step towards understanding transfer learning in RL.

[Cheng et al. \(2022\)](#): A concurrent work of [Cheng et al. \(2022\)](#) also studies the exact same problem as ours. Both works study the setting where the agent performs reward-free exploration in the source tasks for representation learning and used the learned representation for the target task. Both works achieves similar sample complexity in source and target tasks, albeit using very different algorithms. However, in addition to the assumptions that we made, [Cheng et al. \(2022\)](#) has to make an additional assumption.

Assumption B.1 For any two different models in the model class \mathcal{H} , say $P^1(s^0; \mathbf{s}; \mathbf{a}) = h^1(s; \mathbf{a})$; $\phi^1(s^0)$ and $P^2(s^0; \mathbf{s}; \mathbf{a}) = h^2(s; \mathbf{a})$; $\phi^2(s^0)$, there exists a constant C_R such that for all $(s; \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ and $h^2 \in \mathcal{H}$,

$$kP^1(\mathbf{j}; \mathbf{s}; \mathbf{a}) - P^2(\mathbf{j}; \mathbf{s}; \mathbf{a})k_{TV} \leq C_R E_{(s; \mathbf{a}) \sim U(\mathcal{S}; \mathcal{A})} kP^1(\mathbf{j}; \mathbf{s}; \mathbf{a}) - P^2(\mathbf{j}; \mathbf{s}; \mathbf{a})k_{TV}$$

This assumption ensures that the point-wise TV error is bounded, as long as the population-level TV error is bounded. This assumption is used to transfer the MLE error from the source tasks to the target task. This type of assumption is strong in the sense that we typically expect scale with $|\mathcal{S}|$. In contrast, our analysis ([Lemma F.1](#)) shows that this assumption is in fact not necessary assuming online access only to source tasks. The generative access to source task studied here, which enables transfer under weaker reachability assumptions is not studied in their work.

It is worth noting that [Cheng et al. \(2022\)](#) also study of on-policy RL in the target task which we do not cover, while we mainly focus on the setting of generative models in the source tasks and demonstrating a more complete picture by proving generative model access in source tasks is needed without

⁴with a the minor difference that they directly assume the state space is compact with bounded measure

additional assumptions. Comparing to (Cheng et al., 2022), we also further implement and perform experimental evaluations of our algorithm.

C Impossibility Results

Here, we present an interesting result showing that the above assumptions we make are so weak that they do not even permit efficient transfer in supervised learning:

Theorem C.1 (Counter-example in supervised learning) Assume that we want to perform conditional density estimation, where $p_k^?(y|x) = \frac{p_k^?(x) p_k^?(y)}{p_k^?(x)}$. Under [Assumption 3.1](#) (shared representation) and [Assumption 3.4](#) (linear span), and assume that in each source task, one have access to a data generating distribution $p_k(x)$ such that $\min_k E_k [\frac{p_k^?(x)}{p_k^?(x)}] > 0$ (reachability). No algorithm can consistently achieve $E_k [\frac{p_k^?(y|x)}{P_k^?(y|x)}] \leq 1 + \epsilon$ on the target task using the feature learned from the source tasks with probability more than $1 - \epsilon$.

Proof of Theorem C.1 Consider the following example $X = R^2$ and we have the following 3 sets.

$$S_1 = B_{1/2}((1, 1))$$

$$S_2 = B_{1/2}((2, 2))$$

$$S_3 = B_{1/2}((0, 1))$$

where $B_a((x, y))$ stands for the ball with radius a centered at (x, y) . These will be the support of 3 tasks: task 1 and 2 are two source tasks, task 3 is the target task. Let's assume $p_k(x)$ are uniform distribution on S_k .

Suppose that the feature class only contains two functions:

$$f_1 : f(x_1 - 0 \ \& \ x_2 - x_1)g \text{ if } (1, 0), \text{ and } (0, 1) \text{ otherwise}$$

$$f_2 : f(x_2 - 0 \ \& \ x_1 - x_2)g \text{ if } (0, 1), \text{ and } (1, 0) \text{ otherwise}$$

That is, the feature maps from R^2 to the set of binary encoding of dimension 2, $\{f_1(1, 0); (0, 1)g\}$. We further assume that $p_k^? = (p_1(y); p_2(y))$ for some distributions $p_1; p_2$, which is identical for all task k , where $p_1; p_2 \int_{T \setminus V} = 1$. We also assume that $p_k^?$ is known to the learner a priori, i.e. $p_k^? = f_k^?g$ for all $k \in [K]$, so all the learner needs to do is to pick the correct candidate.

Given the above setup, it's easy to verify that both [Assumption 3.1](#) and [Assumption 3.4](#) are satisfied, because the decision boundary of both f_1 and f_2 passes through the support of the source tasks, and all $p_k^?$'s are identical. However, f_1 and f_2 are equivalent in S_1 and S_2 in terms of their representation power, therefore no algorithm can always pick the correct feature function with probability more than $1/2$, regardless of the number of samples. Suppose the true feature and the algorithm incorrectly chooses f_2 . Then, for $x \in S_3$ $f(x_1 - 0)g$ which has probability mass $1/2$, $P_3^?(y|x) = p_2$ whereas $P_3^?(y|x) = p_1$. Thus, the expected total variation distance between $P_3^?$ and $P_3^?$ is $1/2$. \square

The above construction shows that our assumptions are not sufficient to permit reliable representation transfer, even in the supervised learning setting. Yet, surprisingly, these assumptions are sufficient in the RL setting, implying somehow that transfer learning in RL is easier than transfer learning in SL. To understand this phenomenon, observe that in RL, the marginal distribution $p(s)$ is not independent from the conditional density $p(y|s; a)$ we desire to estimate. In particular, if one collects data in the source tasks in an online fashion via running a policy $\pi(s; a)$ is structurally restricted to be an occupancy distribution generated by the ground-truth transition $p(s'; s; a)$. Such a connection can only exist in Markov chains, and our analysis elegantly utilizes this additional structure to establish the soundness of the learned representation. Also note, crucially, that we never learn a representation to capture which would suffer from similar issues as the supervised learning setting, but is not necessary for sample-efficient RL.

Next, we prove the impossibility result in [Theorem 5.1](#), restated below as [Theorem C.3](#). This result shows that one can not achieve online learning in the source tasks without significantly stronger

assumptions such as [Assumption 5.1](#). Before that, we provide a preliminary version, showing that the learned $\hat{\theta}$ is not sufficient to fit the transition model in the target task, which motivates the construction in [Theorem C.3](#).

Theorem C.2 (Impossibility Result: Model Learning) Let M_K be the set of K -task multi-set that satisfies

1. all tasks are Block MDPs;
2. all tasks satisfy [Assumption 3.3](#) and [Assumption 3.4](#);
3. the latent dynamics are exactly the same for all source and target tasks.

For any pre-training algorithm A , there exists $\{P_k^2\}_{k=1:K} \in M_K$ and an occupancy distribution μ_K on the target task, such that with probability at least $\frac{1}{2}$, A will output a feature $\hat{\theta}$ and for any

$$E_{\mu_K} [k^\wedge(s; a)] > \frac{1}{2} \sum_{k=1}^K P_k^2(j; a)_{k_{TV}} \quad (1)$$

Proof of [Theorem C.2](#) Consider a tabular MDP with 2 latent states z_1, z_2 and an observation state space $\mathcal{S} = \{R_1, R_2, B_1, B_2\}$, where in task 1 one can only observe R_1, R_2 and in task 2 one can only observe B_1, B_2 . Correspondingly, $\mu_1(s; z)$ is only supported on $\{R_1, R_2\}$ (i.e., $\mu_1(R_i | z_i) = 1$) and similar for task 2. Let the latent state transition be such that $P(z_1 | z_1; a) = 1$ and $P(z_2 | z_2; a) = 1$, i.e. only self-transition regardless of the actions.

Now, consider a 2-element feature class $\mathcal{F} = \{f_1, f_2\}$ such that

$$\begin{aligned} f_1 &= f(R_1 | z_1; R_2 | z_2; B_1 | z_1; B_2 | z_2) \\ f_2 &= f(R_1 | z_1; R_2 | z_2; B_1 | z_2; B_2 | z_1) \end{aligned} \quad (2)$$

Denote $\mu_i(s; a) = e_{i(s; a)}$ for $i \in \{1, 2\}$. Consider for each task k , a 2-element μ_k class in the form of $\mu_k = \{f(\mu_k(s; z_1); \mu_k(s; z_2)); f(\mu_k(s; z_2); \mu_k(s; z_1))\}$.

Notice that μ_1 and μ_2 are merely permutations of one another and so given any single task data, the two hypothesis will not be distinguishable by any means. Therefore, for any algorithm, there is at least probability $\frac{1}{2}$ that it will choose the wrong hypothesis if the ground truth is sampled between μ_1 and μ_2 uniformly at random. Suppose μ_1 is the correct hypothesis and μ_2 is the one that the algorithm picks (i.e., $\hat{\theta} = \mu_2$). Let task 3 be such that any state emits R_1 or R_2 and B_1 or B_2 each with probability $\frac{1}{2}$ (i.e., $\mu_3(R_i | z_i) = \mu_3(B_i | z_i) = 0.5$). This construction satisfies [Assumption 3.3](#) and [Assumption 3.4](#).

Then, within task 3, one would encounter observations from B_1 and B_2 which should be mapped to latent state z_1 and z_2 respectively by the true decoder f , but are instead both mapped to latent state z_1 by the learned decoder \hat{f} , and thus z_1 and z_2 become indistinguishable. Suppose $\mu_K(z_1) = \mu_K(z_2) = \frac{1}{2}$, then

$$\begin{aligned} E_{\mu_K} [k^\wedge(s; a)] &> \frac{1}{2} \sum_{k=1}^K P_k^2(j; a)_{k_{TV}} \\ &= \frac{1}{4} k^\wedge(R_1)_{k_{TV}} + \frac{1}{4} k^\wedge(B_1)_{k_{TV}} \\ &\quad + \frac{1}{4} k^\wedge(R_2)_{k_{TV}} + \frac{1}{4} k^\wedge(B_2)_{k_{TV}} \\ &= k_{\mu_1}^2 + k_{\mu_2}^2 \\ &= \frac{1}{4} k_{\mu_1}^2 + \frac{1}{4} k_{\mu_2}^2 \\ &= \frac{1}{2}; \end{aligned}$$

where the last second inequality uses triangle inequality, and the last equality comes from the fact that $\mu_3(z_1)$ and $\mu_3(z_2)$ have disjoint support which implies that $k_{\mu_1}^2 + k_{\mu_2}^2 = 1$. \square

Now, we are ready to restate and prove [Theorem 5.1](#).

Theorem C.3 (Impossibility Result: Optimal Policy Identification) Let M_K be the set of K -task multi-set that satisfies

1. all tasks are Block MDPs;
2. all tasks satisfy [Assumption 3.3](#) and [Assumption 3.4](#);
3. the latent dynamics are exactly the same for all source and target tasks.

For any pre-training algorithm A , there exists $P_{k=1:K}^2 \in M_K$, such that with probability at least $1-\epsilon$, A will output a feature \hat{f} , such that for any policy taking the functional form of $\pi(s) = f(\hat{f}(s; a))g_{a2A}; r(s; a)g_{a2A}$, we have

$$V^{\pi} - V^{\pi^*} \leq \epsilon$$

Proof of Theorem C.3 Consider a tabular MDP with $H = 2$, two latent states $z_1; z_2$ for $h = 1$ and two latent states $z_3; z_4$ for $h = 2$.

- For $h = 1$, let there be two actions $a_1; a_2$. Let the observation state space $\mathcal{S} = R_1 \cup R_2 \cup B_1 \cup B_2$, where in task 1 one can only observe $R_1 \cup R_2$ and in task 2 one can only observe $B_1 \cup B_2$. Correspondingly, $\rho_1(s|z)$ is only supported on $R_1 \cup R_2$ (i.e., $\rho_1(R_i|z_i) = 1$) and similar for task 2. Let the latent state transition be such that $P(z_3|z_1; a_1) = P(z_3|z_2; a_2) = 1$, and $P(z_4|z_1; a_2) = P(z_4|z_2; a_1) = 1$. All rewards are 0 for $h = 1$.
- For $h = 2$, in state z_3 , all actions have reward 1, and in state z_4 all actions have reward 0.
- The initial state distribution is $d_0(z_1) = d_0(z_2) = 1/2$.

Now, consider a 2-element feature class $f_1; f_2$ for $h = 1$, such that

$$\begin{aligned} f_1 &= f(R_1 | 1; R_2 | 2; B_1 | 1; B_2 | 2)g \\ f_2 &= f(R_1 | 1; R_2 | 2; B_1 | 2; B_2 | 1)g \end{aligned}$$

Denote $\pi_i(s; a) = e_{\pi_i(s; a)}$ for $i \in [1; 2]$. In addition, define $\pi = f_1; f_2$ where

$$\begin{aligned} \pi_1 &= f(z_3 | (1; 0); z_4 | (0; 1))g \\ \pi_2 &= f(z_4 | (1; 0); z_3 | (0; 1))g \end{aligned}$$

Notice that π_1 and π_2 are merely permutations of one another and so given any single task data, the two hypothesis will not be distinguishable by any means. Therefore, for any algorithm, there is at least probability $1/2$ that it will choose the wrong hypothesis. Suppose $\hat{\pi}$ is the correct hypothesis and π is the one that the algorithm picks (i.e., $\hat{\pi} = \pi_2$). Let task 3 be such that any state emits to R_2 and $B_1 \cup B_2$ each with probability $1/2$ (i.e., $\rho_3(R_i|z_i) = \rho_3(B_i|z_i) = 0.5$). This construction satisfies [Assumption 3.3](#) and [Assumption 3.4](#).

Then, for any policy that only make decision based on $\pi(s; a)$ and $r(s; a)$, would output the same action for observations in R_1 and B_2 , or for B_1 and R_2 . However, notice that the optimal policy, which would try to go to z_3 from either z_1 or z_2 , will pick a_1 at R_1 and B_1 while picking a_2 at R_2 and B_2 , which means that the optimal policy will not agree on R_1 and B_2 , and it also will not agree on R_2 and B_1 . Thus clearly, no such policy as defined above is capable of capturing the optimal policy. From the reward perspective, notice that $d(z_1) = d(z_2) = 1/2$ and $d(R_1) = d(R_2) = d(B_1) = d(B_2) = 1/4$. Since $(R_1) = (B_2)$, the agent will only be able to collect reward at one of the R_1 and B_2 (but not at both). Similarly, since $(R_2) = (B_1)$, the agent will only be able to collect reward at one of the R_2 and B_1 by reaching z_3 (but not at both). This means that will have average reward $1/4$. Since the optimal policy will be able to collect reward at all $R_1; R_2; B_1; B_2$, it will have average reward $1/2$. This concludes the proof. \square

Algorithm 3 Multi-task REPLEARN

- 1: Input: Datasets $\{D_k\}_{k=1:K}$, model class \mathcal{M}_k , $k = 1, \dots, K$.
 - 2: Compute $\hat{\theta}_k := \arg\max_{\theta \in \mathcal{M}_k} E_{D_k} \log(\theta(s; a) \cdot \theta(s^0))$.
 - 3: Return $\hat{\theta}$.
-

Algorithm 4 REWARDFREE REP-UCB

- 1: Input: Regularizer η , bonus scaling β , model class $\mathcal{M} = \bigcup_{h=1:H} \mathcal{M}_h$, number of episodes N .
 - 2: Initialize b_0 as random and $D_{h;0} = \emptyset$.
 - 3: for episode $n = 1; 2; \dots; N$ do
 - 4: Data collection from $b_{h;1}$: for $h = 1; 2; \dots; H$
 - s $\sim D_{h;1}^{b_{h;1}}$; a $\sim \text{Unif}(A)$; $s^0 \sim P_h^?(s; a)$;
 - $\mathbf{e} \sim D_{h;1}^{b_{h;1}}$; $\mathbf{a} \sim \text{Unif}(A)$; $\mathbf{e}^0 \sim P_h^?(s; \mathbf{a})$; $\mathbf{a}^0 \sim \text{Unif}(A)$; $\mathbf{e}^{00} \sim P_h^?(s^0; \mathbf{a}^0)$;
 - $D_{h;n} = D_{h;n-1} \cup \{(s; a; s^0)g; D_{h;n}^0 = D_{h;n-1}^0 \cup \{(\mathbf{e}^0; \mathbf{a}^0; \mathbf{e}^{00})g\}$;

For $h = 0$, only collect $D_{0;n}$.
 - 5: Learn model via MLE: for all $h = 0; 1; \dots; H$

$$\hat{\theta}_{h;n} = (\hat{b}_{h;n}; \hat{b}_{h;n}) = \arg\max_{\theta \in \mathcal{M}_h} E_{D_{h;n}} \log \theta_h(s; a) \theta_h(s^0) :$$
 - 6: Update exploration bonus: for all $h = 0; 1; \dots; H$

$$\hat{b}_{h;n}(s; a) = \eta \hat{b}_{h;n}(s; a) + \frac{\beta}{|D_{h;n}(s; a)|} \sum_{(s; a) \in D_{h;n}} \hat{b}_{h;n}(s; a) \hat{b}_{h;n}(s; a)^\top + \eta I :$$
 - 7: Learn policy $\hat{b}_n = \arg\max_{\theta \in \mathcal{M}} V_{\theta, \hat{b}_n}$ and let \hat{b}_n be its value.
 - 8: Let $\hat{b} = \arg\min_{n=2:N} \hat{b}_n$.
 - 9: Output: $\hat{b}; \hat{\theta}_{\hat{b}}$.
-

[Theorem C.2](#) and [Theorem C.3](#) show that it's impossible to allow online learning in the source tasks without much stronger assumptions. In our paper, we show that our [Assumption 5.1](#), which ensures reachability in the raw states, is sufficient to establish an end-to-end online transfer learning result. However, it is unclear if [Assumption 5.1](#) is necessary for online learning. We leave this as an important direction of future work.

D Reward-free Rep-UCB

In this section, we adapt the Rep-UCB algorithm ([Uehara et al., 2021](#)) for reward-free exploration in a single task. We drop all task subscripts as this section is for a single task only, i.e. think about the task as being each source task. The original Rep-UCB algorithm was for finite-horizon discounted MDPs, so we modify it to work for our undiscounted and finite-horizon setting. Our goal is to prove that Rep-UCB can learn a model that satisfies strong TV guarantees, i.e. [Theorem D.1](#) and [\(5\)](#). Note that FLAMBE ([Agarwal et al., 2020](#), Theorem 2) can be used for this directly, but at a worse (polynomial) sample complexity. Thus, we do a bit more work to derive a new model-learning algorithm for low-rank MDPs, based on Rep-UCB, that is more sample efficient in the source tasks.

A finite-horizon analysis of Rep-UCB was done in BRICE ([Zhang et al., 2022](#)), so here we just need to replace BRICE's RepLearn with that of the MLE, which is how we learn \hat{b} and \hat{b} , as in

Rep-UCB. Recall the notation of (Zhang et al., 2022),

$$\begin{aligned} \phi_{h,n}(s; \mathbf{a}) &= \frac{1}{n} \sum_{i=0}^{X-1} d_h^{b_i}(s) \text{Unif}(\mathbf{a}) \\ \phi_{h,n}(s; \mathbf{a}) &= \frac{1}{n} \sum_{i=0}^{X-1} \mathbb{E}_{\mathbf{e} \sim d_{h-1}^{b_i}; \mathbf{a} \sim \text{Unif}(A)} [P_h^?(s; \mathbf{e}; \mathbf{a}) \text{Unif}(\mathbf{a})] \\ \phi_{h,n}(s; \mathbf{a}) &= \frac{1}{n} \sum_{i=0}^{X-1} d_h^{b_i}(s; \mathbf{a}) \\ \hat{\phi}_{h,n} &= n \mathbb{E}[\phi_{h,n}(s; \mathbf{a}) (s; \mathbf{a})^\top] + \mathbf{1} \end{aligned}$$

By using MLE (Uehara et al., 2021, Lemma 18) to learn models, with probability at least $1 - \delta$ for any $n = 1, 2, \dots; N$ and $h = 0, 1, \dots; H - 1$, we have

$$\max_{\mathbf{a}} \mathbb{E}_{h,n} \|\hat{\phi}_{h,n}(s; \mathbf{a}) - P_h^?(s; \mathbf{a})\|_{TV}^2; \mathbb{E}_{h,n} \|\hat{\phi}_{h,n}(s; \mathbf{a}) - P_h^?(s; \mathbf{a})\|_{TV}^2 \leq \epsilon_n; \quad (3)$$

where

$$\epsilon_n = O\left(\frac{\log(jMj/nH)}{n}\right);$$

and $jMj = \max_{h \in [H]} \sum_j |j_h|$. We also adopt the same choice of ϵ_n parameters as BRICE, which we assume from now on.

$$\begin{aligned} \epsilon_n &= \left(\frac{d \log(jMj/nH)}{n}\right) \\ \epsilon_n &= \frac{p}{njAj^2/n + n d} \end{aligned}$$

As in Rep-UCB, we posit standard assumptions about realizability and normalization, on the (source) task of interest.

Assumption D.1 For any $h = 0, 1, \dots; H - 1$, we have $\|P_h^?\|_2 \leq \frac{1}{d}$ and $\|R_h^?\|_2 \leq \frac{1}{d}$. For any $\mathbf{a} \in A$, $\|k(s; \mathbf{a})\|_2 \leq 1$. For all $\mathbf{a} \in A$ and any function $g: S \rightarrow \mathbb{R}$, we have $\|g\|_s \leq \|g\|_k \leq \|g\|_1 \leq d \|g\|_s$.

Lemma D.1 Let r be any reward function. Suppose we ran Algorithm 4 with line 7 having reward $r + \mathbf{1}_h$ instead of just r . Then, for any $\delta \in (0, 1)$, w.p. at least $1 - \delta$, we have

$$\sum_{n=0}^{X-1} V_{\hat{\phi}_n; r + \mathbf{1}_h}^{b_n} - V_{P^?; r}^{b_n} \leq O\left(H^2 d^2 jAj^{1.5} \frac{p}{N \log(jMj/nH)}\right)$$

Proof. Start from the third equation of Zhang et al. (2022, Theorem A.4). Following their proof until the last page of their proof, we arrive at the following: for any $n = 1, 2, \dots; N$,

$$\begin{aligned} & V_{\hat{\phi}_n; r + \mathbf{1}_h}^{b_n} - V_{P^?; r}^{b_n} \\ & \leq \sum_{h=0}^{X-2} \mathbb{E}_{\mathbf{e}; \mathbf{a}} \sum_{d_{P^?; h}^{b_n}} k_h^?(s; \mathbf{e}; \mathbf{a}) k_{h,n}^? \frac{p}{jAj^2/n + n d} + \frac{q}{jAj^2/n} \\ & + (2H + 1) \sum_{h=0}^{X-2} \mathbb{E}_{\mathbf{e}; \mathbf{a}} \sum_{d_{P^?; h}^{b_n}} k_h^?(s; \mathbf{e}; \mathbf{a}) k_{h,n}^? \frac{p}{njAj/n + n d} + (2H + 1) \frac{p}{jAj/n} \end{aligned}$$

By elliptical potential arguments, we have

$$\sum_{n=0}^{X-1} \mathbb{E}_{\mathbf{e}; \mathbf{a}} \sum_{d_{P^?; h}^{b_n}} k_h^?(s; \mathbf{e}; \mathbf{a}) k_{h,n}^? \leq \frac{s}{dN \log\left(1 + \frac{N}{d}\right)};$$

Thus, summing over n , noting that $n; n; n$ are increasing in n , we can combine the above to get,

$$\sum_{n=0}^{X-1} V_{\hat{\phi}_n; r + \mathbf{1}_h}^{b_n} - V_{P^?; r}^{b_n}$$

$$\begin{aligned}
& \frac{s}{dN \log \left(1 + \frac{N}{d-1}\right)} H^q \frac{1}{|A_j| \frac{2}{N} d + N d + H^{2p} \frac{1}{N |A_j| \frac{1}{N} + N d}} \\
& \frac{s}{dN \log \left(1 + \frac{N}{d-1}\right)} H^p \frac{1}{N |A_j|^3 \frac{1}{N} d + N d^2 + H^{2p} \frac{1}{N |A_j| \frac{1}{N} + N d}} \\
& \frac{s}{dN \log \left(1 + \frac{N}{d-1}\right)} H^{2p} \frac{1}{d |A_j|^3 \log(jMj NH=) + d^3 \log(jMj NH=)} \\
2 O & H^2 d^2 |A_j|^{1.5} \frac{1}{N \log(jMj NH=)} :
\end{aligned}$$

□

This gives the following useful corollary for reward free exploration.

Lemma D.2 For any $\epsilon \in (0, 1)$ w.p. at least $1 - \epsilon$ we have

$$\|v_b - v^*\| \leq H^2 d^2 |A_j|^{1.5} \frac{r \log(jMj NH=)}{N} :$$

Proof. By definition of v_b , we have

$$\frac{N}{2} \|v_b - v^*\| = \sum_{n=N-2}^{N-1} \sum_{b_n, b_{n+1}} V_{b_n, b_{n+1}}^{b_b} - \sum_{n=0}^{N-1} \sum_{b_n, b_{n+1}} V_{b_n, b_{n+1}}^{b_b} ;$$

which is bounded by the previous lemma and the fact that $v_{b_n, b_{n+1}}^{b_b} = 0$, since in Algorithm 4, the reward function is zero. □

Conditioning on this, we now show that the environment \mathcal{P}_b has low TV error for any policy-induced distribution.

Theorem D.1 For any policy π , we have

$$\sum_{h=0}^{H-1} \mathbb{E}_{d_{\pi, h}} \|P_h^{\pi}(s; a) - P_{h; b}(s; a)\|_{TV} \leq H^3 d^2 |A_j|^{1.5} \frac{r \log(jMj NH=)}{N} := \epsilon_{TV} :$$

Proof. In this proof, let $\mathcal{P} = \mathcal{P}_b$, which is the returned environment from the algorithm. Let $r(s; a) = P_h^{\pi}(s; a) - P_{h; b}(s; a) \in [-2, 2]$. Then,

$$\begin{aligned}
& \sum_{h=0}^{H-1} \mathbb{E}_{d_{\pi, h}} \mathbb{E}_{d_{\mathcal{P}, h}} [r(s; a)] \\
& = V_{\pi, r} - V_{\mathcal{P}, r} \\
& = \sum_{h=0}^{H-1} \mathbb{E}_{d_{\mathcal{P}, h}} \mathbb{E}_{P_h^{\pi}(s; a)} \mathbb{E}_{P_{h; b}(s; a)} V_{\pi, r; h+1}(s^0) \quad (\text{Simulation lemma}) \\
& \leq 2H \sum_{h=0}^{H-1} \mathbb{E}_{d_{\mathcal{P}, h}} \|P_h^{\pi}(s; a) - P_{h; b}(s; a)\|_{TV} :
\end{aligned}$$

Thus,

$$\sum_{h=0}^{H-1} \mathbb{E}_{d_{\pi, h}} \|P_h^{\pi}(s; a) - P_{h; b}(s; a)\|_{TV}$$

$$\begin{aligned}
& (2H + 1) \sum_{h=0}^{H-1} E_{d_{\mathbb{P},h}} P_h^?(s; a) \|\mathbb{P}_h(s; a) - P_h^?(s; a)\|_{TV} \quad (\text{by (Zhang et al., 2022, Lemma A.1)}) \\
& \cdot H \sum_{h=0}^{H-1} E_{d_{\mathbb{P},h}} \left(\sum_{i=0}^h \binom{h-i}{i} \frac{r^i}{(jAj)^{N-2}} \right) \\
& H \sum_{\mathbb{P}, \mathbb{B}_b} V_{\mathbb{P}, \mathbb{B}_b} + 2jAj \frac{r \log(jMj NH=)}{N} \\
& H \sum_{\mathbb{P}, \mathbb{B}_b} V_{\mathbb{P}, \mathbb{B}_b}^{b_b} + 2jAj \frac{r \log(jMj NH=)}{N} \\
& \cdot H H^2 d^2 jAj^{1.5} \frac{r \log(jMj NH=)}{N} + jAj \frac{r \log(jMj NH=)}{N} \quad (\text{by Lemma D.2}) \\
& 2 O H^3 d^2 jAj^{1.5} \frac{r \log(jMj NH=)}{N} :
\end{aligned}$$

□

This also gives us a guarantee on the TV distance between the visitation distributions induced by $\mathbb{P}^?$ vs. by \mathbb{P} .

Lemma D.3 Suppose \mathbb{P} satisfies the following for all $h = 0, 1, \dots, H-1$,

$$\delta : E_{d_{\mathbb{P}^?,h}} \|\mathbb{P}_h(s; a) - P_h^?(s; a)\|_{TV} \leq \epsilon_h \quad (4)$$

Then, for any $h = 0, 1, \dots, H-1$, we have

$$\delta : d_{\mathbb{P},h} \leq d_{\mathbb{P}^?,h} + \sum_{t=0}^{h-1} \epsilon_t$$

Note, for $h = 0$, the sum is empty so the right hand side is 0.

Proof. We proceed by induction for $h = 0, 1, \dots, H-1$. For the base case $h = 0$, no transition has been taken, so that $d_{\mathbb{P},0} = d_{\mathbb{P}^?,0}$. Now let $h = 0, 1, \dots, H-2$ be arbitrary, and suppose that the claim is true for h (IH). We want to show the claim holds for $h+1$. One key fact we'll use is that, for any measure ν , $k_{TV} = \sup_{k \in \mathcal{K}_1} \int k \, d\nu$. Below we use the notation that $f(s; a) = E_{a \sim (s)} f(s; a)$.

$$\begin{aligned}
& k d_{\mathbb{P},h+1} \leq d_{h+1} k_{TV} \\
& = \sup_{k \in \mathcal{K}_1} E_{d_{\mathbb{P},h+1}} [k f(s; a)] - E_{d_{h+1}} [k f(s; a)] \\
& = \sup_{k \in \mathcal{K}_1} E_{(\mathbf{e}; a) \sim d_{\mathbb{P},h}; (s; a)} \mathbb{P}_h(\mathbf{e}; a) [k f(s; h+1)] - E_{(\mathbf{e}; a) \sim d_h; (s; a)} P_h^?(s; a) [k f(s; h+1)] \\
& \quad \sup_{k \in \mathcal{K}_1} E_{(\mathbf{e}; a) \sim d_{\mathbb{P},h}} E_{(\mathbf{e}; a) \sim d_h} E_{\mathbb{P}_h(\mathbf{e}; a)} f(s; h+1) \\
& \quad + \sup_{k \in \mathcal{K}_1} E_{(\mathbf{e}; a) \sim d_h} E_{\mathbb{P}_h(\mathbf{e}; a)} f(s; a) - E_{P_h^?(s; a)} f(s; h+1) \\
& \leq \sum_{t=0}^{h-1} \epsilon_t + E_{(\mathbf{e}; a) \sim d_h} \sup_{k \in \mathcal{K}_1} E_{\mathbb{P}_h(\mathbf{e}; a)} E_{P_h^?(s; a)} f(s; h+1) \quad (\text{by IH and Jensen}) \\
& \leq \sum_{t=0}^{h-1} \epsilon_t + \epsilon_h; \quad (\text{by (4) and } k f(s; h+1) \leq 1)
\end{aligned}$$

as desired. □

Thus, when combined with [Theorem D.1](#), we have for $0 \leq h \leq H-1$ and any policy π ,

$$k_{\pi,h} - d_{\pi,h} \leq O\left(\frac{H^3 d^2 j A_j^{1.5} \log(jMj NH)}{N}\right) = \epsilon_{TV} \quad (5)$$

In other words, the sample complexity needed for a model-error of ϵ is

$$O\left(\frac{H^6 d^4 j A_j^3 \log(jMj NH)}{\epsilon^2}\right)$$

Note this is much better than FLAMBE's guarantee ([Agarwal et al., 2020](#), Theorem 2) which requires,

$$O\left(\frac{H^{22} d^7 j A_j^9 \log(jMj NH)}{\epsilon^{10}}\right)$$

E Reward-free Exploration

In this section, we show that the mixture policy returned by [Algorithm 2](#) has good coverage. Recall that [Algorithm 2](#) contains two main steps:

Step 1 Learn a model \hat{P} . This was the focus of the previous section, where our model-free UCB method obtained a strong TV guarantee ((5)) by requiring number of episodes at most,

$$N_{\text{REWARDFREE}} = O\left(\frac{H^6 d^4 j A_j^3 \log(jMj NH)}{\epsilon_{TV}^2}\right)$$

Step 2 Run LSVI-UCB ([Algorithm 6](#)) in the learned model \hat{P} with reward at the e -th episode being $b_{h,e}$ and `UNIFORMACTIONS = TRUE`. The optimistic bonus pushes the algorithm to explore directions that are not well-covered yet by the mixture policy up to this point. With elliptical potential, we can establish that this process will terminate in polynomial number of steps.

We now focus on Step 2. Let π_h^{+1} denote rolling-in for h steps and taking uniform actions on the $h+1$ step, thus inducing a distribution over $s_{h+1}; a_{h+1}$. We abuse the notation a little and use π_h^{+1} for a policy that just takes one uniform action from the initial distribution.

Lemma E.1 Let $\epsilon \in (0, 1)$ and run REWARDFREE ([Algorithm 2](#)). Let $\Sigma_{h,N}$ be the empirical covariance at the N -th iteration of LSVI-UCB ([Algorithm 6](#)). Then, w.p. at least $1-\epsilon$ we have,

$$\sup_{h=0}^{H-1} \mathbb{E}_{\pi_h^{+1}; \mathcal{D}_h} \left| \frac{1}{N} \sum_{h=1}^N b_h(s_h; a_h) - \mathbb{E}_{\pi_h^{+1}} \left[\frac{1}{N} \sum_{h=1}^N b_h(s_h; a_h) \right] \right| \leq \epsilon$$

Proof. In this proof, we'll treat the empirical MDP $\mathcal{M}^?$, as that is the environment we're running in. Thus, we abuse notation and $\mathcal{M}_{h,e}$ is the model-based perspective of the linear MDP, $\mathcal{M}_{h,e}$ where $b_{h,e}$ is $\frac{1}{N} \sum_{k=1}^N b_h(s_h^k; a_h^k) (s_{h+1}^k)$. Also, in [Algorithm 2](#), we set reward to be zero, but for the purpose of this analysis, suppose the reward function is precisely the (unscaled) bonus in LSVI-UCB, i.e. $r_{h,e}(s_h; a_h) = b_{h,e}(s_h; a_h)$. This does not change the algorithm at all since the ϵ -scaling of the bonus dominates this reward in the definition of $\mathcal{M}_{h,e}$, but thinking about the reward in this way will make our analysis simpler.

Recall the high-level proof structure of reward free guarantee of linear MDP (with known features) ([Wang et al., 2020](#), Lemma 3.2).

Step 1 Show that $\frac{1}{N} \sum_{h=1}^N V_h$ and w.p. $1-\epsilon$, for all $h; e$,

$$\left| \frac{1}{N} \sum_{h=1}^N V_h - \mathbb{E}_{\pi_h^{+1}} \left[\frac{1}{N} \sum_{h=1}^N V_h \right] \right| \leq \epsilon$$

This step only uses self-normalized martingale bounds. So, [line 9](#) can use any martingale sequence of states and actions, and this claim still holds, with bonus using the appropriate covariance under the data.

Step 2 Show optimism conditioned on Step 1. Specifically, for $e = 1; 2; \dots; N$, we have $E_{d_0} V_0^?(s_0; r_e) - v_{0,e}(s_0) \geq 0$. To show this, we need that $b_{h,e}(s_h; a_h) = Q_{h,e}(s_h; a_h) - Q_{h,e}(s_h; a_h^?)$ (this is for the unclipped case of optimism), which we have satisfied in the algorithm, i.e. e_h is greedy w.r.t. $Q_{h,e}$.

Step 3 Bound the sum $\sum_{e=1}^P v_{h,e}$, where we decompose it as a sum of expected bonuses with the expectation is under e .

Step 3 is the only place where we use the fact that $s_{h,h}$ are data sampled from rolling out. For Step 1 and 2, please refer to existing proofs in (Agarwal et al., 2019; Jin et al., 2020b; Wang et al., 2020).

Now we show Step 3 for our modified algorithm with uniform actions. First, let us show a simulation lemma. For any episode $e = 1; 2; \dots; N$, for any s , recalling the definition of reward being $g_{h,e}$, we have

$$v_{0,e}(s_0) \leq (1 + \gamma) b_{0,e}(s_0; g_0^e(s_0)) + \gamma v_{1,e} \\ (1 + 2\gamma) b_{0,e}(s_0; g_0^e(s_0)) + \gamma^2 v_{1,e} + \dots$$

where the first inequality is due to the thresholding $v_{h,e}$'s and the second inequality is due to Step 1. Continuing in this fashion, we have

$$E_{d_0} v_{0,e}(s_0) \leq (1 + 2\gamma)^{H-1} E_e [b_{h,e}(s_h; a_h)]:$$

Summing over $e = 1; 2; \dots; N$, we have

$$\sum_{e=1}^N E_{d_0} v_{0,e}(s_0) \leq \sum_{h=0}^{H-1} \sum_{e=1}^N E_e [b_{h,e}(s_h; a_h)] \\ \leq A \sum_{h=0}^{H-1} \sum_{e=1}^N E_{(\frac{e}{H-1})^{+1}} [b_{h,e}(s_h; a_h)]$$

For each $h = 0; 1; \dots; H-1$, apply Azuma's inequality to the martingale difference sequence $e = E_{(\frac{e}{H-1})^{+1}} [b_{h,e}(s_h; a_h)] - b_{h,e}(s_h^e; a_h^e)$. The envelope is at most Δ so, w.p. $1 - \delta$,

$$\sum_{h=0}^{H-1} \sum_{e=1}^N b_{h,e}(s_h^e; a_h^e) + A \sqrt{p \sum_{h=0}^{H-1} \sum_{e=1}^N \Delta^2} \leq \sum_{h=0}^{H-1} \sum_{e=1}^N E_{(\frac{e}{H-1})^{+1}} [b_{h,e}(s_h; a_h)] + A \sqrt{p N \log(H/\delta)}$$

Now apply a self-normalized elliptical potential bound to the first term, giving that

$$\sum_{h=0}^{H-1} \sum_{e=1}^N b_{h,e}(s_h^e; a_h^e) \leq \sum_{h=0}^{H-1} \sqrt{p \sum_{e=1}^N \frac{1}{N} \sum_{e=1}^N b_{h,e}(s_h^e; a_h^e)^2} + H \sqrt{p d N \log(N)}$$

Thus, we finally have

$$\sum_{e=1}^N E_{d_0} v_{0,e}(s_0) \leq A H \sqrt{p d N \log(NH/\delta)}$$

Consider any episode $e = 1; 2; \dots; N$. By definition, $v_{h,N}(s; a) = v_{h,e}(s; a)$, so for all $s; a$ we have pointwise that $v_{h,N}(s; a) = v_{h,e}(s; a)$. Hence, for all s , we have $V_0^?(s; r^N) = V_0^?(s; r^e)$, and further using optimism, we have

$$N E_{d_0} [V_0^?(s_0; r^N)] \leq \sum_{e=1}^N E_{d_0} [V_0^?(s_0; r_e)] \leq \sum_{e=1}^N E_{d_0} v_{0,e}(s_0) \leq A H \sqrt{p d N \log(NH/\delta)}$$

Now consider any π and policy π , and consider rolling it out for H steps and taking a random action. Then we have

$$E_{\pi} [g(s_{h+1}; a_{h+1})] = E_{d_0} [V_0^{h+1}(s_0; r_N) - A^H P^H d \log(NH) = N]:$$

Summing over h incurs an extra H factor on the right. This concludes the proof. \square

Lemma E.2 (One-step back for Linear MDP) Suppose $P_h = (P_{h,s}; P_{h,a})$ is a linear MDP. Suppose π is any mixture of policies, and let $\Sigma_h := n E_{\pi} [k_h(s_h; a_h) k_h(s_h; a_h)^T] + I$ denote the unnormalized covariance. For any $g: S \rightarrow \mathbb{R}$, policy π , and $h = 0, 1, \dots, H-2$, we have

$$E_{\pi} [g(s_{h+1}; a_{h+1})] = E_{\pi} [k_h(s_h; a_h) k_h^{-1} \frac{1}{n} E_{\pi} [g(s_{h+1}; a_{h+1})^2] + d \text{kgk}_1^2]$$

Proof.

$$E_{\pi} [g(s_{h+1}; a_{h+1})] = E_{\pi} [k_h(s_h; a_h) \int_{S_{h+1}} g(s_{h+1}; a_{h+1}) d_{\pi}(s_{h+1}) + E_{\pi} [k_h(s_h; a_h) k_h^{-1} \int_{S_{h+1}} g(s_{h+1}; a_{h+1}) d_{\pi}(s_{h+1})];$$

where

$$\int_{S_{h+1}} g(s_{h+1}; a_{h+1}) d_{\pi}(s_{h+1}) = n E_{\pi} [E_{S_{h+1}} [P_h(s_h; a_h) [g(s_{h+1}; a_{h+1})]^2] + \int_{S_{h+1}} g(s_{h+1}; a_{h+1}) d_{\pi}(s_{h+1})^2]$$

$$= n \int_{S_{h+1}} g(s_{h+1}; a_{h+1})^2 + d \text{kgk}_1^2:$$

\square

Under reachability, we can show that small (squared) bonuses and spectral coverage, in the sense of having lower bounded eigenvalues, are somewhat equivalent.

Lemma E.3 Let Σ_h be a symmetric positive definite matrix and define the bonus $b_h(s; a) = k_h^{-1}(s; a) \Sigma_h^{-1} k_h(s; a)$. Then we have

1. For any policy π , $E_{d_h} [b_h^2(s; a)] \geq \frac{1}{\min(\cdot)}$. That is, coverage implies small squared bonus.
2. Suppose reachability under π (Assumption 3.3), then we have the converse: there exists ϵ such that for any policy π , $E_{d_h} [b_h^2(s; a)] \geq \frac{\epsilon}{\min(\cdot)}$. That is, small squared bonus implies coverage.

Proof. The first claim follows directly from Cauchy-Schwartz. Indeed, for any policy π we have

$$E_{d_h} [b_h^2(s; a)] = E_{d_h} [k_h^{-1}(s; a) k_h^{-1} k_h^{-1} k_h(s; a)] \geq \frac{1}{\min(\cdot)}:$$

For the second claim, Assumption 3.3 implies that there exist a policy π such that for all vectors $v \in \mathbb{R}^d$ with $\|v\|_2 = 1$, we have $E_{d_h} [(k_h^{-1}(s; a) v)^2] \geq \epsilon$. Now decompose $v = \sum_{i=1}^d v_i v_i^>$, where (λ_i, v_i) are eigenvalue/vector pairs with $\|v_i\|_2 = 1$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$. Then substituting this into the definition of the bonus, we have

$$E_{d_h} [b_h^2(s; a)] = \sum_{i=1}^d \frac{1}{\lambda_i} E_{d_h} [(k_h^{-1}(s; a)^T v_i)^2]$$

$$\frac{1}{d} \mathbb{E}_{d_h^e} \left(\sum_{h=1}^H (s_h; a_h)^T v_d \right)^2$$

$$\frac{1}{\min(\cdot)}:$$

□

We now prove our main lemma for reward-free exploration, [Lemma 4.2](#).

Lemma 4.2 (Reward-free Exploration) Fix any source task $\mathcal{K} \in \mathcal{K}$. Suppose Assumptions [3.2](#) and [3.3](#). Then, for any $\epsilon \in (0, 1)$, w.p. $1 - \epsilon$, `REWARDFREE` ([Algorithm 2](#)) with $N_{\text{LSVI-UCB}} = \lceil e A^3 d^6 H^8 \rceil$ and $N_{\text{REWARDFREE}} = \lceil e A^3 d^4 H^6 \log \frac{1}{\epsilon} \rceil N_{\text{LSVI-UCB}}^2$ returns a ϵ_{\min} -exploratory policy π_k where $\epsilon_{\min} = e^{-1} A^{-3} d^{-5} H^{-7} \epsilon^2$. The sample complexity here is $N_{\text{REWARDFREE}}$ episodes in the source task.

Proof of [Lemma 4.2](#) In this proof, let

$$\tilde{b}_h = N_{\text{LSVI-UCB}} \mathbb{E}_{h-1} \left[\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_h(s_h; a_h) b_h(s_h; a_h)^T + I \right];$$

$$\hat{b}_h = N_{\text{LSVI-UCB}} \mathbb{E}_{h-1} \left[\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_h(s_h; a_h) b_h(s_h; a_h)^T + I \right];$$

where $\tilde{b}_h = dH \log(N_{\text{LSVI-UCB}}) \cdot I$. This setting of \tilde{b}_h satisfies the precondition for the Concentration of Inverse Covariances [Zanette et al. \(2021, Lemma 39\)](#), which implies w.p. at least that

$$\hat{b}_h^{-1} \leq 2 \sum_{e=1}^{N_{\text{LSVI-UCB}}} b_h(s_h^e; a_h^e) b_h(s_h^e; a_h^e)^T + I \leq 2 \tilde{b}_h^{-1};$$

where we've also used the fact that $\tilde{b}_h \geq I$, so $(A + I)^{-1} \leq (A + I)^{-1}$.

Under this event, for any, we have,

$$\sum_{h=1}^H \mathbb{E}_{\pi} \left[\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_h(s_h; a_h) b_h(s_h; a_h)^T \right] \leq \sum_{h=1}^H \mathbb{E}_{\pi} \left[\sum_{e=1}^{N_{\text{LSVI-UCB}}} \hat{b}_h(s_h; a_h) b_h(s_h; a_h)^T \right]; \quad (6)$$

Now let $h = 0; 1; \dots; H-2$ be arbitrary. By [Assumption 3.3](#) (there exists some policy with coverage) such that,

$$\begin{aligned} & \frac{1}{\min_{h+1} \left(\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T \right)} \\ & \mathbb{E}_e \left[\sum_{h+1}^H \left(\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T \right)^{-1} \right] \quad (\text{by [Lemma E.3](#)}) \\ & \mathbb{E}_e \left[\sum_{h+1}^H \left(\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T \right)^{-1} \right] \quad (\text{by (6)}) \\ & \mathbb{E}_{\pi} \left[\sum_{h+1}^H \frac{1}{\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T} \right] \leq \frac{1}{\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T} + \epsilon_{\text{TV}} \quad (\text{by [Corollary E.1](#)}) \\ & \mathbb{E}_{\pi} \left[\sum_{h+1}^H \frac{1}{\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T} \right] \leq \frac{1}{\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T} + \epsilon_{\text{TV}} \quad (\text{by } \epsilon_{\text{TV}} = 1 - \epsilon_{\text{LSVI-UCB}}) \\ & \leq A^{1.5} d^2 H^3 \frac{1}{\log(dHN_{\text{LSVI-UCB}}) - N_{\text{LSVI-UCB}}} + 1 = N_{\text{LSVI-UCB}} \quad (\text{by (6) and [Lemma E.1](#)}) \\ & \leq A^{1.5} d^2 H^3 \frac{1}{\log(dHN_{\text{LSVI-UCB}}) - N_{\text{LSVI-UCB}}}. \end{aligned}$$

Recall that $\tilde{b}_h = dH \log(N_{\text{LSVI-UCB}}) \cdot I$, we have,

$$\begin{aligned} & \min_{h+1} \mathbb{E}_{h+1} \left[\sum_{e=1}^{N_{\text{LSVI-UCB}}} b_{h+1}(s_{h+1}; a_{h+1}) b_{h+1}(s_{h+1}; a_{h+1})^T \right] \\ & = \frac{1}{N_{\text{LSVI-UCB}}} \end{aligned}$$

$$\frac{1}{N_{\text{LSVI-UCB}}} \frac{C}{A^{1.5} d^2 H^3} \frac{dH}{N_{\text{LSVI-UCB}}} \log(dHN_{\text{LSVI-UCB}}) = N_{\text{LSVI-UCB}} \log(N_{\text{LSVI-UCB}})$$

$$\& \frac{C}{A^{1.5} d^2 H^3} \frac{dH}{N_{\text{LSVI-UCB}}};$$

where we've omitted the log terms for simplicity in the &. Now we optimize $N_{\text{LSVI-UCB}}$ to maximize this bound. For $a, b > 0$, to maximize a function of the form $f(x) = \frac{a}{x} - \frac{b}{x^2}$, it's best to seek x^* such that $f'(x^*) = 0$, resulting in value $f(x^*) = \frac{a^2}{4b}$. Setting,

$$x = N_{\text{LSVI-UCB}};$$

$$a = \frac{C}{A^{1.5} d^2 H^3};$$

$$b = dH;$$

Hence, we need to set

$$N_{\text{LSVI-UCB}} = e \frac{a^2}{b} = e \frac{A^3 d^6 H^8}{2};$$

which results in a \min lower bound of

$$\min E_{\pi_h} \sum_{h+1}^T (s_{h+1}; a_{h+1})^T (s_{h+1}; a_{h+1}) = e \frac{a^2}{b} = e \frac{2}{A^3 d^5 H^7};$$

Finally, we used the fact that $\gamma = 1 = N_{\text{LSVI-UCB}}$, which is set by the choice of $N_{\text{REWARDFREE}}$ in the lemma statement to satisfy (5).

The above proves coverage of $\hat{\mu}_h$ for $h = 0; 1; \dots; H-2$. Finally to argue for $h = H-1$, which is simply taking a random action at time t , we can simply invoke [Assumption 3.3](#) for $\gamma = 0$ to get a policy that

$$E_{\pi_{H-1}} \sum_{t=0}^T (s_t; a_t)^T (s_t; a_t) > \frac{1}{A} E_{\pi_{H-1}} \sum_{t=0}^T (s_t; a_t)^T (s_t; a_t) > \frac{1}{A};$$

□

Corollary E.1 Let $\hat{\mu}_h; \hat{\Sigma}_h$ be defined as in the proof of [Lemma 4.2](#). For any $h = 0; 1; \dots; H-2$ and any policy π_h , we have

$$E_{\pi_h} \sum_{h+1}^T (s_{h+1}; a_{h+1})^T (s_{h+1}; a_{h+1}) \leq E_{\pi_h} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h (s_{h+1}; a_{h+1}) (s_{h+1}; a_{h+1})^T) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h);$$

Intuitively, this means that coverage in the learned features implies coverage in the true features.

Proof. For shorthand, let $N = N_{\text{LSVI-UCB}}$. Apply [Lemma E.2](#) (one-step back) to the learned model and the function $f(s; a) = k_{h+1}^T (s; a) k_{h+1}$, which is bounded by $\|k_{h+1}\|^2 = 1$. We have,

$$E_{\pi_h} \sum_{h+1}^T (s_{h+1}; a_{h+1})^T (s_{h+1}; a_{h+1}) \leq \frac{1}{r} \frac{E_{\pi_h} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h (s_{h+1}; a_{h+1}) (s_{h+1}; a_{h+1})^T) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h)}{N \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h (s_{h+1}; a_{h+1}) (s_{h+1}; a_{h+1})^T) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h)} + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h) + \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h);$$

where we used the fact that

$$E_{\pi_h} \sum_{h+1}^T k_{h+1}^T (s_{h+1}; a_{h+1}) k_{h+1} \leq \frac{1}{A} \sum_{h+1}^T \text{tr}(\hat{\Sigma}_h);$$

$$\begin{aligned}
&= \text{Tr} \left(E_{h+1} \left[\frac{1}{N} \sum_{i=1}^N (s_{h+1}^i; a_{h+1}^i) \frac{1}{N} \sum_{i=1}^N (s_{h+1}^i; a_{h+1}^i) \right] - N E_{h+1} \left[\frac{1}{N} \sum_{i=1}^N (s_h; a_h) \frac{1}{N} \sum_{i=1}^N (s_h; a_h) \right] + I \right) \\
&= \frac{1}{N} \text{Tr}(I - M) \frac{d}{N};
\end{aligned}$$

where M is a positive definite matrix. Thus, doing an initial change from d_{h+1} to $d_{p,h+1}$ concludes the proof. \square

F Representation Transfer

First, we prove [Lemma 4.1](#), restated below.

Lemma 4.1 Suppose [Assumption 3.4](#) and that for all source tasks $k \in [K-1]$ we have a \min -exploratory policy π_k . Then, for any $\epsilon \in (0, 1)$, learning features \hat{b} using the cross-sampling procedure of [Algorithm 1](#) satisfies ϵ w.p. $1 - \epsilon$. Furthermore, for any $h = 0, 1, \dots, H-1$, there exist $e_h : S \rightarrow \mathbb{R}^d$ such that for any function $g : S \rightarrow [0, 1]$, $\|g(s)d - e_h(s)\|_2 \leq \frac{\epsilon}{\sqrt{d}}$ and

$$\sup_{g: S \rightarrow [0,1]} E_{\mathcal{P}_K} \left[\frac{1}{N} \sum_{i=1}^N b_h(s_h; a_h)^i - e_h(s) \right] \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \leq \frac{\epsilon}{\sqrt{d}} \frac{1}{\sqrt{K}} \frac{1}{\sqrt{n}} := \frac{\epsilon}{\sqrt{d}} \frac{1}{\sqrt{K}} \frac{1}{\sqrt{n}}.$$

Proof of Lemma 4.1 Fix an arbitrary ϵ . Denote $e_h(s^0) = \frac{1}{K} \sum_{k=0}^{K-1} b_{k,h}(s^0) b_{k,h}(s^0)$. First, note that

$$\begin{aligned}
&\max_{g: S \rightarrow [0,1]} \frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - e_h(s) \Big|_{TV} \\
&\leq \max_{g: S \rightarrow [0,1]} \frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \\
&\leq \max_{s \in S} \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \quad (\text{Since } \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \leq \frac{\epsilon}{\sqrt{d}} \text{ by 3.2}) \\
&= \frac{\epsilon}{\sqrt{d}}
\end{aligned}$$

For any $h = 0, 1, \dots, H-1$, we have

$$\begin{aligned}
&E_{\mathcal{P}_K} \left[\frac{1}{N} \sum_{i=1}^N b_h(s_h; a_h)^i - e_h(s) \right] \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \tag{3} \\
&= E_{\mathcal{P}_K} \left[\frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N (s_{h+1}^i; a_{h+1}^i) - b_{k,h}(s_{h+1}^i) \right] \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \tag{5} \\
&\leq \frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N (s_{h+1}^i; a_{h+1}^i) - b_{k,h}(s_{h+1}^i) \Big|_{TV} \tag{3} \\
&= E_{\mathcal{P}_K} \left[\frac{1}{N} \sum_{i=1}^N \frac{1}{N} \sum_{i=1}^N (s_{h+1}^i; a_{h+1}^i) - b_{k,h}(s_{h+1}^i) \right] \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} \tag{5} \\
&\leq \max_{k=1}^{K-1} E_{\mathcal{P}_K} \left[\frac{1}{N} \sum_{i=1}^N b_h(s_h; a_h)^i - b_{k,h}(s) \right] \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV} :
\end{aligned}$$

First consider the case when $h=0$. At $h=0$, the distribution under \mathcal{P}_K is the same as \mathcal{P}_k , and so, we directly get that the above quantity is at most $\frac{\epsilon}{\sqrt{d}}$, which proves the $h=0$ case.

Now consider any $h = 1, 2, \dots, H-1$. To simplify notation, let us denote

$$\begin{aligned}
\text{err}_{k,h}(s_h; a_h) &= \frac{1}{N} \sum_{i=1}^N b_h(s_h; a_h)^i - b_{k,h}(s) \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i \Big|_{TV}; \\
w_{k,h} &= \frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i E_{a_h} \left[\frac{1}{N} \sum_{i=1}^N (s_h; a_h)^i - b_{k,h}(s) \right] \text{err}_{k,h}(s_h; a_h);
\end{aligned}$$

$$V_{k,h} = \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_h; a_h) - \mathbb{E}_h(s_h; a_h) \rangle :$$

Note that $\min_{k,h} (V_{k,h})$ by assumption. Now continuing from where we left off, we take a one-step back as follows,

$$\begin{aligned} & \max_{k=1}^K \mathbb{E}_{\mathcal{P}_k^?} \text{err}_{k,h}(s_h; a_h) \\ = & \max_{k=1}^K \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \\ & \max_{k=1}^K \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \end{aligned}$$

By \min guarantee of $V_{k,h}$, and Jensen's inequality to push the square inside,

$$\begin{aligned} & \frac{1}{\min_{k=1}^K} \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \\ & \frac{1}{\min_{k=1}^K} \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \end{aligned}$$

By Assumption 3.4, the expectation over $\mathcal{P}_k^?$ is a linear combination of expectations over $\mathcal{P}_{j,h}^?$,

$$\begin{aligned} & \frac{1}{\min_{k=1}^K} \sum_{j=1}^J \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \\ & \frac{1}{\min_{k=1}^K} \sum_{j=1}^J \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \\ & \frac{1}{\min_{k=1}^K} \sum_{j=1}^J \mathbb{E}_{\mathcal{P}_k^?} \langle \mathbb{E}_h(s_{h-1}; a_{h-1}) - \mathbb{E}_h(s_{h-1}; a_{h-1}) \rangle w_{k,h} \end{aligned}$$

where we used the MLE guarantee (2) in the last step. □

Algorithm 5 Transfer learning with online access

PRE-TRAINING PHASE

Input: num. LSVI-UCB episodes $N_{\text{LSVI-UCB}}$, num. model-learning episodes $N_{\text{REWARDFREE}}$, size of cross-sampled datasets, failure probability δ .

- 1: for source task $k = 1; \dots; K - 1$ do
- 2: Find policy cover $\pi_k = \text{REWARDFREE}(\mathcal{P}_k^?; N_{\text{LSVI-UCB}}; N_{\text{REWARDFREE}}; \delta)$. (Algorithm 2)
- 3: for source task $k = 1; \dots; K - 1$ do
- 4: For each $h = 0; 1; \dots; H - 1$, sample \mathcal{D}_k as n i.i.d. $(s_h; a_h; s_{h+1})$ tuples from π_k .
- 5: For each $h = 0; 1; \dots; H - 1$, learn $b_h = \text{Multi-task REPLEARN}(\mathcal{D}_{k,h}; g_{k,2[K-1]})$. (Algorithm 3)

DEPLOYMENT PHASE

Additional Input: number of deployment episodes \bar{s}

- 1: Set $\bar{d} = H \bar{s} + dH \frac{1}{\log(dHT)}$.
- 2: Run LSVI-UCB \bar{d} in the target task $\mathcal{P}_K^?$ (Algorithm 6).

Next we state an analogous lemma for when we don't need generative access to the source task, but instead assume Assumption 5.1, and Assumption 5.2.

Lemma F.1 Suppose [Assumption 5.1](#), and [Assumption 5.2](#). Now take the setup of [Lemma 4.1](#) with the only difference being that θ is learned as in [Algorithm 5](#). Then, the same guarantee of [Lemma 4.1](#) holds with a slightly different right hand side for the bound on the TV-error,

$$\sup_{\mathbb{P}_K} \mathbb{E} \left[\int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_h(s_h; a_h) \pi_{K;h}(\cdot) \, \text{TV} \left(\frac{\max_{\text{raw}} K^{1=2} \frac{1=2}{n}}{\min} \right) \right]$$

Proof of [Lemma F.1](#) Fix an arbitrary \cdot . Denote $\pi_h(s^0) = \prod_{k=0}^{K-1} \pi_{k;h}(s^0) \wedge \pi_{k;h}(s^0)$. Then, some algebra with importance sampling gives us the bound,

$$\begin{aligned} & \mathbb{E}_{\mathbb{P}_K} \int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_h(s_h; a_h) \pi_{K;h}(\cdot) \, \text{TV} \tag{3} \\ & \mathbb{E}_{\mathbb{P}_K} \int_{\mathcal{S}_{h+1}} \prod_{k=1}^{K-1} \pi_{k;h}(s_{h+1}) b_h(s_h; a_h) \pi_{k;h}(s_{h+1}) \pi_h(s_h; a_h) \pi_{k;h}(s_{h+1}) \, \text{TV} \tag{5} \\ & \max_{k=1}^{K-1} \mathbb{E}_{\mathbb{P}_K} \int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_{k;h}(\cdot) \pi_h(s_h; a_h) \pi_{k;h}(\cdot) \, \text{TV} \\ & \max_{k=1}^{K-1} K^{1=2} \int_{\mathcal{S}_h} \mathbb{E}_{\mathbb{P}_K} \left[\int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_{k;h}(\cdot) \pi_h(s_h; a_h) \pi_{k;h}(\cdot) \, \text{TV} \right]^2 \end{aligned}$$

By [Assumption 5.1](#), [Assumption 5.2](#), for any $s; a$, we have $\frac{d_{K;h}(s;a)}{d_{k;h}(s;a)}$ $\frac{1}{\min_{\text{raw}}}$, where we used the coverage-under- π -assumption in the last inequality. In other words, for each $k=1; 2; \dots; K-1$, we have

$$\begin{aligned} & \frac{d_{K;h}}{d_{k;h}} \frac{1}{\min_{\text{raw}}} \int_{\mathcal{S}_h} \mathbb{E}_{\mathbb{P}_K} \left[\int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_h(s_h; a_h) \pi_{k;h}(\cdot) \, \text{TV} \right]^2 \\ & \frac{\max_{k=1}^{K-1} K^{1=2} \int_{\mathcal{S}_h} \mathbb{E}_{\mathbb{P}_K} \left[\int_{\mathcal{S}_h} b_h(s_h; a_h) \pi_{k;h}(\cdot) \pi_h(s_h; a_h) \pi_{k;h}(\cdot) \, \text{TV} \right]^2}{\left(\min_{\text{raw}} \right)^{1=2} \frac{1=2}{n}} \\ & \frac{\max_{k=1}^{K-1} K^{1=2} \frac{1=2}{n}}{\left(\min_{\text{raw}} \right)^{1=2}} \end{aligned}$$

□

G Proofs for Deployment Phase

G.1 Online RL Lemmas

Lemma G.1 (Self-normalized Martingale) Consider iterations $F_i, g_{i=1; 2; \dots}$, so that $\mathbb{E}[F_i | F_{i-1}] = 0$ and $F_i | F_{i-1}, g_{i=1; 2; \dots}$ are sub-Gaussian with parameter σ . Let $X_i | g_{i=1; 2; \dots}$ be random variables in a Hilbert space \mathcal{H} . Suppose a linear operator $\Sigma_0 : \mathcal{H} \rightarrow \mathcal{H}$ is positive definite. For any, denote $\Sigma_t = \Sigma_0 + \sum_{i=1}^t X_i X_i^T$. Then w.p. at least $1 - \delta$, we have,

$$\| \Sigma_t^{-1} \sum_{i=1}^t X_i \|^2 \leq 2 \log \frac{\det(\Sigma_t) \det(\Sigma_0)^{-1}}{\delta}$$

Proof. Lemma A.8 of ([Agarwal et al., 2019](#)). □

Lemma G.2 Let $\Sigma_t = I + \sum_{i=1}^t x_i x_i^T$ for $x_i \in \mathbb{R}^d$ and $\|x_i\| \leq 1$. Then $\Sigma_t^{-1} \sum_{i=1}^t x_i \leq d$.

Proof. Lemma D.1 of ([Jin et al., 2020b](#)). □

G.2 Proof of main results

Let $(x)_y$ refer to the clamping operator, i.e. $(x)_y = \min\{x, y\}$; $y \cdot g$. Let M_V be the maximum possible value in the MDP with the given reward function.

Algorithm 6 LSVI-UCB

- 1: Input: Features $\{b_h\}_{h=0;1;\dots;H-1}$, reward $\{r_h\}_{h=0;1;\dots;H-1}$, number of episodes N , bonus scaling parameter, UNIFORMACTIONS = FALSE.
- 2: for episode $e = 1; 2; \dots; N$ do
- 3: Initialize $v_{h,e}(s) = 0; \forall s$
- 4: for step $h = H-1; H-2; \dots; 0$ do
- 5: Learn best predictor for v_{h+1}^e ,

$$w_{h,e} = \sum_{k=1}^K b_h(s_h^k; a_h^k) v_{h+1}^e(s_{h+1}^k) + I;$$

$$w_{h,e} = \sum_{k=1}^K b_h(s_h^k; a_h^k) v_{h+1}^e(s_{h+1}^k);$$

- 6: Set bonus and value functions,

$$b_{h,e}(s; a) = b_h(s; a) \cdot \frac{1}{M_V};$$

$$Q_{h,e}(s; a) = w_{h,e} + b_{h,e}(s; a) + r_h(s; a) + v_{h,e}(s; a);$$

$$v_{h,e}(s) = \max_a Q_{h,e}(s; a);$$

- 7: Set $a_h^e(s) = \arg\max_a Q_{h,e}(s; a)$.
 - 8: Execute a_h^e to collect a trajectory $(s_h^e; a_h^e)_{h=0}^{H-1}$.
 - 9: If UNIFORMACTIONS = TRUE, discard a_h^e and draw freshly sampled uniform actions independently for all h , i.e. $a_h^e \sim \text{Unif}(A)$.
 - 10: Return: uniform mixture $\mu = \text{Uniform}(\{a_h^e\}_{e=1}^N)$.
-

Previously, [Jin et al. \(2020b\)](#) analyzed LSVI-UCB under point-wise model-misspecification. Here, we show that similar guarantees hold under a more general μ -distribution model-misspecification μ_{ms} , captured by [Assumption G.1](#).

Assumption G.1 Suppose for every $h = 0; 1; \dots; H-1$, there exists ϵ_h such that for any policy π ,

$$\mathbb{E} \left[\epsilon_h \left(\sum_{i=1}^I \pi_i(s_h; a_h) \right)^T b_h(s_h; a_h) \right] \leq \epsilon_h \mu_{ms};$$

We further assume that $\sup_{s,a,h} \epsilon_h \left(\sum_{i=1}^I \pi_i(s; a) \right)^T b_h(s; a) \leq M$ and $\| \sum_{i=1}^I \pi_i(s; a) \|^2 \leq M$ for $\forall s, a$.
 $S \subseteq \mathbb{R}^d$, for some positive constant M .

In other words, we only need the model to be accurate on average under the occupancy distributions realizable by policies. We also make a slight generalization on the regularization constant λ which is set to 1 in the original linear MDP definition ([Jin et al., 2020b](#)). Later, we will later instantiate the above assumption with our transferred model $\mathbb{P}(s^0) = \sum_{k=1}^K \pi_{k;h}(s^0) b_{k;h}(s^0)$, then for any $s; a$, we have

$$\| \sum_{k=1}^K \pi_{k;h}(s^0) b_{k;h}(s^0) \|^2 \leq M$$

$$\| \sum_{k=1}^K \pi_{k;h}(s^0) b_{k;h}(s^0) \|^2 \leq M$$

$$\sum_{k=1}^{K-1} \max_{s^0} \sum_{j: k;h(s^0)} b_{k;h}(s^0) f(s^0) \quad (\text{by } b_{k;h}(s; a) \geq 0) \quad (1)$$

Also,

$$\begin{aligned} \|k f^\top e_h\|_2 &= \sum_{k=1}^{K-1} \sum_{s^0} \sum_{j: k;h(s^0)} b_{k;h}(s^0) f(s^0) \\ &= \sum_{k=1}^{K-1} \max_{s^0} \sum_{j: k;h(s^0)} b_{k;h}(s^0) f(s^0) \\ &= \sum_{k=1}^{K-1} \|k f^\top b_{k;h}\|_2 \quad (\text{by } \|k f^\top b_{k;h}\|_2 = \sum_{s^0} \sum_{j: k;h(s^0)} b_{k;h}(s^0) f(s^0)) \end{aligned}$$

So we will set $M = \dots$

Note that we only need the existence of $b_{k;h}(s; a)$ here, and $b_{k;h}(s; a)$ need not be a valid probability kernel. In fact, it may even be negative valued.

In this section, we make a model-based analysis of LSVI. Similar approaches have been used in prior works, e.g. [Lykouris et al. \(2021\)](#); [Agarwal et al. \(2019\)](#); [Zhang et al. \(2022\)](#). For simplicity, we suppose that S is finite, but may be exponentially large, as we suffer no dependence on $|S|$. The proof can be easily extended to infinite state spaces by replacing inner products with integrals.

Consider the following quantity,

$$b_{h,e} = \sum_{k=1}^{K-1} (s_{h+1}^k) b_h(s_h^k; a_h^k)^\top \left(\sum_{j: k;h(s_h^k)} b_{k;h}(s_h^k) f(s_h^k) \right) \quad \text{argmin}_{s^0 \in \mathcal{R}^{S^d}} \sum_{k=1}^{K-1} \|b_h(s_h^k; a_h^k) - (s_{h+1}^k) b_{k;h}(s_h^k)\|_2 + k \|f\|_2^2;$$

where (s) is a one-hot encoding of the state. In words, this is the best choice for linearly (in $b_h(s; a)$) predicting $E_{s^0 \sim P_h^?(s; a)} [f(s^0)] = P_h^?(s^0; s; a)$. We highlight that this is just a quantity for analysis and not computed in the algorithm. Finally, denote

$$\begin{aligned} \tilde{b}_{h,e} &= b_{h,e} b_h; \\ \tilde{e}_h &= e_h b_h; \end{aligned}$$

We will also sometimes use the shorthand $\tilde{b}_h(s; a)$ for $E_{s^0 \sim P_h^?(s; a)} [f(s^0)]$.

For each $h = 0; 1; \dots; H-1$, let V_h denote the class of functions

$$\mathcal{V}_h = \left\{ \sum_{a \in \mathcal{A}} w_a^\top b_h(s; a) + r_h(s; a) + \sum_{k=1}^{K-1} \|b_{k;h}(s; a)\|_2 \mid w \in \mathbb{R}^M; e \in [0; B]; I \text{ symmetric} \right\}$$

The motivation behind this construction is that it satisfies the key property that all of the learned value functions $\tilde{b}_{h,e}$ during [Algorithm 6](#) are captured in this class.

Lemma G.3 For any $h = 0; 1; \dots; H-1$,

1. $\sup_s \|\tilde{b}_{h,e}(s)\|_2 \leq M_V$.
2. For any $e = 1; 2; \dots; N$, we have $\tilde{b}_{h,e} \in \mathcal{V}_h$.
3. If $\tilde{b}_{h,e} \in \mathcal{V}_h$, we have $\sup_s \|f(s)\|_2 \leq M_V$.

Proof. Recall that

$$\tilde{b}_{h,e}(s) = \sum_{a \in \mathcal{A}} w_a^\top b_h(s; a) + r_h(s; a) + \sum_{k=1}^{K-1} \|b_{k;h}(s; a)\|_2 \quad M_V$$

$$\text{where } v_{h,e} = \sum_{k=1}^K b_h(s_h^k; a_h^k) v_{h+1,e}(s_{h+1}^k):$$

From the thresholding, we have

$$v_{h,e}(s) \leq M_V:$$

We can bound the norm of $v_{h,e}$ as follows,

$$\|v_{h,e}\|_k \leq \sum_{k=1}^K b_{h+1,e}(s_{h+1}^k) \leq N \sup_s v_{h+1,e}(s) \leq N M_V:$$

We also required B , and we regularized the covariance with λ_{\min} is at least 1. Hence $v_{h,e}$ satisfies all the conditions to be M_h . \square

Now we control the metric entropy of \mathcal{F}_h in \mathcal{L}_1 , i.e. $d(f_1; f_2) = \sup_s |f_1(s) - f_2(s)|$ for $f_i \in \mathcal{F}_h$. Lemma G.4 Let $\epsilon > 0$ be arbitrary and let \mathcal{N}_ϵ be the smallest-net with ϵ of \mathcal{F}_h . Then,

$$\log |\mathcal{N}_\epsilon| \leq d \log(1 + 6L/\epsilon) + \log(1 + 6B/\epsilon) + d^2 \log(1 + 18B^2/\epsilon^2):$$

Proof. Let $f_1, f_2 \in \mathcal{F}_h$. Then,

$$\begin{aligned} & |f_1(s) - f_2(s)| \\ & \leq \max_a (w_1 - w_2)^T b_h(s; a) + \sum_{i=1}^2 |b_h(s; a_i) - b_h(s; a_{i+1})| \\ & \leq k |w_1 - w_2|_2 + \max_a (|b_h(s; a) - b_h(s; a_{i+1})| + |b_h(s; a_{i+1}) - b_h(s; a_{i+2})|) \\ & \leq k |w_1 - w_2|_2 + j |a_{i+1} - a_{i+2}| + B \max_a |b_h(s; a) - b_h(s; a_{i+1})| + (|a_{i+1} - a_{i+2}|) \\ & \leq k |w_1 - w_2|_2 + j |a_{i+1} - a_{i+2}| + B \sqrt{\frac{1}{d} \sum_{i=1}^d |a_{i+1} - a_{i+2}|^2}; \end{aligned}$$

where we used for any $a, b \in \mathbb{R}^d$, we have $\|a - b\|_1 \leq \sqrt{d} \|a - b\|_2$. Now proceeding like the Lemma 8.6 in the RL Theory Monograph (Agarwal et al., 2019), we have the result. \square

In this section, we'll use the following bonus scaling parameter,

$$:= O \left(\sqrt{\frac{1}{N d^m}} M_V + M_V M \sqrt{\frac{1}{d \log(d N M_V)}} \right): \quad (7)$$

The following high probability event ($\mathcal{E}_{\text{model}}$) is a key step in our proof. Essentially, Theorem G.1 guarantees that, for all functions v_h , the model we learn is an accurate predictor of the expectation, up to a bonus and some vanishing terms.

For all the following lemmas and theorems, suppose Assumption G.1 and the bonus scaling as in (7). Throughout the section, $v_h(\cdot)$ refers to indicator functions of the trajectory, where $\tau_h = (s_0; s_1; \dots; s_h)$. As before, the expectations $\mathbb{E}[g(\tau_h)]$ are with respect to the distribution of trajectories when π is executed in the environment \mathcal{P} .

Theorem G.1 Let $\epsilon \in (0, 1)$. Then, w.p. $1 - \epsilon$, for any time h , episode e , indicator functions $v_1; \dots; v_H$, and policy π , we have

$$\sup_{f \in \mathcal{F}_V} \mathbb{E} \left[\sum_{h=1}^H \sum_{e=1}^E |b_{h,e}^e(s_h; a_h) - P_h^e(s_h; a_h) f_h(\tau_h)| \right] \leq \mathbb{E} \left[\sum_{h=1}^H |b_{h,e}^e(s_h; a_h) - v_h(\tau_h)| \right] + k V_h k_1 \epsilon^m: \quad (\mathcal{E}_{\text{model}})$$

Proof. Condition on the outcome of Lemma G.5, which implies that w.p. 1, for any $h; e; \epsilon; \delta; \eta$, we have

$$\sup_{f \in \mathcal{F}_{2V_h}} E \left[\mathbb{P}_{h;e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h) \right] f_h(s_h) \leq E \left[b_{h;e}(s_h; a_h) \right]$$

Also, for any $h; e; \epsilon; \delta; \eta$, by Assumption G.1, we have (w.p. 1) that

$$\begin{aligned} \sup_{f \in \mathcal{F}_{2V_h}} E \left[\mathbb{P}_h(s_h; a_h) - P_h^?(s_h; a_h) \right] f_h(s_h) &\leq E \left[\sup_{f \in \mathcal{F}_{2V_h}} \mathbb{P}_h(s_h; a_h) - P_h^?(s_h; a_h) \right] f_h(s_h) \\ &\leq E \left[\sup_{f \in \mathcal{F}_{2V_h}} \mathbb{P}_h(s_h; a_h) - P_h^?(s_h; a_h) \right] \\ &\leq \epsilon \end{aligned}$$

Combining these two yields the result, as

$$\begin{aligned} \sup_{f \in \mathcal{F}_{2V_h}} E \left[\mathbb{P}_{h;e}(s_h; a_h) - P_h^?(s_h; a_h) \right] f_h(s_h) \\ \leq \sup_{f \in \mathcal{F}_{2V_h}} E \left[\mathbb{P}_h(s_h; a_h) - P_h^?(s_h; a_h) \right] f_h(s_h) + \sup_{f \in \mathcal{F}_{2V_h}} E \left[\mathbb{P}_{h;e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h) \right] f_h(s_h) \end{aligned}$$

□

Lemma G.5 Suppose Assumption G.1 and the bonus scaling is set as in (7). For any $\epsilon \in (0, 1)$, w.p. at least $1 - \epsilon$, we have for any time t , episode e , and policy π ,

$$b_{h;e}(s_h; a_h) := \sup_{f \in \mathcal{F}_{2V_h}} \mathbb{P}_{h;e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h) + b_{h;e}(s_h; a_h)$$

Proof. Consider any $h; e; \epsilon; \delta; \eta$. Define $\mu_h^k := \mathbb{P}_h(s_{h+1}^k) + P_h^?(s_{h+1}^k | s_h^k; a_h^k)$, so that $E[\mu_h^k | H_{k-1}] = 0$, where H_{k-1} contains the states and actions before episode k . In what follows, we slightly abuse notation, as $\mathbb{P}(s; a) b^T(s; a)$ will denote the outer product, and hence a $d \times d$ quantity.

$$\begin{aligned} b_{h;e} &= \sum_{k=1}^K (s_{h+1}^k) b_h(s_h^k; a_h^k)^T \\ &= \sum_{k=1}^K P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k) b_h(s_h^k; a_h^k)^T + \sum_{k=0}^K \mathbb{P}_h(s_h^k; a_h^k) - \mu_h^k b_h(s_h^k; a_h^k)^T \\ &= \sum_{k=1}^K P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k) b_h(s_h^k; a_h^k)^T + e_h(s_{h;e}) + \sum_{k=0}^K \mu_h^k b_h(s_h^k; a_h^k)^T \end{aligned}$$

Rearranging, we have

$$\begin{aligned} b_{h;e} - e_h &= \sum_{k=0}^K P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k) b_h(s_h^k; a_h^k)^T - \sum_{k=0}^K \mu_h^k b_h(s_h^k; a_h^k)^T \\ &= e_h(s_{h;e}) + \sum_{k=0}^K \mu_h^k b_h(s_h^k; a_h^k)^T \end{aligned}$$

Now let $f \in \mathcal{F}_{2V_h}$ be arbitrary. For any $s_h; a_h$, multiply the above with $b_h(s_h; a_h)$ and multiply with f , we have

$$\begin{aligned} &\mathbb{P}_{h;e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h) f \\ &= f^T (b_{h;e} - e_h) b_h(s_h; a_h) \end{aligned}$$

$$\begin{aligned}
& \left| \sum_{k=1}^{\infty} f^T P_h^k(s_h^k; a_h^k) \mathbb{P}_h(s_h^k; a_h^k) b_h(s_h^k; a_h^k)^T \right|_{h,e} \mathbb{1} b_h(s_h; a_h) \\
& \text{Term(a)} \\
& + \left| f^T e_h \right|_{h,e} \mathbb{1} b_h(s_h; a_h) \\
& \text{Term(b)} \\
& + \left| \sum_{k=1}^{\infty} \mathbb{1}^k b_h(s_h^k; a_h^k)^T \right|_{h,e} \mathbb{1} b_h(s_h; a_h) : \\
& \text{Term(c)}
\end{aligned}$$

We can deterministically bound Term (b) as follows,

$$\begin{aligned}
& \sup_{f \in 2V_h} f^T e_h \mathbb{1} b_h(s_h; a_h) \\
& = \sup_{f \in 2V_h} (\mathbb{1}^2 f^T e_h)^T \mathbb{1}^2 b_h(s_h; a_h) \\
& \sup_{f \in 2V_h} \mathbb{1}^2 k f^T e_h \mathbb{1} b_h(s_h; a_h) \\
& k V_h k_1 M \mathbb{1}^2 b_h(s_h; a_h) : \quad (\text{by Assumption G.1})
\end{aligned}$$

This term will be lower order compared to the other two.

We now derive the bound for Term (c) for any fixed $f \in 2V_h$. Observe that

$$\begin{aligned}
f^T \sum_{k=1}^{\infty} \mathbb{1}^k b_h(s_h^k; a_h^k)^T \mathbb{1} b_h(s_h; a_h) & = \sum_{k=1}^{\infty} \mathbb{1}^2 b_h(s_h^k; a_h^k) (f^T \mathbb{1}^k) \mathbb{1}^2 b_h(s_h; a_h) \\
& \sum_{k=1}^{\infty} \mathbb{1}^k b_h(s_h^k; a_h^k) (f^T \mathbb{1}^k) \mathbb{1}^2 b_h(s_h; a_h) :
\end{aligned}$$

Now we argue w.p.1, for any $e; h$ we have

$$\sum_{k=1}^{\infty} \mathbb{1}^k b_h(s_h^k; a_h^k) (f^T \mathbb{1}^k) \mathbb{1}^2 b_h(s_h; a_h) \leq 2kV_h k_1 \mathbb{1}^2 \sqrt{2 \log(1/\delta) + d \log(N+1)} ;$$

which implies the claim about $\mathbb{1}^2 b_h(s_h; a_h)$. Indeed, we can apply Lemma G.1. Checking the preconditions, $E_{P_h^k(s_h^k; a_h^k)} f^T \mathbb{1}^k j H_{k-1} = 0$, $\mathbb{1}^2 f^T \mathbb{1}^k j k f k_1 k \mathbb{1}^k k_1 \leq 2kV_h k_1$, $\det(\mathbb{1}^2) = \det I = 1$, and $\det(\mathbb{1}^2) = \det(\mathbb{1}^2) (e+1)^d$ since the largest eigenvalue is 1. So, w.p. at least $1-\delta$, for all e , we have the above inequality.

Thus, for any fixed $f \in 2V_h$, w.p. 1, for all $e; h$ we have,

$$\begin{aligned}
& \mathbb{1}^2 b_h(s_h; a_h) \mathbb{1}^2 b_h(s_h; a_h) f \\
& \text{Term(a)} + \text{Term(b)} + \text{Term(c)} \\
& 4kV_h k_1 (1 + M) \mathbb{1}^2 \sqrt{2 \log(1/\delta) + d \log(N)} + \mathbb{1}^2 \sqrt{2N} kV_h k_1 \mathbb{1}^2 b_h(s_h; a_h) \\
& + kV_h k_1 M \mathbb{1}^2 b_h(s_h; a_h) \\
& + 4kV_h k_1 \mathbb{1}^2 \sqrt{2 \log(1/\delta) + d \log(N)} \mathbb{1}^2 b_h(s_h; a_h) \\
& \mathbb{1}^2 \sqrt{2N} kV_h k_1 \mathbb{1}^2 b_h(s_h; a_h) + kV_h k_1 M \mathbb{1}^2 \sqrt{2 \log(1/\delta) + d \log(N)} \mathbb{1}^2 b_h(s_h; a_h) :
\end{aligned}$$

Now we apply a covering argument. Namely, union bound the above argument to every element in an ϵ -net of V_h . For any $s \in V_h$, let s^e be its neighbor in the net such that $\|s - s^e\| \leq \epsilon$, so we have

$$|b_{h,e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h)| \leq |b_{h,e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h)| + |b_{h,e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h)| \epsilon$$

and

$$|b_{h,e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h)| \leq \epsilon + \epsilon k_1 (N+1) \epsilon$$

Setting $\epsilon = \sqrt{\frac{1}{N}}$, the metric entropy is of the order $\log(N(M_V + B)) + \log(BN) + d^2 \log(BdN)$. The error incurred with this epsilon net is a constant, which is lower order.

Thus, we have

$$\begin{aligned} \delta_{s_h; a_h} &: \sup_{f \in \mathcal{F}_h} |b_{h,e}(s_h; a_h) - \mathbb{P}_h(s_h; a_h)| \\ &\leq \sqrt{\frac{1}{N}} \sqrt{M_V + B} + \sqrt{\frac{1}{N}} \sqrt{M_V + B} \sqrt{\log(1 + d \log(M_V) + d^2 \log(BdN))} \\ &\leq \sqrt{\frac{1}{N}} \sqrt{M_V + B} \sqrt{\log(1 + d \log(M_V) + d^2 \log(BdN))} \end{aligned}$$

Note that δ scales as $\frac{1}{\sqrt{N}}$, so one can find a valid B by solving B for B .

□

Lemma G.6 Let $f \in \mathcal{F}_h$. For any $\epsilon \in (0, 1)$, w.p. at least $1 - \epsilon$, for any time h , episode e , we have

$$\begin{aligned} \delta_{s_h; a_h} &: \mathbb{E} \left[\sum_{k=1}^K |P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k)| \right] \\ &\leq \sqrt{\frac{1}{N}} \sqrt{M_V + B} \sqrt{\log(1 + d \log(N))} + \sqrt{\frac{1}{N}} \sqrt{M_V + B} \end{aligned}$$

Proof. First observe that

$$\begin{aligned} &\mathbb{E} \left[\sum_{k=1}^K |P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k)| \right] \\ &= \mathbb{E} \left[\sum_{k=1}^K |b_h(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k)| \right] \\ &\leq \sqrt{\frac{1}{N}} \sqrt{M_V + B} \sqrt{\log(1 + d \log(N))} + \sqrt{\frac{1}{N}} \sqrt{M_V + B} \end{aligned}$$

where $b_k = |P_h^?(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k)|$.

Now we will argue that w.p.1, for all $e; h$,

$$\sum_{k=1}^K |b_h(s_h^k; a_h^k) - \mathbb{P}_h(s_h^k; a_h^k)| \leq \sqrt{\frac{1}{N}} \sqrt{M_V + B} \sqrt{\log(1 + d \log(N))} + \sqrt{\frac{1}{N}} \sqrt{M_V + B}$$

which will imply the claim for all $s_h; a_h$.

Apply self-normalized martingale concentration (Lemma G.1) to $b_i = b_h(s_h^i; a_h^i)$ and $\sigma_i = \sqrt{b_h(s_h^i; a_h^i)}$, where the expectation is over $(s_h^i; a_h^i)$ in the definition of b_h . To see sub-Gaussianity, bound the envelope $|b_h(s_h^i; a_h^i)| \leq \sqrt{M_V + B}$, and thus

$j^k - j^k \leq 2kV_h k_1 (1 + M)$. Now compute the determinant $\det(\mathbf{A}_{h,e}) = 1$ and since $\max(\mathbf{A}_{h,e}) \leq e + 1$, we have that $\log \det(\mathbf{A}_{h,e}) \leq d \log(e + 1)$. Hence, w.p. at least $1 - \delta$, we have

$$\mathbb{E} \left[\sum_{k=1}^K b_h(s_h^k; a_h^k) \left(\mathbb{E} [b_k | H_{k-1}] - 2kV_h k_1 (1 + M) \right)^2 \right] \leq \frac{d \log(e + 1) + d \log(N + 1)}{2 \log(1 - \delta)}$$

By [Assumption G.1](#) applied to k (the data-generating policy for episode k), we have $\mathbb{E} [b_k | H_{k-1}] \leq \frac{kV_h k_1}{\sigma^2} \mu_{ms}$. Recall for any scalars α and vectors x_i , we have $\sum_{i=1}^d \alpha^2 x_i^2 \leq \alpha^2 \sum_{i=1}^d x_i^2$. Thus,

$$\begin{aligned} & \sum_{k=1}^K \frac{\sum_{i=1}^d b_h(s_h^k; a_h^k) \mathbb{E} [b_k | H_{k-1}]^2}{k^2 b_h(s_h^k; a_h^k)} \leq \sum_{k=1}^K \frac{\sum_{i=1}^d \mathbb{E} [b_k | H_{k-1}]^2}{k^2} \\ & \leq \frac{d \sum_{k=1}^K \frac{1}{k^2}}{(\sigma^2)^2} \mu_{ms}^2 \leq \frac{d \mu_{ms}^2}{(\sigma^2)^2} \end{aligned} \quad (\text{by Lemma G.2})$$

Combining these two bounds concludes the proof. \square

Lemma G.7 (Optimism) Suppose [\(E_{model}\)](#) holds. Let $\mu = kV_h k_1 \mu_{ms}$. Then, for any episode $e = 1; 2; \dots; N$, we have

$$\mathbb{E} \left[\sum_{h=0}^H \sum_{i=1}^d \left(Q_{h,e}^2(s_h; a_h) - Q_{h,e}(s_h; a_h) \right)^2 \right] \leq \mu^2 \sum_{h=0}^H \sum_{i=1}^d (H - h);$$

and

$$\mathbb{E} \left[\sum_{h=0}^H \sum_{i=1}^d \left(V_h^2(s_h) - V_{h,e}(s_h) \right)^2 \right] \leq \mu^2 \sum_{h=0}^H \sum_{i=1}^d (H - h);$$

where

$$\begin{aligned} Q_{h,e}(s_h) &:= \mathbb{E} \left[Q_{h,e}(s_h; b_h^e(s_h)) \right] \\ V_h(s_h) &:= \mathbb{E} \left[V_h(s_h) \right] \\ V_{h,e}(s_h) &:= \mathbb{E} \left[V_{h,e}(s_h) \right] \end{aligned}$$

Abusing notation, $\mathbb{1}(\cdot)$ is the constant function 1.

In particular, we have that

$$\mathbb{E}_{d_0} \left[\sum_{h=0}^H \sum_{i=1}^d V_0^2(s_0) - V_{0,e}(s_0) \right] \leq \mu^2 H;$$

Proof. Fix any episode e . We prove both claims via induction on $n = H; H - 1; H - 2; \dots; 1; 0$. The base case holds trivially since $Q_{H,e}$ and V_H are zero at every state by definition. Indeed, we have that for any s , including s_0 , that

$$\begin{aligned} & \mathbb{E} \left[\sum_{h=0}^H \sum_{i=1}^d \mathbb{1}(s_{H-1}; a_{H-1}) (V_{H-1}^2 - V_{H-1,e}) \right] \\ & = \mathbb{E} \left[\sum_{h=0}^H \sum_{i=1}^d (0 - 0) \right] = 0; \end{aligned}$$

Now let's show the inductive step. Let $\mathcal{I} = \{H - 1; H - 2; \dots; 1; 0\}$ be arbitrary and suppose the inductive hypothesis. So suppose that [\(E_{model}\)](#) holds at $t + 1$ (we don't even need [\(E_{model}\)](#) in the future), i.e.

$$\mathbb{E} \left[\sum_{h=0}^H \sum_{i=1}^d \left(V_{h+1}^2(s_{h+1}) - V_{h+1,e}(s_{h+1}) \right)^2 \right] \leq \mu^2 \sum_{h=0}^H \sum_{i=1}^d (H - h);$$

$$(H - h - 1) \tag{IH}$$

Recalling that $Q_{h,e}(s_h; a_h) = r_h(s_h; a_h) + \beta_{h,e}(s_h; a_h) v_{h+1,e} + b_{h,e}(s_h; a_h)$, we have

$$\begin{aligned} & E_{\mathcal{D}_h} [Q_h^2(s_h; a_h) - Q_{h,e}(s_h; a_h)]_{h(h)} \\ &= E_{\mathcal{D}_h} [P_h^2(s_h; a_h) V_{h+1}^2 - \beta_{h,e}(s_h; a_h) v_{h+1,e} - b_{h,e}(s_h; a_h)]_{h(h)} \\ &= E_{\mathcal{D}_h} [P_h^2(s_h; a_h) - \beta_{h,e}(s_h; a_h) v_{h+1,e} - b_{h,e}(s_h; a_h)]_{h(h)} + (H - h - 1) \tag{by (IH)} \\ &= E_{\mathcal{D}_h} [\beta_{h,e}(s_h; a_h) - P_h^2(s_h; a_h) v_{h+1,e}]_{h(h)} + E_{\mathcal{D}_h} [b_{h,e}(s_h; a_h)]_{h(h)} + (H - h - 1) \\ &+ (H - h - 1) = (H - h); \tag{by (E_{model}) and v_{h+1,e} \leq 2 V_h (Lemma G.3)} \end{aligned}$$

which proves the Q-optimism claim.

Now let's prove V-optimism.

$$\begin{aligned} & E_{\mathcal{D}_h} [V_h^2(s_h) - v_{h,e}(s_h)]_{h-1(h-1)} \\ &= E_{\mathcal{D}_h} [Q_h^2(s_h; a_h) - Q_{h,e}(s_h; b_h^e(s_h))]_{M_V, h-1(h-1)} \\ &= E_{\mathcal{D}_h} [(Q_h^2(s_h; a_h) - M_V)_{h-1(h-1)} (1 - \beta_h(s_h))] \\ &+ E_{\mathcal{D}_h} [Q_h^2(s_h; a_h) - Q_{h,e}(s_h; b_h^e(s_h))]_{h-1(h-1), h(s_h)} \\ &= E_{\mathcal{D}_h} [Q_h^2(s_h; a_h) - Q_{h,e}(s_h; b_h^e(s_h))]_{h(h)} \\ &= E_{\mathcal{D}_h} [Q_h^2(s_h; a_h) - Q_{h,e}(s_h; \beta_h^2(s_h))]_{h(h)} \\ &(H - h); \end{aligned}$$

by Q-optimism. □

Remark G.1 We did not require $\beta_{h,e}$ to be a valid transition! It is in general unbounded and can even have negative entries!

Lemma G.8 (Simulation) For any episode $e = 1; 2; \dots; N$, we have

$$E_{\mathcal{D}_0} [v_{0,e}(s_0) - V_0^e(s_0)] \leq \sum_{h=0}^{H-1} E_{\mathcal{D}_e} [b_{h,e}(s_h; a_h) + (\beta_h(s_h; a_h) - P_h^2(s_h; a_h)) v_{h+1,e}]$$

Proof. We progressively unravel the left hand side. For any

$$\begin{aligned} & v_{0,e}(s_0) - V_0^e(s_0) \\ &= Q_{0,e}(s_0; \beta_0^e(s_0)) - Q_0^e(s_0; \beta_0^e(s_0)) \\ &= b_{0,e}(s_0; \beta_0^e(s)) + \beta_{0,e}(s_0; \beta_0^e(s)) - P_0^2(s_0; \beta_0^e(s_0)) v_{1,e} + P_0^2(s_0; \beta_0^e(s_0)) v_{1,e} - V_1^e; \end{aligned}$$

where the inequality is due to the thresholding on the value function. Now, perform this recursively on the $P_0^2(s_0; \beta_0^e(s_0)) v_{1,e} - V_1^e$ term. Doing this unravelling h times gives the result. □

Theorem G.2 Suppose [Assumption G.1](#) and let α and β be defined as [\(7\)](#). Let $\gamma \in (0, 1)$. Then w.p. at least $1 - \gamma$, we have that the regret of LSVI is sublinear,

$$NV \leq \sum_{e=0}^{H-1} V^{b^e} \leq \mathcal{O}(d \text{HNM}_V^p \frac{1}{\log(\text{HNM}_V)} + d^{1.5} H^p \overline{\text{NM}}_V M \log(d \text{HNM}_V))$$

where \mathcal{O} hides log dependence.

Proof. We first condition on the high-probability event $\mathcal{E}_{\text{model}}$, which occurs w.p. at least $1 - \delta$. Fix any arbitrary episode e . By optimism [Lemma G.7](#) and the simulation lemma [Lemma G.8](#),

$$\begin{aligned} E_{d_0} V_0^?(s_0) - V_0^{b^e}(s_0) &= E_{d_0} \sum_{h=0}^{H-1} (V_h^?(s_h) - V_h^{b^e}(s_h)) + H \\ &\leq \sum_{h=0}^{H-1} E_{b^e} [b_{h,e}^e(s_h; a_h) + \sum_{i=h+1}^H P_h^?(s_h; a_h) V_{h+1,e}^i] + H \end{aligned}$$

Applying [\(E_{model}\)](#) with no indicators, i.e. $\mathbb{1}_{\mathcal{E}_{\text{model}}}(s_h) = 1$ always, gives,

$$\sum_{h=0}^{H-1} E_{b^e} [2 b_{h,e}^e(s_h; a_h)] + 2H :$$

Now, summing over $e = 1; 2; \dots; N$, we have

$$\begin{aligned} \sum_{e=1}^N E_{d_0} \sum_{h=0}^{H-1} (V_h^?(s_h) - V_h^{b^e}(s_h)) & \\ \leq 2HN + 2 \sum_{h=0}^{H-1} \sum_{e=1}^N E_{b^e} [b_{h,e}^e(s_h; a_h)] \end{aligned}$$

By Azuma's inequality applied to the martingale difference $\xi_{h,e} = E_{b^e} [b_{h,e}^e(s_h; a_h) | \mathcal{F}_{h-1}^e] - b_{h,e}^e(s_h^e; a_h^e)$, which has envelope bounded by β , implies w.p. $1 - \delta$,

$$2HN + 2 \sum_{h=0}^{H-1} \sum_{e=1}^N b_{h,e}^e(s_h^e; a_h^e) + 4 \sqrt{p \sum_{h=0}^{H-1} \sum_{e=1}^N \beta^2} :$$

It remains to bound the sum of expected bonuses. By [Lemma G.2](#), we know that almost surely,

$$\sum_{h=0}^{H-1} \sum_{e=1}^N b_{h,e}^e(s_h^e; a_h^e) \leq H \sqrt{dN \log(N)} :$$

So, putting everything together,

$$\begin{aligned} \sum_{e=1}^N E_{d_0} \sum_{h=0}^{H-1} (V_h^?(s_h) - V_h^{b^e}(s_h)) & \\ \leq HN + H \sqrt{dN \log(N)} + \sqrt{p \sum_{h=0}^{H-1} \sum_{e=1}^N \beta^2} & \\ \leq HN + \sqrt{p d N M_V^2 m_s} + M_V M \sqrt{d \log(dHN)} + H \sqrt{p d N \log(HN)} & \\ = HN + dHN M_V \sqrt{\log(HN)} m_s + d^{1.5} H \sqrt{p N M_V M \log(dHN)} & \end{aligned}$$

Note that $H \sqrt{dN \log(N)} = H \sqrt{dN} \sqrt{\log(N)} = H M_V \sqrt{m_s}$ is of lower order (with respect to N), we can simply drop it. This concludes the proof. \square

Corollary G.1 By setting $\beta = 1/N$, we have that expected regret also has the same rate as above.

Proof. The expected regret by law of total probability, since regret is at most

$$\begin{aligned} E[\text{Reg}_N] &= E[\text{Reg}_N | \mathcal{E}_{\text{model}}] + NH(1 - P(\mathcal{E}_{\text{model}})) \\ &= E[\text{Reg}_N | \mathcal{E}_{\text{model}}] + H : \end{aligned}$$

Since H is lower-order, we have the same rate. \square

H Proof of Main Theorems

First we prove [Theorem 4.1](#) and [Theorem 4.2](#).

Theorem 4.1 (Regret under generative source access) Suppose Assumptions 3.1-3.4 and 4.1, and suppose the input policies π_k are ϵ -exploratory. Then, for any $\delta \in (0, 1)$, w.p. $1 - \delta$, REPTRANSFER when deployed in the target task has regret at most $\mathcal{O}(H^2 d^{1.5} \sqrt{T \log(1/\delta)})$, with at most Kn generative accesses per source task, with $\mathcal{O}(\frac{1}{\epsilon} A^3_{\max} K T \log \frac{j}{j'} + K \log j)$.

Theorem 4.2 (Regret for REPTRANSFER with REWARDFREE subroutine) Suppose Assumptions 3.1-3.4 and 4.1 hold, and let $\delta \in (0, 1)$. Let $f_{k, g_{k=1}^K}$ be exploratory policies learned from running REWARDFREE on each source task with $N_{\text{LSVI-UCB}}$ and $N_{\text{REWARDFREE}}$ set as in Lemma 4.2. Then, w.p. $1 - \delta$, running REPTRANSFER with $f_{k, g_{k=1}^K}$ has regret in the target task of $\mathcal{O}(H^2 d^{1.5} \sqrt{T \log(1/\delta)})$, with at most $\mathcal{O}(A^4 \frac{3}{\max} d^5 H^7 K^2 T^2 (\log(j/j') + K \log j))$ generative accesses per source task.

Proof of Theorem 4.1 and Theorem 4.2 For the regret bound, set $M_V = H$ and $M = \frac{1}{\epsilon}$ and apply Theorem G.2. This choice of M is valid by the argument following Assumption G.1. This gives us a regret bound of

$$\mathcal{O}(dH^2 T \epsilon_{\text{ms}} + d^{1.5} H^2 \sqrt{T \log(1/\delta)});$$

where ϵ_{ms} can be made smaller than $\frac{\delta}{T}$, in which the second term dominates.

Now, we calculate the pre-training phase sample complexity in a source task.

First, let's calculate the reward-free model learning sample complexity, i.e. this is the number of samples required for learning θ_k . Recall that we need this to be sufficiently large such that $\|T_V - 1\| = N_{\text{LSVI-UCB}}$. As required by Lemma 4.2, we need,

$$\begin{aligned} N_{\text{REWARDFREE}} &= \mathcal{O}(A^3 d^4 H^6 \log(j/j')) N_{\text{LSVI-UCB}}^2 \\ &= \mathcal{O}(A^3 d^4 H^6 \log(j/j')) A^3 d^6 H^8 \sqrt{2}^2 \\ &= \mathcal{O}(A^9 d^{16} H^{22} \sqrt{4} \log(j/j')) : \end{aligned}$$

Second, we calculate the cross-sampling sample complexity. Recall that the number of samples in each pairwise dataset. In order to reduce to $\frac{\delta}{T}$, by Lemma 4.1, we need

$$\begin{aligned} \frac{\delta}{T} \epsilon_{\text{ms}} &= A^3_{\max} K = \min_{s=1,2} \frac{\delta}{n} \\ &= A^3_{\max} K = \min_{s=1,2} \frac{1}{n} \log \frac{j}{j'} + K \log j/j' \quad (\text{by (2)}) \end{aligned}$$

which implies that we need

$$n \geq \frac{1}{\epsilon} A^3_{\max} K T \log \frac{j}{j'} + K \log j/j' :$$

Incorporating the coverage result from Lemma 4.2 gives,

$$n \geq A^4 \frac{3}{\max} d^5 H^7 K T^2 \log \frac{j}{j'} + K \log j/j' :$$

Since each task is in at most $\sqrt{2}$ pairwise datasets, each of size n , the total pre-training sample complexity per task is at most,

$$\begin{aligned} &N_{\text{REWARDFREE}} + (\sqrt{2} - 1) n \\ &= \mathcal{O}(A^9 d^{16} H^{22} \sqrt{4} \log(j/j')) + A^4 \frac{3}{\max} d^5 H^7 K^2 T^2 \log \frac{j}{j'} + K \log j/j' : \end{aligned}$$

□

Now we prove [Theorem 5.2](#), restated below.

[Theorem 5.2](#) (Regret with online access) Suppose Assumptions [3.1-3.4, 5.1, 5.2](#) hold. For any $\epsilon \in (0, 1)$, w.p. $1 - \epsilon$, [Algorithm 5](#) with appropriately set parameters achieves a regret in the target task at most $\epsilon^{-1} d^{1.5} H^2 \sqrt{T \log(1/\epsilon)}$, with at most $\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))$ online queries in the source tasks.

[Proof of Theorem 5.2](#) We follows the same format as the proof of [Theorem 4.1](#). The regret bound is identical.

Now let's compute the pre-training sample complexity. The regret bound requires us to set $\epsilon = \frac{1}{T}$. Here, our n_{ms} comes from [Lemma F.1](#), so

$$n_{\text{ms}} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\epsilon} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\frac{1}{T}}$$

which implies we need

$$n_{\text{ms}} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\epsilon} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\frac{1}{T}}$$

Plugging in the coverage of [Lemma 4.2](#),

$$n_{\text{ms}} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\epsilon} = \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\frac{1}{T}}$$

Here, we only collect one dataset, so the total pre-training sample complexity is

$$N_{\text{REWARDFREE}} + n_{\text{ms}} = \text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =)) + \frac{\text{poly}(A; \max_{\text{raw}} d; H; K; T; \epsilon^{-1}; \frac{1}{\text{raw}}; \log(j \cdot j \cdot j =))}{\epsilon}$$

□

I Experiment Details

I.1 Construction of Comblock

In this section we first introduce the vanilla Combination lock (Comblock) environment that is widely used as the benchmark for algorithms for Block MDPs. We provide a visualization of the comblock environment in [Fig. 1\(a\)](#). Concretely, the environment has a horizon of 3 latent states $z_{i,h}; i \in \{0, 1, 2\}$ for each timestep h and 10 actions. Among the three latent states, we denote $z_{0,h}$ and $z_{1,h}$ as the good states which leads to the final reward and the bad states. At the beginning of the task, the environment will uniformly and independent sample 1 out of the 10 actions for each good state $z_{0,h}$ and $z_{1,h}$ for each timestep h , and we denote these actions $a_{0,h}$ and $a_{1,h}$ as the optimal actions (corresponding to each latent state). These optimal actions, along with the task itself, determines the dynamics of the environment. At each good latent state $s_{i,h}$, if the agent takes the correct action, the environment transits to the either good state at the next timestep (i.e., $s_{i,h+1}$) with equal probability. Otherwise, if the agent takes any 9 of the bad actions, the environment will transition to the bad state $s_{2,h+1}$ deterministically, and the bad states transit to only bad states at the next timestep deterministically. There are two situations where the agent receives a reward: one is upon arriving the good states at the last timestep, the agent receives a reward of 1. The other is upon the first ever transition into the bad state, the agent receives an “anti-shaped” reward of 0.1 with

probability 0.5. Such design makes greedy algorithms without strategic exploration such as policy optimization methods easily fail. For the initial state distribution, the environment starts for $s_{1,0}$ with equal probability. The dimension of the observation is $d \times (H + jS_j + 1) \times e$. For the emission distribution, given a latent state $s_{t,h}$, the observation is generated by first concatenate the one-hot vectors of the state and the horizon, adding $\text{N}(0,0.1)$ noise at each entry, appending 0 at the end if necessary. Then finally we apply a linear transformation on the observation with a Hadamard matrix. Note that without a good feature or strategic exploration, it takes actions to reach the final goal with random actions.

I.2 Construction of transfer setup in Section. 6.1

In this section we introduce the detailed construction of the experiment in Section. 6.1. For the source environment, we simply generate 5 random vanilla comblock environment described in Section. I.1. Note that in this way we ensure that the emission distribution shares across the sources, but the latent dynamics are different because the optimal actions are independently randomly selected. For the target environment, for each timestep, we randomly acquire the optimal actions from one of the sources and set it to be the optimal action of the target environment at time t if the selected optimal actions are different for the two good states. Otherwise we keep sampling until they are different. Note that under such construction, since we fix the emission distribution, Assumption 3.4 is satisfied if we set $\beta = 1$ for the source environment where we select the optimal action and 0 for the other sources, at each timestep. To see how Assumption 5.1 is satisfied, recall that comblock environment naturally satisfies Assumption 3.3, and identical emission implies that the conditional ratio of all observations between source and target is 1.

I.3 Construction of transfer setup in Section. 6.2

Now we introduce the construction of the Comblock with Partitioned Observation (Comblock-PO) environment. Comparing with the vanilla comblock environment, the major difference is in the observation space. In this setting, the size of the observation depends on the number of source environments K . Let the size of the original observation space be $|O|$, the size of the observation for comblock-PO is $K|O|$. For the k -th source environment, where $k \in [K]$, the environment first generates the $|O|$ -dimensional observation vector as in the original comblock, and then embed it to the $(k-1)|O|$ -th to $k|O|$ -th entries of the $K|O|$ -dimensional observation vector, where it is 0 everywhere else. Thus we can see that the observation space for each source environment is disjoint (and thus the name partitioned observations). For the target environment, since the latent dynamics are the same, we only need to design the emission distribution: for each latent state, we assign the emission distribution uniformly at random from one of the sources.

I.4 Implementation details

Our implementation builds on BRIEE (Zhang et al., 2022)⁵. In the Multi-task REPLEARN stage, we require our learned feature to predict the Bellman backup of all the sources simultaneously. Therefore, in each iteration we have k discriminators and k sets of linear weights (instead of 1 in BRIEE), where k is the number of source environments. For the deployment stage we implement LSVI following Algorithm 6.

To create the training dataset for Multi-task REPLEARN, for each $(i; j)$ environment pairs where $i \in j$, we collect 500 samples for each timestep. For each $(i; i)$ environment pairs, we collect 500 $(k-1) \times k$ samples for each timestep where k denotes the number of sources. Thus we ensure that the number of samples from cross transition of different environments is the same as the number of samples from cross transition of the same environment. For the online setting, we simply sample 1000 $(k-1) \times k$ samples for each $(i; i)$ cross transition to ensure that the total number of samples is the same for GERTTRANSFER and O-REPTRANSFER.

⁵Code based on public repository <https://github.com/yudasong/briee>.

To sample the initial state action pair $(s; a)$ pair as in (1), for 90% of the samples, we follow the policy from each source environment trained using REE. For the remaining 10%, we follow the same policy to state, and then take a uniform random policy to get s . With this sampling scheme we ensure that Assumption 3.3 is satisfied. In the setting of Section. 6.2, we follow a more simple procedure to ensure that the samples are more balanced among the three states: we skip the first sampling step from environment i , i.e., samples given $(s; a)$, and simply reset environment to s , where s is one of the three states with equal probability, and generate the observation accordingly. Note that such visitation distribution is also possible in the online setting with a more nuanced sampling procedure, and in the experiment we use the same sampling procedure for both G-REPTTRANSFER and O-REPTTRANSFER for a fair comparison.

1.5 Hyperparameters

In this section, we record the hyperparameters we try and the final hyperparameter we use for each baseline. The hyperparameters for REPTTRANSFER in Section. 6.1 is in Table. 3. The hyperparameters for REPTTRANSFER in Section. 6.2 is in Table. 4. The hyperparameters for REE is in Table. 5. We use the same set of hyperparameters for G-REPTTRANSFER and O-REPTTRANSFER.

Table 3: Hyperparameters for REPTTRANSFER in Comblock.

	Value Considered	Final Value
Decoder learning rate	{1e-2}	1e-2
Discriminator learning rate	{1e-2}	1e-2
Discriminator hidden layer size	{256}	256
RepLearn Iteration	{30}	30
Decoder number of gradient steps	{64}	64
Discriminator number of gradient steps	{64}	64
Decoder batch size	{256}	256
Discriminator batch size	{512}	512
RepLearn regularization coefficient	{0.01}	0.01
Decoder softmax temperature	{1}	1
Decoder β_0 softmax temperature	{0.1}	0.1
LSVI bonus coefficient	$\{1, \frac{H}{5}\}$	1
LSVI regularization coefficient	{1}	1
Buffer size	{1e5}	1e5
Update frequency	{50}	50
Optimizer	{SGD}	SGD

Table 4: Hyperparameters for REPTTRANSFER in Comblock-PO.

	Value Considered	Final Value
Decoder learning rate	{1e-2}	1e-2
Discriminatorf learning rate	{1e-2}	1e-2
Discriminatorf hidden layer size	{256,512}	256
Discriminatorf hidden layer number	{2,3}	3
RepLearn IteratioT	{30,40,50,100,150}	50
Decoder number of gradient steps	{64,80,128,256}	64
Discriminatorf number of gradient steps	{64,80,128,256}	64
Decoder batch size	{256,512}	512
Discriminatorf batch size	{256,512}	512
RepLearn regularization coef cient	{0.01}	0.01
Decoder softmax temperature	{1}	1
Decoder ₀ softmax temperature	{0.1,1}	1
LSVI bonus coef cient	{1, $\frac{H}{5}$ }	1
LSVI regularization coef cient	{1}	1
Buffer size	{1e5}	1e5
Update frequency	{50}	50
Optimizer	{SGD, Adam}	Adam

Table 5: Hyperparameters for REE in Comblock and Comblock-PO.

	Value Considered	Final Value
Decoder learning rate	{1e-2}	1e-2
Discriminatorf learning rate	{1e-2}	1e-2
Discriminatorf hidden layer size	{256}	256
RepLearn IteratioT	{30}	30
Decoder number of gradient steps	{64}	64
Discriminatorf number of gradient steps	{64}	64
Decoder batch size	{512}	512
Discriminatorf batch size	{512}	512
RepLearn regularization coef cient	{0.01}	0.01
Decoder softmax temperature	{1}	1
Decoder ₀ softmax temperature	{0.1}	0.1
LSVI bonus coef cient	{ $\frac{H}{5}$ }	$\frac{H}{5}$
LSVI regularization coef cient	{1}	1
Buffer size	{1e5}	1e5
Update frequency	{50}	50

I.6 Visualizations

In this section we provide a comprehensive visualization of the decoders for all baselines in the target environment. We observe that the behaviors of all baselines are similar across the 5 random seeds. Thus to avoid redundancy, we only show the visualization from 1 random seed.

I.6.1 Visualizations from Section. 6.1

We record the visualization of the 5 sources from Fig. 3 to Fig. 7, REPTTRANSFER in Fig. 8; G-REPTTRANSFER in Fig. 9; running REE on target in Fig. 10.

I.6.2 Visualizations from Section. 6.2

We record the visualization of the 2 sources from Fig. 11 and Fig. 12, REPTTRANSFER in Fig. 13; G-REPTTRANSFER in Fig. 14; running REE on target in Fig. 15. Note that the features collapse at some timesteps in Fig. 14 and Fig. 15, but this is acceptable because the optimal actions at those timesteps are the same for the collapsed states.

Figure 3: Visualization of decoders from source 1. Note the collapse happens at timestep 5, 9 and 17.

Figure 4: Visualization of decoders from source 2. Note the collapse happens at timestep 1 and 10.

Figure 5: Visualization of decoders from source 3. Note the collapse happens at timestep 14 and 15.

Figure 6: Visualization of decoders from source 4. Note the collapse happens at timestep 7, 16, 24.

Figure 7: Visualization of decoders from source 3. Note the collapse happens at timestep 13 and 16.

Figure 8: Visualization of decoders from `CERTTRANSFER`

Figure 9: Visualization of decoders from GERT RANSFER

Figure 10: Visualization of decoders from running BE on target.

Figure 11: Visualization of decoders from source environment 1.

Figure 12: Visualization of decoders from source environment 2.

Figure 13: Visualization of decoders from $\text{C}\overline{\text{E}}\text{R}\text{T}\text{R}\text{A}\text{N}\text{S}\text{F}\text{E}\text{R}$

Figure 14: Visualization of decoders from $\text{G}\overline{\text{E}}\text{R}\text{T}\text{R}\text{A}\text{N}\text{S}\text{F}\text{E}\text{R}$

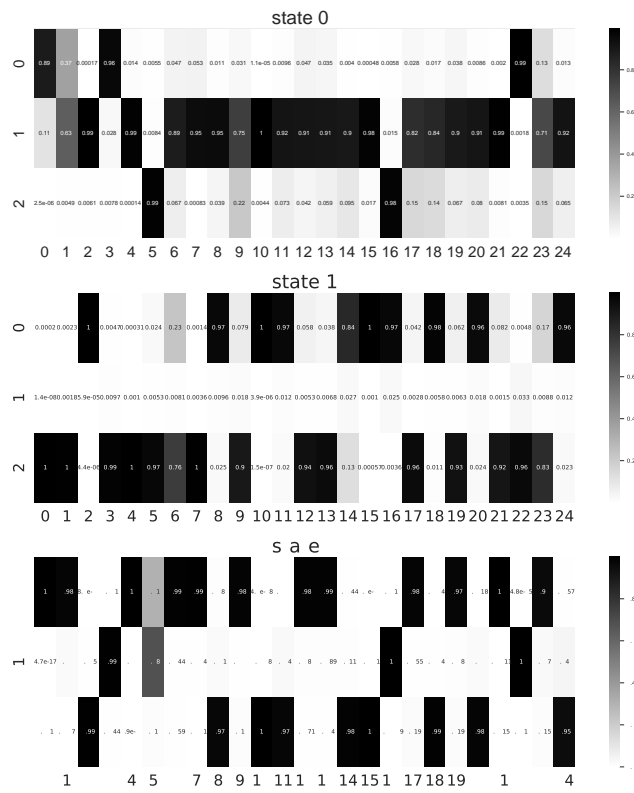


Figure 15: Visualization of decoders running BRIEE in the target.