
Studying Exploration in RL: An Optimal Transport Analysis of Occupancy Measure Trajectories

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The rising successes of RL are propelled by combining smart algorithmic strategies
2 and deep architectures to optimize the distribution of returns and visitations over the
3 state-action space. A quantitative framework to compare the learning processes of
4 these eclectic RL algorithms is currently absent but desired in practice. We address
5 this gap by representing the learning process of an RL algorithm as a sequence of
6 policies generated during training, and then studying the policy trajectory induced
7 in the manifold of state-action occupancy measures. Using an optimal transport-
8 based metric, we measure the length of the paths induced by the policy sequence
9 yielded by an RL algorithm between an initial policy and a final optimal policy.
10 Hence, we first define the *Effort of Sequential Learning (ESL)*. ESL quantifies
11 the relative distance that an RL algorithm travels compared to the shortest path
12 from the initial to the optimal policy. Furthermore, we connect the dynamics of
13 policies in the occupancy measure space and regret (another metric to understand
14 the suboptimality of an RL algorithm), by defining the *Optimal Movement Ratio*
15 (OMR). OMR assesses the fraction of movements in the occupancy measure space
16 that effectively reduce an analogue of regret. Finally, we derive approximation
17 guarantees to estimate ESL and OMR with a finite number of samples and without
18 access to an optimal policy. Through empirical analyses across various environ-
19 ments and algorithms, we demonstrate that ESL and OMR provide insights into
20 the exploration processes of RL algorithms and the hardness of different tasks in
21 discrete and continuous MDPs.

1 Introduction

23 In recent years, significant advancements in Reinforcement Learning (RL) have been achieved
24 in developing exploration techniques that improve learning [10, 14, 20] along with new learning
25 methods [46, 59, 49]. With growing computational resources, these techniques have led to various
26 successful applications of RL, such as playing games up to human proficiency [72, 34], controlling
27 robots [33, 39], tuning databases and computer systems [80, 9], etc. However, there remains a lack of
28 consensus over approaches that can quantitatively compare these exploratory processes across RL
29 algorithms and tasks [71, 6, 42]. This is attributed to some methods being algorithm-specific [78],
30 while others provide theoretical guarantees for very specific settings [45, 2]. Thus, comparing the
31 exploratory processes of these eclectic algorithms across the multi-directional space of RL algorithm
32 design emerges as a natural question. However, the present literature lacks a metric to compare them
33 except regret, which is often hard to estimate [69, 68].

34 This paper aims to address this gap based on two key observations. *First*, we observe from the linear
35 programming formulation of RL that solving the value maximisation problem is equivalent to finding
36 an optimal occupancy measure [77, 60, 37]. Occupancy measure is the distribution of state-action pair

visits induced by a policy [4, 43]. Under mild assumptions, a policy maps uniquely to an occupancy measure. *Second*, we observe that any RL algorithm learns by sequentially updating policies starting from an initial policy to reach an optimal policy. The search for an optimal policy is influenced by the exploration-exploitation strategy and functional approximators, both of which impact the overall performance of the agent by determining the quality of experiences from which it learns [84, 42]. Hereby, we term collectively the learning strategy and the exploration-exploitation interplay as the *exploratory process*.

Contributions. *1. A Framework.* Motivated by our observations, we abstract any RL algorithm as a trajectory of occupancy measures induced by a sequence of policies between an initial and a final (optimal) policy. The occupancy measure of a policy given an environment corresponds to the data-generating distribution of state-actions. Thus, we can quantify the effort of each policy update, i.e. the effort to shift the state-action data distributions, as the transportation distance between their occupancy measures. The total effort of learning by the algorithm can be measured as the total distance covered by its occupancy measure trajectory. We provide a mathematical basis for this quantification by proving that the space of occupancy measures is a differentiable manifold for smoothly parameterized policies (Section 3). Hence, we can compute the length of the occupancy measure trajectory on this manifold using Wasserstein distance as the metric [79].

2. Effort of Sequential Learning. In contrast to RL, if we knew the optimal policy we could update our initial policy directly via supervised or imitation learning. Effort of this learning is represented by a direct, shortest (geodesic) path from initial to optimal policy on the occupancy measure manifold. To quantify the cost of the exploratory process to learn the environment, we define the *Effort of Sequential Learning* (ESL) as the ratio of the (indirect) path traversed by an RL algorithm in the occupancy measure space to the direct distance between the initial and optimal policy (Section 3.1). Lower ESL implies more efficient exploratory process.

3. Efforts to learn that lead to Regret-analogue minimisation. Regret is a widely used optimality measure for reward-maximizing RL algorithms [35]. It measures the total deviation in the value functions achieved by a sequence of policies learned by an RL algorithm with respect to the optimal algorithm that always uses the optimal policy [73]. We show that regret is related to the sum of distances between the optimal policy and each policy in the sequence learned by the RL algorithm, in the occupancy measure space. We can define an analogue of instantaneous regret (at any one step during learning rather than cumulative), in the occupancy measure space, as the geodesic distance between the occupancy measure of the policy at this step in the learning sequence, and the optimal one. We find that not all policy updates lead to a reduction in this analogue of immediate regret, and thus define another index *Optimal Movement Ratio* (OMR) that measures the fraction that do (Section 3.2).

4. Computational and Numerical Insights. We prove sample complexity guarantees to approximate ESL and OMR in practice as we do not have access to the occupancy measures but collection of rollouts from the corresponding policies (Section 4). We show the relation of empirical OMR and ESL to the true ones if the optimal policy is never reached by an algorithm. We conduct experiments on multiple environments, both discrete and continuous, with sparse and dense rewards, comparing state-of-the-art algorithms. We observe that by visualizing aspects of the path traversed (and by comparing ESL and OMR), we are able to compare and provide insights into their exploratory processes and the impact of task hardness on them (Section 5).

2 Preliminaries

Markov Decision Processes. Consider an agent interacting with an environment in discrete timesteps. At each timestep $t \in \mathbb{N}$, the agent observes a state s_t , executes an action a_t , and receives a scalar reward $\mathcal{R}(s_t, a_t)$. The behaviour of the agent is defined by a policy $\pi(a_t|s_t)$, which maps the observed states to actions. The environment is modelled as a Markov Decision Process (MDP) \mathbb{M} with a state space \mathcal{S} , action space \mathcal{A} , transition dynamics $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$, and reward function $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. During task execution, the agent issues actions in response to visited states, and hence a sequence of states and actions $h_t = (s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t)$, here called a rollout, is observed. In infinite-horizon settings, the state value function for a given policy π is the expected discounted cumulative reward over time $V_\pi(s) \triangleq \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t) \mid s_0 = s]$, where $\gamma \in [0, 1)$

is the discount rate. The goal is to learn a policy that maximises the objective $J_\mu^\pi \triangleq \mathbb{E}_{s \sim \mu}[V_\pi(s)]$, where $\mu(s)$ is the initial state distribution.

Occupancy Measure. The state-action occupancy measure is a distribution over the $\mathcal{S} \times \mathcal{A}$ space that represents the discounted frequency of visits to each state-action pair when executing a policy π in the environment [77]. Formally, the occupancy measure of π is $v_\pi(s, a) \triangleq \rho \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a \mid \pi, \mu)$, where $\rho = 1 - \gamma$ is the normalizing factor.

Stationary Markovian policies allow a bijective correspondence with their state-action occupancy measures [26]. We express the objective J_μ^π in terms of the occupancy measure as

$$J_\mu^\pi = \frac{1}{\rho} \mathbb{E}_{(s,a) \sim v_\pi} [\bar{\mathcal{R}}(s, a)], \quad (1)$$

where $\bar{\mathcal{R}}(s, a)$ is the expected immediate reward for the state-action pair (s, a) .

Wasserstein Distance. Let $\mu, \nu \in \mathcal{P}(\mathcal{X})$ be probability measures on a complete and separable metric (Polish) space $(\mathcal{X}, d_\mathcal{X})$. The p-Wasserstein distance between μ and ν is [79]

$$\mathcal{W}_p(\mu, \nu) \triangleq \left(\min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} c(x, x') d\pi(x, x') \right)^{1/p}, \quad (2)$$

where the cost function is given by the metric as $c(x, x') = (d_\mathcal{X}(x, x'))^p$ for some $p \geq 1$. $\Pi(\mu, \nu)$ is a set of all admissible transport plans between μ and ν , i.e. probability measures on $\mathcal{X} \times \mathcal{X}$ space with marginals μ and ν . Wasserstein distances induce geodesic in well-behaved spaces of probability measures. For more discussion, we refer to Appendix A.9. For this work, we consider 1-Wasserstein distance, i.e. $p = 1$, though the results are generalisable to $p > 1$.

MDPs with Lipschitz Rewards. Following [65] and [38], we assume an MDP with $L_\mathcal{R}$ -Lipschitz rewards (ref. Appendix A.1 for elaboration) that satisfies $|\bar{\mathcal{R}}(s, a) - \bar{\mathcal{R}}(s', a')| \leq L_\mathcal{R} d_{\mathcal{S}\mathcal{A}}((s, a), (s', a'))$ for all $s, s' \in \mathcal{S}$ and $a, a' \in \mathcal{A}$. Here, $d_{\mathcal{S}\mathcal{A}}((s, a), (s', a')) = d_\mathcal{S}((s, s')) + d_\mathcal{A}((a, a'))$ is the metric defined on the joint state-action space $\mathcal{S} \times \mathcal{A}$. This is a weaker condition than assuming a completely Lipschitz MDP. [65] showed that for any pair of stationary policies π and π' , the absolute difference between their corresponding objectives is

$$|J_\mu^\pi - J_\mu^{\pi'}| \leq \frac{L_\mathcal{R}}{\rho} \mathcal{W}_1(v_\pi, v_{\pi'}), \quad (3)$$

where $\mathcal{W}_1(v_\pi, v_{\pi'})$ is the 1-Wasserstein distance between the occupancy measures of the policies (ref. Appendix A.2 for details).

3 RL Algorithms as Trajectories of Occupancy Measures

The exploration process (i.e. the exploration-exploitation interplay and learning strategy) of an RL algorithm influences how the policy model updates its policies [36, 76]. During training, a *policy trajectory*, i.e. sequence of policies $(\pi_0, \pi_1, \dots, \pi_N)$, is generated over policy updates due to the exploratory process. We assume these policies belong to a set of stationary Markov policies parameterized by $\theta \in \Theta$. For policies in this set $\pi_\theta \in \Gamma_\Theta$, we define the space of occupancy measures corresponding to Γ_Θ as $\mathcal{M} = \{v_{\pi_\theta}(s, a) \mid \pi_\theta \in \Gamma_\Theta, \theta \in \Theta\}$.

Proposition 1 (Properties of \mathcal{M}). *If the policy π has a smooth parameterization θ , then the space of occupancy measures \mathcal{M} is a differentiable manifold. (Proof in Appendix A.3)*

We can endow the manifold \mathcal{M} with a 1-Wasserstein metric \mathcal{W}_1 to compute the length of any path on \mathcal{M} since $(\mathcal{M}, \mathcal{W}_1)$ is a geodesic space (ref. Appendix A.9 for details). The path distance between occupancy measures corresponding to policies parameterised by $\theta, \theta + d\theta \in \mathcal{M}$ is $ds = \mathcal{W}_1(v_{\pi_\theta}, v_{\pi_{\theta+d\theta}})$. Additionally, in imitation learning, the 1-Wasserstein distance between the occupancy measures of the learner and expert can be used as a minimisable loss function to learn the expert's policy [85]. Hence, the 1-Wasserstein distance reflects the effort required to achieve this imitation learning. Similarly, we propose the following quantification of the effort to update from one policy to another.

Definition 1 (Effort of Learning). *We define the 1-Wasserstein metric between occupancy measures of two policies π and π' , i.e. $\mathcal{W}_1(v_\pi, v_{\pi'})$, as the effort required to learn or update from one policy to the other.*

When a learning process causes an update between occupancy measures in \mathcal{M} , we attribute the resulting update effort to the learning process and refer to it as the effort of learning. In a learning process, first the initial policy π_0 is obtained typically by randomly sampling the model parameters, then these parameters θ undergo updates until a predefined convergence criterion is satisfied, yielding the final optimal policy $\pi_N = \pi^*$. Since each policy has a corresponding occupancy measure, this process yields a sequence of points on \mathcal{M} , which can be connected by geodesics between successive points, producing a curve. The length of the curve is computed by the summation of the finite geodesic distances between consecutive policies along it [53],

$$C \triangleq \sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_{\theta_k}}, v_{\pi_{\theta_{k+1}}}), \quad (4)$$

where θ_0 and θ_N are respectively the initial and final parameter values before and after learning.

3.1 Effort of Sequential Learning (ESL)

As we saw above, RL generates a trajectory in the occupancy measure manifold \mathcal{M} , whose length is given by Equation (4). Compared to the long trajectory of sequential policies generated by the exploratory process, the geodesic L is the ideal shortest path to the optimal policy $\pi_N = \pi^*$ from π_0 , whose length is $L = \mathcal{W}_1(v_{\pi_0}, v_{\pi_N})$. This path would be taken by an imitation-learning oracle algorithm that knows π^* . Both these paths are schematically depicted in Figure 1.

Definition 2 (Effort of Sequential Learning (ESL)). *We define the effort of sequential learning incurred by a trajectory of the exploratory process of an RL algorithm, relative to the oracle that knows π^* ($= \pi_N$) as,*

$$\eta \triangleq \frac{\sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \quad (5)$$

Due to the stochasticity of the exploratory process, we introduce an expectation to obtain $\bar{\eta} = \mathbb{E}_{\pi_0, \mu} [\eta]$. We refer to $\bar{\eta}$ as the effort of sequential learning (ESL).

$\bar{\eta} \geq 1$ and a larger $\bar{\eta}$ corresponds to a less efficient exploratory process of the RL algorithm. Hence, an RL algorithm with $\bar{\eta} \approx 1$ closely mimics the oracle and has an efficient exploratory process.

3.2 Optimal Movement Ratio (OMR)

Regret measures the total deviation in value functions incurred by a sequence of policies learned by an RL algorithm with respect to the optimal algorithm that always uses the optimal policy [73]. We show that regret is connected to the sum of distances from each policy (in the sequence learned by an RL algorithm) to the optimal policy in the occupancy measure space.

Proposition 2 (Regret and Occupancy Measures). *Given an MDP with $L_{\mathcal{R}}$ -Lipschitz rewards, we obtain $\text{Regret} \triangleq \sum_{k=1}^N (J_{\mu}^{\pi^*} - J_{\mu}^{\pi_k}) \leq \frac{L_{\mathcal{R}}}{\rho} \sum_{k=1}^N \mathcal{W}_1(v_{\pi_k}, v_{\pi^*})$. (Proof in Appendix A.4)*

We refer to $\mathcal{W}_1(v_{\pi_k}, v_{\pi^*})$ as the *distance-to-optimal*, and analogously use it as the expected immediate regret in the occupancy measure space. Furthermore, we refer to $\mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})$ as *stepwise-distance*. Interestingly, during training, the *distance-to-optimal* and *stepwise-distance* share a relationship illustrated in Figure 2. From Figure 2, we observe that if the change in *distance-to-optimal*, $\delta_k \triangleq \mathcal{W}_1(v_{\pi_k}, v_{\pi^*}) - \mathcal{W}_1(v_{\pi_{k+1}}, v_{\pi^*}) > 0$, it indicates that the agent got closer to the optimal. We define the set K^+ as containing indices k for which $\delta_k > 0$, while K^- contains the rest.

Definition 3 (Optimal Movement Ratio (OMR)). *We define the proportion of policy transitions that effectively reduce the distance-to-optimal, in a learning trajectory, as*

$$\kappa \triangleq \frac{\sum_{k \in K^+} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}{\sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}. \quad (6)$$

Due to the stochasticity of the exploratory process, we introduce an expectation to obtain $\bar{\kappa} = \mathbb{E}_{\pi_0, \mu} [\kappa]$. We refer to $\bar{\kappa}$ as the optimal movement ratio (OMR).

Note that $\bar{\kappa} \in [0, 1]$, and $\bar{\kappa} \rightarrow 1$ indicates that nearly all the policy updates reduce the *distance-to-optimal*, thus showing high efficiency. $\bar{\kappa} \rightarrow 0$ implies low efficiency, since only a small fraction of the policy updates contribute towards the reduction of the *distance-to-optimal*.

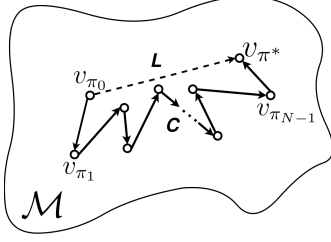


Figure 1: Schematic of the policy trajectory C in the space of occupancy measures \mathcal{M} during RL training (solid line) vs the geodesic L (shortest path, dashed line) between the initial and final points (i.e. π_0 and $\pi_N = \pi^*$).

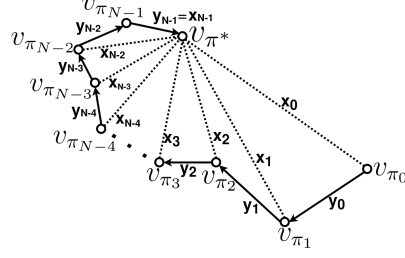


Figure 2: Schematic of how *distance-to-optimal* (denoted by x_k) and *stepwise-distance* (denoted by y_k) on the occupancy measure space describe exploratory process of an RL algorithm during training.

176 The definitions of ESL and OMR assume that the policy at the end of learning is optimal. In
 177 Section 4.2, we define a version of ESL that is useful for the cases where an optimal policy is not
 178 reached. While this is not an empirical proxy, we show in Section 5.3 and Appendix B.5 that it is
 179 useful when the final policy is closer to optimal than the initial one. While regret also depends on an
 180 optimal policy, it is related to cumulative rewards, whereas our metrics do not explicitly depend on
 181 rewards. Still, we show a bound with regret in Proposition 2, and further discuss the possibility of
 182 extending our metrics to be reward-aware in Section 6. We show empirically that our metrics are
 183 complementary to regret in Section 5.2, and discuss other connections with regret in Section 7.

184 3.3 Extension to Finite-Horizon Episodic Setting

185 In the episodic finite-horizon MDP formulation of RL, in short *Episodic RL* [61, 7, 62], the agent
 186 interacts with the environment in multiple episodes of H steps. An episode starts by observing
 187 state s_1 , then for $t = 1, \dots, H$, the agent draws action a_t from a (possibly time-dependent) pol-
 188 icy $\pi_t(\cdot | s_t)$, observes the reward $r(s_t, a_t)$, and transits to a state $s_{t+1} \sim T(\cdot | s_t, a_t)$. Here,
 189 the value function and the state-action value functions at step $h \in [H]$ are respectively defined
 190 as $V_h^\pi(s) \triangleq \mathbb{E}_\pi \left[\sum_{t=h}^H r(s_t, a_t) | s_h = s \right]$ and $Q_h^\pi(s, a) \triangleq \mathbb{E}_\pi \left[\sum_{t=h}^H r(s_t, a_t) | s_h = s, a_h = a \right]$.
 191 Following [4], we can define a finite-horizon version of occupancy measures as

$$v_\pi^H(s, a) \triangleq \frac{1}{H} \sum_{t=1}^H \mathbb{P}(s_t = s, a_t = a | \pi, \mu). \quad (7)$$

192 Following [77], work by [37] shows that v_π^H can be used in the linear programming formulation for
 193 solving MDPs and satisfies the *Bellman Flow Constraints* (in Equation 19 from Appendix A.3). We
 194 prove that under some assumptions, the finite-horizon occupancy measures also construct a manifold,
 195 referred to as \mathcal{M}^H .

196 **Proposition 3** (Properties of \mathcal{M}^H). *If the policy π has a smooth parametrization θ , then the space of*
 197 *finite-horizon occupancy measures \mathcal{M}^H is a differentiable manifold. (Proof in Appendix A.5)*

198 This allows us to similarly define a Wasserstein metric on this manifold, which in turn, allows us to
 199 compute ESL and OMR for evaluating different RL algorithms.

200 4 Computational Challenges and Solutions

201 Similar to regret, our method requires knowing the optimal policy. This is because the efficiency and
 202 effectiveness of exploratory processes of RL algorithms are highly coupled with their ability to reach
 203 optimal policy. ESL and OMR depend on the policies being stationary and Markovian.

204 4.1 Policy datasets for computing occupancy measures

205 We consider approximations of occupancy measures using datasets assumed to be drawn from these
 206 measures. We estimate the Wasserstein distance between the occupancy measures using a method

introduced by [5] known as the *optimal transport dataset distance* (OTDD). OTDD uses datasets to estimate the Wasserstein distance between the underlying distributions. See Appendix A.6 for a detailed account of OTDD.

Definition 4 (Policy dataset). *A dataset of a policy \mathcal{D}_π is a set of state-action pairs drawn from the policy’s occupancy measure, i.e. $\mathcal{D}_\pi = \{(s(i), a(i))\}_{i=1}^m \sim v_\pi$. These can be constituted from the rollouts generated by the policy during task execution.*

We know from imitation learning that if we are given \mathcal{D}_π , generated by an expert policy, we can train a policy model on it in a supervised manner via behaviour cloning [32]. Thus, knowing \mathcal{D}_π can allow converting an RL task into a Supervised Learning (SL) task. Consider a scenario when we have access to a sequence of datasets $(\mathcal{D}_{\pi_0}, \dots, \mathcal{D}_{\pi_N})$, each corresponding to policy π_t for $t \geq 0$. If we train (in a supervised manner) a policy model sequentially on these datasets, the model will undergo a similar policy evolution as the RL algorithm that generated the policy trajectory $(\pi_t)_{t \geq 0}$. This allows us to conceptualise learning in RL as a sequence of SL tasks with sequential transfer learning across the datasets $(\mathcal{D}_{\pi_0}, \dots, \mathcal{D}_{\pi_N})$. We employ OTDD to estimate $\mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})$ using these datasets, i.e. $d_{OT}(\mathcal{D}_{\pi_k}, \mathcal{D}_{\pi_{k+1}}) \approx \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})$, based on Proposition 4.

Proposition 4 (Upper Bound on Estimation Error). *Let an RL algorithm yield a sequence of policies π_0, \dots, π_N while training. Now, we construct N datasets $\mathcal{D}_{\pi_0}, \dots, \mathcal{D}_{\pi_N}$, each consisting of M rollouts of the corresponding policies. Then, we can use these datasets to approximate $\sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})$ by $\sum_{k=0}^{N-1} d_{OT}(\mathcal{D}_{\pi_k}, \mathcal{D}_{\pi_{k+1}})$ with an expected error upper bound $\frac{2N\mathcal{E}_2}{\sqrt{M}} + N\gamma^{T+1} \text{diam}(SA)$. Here, T is the total number of steps per episode, $\text{diam}(SA)$ is the diameter of the state-action space, and \mathcal{E}_2 is a positive-valued and polylogarithmic function of S and A . For finite horizon case, we can further reduce the error bound to $\frac{2N\mathcal{E}_2}{\sqrt{M}}$.*

Proof of Proposition 4 is in Appendix A.7. The results support that ESL and OMR can be estimated as

$$\bar{\eta} = \mathbb{E}_{\pi_0, \mu} \left[\frac{\sum_{k=0}^{N-1} d_{OT}(\mathcal{D}_{\pi_k}, \mathcal{D}_{\pi_{k+1}})}{d_{OT}(\mathcal{D}_{\pi_0}, \mathcal{D}_{\pi_N})} \right], \text{ and } \bar{\kappa} = \mathbb{E}_{\pi_0, \mu} \left[\frac{\sum_{k \in K^+} d_{OT}(\mathcal{D}_{\pi_k}, \mathcal{D}_{\pi_{k+1}})}{\sum_{k=0}^{N-1} d_{OT}(\mathcal{D}_{\pi_k}, \mathcal{D}_{\pi_{k+1}})} \right]. \quad (8)$$

4.2 When an optimal policy is not reached

So far we have assumed that the algorithms converge at the optimal policy, i.e. $\pi_N = \pi^*$. However, this is not always true. We consider a scenario when $\pi_N \neq \pi^*$, and define

$$\eta_{sub} = \frac{\sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})}, \pi_N \neq \pi^*. \quad (9)$$

Proposition 5. *Given $N \geq 2$ and $\pi_0 \neq \pi_N \neq \pi^*$, we obtain*

$$\frac{\eta - \eta_{sub}}{\eta} \leq \frac{2\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})}. \quad (10)$$

This is true due to the triangle inequalities: $\mathcal{W}_1(v_{\pi_0}, v_{\pi^*}) + \mathcal{W}_1(v_{\pi_N}, v_{\pi^*}) \geq \mathcal{W}_1(v_{\pi_0}, v_{\pi_N})$ and $\mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi_N}) + \mathcal{W}_1(v_{\pi_N}, v_{\pi^*}) \geq \mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi^*})$. The proof is provided in Appendix A.8. Note that Equation (10) shows that when π_N is close to π^* , then η_{sub} is a good approximation of η , and thus a good quantifier to determine the efficiency of the algorithm’s exploratory process. However, $\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})$ is dependent on the RL algorithm and hence a bound cannot be provided here. Still, η_{sub} might be useful when π_N is closer to π^* than π_0 . A fallible proxy for this could be when the performance of π_N is better than that of π_0 , i.e. $J_\mu^{\pi_N} > J_\mu^{\pi_0}$. We show the usefulness of η_{sub} in our experimental results in Section 5.3 and Appendix B.5 for simple environments. It remains to be seen how useful η_{sub} is in complex environments.

5 Experimental Evaluation

In this section, we evaluate the proposed methods in the *2D-Gridworld* and *Mountain Car* [58, 13] environments, to analyze our methods in discrete and continuous state-action spaces respectively. The 2D-Gridworld environment is of size 5×5 with actions: {up, right, down, left}. In the gridworld,

we perform experiments on 3 settings namely:- A) deterministic with dense rewards, B) deterministic with sparse rewards, and C) stochastic with dense rewards. Further details about these settings are provided in Appendix B.1. The Mountain Car environment, in our experimentation, is a deterministic MDP with dense rewards that consists of both continuous states and actions - described in detail in [13]. Note that we used $L1$ distance and $L2$ distance as metrics (d_X) for the state spaces of the 2D-Gridworld and Mountain Car, respectively, which underpin the \mathcal{W}_1 .

Our experiments aim to address the following questions:

1. What information can the visualization of the policy evolution during RL training provide about the exploratory process of the algorithm?
2. How do ESL and OMR allow us to analyze the exploratory processes of RL algorithms?
3. Does ESL scale proportionally with task difficulty?

Summary of Results. In Section 5.1, we demonstrate that visualizing the evolution of *distance-to-optimal* and *stepwise-distance* of different RL algorithms during training reveals: 1) whether the agent is stuck in suboptimal policies, 2) the coverage area of the exploration processes, and 3) their varied characteristics over time. We further compare ESL and OMR of different algorithms on a few environments in Section 5.2. Finally, we show in Section 5.4 that ESL scales proportionally with task difficulty, and thus, reflects the effects of task difficulty on exploration and learning.

5.1 Exploration Trajectories of RL Algorithms

(I) **DISCRETE MDP.** To understand the utility of visualizing exploratory processes, we use the following RL algorithms: 1) Tabular Q-learning with a) ϵ -greedy ($\epsilon = 0$) and b) ϵ -greedy ($\epsilon = 1$) strategies; 2) UCRL2 [35]; 3) PSRL [61]; 4) SAC [28, 17]; and 5) DQN [57] with ϵ -decay. The algorithms solve a simple 5×5 gridworld with dense rewards, starting from top-left (0,0) to reach bottom-right (4,4). Figure 3 presents exploratory behaviours of the algorithms in both the occupancy measure space and state space.

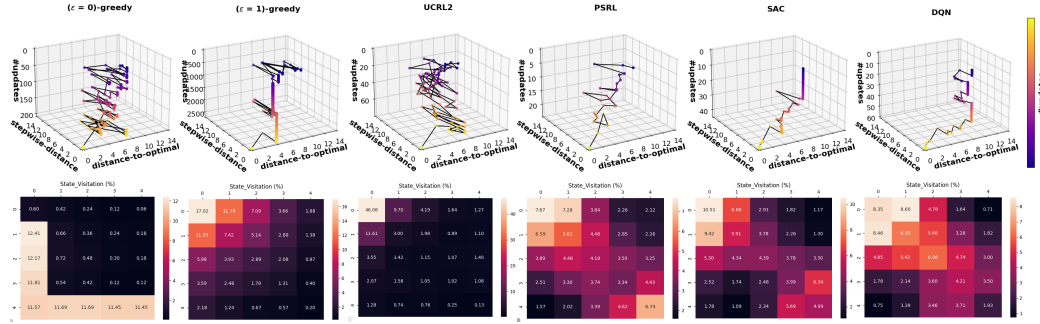


Figure 3: Top row: 3D plots of *distance-to-optimal* (x-axis) and *stepwise-distance* (y-axis) across number of updates (z-axis), illustrating policy evolution in the occupancy measure space for algorithms: $\epsilon(=0)$ -greedy and $\epsilon(=1)$ -greedy Q-learning, UCRL2, PSRL, SAC, and DQN (left to right). Bottom row: Corresponding state visitation frequencies over the full training. The problem setting is deterministic with dense rewards and 15 maximum number of steps per episode. (NB. Larger versions of these plots are presented in Appendix D.1, while their 2D projections are in Figure 15 at Appendix D.3, and corresponding performance plots are in Figure 12 at Appendix D.2.)

Q-learning: $\epsilon = 0$ vs $\epsilon = 1$. Note that $\epsilon = 0$ updates the Q-table by only exploiting, while $\epsilon = 1$ by exploring. From the state visitations, we observe expected characteristics, like a preferred visit path for $\epsilon = 0$, versus $\epsilon = 1$ with visitation frequencies that are similar at states equidistant from the start-state and gradually decreasing as the distance from the start-state increases. From the policy evolution, we see how scattered and erratic the policy transitions are for $\epsilon = 0$. Whereas $\epsilon = 1$ is dominated by unchanging or little-changing policies seen by straight vertical line segments (indicating being *stuck in suboptimality*). In this setting, $\epsilon = 0$ is characterised by transitioning between diverse policies (i.e. being aggressive with larger coverage area) while $\epsilon = 1$ is likely to be stuck in suboptimality. The stuckness is due to high action randomness in $\epsilon = 1$ that causes the agent to select suboptimal actions, slowing the convergence of the Q-table and not changing the learning policy until the best actions are discovered.

283 **UCRL2 vs PSRL.** UCRL2 has nearly uniform state visits (with the exception of the start-state
 284 because the initial state distribution is 1 at state (0,0)), thus being consistent with literature since the
 285 algorithm selects exploratory state-action pairs more uniformly [35]. In contrast, PSRL has high
 286 visit frequencies along the diagonal states, because it selects actions according to the probability that
 287 they are optimal [61]. We observe from the policy evolution plots that PSRL has smoother policy
 288 transitions that are orientated towards optimality, while UCRL2 behaves more aggressively with
 289 policy transitions that do not taper as it approaches optimality. [61] highlighted that exploration in
 290 PSRL is guided by the variance of sampled policies as opposed to optimism in UCRL2. We observe
 291 in Figure 3 that the guiding variance in PSRL reduces after every policy update until optimality is
 292 reached, while UCRL2 maintains high variance.

293 **SAC vs DQN.** The state visits of both the algorithms appear to be similar. SAC has higher visitation
 294 frequencies at the corners than DQN. Surprisingly from the policy evolution plots, we learn that
 295 both algorithms have a reluctance to transition between policies - hence the *stuck in suboptimality*
 296 vertical line segments, especially initially. This reluctance is due to the slow ‘soft updates’ of target
 297 networks [51] in the algorithms. We also observe that SAC approaches optimality more gradually
 298 than DQN.

299 **All algorithms.** Figure 3 shows that UCRL2 was more meandering (with larger coverage area)
 300 towards optimality than the rest. SAC and DQN approached optimality more directly and smoothly
 301 (with smaller coverage area) than the rest. These characteristics are intuitively revealed by policy visu-
 302 alisation plots, and are aligned with literature, hence enhancing our understanding of the exploratory
 303 processes.

304 **(II) CONTINUOUS MDP.** We use DDPG [51] and SAC to solve the Mountain Car. The policy
 305 evolutions of these algorithms are presented in Figure 4.

306 **DDPG vs SAC.** Both exhibit short-distances (< 1) between policy updates (i.e. small coverage
 307 area). They depict no sign of being stuck or settling early on any particular policy, which shows their
 308 continuously exploratory nature. While they begin with almost constant mean *distances-to-optimal*
 309 and *stepwise-distances*, SAC drops its mean *distance-to-optimal* earlier than DDPG.

310 Figure 4 illustrates how OMR changes with update number k . $OMR(k)$ represents OMR starting
 311 with the k^{th} policy as the initial policy, while OMR starts from the 0^{th} policy (details of computing
 312 $OMR(k)$ are in Appendix B.2). For both algorithms, $OMR(k)$ remains near chance level (~ 0.5)
 313 initially, then sharply increases near the final updates. This suggests that early policy updates are
 314 purely exploratory and oblivious to policy improvement but align with the optimal policy just before
 315 convergence. The efficiency of the algorithm depends on how early this transition occurs, e.g. starts
 earlier for SAC than DDPG, rendering SAC more efficient.

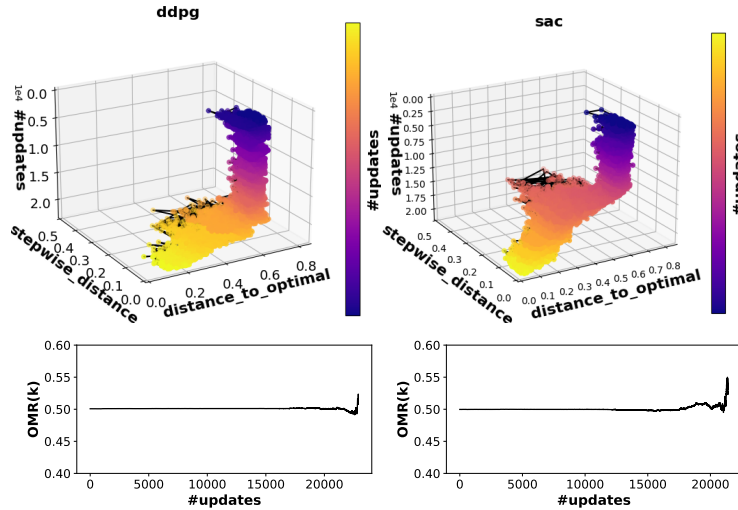


Figure 4: Top row: 3D plots of *distance-to-optimal* and *stepwise-distance* vs. number of updates for DDPG and SAC. Bottom row: $OMR(k)$ vs. #update, k , for the corresponding algorithms. Note that corresponding performance plots are in Figure 13 in Appendix D.2.

5.2 Comparison of ESL and OMR across RL Algorithms and Environments, and their complementarity to number of updates and regret

Tables 1-3 showcase how ESL and OMR are summary indices of the policy trajectories during learning by evaluating the algorithms in various settings.

| Algo. | ESL | OMR | UC | SR% |
|-------------------------|------------------|------------------|------------------|-----|
| SAC | 9.26±5.54 | 0.58±0.14 | 980±670 | 100 |
| UCRL2 | 47.2±8.20* | 0.49±0.04 | 60.7±11 | 100 |
| PSRL | 23.2±11.5 | 0.52±0.06 | 34.1±9.34 | 100 |
| DQN | 12.4±7.13 | 0.54±0.11 | 161±93 | 98 |
| $\epsilon(=1)$ -greedy | 6.27±2.22 | 0.61±0.09 | 672±385 | 100 |
| $\epsilon(=0.9)$ -decay | 8.10±3.43 | 0.61±0.10 | 389±138 | 100 |
| $\epsilon(=0)$ -greedy | 15.5±5.28 | 0.53±0.06 | 176±37.9 | 84 |

Table 1: Evaluation of RL algorithms (over 40 runs) in the **deterministic, dense-rewards setting** for 5×5 gridworld, including Effort of Sequential Learning (ESL), Optimal Movement Ratio (OMR), number of updates to convergence (UC), and success rate (SR). Lowest ESL, lowest UC, and highest OMR values are in **bold**. The highest ESL value is starred (*).

| Algo. | ESL | OMR | UC | SR% |
|------------------------------|------------------|------------------|------------------|-----------|
| Deterministic, sparse | | | | |
| SAC | 27.8±21.9 | 0.57±0.13 | 4385±3274 | 100 |
| UCRL2 | 73.3±0.0 | 0.45±0.0 | 93.0±0.0 | 100 |
| PSRL | 73.2±54.1 | 0.52±0.076 | 100±67.3 | 100 |
| DQN | 137±154* | 0.49±0.08 | 12638±4431 | 80 |
| Stochastic, dense | | | | |
| SAC | 445±245 | 0.501±0.004 | 2463±2043 | 92 |
| UCRL2 | 198±121 | 0.502±0.027 | 268±155 | 32 |
| PSRL | 55.4±33.6 | 0.52±0.04 | 76.1±50.6 | 92 |
| DQN | 458±311* | 0.502±0.01 | 1586±1077 | 24 |

Table 2: Evaluation of RL algorithms (over 40 runs) in the **deterministic, sparse-rewards** and **stochastic, dense-rewards** settings for 5×5 gridworld. Lowest ESL, highest OMR and lowest UC values are in **bold**. The highest ESLs are starred.

Dense Rewards. We observe, in Table 1, that PSRL took the lowest number of updates (UC) to reach the optimal policy in contrast with SAC. Yet, PSRL was meandering more than SAC (see Figure 3). The relative directness of SAC is captured by lower ESL and higher OMR compared to PSRL. Even though SAC has larger UC than PSRL, it took a shorter path to optimality than PSRL. This shows that UC does not necessarily correlate with ESL and OMR, and it provides incomplete information about the exploratory processes. Indeed, two algorithms may have the same UC, but different ESL and/or OMR due to different *step-wise distances* and varied movement towards optimality.

Sparse Rewards. In the sparse rewards setting (Table 2), low performance of DQN is observed in both our indices and UC. However, SAC is more efficient with lowest ESL and highest OMR, yet UCRL2 has the lowest number of updates (UC). Note that *UCRL2 is provably regret-optimal, while SAC does not have such rigorous theoretical guarantees but is known to be practically efficient, and this is well captured by ESL and OMR*. So far, Tables 1 and 2 demonstrate how our indices provide a clearer picture of exploratory processes than the number of updates.

Stochastic Transitions. In the stochastic setting (Table 2), by observing only successful cases, we notice that the meandering characteristic of PSRL and UCRL2 is more suitable for this setting than SAC and DQN (based on better ESL and OMR values). PSRL and UCRL2 have similar regret bounds [61]; yet in Tables 1 and 2, PSRL has better ESL and OMR (along with a higher success rate). This aligns with the regret analysis presented in [61].

Table 3 corroborates with the policy evolution plots in Figure 4, in that due to SAC dropping its mean *distance-to-optimal* earlier than DDPG it exhibits a lower ESL. Additionally, we notice a trend of increasing ESL and decreasing OMR across algorithms when shifting from dense-rewards to sparse-rewards settings, from deterministic to stochastic transitions, from discrete to continuous environments, indicating an increase in the effort of the exploratory processes. We have shown that ESL and OMR enhance the understanding of exploratory processes by effectively summarizing the policy trajectories of algorithms during learning. They offer more insight than the number of updates, and align with regret when algorithms reach an optimal policy. In the next Section, we demonstrate the utility of ESL when an optimal policy is not reached.

5.3 Usefulness of ESL when optimal policy is not reached

When an optimal policy is not reached at the end of learning, ESL cannot be computed exactly. At this point, we propose to use an approximation of ESL, i.e. η_{sub} (Equation 9). A natural question

| Algo. | ESL | OMR | UC | SR% |
|-------|----------|-------|------------|-----|
| DDPG | 1881±500 | 0.501 | 23500±5268 | 100 |
| SAC | 1619±189 | 0.5 | 22700±2971 | 100 |

Table 3: Evaluation of RL algorithms in the Mountain Car continuous MDP (over 5 runs). The variances for OMR are negligible.

350 arises: *When the optimal policy is not reached, does the η_{sub} still yield insights about the exploratory*
351 *process of RL?* In Table 4, we compare ESL when the optimal policy was reached, i.e. η , versus when it was not, i.e. η_{sub} .

| Algo. | η | η_{sub} | d | c |
|-------|-----------|--------------|-----------|-----------|
| SAC | 445±246 | 853±127 | 5.63±1.23 | 7.26±1.45 |
| UCRL2 | 198±121 | 510±274 | 5.36±0.84 | 4.58±1.90 |
| PSRL | 55.4±33.6 | 361±43.6 | 4.97±1.34 | 3.91±0.48 |
| DQN | 458±311 | 1971±250 | 4.88±1.06 | 6.52±0.31 |

Table 4: Evaluation of algorithms in the **stochastic, dense-rewards setting** for 5×5 gridworld. When the algorithm converged to optimality, η is the Effort of Sequential Learning, and $d = \mathcal{W}_1(v_{\pi_0}, v_{\pi^*})$ is the distance between the initial and optimal policies. When the algorithm did not converge to the optimal policy but some π_N , we used η_{sub} and $c = \mathcal{W}_1(v_{\pi_0}, v_{\pi_N})$ to denote the aforementioned quantities. 40 training trials were used.

352

353 We observe that η_{sub} values are always greater than η values. However, they both yield the same
354 efficiency ranking (e.g. PSRL, UCRL2, SAC and DQN). This indicates that η_{sub} reliably predicts
355 results provided by η for relative comparison of algorithms.

356 5.4 ESL Increases with Task Difficulty

357 Figure 5 depicts the ESL values for Q-learning with ϵ -decay strategy (for $\epsilon = 0.9$) across tasks with
358 varying hardness. These tasks are deterministic 2D-Gridworld of sizes 5×5 and 15×15 matched
359 with either dense or sparse rewards (as specified in Appendix B.1). We chose to assess the ϵ -decay
360 Q-learning algorithm because it is simple and yet completes all these tasks. We observe that the
361 ESL is lowest for $[5 \times 5]$ dense (5×5 grid, dense rewards) and highest for $[15 \times 15]$ sparse (15×15
362 grid, sparse rewards) as anticipated. The results demonstrate that ESL scales proportionally with task
363 difficulty, matching expectations that more difficult tasks demand greater effort of the exploratory
364 process.

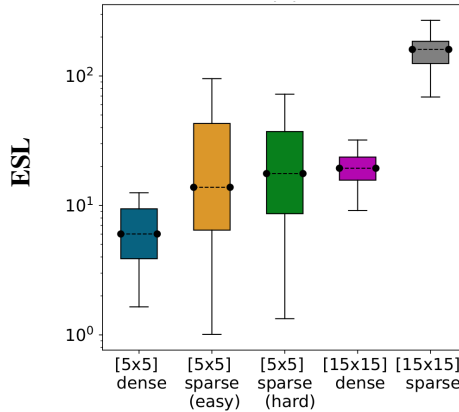


Figure 5: Q-learning with ϵ -greedy ($\epsilon = 0.9$ decaying, averaged over 40 runs) across deterministic 2D-Gridworld (5×5 and 15×15) tasks. The 1st and 4th (from left to right) have dense rewards, while the rest have sparse rewards (details in Appendix B.1).

365 **Remarks.** When the optimal policy is reached, we can use the visualisation of policy trajectories, as
366 well as ESL (η) and OMR to study exploratory processes of an RL algorithm. When the optimal policy

is not reached, we can still use the visualisation of policy trajectories and ESL-sub (η_{sub}) to study exploratory processes because they still capture characteristics of exploratory processes (e.g. high coverage, smooth exploration). Additionally, ESL (η or η_{sub}) captures the hardness of the task that we are solving. Thus, studying the occupancy measure trajectories and their corresponding indices can aid in making a knowledgeable choice of an RL algorithm that exhibits desired characteristics.

6 Related Works

Several prior works have utilized various components leveraged in our work, namely Wasserstein distance, occupancy measures, and the trajectory of RL on a manifold, but for different purposes. Here, we summarise them and elucidate the connections.

In supervised learning, [5] proposed an optimal transport approach, namely Optimal Transport Dataset Distance (OTDD), to quantify the transferability between two supervised learning tasks by computing the similarity (i.e. distance) between the task datasets. Here, we conceptualise and define the effort of learning for RL, as a sequence of such supervised learning tasks. We observe that *the total effort of sequential learning can be computed as the sum of OTDD distances between consecutive occupancy measures*. Recently, [86] have developed generalized occupancy models by defining cumulative features that are transferable across tasks. In future, one can generalize our indices for the cumulative features constructed from some invertible functions of the step-wise occupancy measures.

Optimal transport-based approaches are also explored in RL literature. These works broadly belong to two families. First line of works uses Wasserstein distance over a posterior distribution of Q-values [55, 50] or return distributions [75] to quantify uncertainty, and then to use this Wasserstein distance as a loss to learn better models of the posterior distribution of Q-values or return distributions, respectively. The second line of works uses Wasserstein distance between a feasible family of MDPs as an additional robustness constraint to design robust RL algorithms [1, 18, 30]. Here, *we introduce a concept of using Wasserstein distance between occupancy measures to understand the exploratory dynamics*. Incorporating this insight into better algorithm design would be an interesting future work. Recently, [15] related Wasserstein distance between reward-labelled Markov chains to bisimulation metrics which abstract state spaces. In the same spirit, we could use reward as the cost-function in computing our nested Wasserstein distance (OTDD) to obtain a reward- or value-aware OTDD to define broader bisimulation metrics with abstract state-action spaces, instead of just state spaces.

As a parallel approach to optimal transport, the information geometries of the trajectory of an RL algorithm under different settings are studied. These approaches use mutual information as a metric instead of Wasserstein distance. [8] studied the information geometry of Bayesian multi-armed bandit algorithms. They considered a bandit algorithm as a trajectory on a belief-reward manifold, and proposed a geometric approach to design a near-optimal Bayesian bandit algorithm. [21, 44] studied information geometry of unsupervised RL and proposed mutual information maximisation schemes over a set of tasks and their marginal state distributions. [83] extended this approach with Wasserstein distance and demonstrated benefits of using Wasserstein distance than mutual information. *We use Wasserstein distance as a natural metric in occupancy manifold which aligns with the hardness of different tasks*. It would be interesting to extend our framework to understand the dynamics of unsupervised RL algorithms.

7 Discussion

Our work introduces methods to theoretically and quantitatively understand and compare the exploratory strategies of different RL algorithms. Since learning in a typical RL algorithm happens through a sequence of policy updates, we propose to understand the exploratory process by visualizing and analyzing the path traversed by an RL algorithm in the space of occupancy measures. We show the usefulness of this approach by conducting experiments on various environments and RL algorithms.

Our results show that ESL and OMR provide insight into the evolution of the agent’s policy, revealing whether it is approaching the optimal policy in a steady or meandering manner. Additionally, they allow us to understand how the learning process of the same algorithm changes with different rewards and transitions structures, and task hardness. We emphasize that ESL and OMR complement the number of updates to converge and regret (see Appendix B.7) rather than replace them.

We now discuss various practical aspects and future possibilities.

Computational complexity. While efficiency is an important aspect, the primary focus of this work is to introduce a framework for analyzing exploration in RL using occupancy measures. It is of interest to utilize our approach to benchmark and compare the learning dynamics of various RL algorithms in more environments, especially large-scale or high-dimensional ones. Computing Wasserstein distances in such environments would incur high computational costs, however in recent years several methods such as greedy computation [16], hierarchical methods [47], and inexact proximal point methods [82] have been introduced to handle large-scale OT problems. For example, [23] leverages a block-diagonal approximation method to deal with high-dimensional probability distributions similar to ours, while anchor space OT [31] specifically addresses multiple OT problems with multiple distributions.

Use of 1-Wasserstein metric \mathcal{W}_1 . The \mathcal{W}_1 satisfies the Kantorovich-Rubinstein duality, making Equation (3) applicable and providing a basis for bounding regret in Proposition 2. In contrast, $\mathcal{W}_{p>1}$ does not have such duality. Since $\mathcal{W}_1(P, Q) \leq \mathcal{W}_{p>1}(P, Q)$ [79], \mathcal{W}_1 yields tighter regret bounds. Moreover, \mathcal{W}_1 is less sensitive to outliers and sampling discrepancies than $\mathcal{W}_{p>1}$ [67], making it well-suited for our setting. Nevertheless, our approach extends to $\mathcal{W}_{p>1}$. Compared to KL-divergence, \mathcal{W}_1 is a metric that satisfies symmetry and triangle inequality, which have been instrumental in our proofs and guarantees, while KL-divergence is not a metric. Additionally, $\mathcal{W}_{p>1}$ leverages the geometry of the underlying support space [64]. This allows it to capture distances between distributions with disjoint support [64], which KL-divergence cannot capture. Hence, we use \mathcal{W}_1 in our approach.

Choice of distance metric $d_{\mathcal{X}}$. The choice of the distance metric impacts the geometry of the support space [48], consequently the Wasserstein distances, and thus ESL and OMR. In our case, the support space (i.e. state-action space) is reduced to the state space, as the action space maps back to the state space via OTDD. Thus, the choice of distance metric should reflect the effort of moving in the state space. For example, in the Gridworld, we used L1 distance because only vertical and horizontal displacements are allowed, and L2 distance in the Mountain Car, as applicable to real-world spaces.

Regret and Sum-of-distances-to-optimal ($\sum_{k=1}^N \mathcal{W}_1(v_{\pi_k}, v_{\pi^*})$). Comparing regret across policies measures the similarity of returns (with respect to optimal), disregarding behavioural differences, like variations in actions at the same states. This makes regret advantageous in settings without critical safety and physical constraints, e.g. games, due to its computational efficiency. In contrast, the *sum-of-distance-to-optimal* focuses on behavioural differences between policies, is reward-agnostic, and is well-suited for environments where safety and physical constraints are critical, e.g. robotics. While minimizing regret prioritizes matching the performance of the optimal policy, minimizing the *sum-of-distance-to-optimal* focuses on mimicking its behaviour. Thus, the *sum-of-distance-to-optimal* can be used similarly to regret, especially where the behaviour of achieving good performance is essential.

Algorithm Selection and Design. Depending on the environment, we can select an algorithm with promising characteristics and spend more time optimising it to improve performance. For example, rather than tuning the hyper-parameters of multiple competing algorithms to find the best one, it may be more effective to first identify an algorithm best suited to an environment based on ESL and OMR, and then fine-tune it. The chosen algorithm might remain suitable across similar environments as well.

Furthermore, we could incorporate an online adaptation of the exploratory process of the RL algorithm itself, based on recent estimates of the ESL and OMR. For example, if the current policy gives a better return than the initial policy, then we could adapt the exploratory parameters (at a slow rate) to optimise a running estimate of η_{sub} and a suitable approximation of OMR, thus enabling better exploration. However, the feasibility and convergence of such a scheme remain open.

References

- [1] M. A. Abdullah, H. Ren, H. B. Ammar, V. Milenkovic, R. Luo, M. Zhang, and J. Wang. Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*, 2019.
- [2] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun. Reinforcement learning: Theory and algorithms, 2022.

- [3] M. Aleksandrowicz and J. Jaworek-Korjakowska. Metrics for assessing generalization of deep reinforcement learning in parameterized environments. *JAISCR*, 14(1):45–61, 2023.
- [4] E. Altman. *Constrained Markov Decision Processes*. Routledge, 1999.
- [5] D. Alvarez-Melis and N. Fusi. Geometric dataset distances via optimal transport. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [6] S. Amin, M. Gomrokchi, H. Satija, H. van Hoof, and D. Precup. A survey of exploration methods in reinforcement learning, 2021.
- [7] M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In *International conference on machine learning*, pages 263–272. PMLR, 2017.
- [8] D. Basu, P. Senellart, and S. Bressan. Belman: An information-geometric approach to stochastic bandits. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2019, Würzburg, Germany, September 16–20, 2019, Proceedings, Part III*, pages 167–183. Springer, 2020.
- [9] D. Basu, X. Wang, Y. Hong, H. Chen, and S. Bressan. Learn-as-you-go with megh: Efficient live migration of virtual machines. *IEEE Transactions on Parallel and Distributed Systems*, 30(8):1786–1801, 2019.
- [10] M. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 1471–1479. Curran Associates, Inc., 2016.
- [11] H. Bojun. Steady state analysis of episodic reinforcement learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9335–9345. Curran Associates, Inc., 2020.
- [12] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [13] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- [14] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.
- [15] S. Calo, A. Jonsson, G. Neu, L. Schwartz, and J. Segovia-Aguas. Bisimulation metrics are optimal transport distances and can be computed efficiently, 2024.
- [16] G. Carlier, A. Galichon, and F. Santambrogio. From knothe’s transport to brenier’s map and a continuation method for optimal transport. *SIAM Journal on Mathematical Analysis*, 41(6):2554–2576, 2010.
- [17] P. Christodoulou. Soft actor-critic for discrete action settings, 2019.
- [18] E. Derman and S. Mannor. Distributional robustness and regularization in reinforcement learning. *arXiv preprint arXiv:2003.02894*, 2020.
- [19] M. P. Drazin. Pseudo-inverses in associative rings and semigroups. *The American mathematical monthly*, 65(7):506–514, 1958.
- [20] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.
- [21] B. Eysenbach, R. Salakhutdinov, and S. Levine. The information geometry of unsupervised reinforcement learning. *arXiv preprint arXiv:2110.02719*, 2021.
- [22] R. Flamary, N. Courty, A. Gramfort, M. Z. Alaya, A. Boisbunon, S. Chambon, L. Chapel, A. Corenflos, K. Fatras, N. Fournier, L. Gautheron, N. T.H. Gayraud, H. Janati, A. Rakotomamonjy, I. Redko, A. Rolet, A. Schutz, V. Seguy, D. J. Sutherland, R. Tavenard, A. Tong, and T. Vayer. Pot: Python optimal transport. *Journal of Machine Learning Research*, 22(78):1–8, 2021.

- [23] Y. Gao and P. Chaudhari. An information-geometric distance on the space of tasks. In *Proceedings of the 38th International conference on machine learning*, pages 3553–3563. PMLR, 2021.
- [24] C. Gelada, S. Kumar, J. Buckman, O. Nachum, and M. G. Bellemare. Deepmdp: Learning continuous latent space models for representation learning. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2170–2179. PMLR, 2019.
- [25] A. Gibbs and F. E. Su. On choosing and bounding probability metrics. *International Statistical Review / Revue Internationale de Statistique*, 70(3):419–435, 2002.
- [26] A. Givchi. *Optimal Transport in Reinforcement Learning*. PhD thesis, Graduate School-Newark Rutgers, The State University of New Jersey, 2021.
- [27] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *International Conference on Artificial Intelligence and Statistics*, 2010.
- [28] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [29] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine. Soft actor-critic algorithms and applications, 2019.
- [30] L. Hou, L. Pang, X. Hong, Y. Lan, Z. Ma, and D. Yin. Robust reinforcement learning with wasserstein constraint. *arXiv preprint arXiv:2006.00945*, 2020.
- [31] J. Huang, X. Su, Z. Fang, and H. Kasai. Anchor space optimal transport as a fast solution to multiple optimal transport problems. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [32] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys*, 50:1–35, 2017.
- [33] J. Ibarz, J. Tan, C. Finn, M. Kalakrishnan, P. Pastor, and S. Levine. How to train your robot with deep reinforcement learning: lessons we have learned. *The International Journal of Robotics Research*, 40(4-5):698–721, 2021.
- [34] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, A. G. Castaneda, C. Beattie, N. C. Rabinowitz, A. S. Morcos, A. Ruderman, N. Sonnerat, T. Green, L. Deason, J. Z. Leibo, D. Silver, D. Hassabis, K. Kavukcuoglu, and T. Graepel. Human-level performance in 3d multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019.
- [35] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010.
- [36] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Journal of artificial intelligence research. *Reinforcement Learning: A Survey*, 4:237–285, 1996.
- [37] K. Kalagarla, R. Jain, and P. Nuzzo. A sample-efficient algorithm for episodic finite-horizon mdp with constraints. In *The 35th AAAI Conference on Artificial Intelligence*, pages 8030–8037, 2021.
- [38] M. Kallel, D. Basu, R. Akrou, and C. D’Eramo. Augmented bayesian policy search. In *The 12th International Conference of Learning Representations*, volume 139, 2024.
- [39] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza. Champion-level drone racing using deep reinforcement learning. *Nature*, 620(7976):982–987, 2023.
- [40] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017.

- [41] S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing and machine-learning applications. *IEEE Signal Processing Magazine*, 34:43–59, 2017.
- [42] P. Ladosz, L. Weng, M. Kim, and H. Oh. Exploration in deep reinforcement learning: A survey. *Information Fusion*, 85:1–22, 2022.
- [43] R. Larocche and R. T. des Combes. On the occupancy measure of non-markovian policies in continuous mdps. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, 2023.
- [44] M. Laskin, H. Liu, X. Bin Peng, D. Yarats, A. Rajeswaran, and P. Abbeel. CIC: contrastive intrinsic control for unsupervised skill discovery. *CoRR*, abs/2202.00161, 2022.
- [45] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [46] A. Lazaridis, A. Fachantidis, and I. Vlahavas. Deep reinforcement learning: A state-of-the-art walkthrough. *Journal of Artificial Intelligence Research*, 69:1421–1471, 2020.
- [47] J. Lee, M. Dabagia, E. Dyer, and C. Rozell. Hierarchical optimal transport for multimodal distribution alignment. *Advances in Neural Information Processing Systems*, 32, 2019.
- [48] J. M. Lee. *Manifolds and differential geometry*, volume 107. American Mathematical Society, 2009.
- [49] S. E. Li. Deep reinforcement learning. In *Reinforcement learning for sequential decision and optimal control*, pages 365–402. Springer, 2023.
- [50] A. Likmeta, M. Sacco, A. M. Metelli, and M. Restelli. Wasserstein actor-critic: directed exploration via optimism for continuous-actions control. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8782–8790, 2023.
- [51] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016.
- [52] S. Liu. An evaluation of ddpg, td3, sac, and ppo: Deep reinforcement learning algorithms for controlling continuous system. In *International Conference on Data Science, Advanced Algorithm and Intelligent Computing*, 2023.
- [53] J. Lott. Some geometric calculations on wasserstein space. *Communications in Mathematical Physics*, 277:423–437, 2008.
- [54] M. Memmel, P. Liu, D. Tateo, and J. Peters. Dimensionality reduction and prioritized exploration for polio search. In *25th International Conference on Artificial Intelligence and Statistics*, 2022.
- [55] A. M. Metelli, A. Likmeta, and M. Restelli. Propagating uncertainty in reinforcement learning via wasserstein barycenters. *Advances in Neural Information Processing Systems*, 32, 2019.
- [56] C. D. Meyer. The role of the group generalized inverse in the theory of finite markov chains. *SIAM Review*, 17(3):443–464, 1975.
- [57] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning, 2013.
- [58] A. W. Moore. *Efficient Memory-based Learning for Robot Control*. PhD thesis, University of Cambridge, 1990.
- [59] S. Müller, A. von Rohr, and S. Trimpe. Local policy search with bayesian optimization. *Advances in Neural Information Processing Systems*, 34:20708–20720, 2021.
- [60] G. Neu and C. Pike-Burke. A unifying view of optimism in episodic reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1392–1403, 2020.

- [61] I. Osband, B. V. Roy, and D. Russo. (more) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems*, pages 3003–3011, 2013.
- [62] R. Ouhamma, D. Basu, and O. Maillard. Bilinear exponential family of mdps: frequentist regret bound with tractable exploration & planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 9336–9344, 2023.
- [63] V. M. Panaretos and Y. Zemel. Statistical aspects of wasserstein distances. *Annual Review of Statistics and its Applications*, 6:405–431, 2019.
- [64] G. Peyré. Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- [65] M. Pirodda, M. Restelli, and L. Bascetta. Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2-3):255–283, 2015.
- [66] M. L. Puterman. *Markov decision processes: Discrete stochastic dynamic programming*. John Wiley and Sons, 1994.
- [67] S. Raghvendra, P. Shirzadian, and K. Zhang. A new robust partial p -wasserstein-based metric for comparing distributions, 2024.
- [68] G. O. Ramos, A. L. C. Bazzan, and B. C. da Silva. Analysing the impact of travel information for minimising the regret of route choice. *Transportation Research Part C: Emerging Technologies*, 88:257–271, 2018.
- [69] G. O. Ramos, B. C. da Silva, and A. L. C. Bazzan. Learning to minimise regret in route choice. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, pages 846–855, 2017.
- [70] F. Santambrogio. *Optimal Transport for Applied Mathematicians*. Birkhäuser Cham, 2015.
- [71] H. Van Seijen, H. Nekoei, E. Racah, and S. Chandar. The loca regret: a consistent metric to evaluate model-based behavior in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:6562–6572, 2020.
- [72] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [73] S. R. Sinclair, S. Banerjee, and C. L. Yu. Adaptive discretization in online reinforcement learning. *Operations Research*, 71(5):1636–1652, 2023.
- [74] M. Sommerfeld, J. Schrieber, Y. Zemel, and A. Munk. Optimal transport: Fast probabilistic approximation with exact solvers. *Journal of Machine Learning Research*, 20:1–23, 2019.
- [75] K. Sun, Y. Zhao, Y. Liu, B. Jiang, and L. Kong. Distributional reinforcement learning via sinkhorn iterations. *arXiv preprint arXiv:2202.00769*, 2022.
- [76] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction, 2nd Edition*. MIT press, 2018.
- [77] U. Syed, M. Bowling, and R. E. Schapire. Apprenticeship learning using linear programming. In *Proceedings of the 25th International Conference on Machine Learning*, pages 1032–1039. PMLR, 2008.
- [78] H. Tang, R. Houthoofd, D. Foote, X. Chen, A. Stooke, Y. Duan, J. Schulman, F. De Turck, and P. Abbeel. #exploration: A study of count-based exploration for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2753–2762. Curran Associates, Inc., 2017.
- [79] C. Villani. *Optimal Transport Old and New*. Springer Berlin, Heidelberg, 2009.

- 653 [80] J. Wang, I. Trummer, and D. Basu. Udo: universal database optimization using reinforcement
654 learning. *Proceedings of the VLDB Endowment*, 14(13):3402–3414, 2021.
- 655 [81] P. N. Ward. Linear programming in reinforcement learning, 2021. MSc thesis.
- 656 [82] Y. Xie, X. Wang, R. Wang, and H. Zha. A fast proximal point method for computing exact
657 wasserstein distance. In *Uncertainty in artificial intelligence*, pages 433–453. PMLR, 2020.
- 658 [83] Y. Yang, T. Zhou, Q. He, L. Han, M. Pechenizkiy, and M. Fang. Task adaptation from skills:
659 Information geometry, disentanglement, and new objectives for unsupervised reinforcement
660 learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- 661 [84] L. Zhang, K. Tang, and X. Yao. Explicit planning for efficient exploration in reinforcement
662 learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett,
663 editors, *Advances in Neural Information Processing Systems*, volume 32, pages 7488–7497.
664 Curran Associates, Inc., 2019.
- 665 [85] M. Zhang, Y. Wang, X. Ma, L. Xia, J. Yang, Z. Li, and X. Li. Wasserstein distance guided
666 adversarial imitation learning with reward shape exploration. In *IEEE 9th Data Driven Control
667 and Learning Systems Conference*, pages 1165–1170, 2020.
- 668 [86] C. Zhu, X. Wang, T. Han, S. S. Du, and A. Gupta. Transferable reinforcement learning via
669 generalized occupancy models, 2024.

670 A Theoretical Analysis

671 A.1 MDP with Lipschitz Rewards

672 Given two metric spaces (\mathcal{X}, d_X) and (\mathcal{Y}, d_Y) , a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ is called 1-Lipschitz continuous
673 if [79]:

$$d_Y(f(x), f(x')) \leq d_X(x, x'), \forall (x, x') \in X \quad (11)$$

674 This implies that the Lipschitz semi-norm over the function space $\mathcal{F}(X, Y)$, defined as

$$\|f\|_L = \sup_{x \neq x'} \left\{ \frac{d_Y(f(x), f(x'))}{d_X(x, x')} \mid \forall (x, x') \in \mathcal{X} \right\}, \quad (12)$$

675 is ≤ 1 . When (\mathcal{X}, d_X) is a Polish space and $\mu, \nu \in \mathcal{P}(\mathcal{X})$, the **Kantorovich-Rubinstein** formula
676 states that [79]:

$$\begin{aligned} \mathcal{W}_1(\mu, \nu) &= \sup_{\|f\|_L \leq 1} \left\{ \int_{\mathcal{X}} f d\mu - \int_{\mathcal{X}} f d\nu \right\} \\ &= \sup_{\|f\|_L \leq 1} \{ \mathbb{E}_{\mu}[f(X)] - \mathbb{E}_{\nu}[f(X)] \}, \end{aligned} \quad (13)$$

677 where $\mathcal{W}_1(\mu, \nu)$ is the 1-Wasserstein distance between μ and ν with f as the cost function.

678 Note that when $\|f\|_L \leq L_{\mathcal{R}}$ for any $L_{\mathcal{R}} > 0$, then the function f is called $L_{\mathcal{R}}$ -Lipschitz continuous,
679 and Eq. (13) becomes [24],

$$\mathcal{W}_1(\mu, \nu) = \frac{1}{L_{\mathcal{R}}} \sup_{\|f\|_L \leq L_{\mathcal{R}}} \{ \mathbb{E}_{\mu}[f(X)] - \mathbb{E}_{\nu}[f(X)] \}. \quad (14)$$

680 Now, we consider $\mathcal{X} = \mathcal{S} \times \mathcal{A}$, i.e. the state-action space, $\mathcal{Y} = \mathbb{R}$, i.e. the real line, and the function
681 f to be the reward function \mathcal{R} . Then, we can call the reward function \mathcal{R} to be $L_{\mathcal{R}}$ -Lipschitz if

$$|\bar{\mathcal{R}}(s, a) - \bar{\mathcal{R}}(s', a')| \leq L_{\mathcal{R}} d_{\mathcal{S}\mathcal{A}}((s, a), (s', a'))$$

682 for all $s, s' \in \mathcal{S}$, and $a, a' \in \mathcal{A}$, and $d_{\mathcal{S}\mathcal{A}}((s, a), (s', a')) = d_{\mathcal{S}}((s, s')) + d_{\mathcal{A}}((a, a'))$ being the
683 metric on the state-action space $\mathcal{S} \times \mathcal{A}$. If the reward function \mathcal{R} of an MDP is $L_{\mathcal{R}}$ -Lipschitz, we
684 refer it as an MDP with Lipschitz rewards.

685 A.2 Performance Difference and Occupancy Measures

686 We know that

$$J_{\mu}^{\pi} = \frac{1}{\rho} \mathbb{E}_{(s,a) \sim v_{\pi}} [\bar{\mathcal{R}}(s, a)]. \quad (15)$$

687 Using Eq. (15), we write for two policies π and π' , with $\mu(s)$ as the initial state distribution,

$$\left| J_{\mu}^{\pi} - J_{\mu}^{\pi'} \right| = \frac{1}{\rho} \left| \mathbb{E}_{(s,a) \sim v_{\pi}} [\bar{\mathcal{R}}(s, a)] - \mathbb{E}_{(s,a) \sim v_{\pi'}} [\bar{\mathcal{R}}(s, a)] \right| \quad (16)$$

688 Given an MDP with $L_{\mathcal{R}}$ -Lipschitz rewards, the **Kantorovich-Rubinstein** formula dictates that [24]:

$$\sup_{\|\bar{\mathcal{R}}\|_L \leq L_{\mathcal{R}}} \left| \mathbb{E}_{(s,a) \sim v_{\pi}} [\bar{\mathcal{R}}(s, a)] - \mathbb{E}_{(s,a) \sim v_{\pi'}} [\bar{\mathcal{R}}(s, a)] \right| = L_{\mathcal{R}} \mathcal{W}_1(v_{\pi}, v_{\pi'}) \quad (17)$$

689 By dividing both sides of Eq. (17) by ρ , and due to an upper bound by the supremum, this inequality
690 follows:

$$\left| J_{\mu}^{\pi} - J_{\mu}^{\pi'} \right| \leq \frac{L_{\mathcal{R}}}{\rho} \mathcal{W}_1(v_{\pi}, v_{\pi'}) \quad (18)$$

691 A.3 Proof of proposition 1

692 The Linear Programming formulation for solving MDPs, assuming discrete state and action spaces,
693 is [66]:

$$\begin{aligned} & \max_{v_\pi} \sum_{s,a} r(s,a) v_\pi(s,a) \\ & \text{subject to } \sum_a v_\pi(s,a) = p_0(s) + \gamma \sum_{s',a} T(s | s',a) v_\pi(s',a) \\ & v_\pi(s,a) \geq 0 \quad \forall (s,a) \in \mathcal{S} \times \mathcal{A}, \end{aligned} \quad (19)$$

694 where $p_0(s)$ is the initial state distribution and $T(s | s',a)$ is the transition probability. The constraints
695 of this optimization problem are often referred to as *Bellman Flow Constraint*.

696 A stationary policy π has a corresponding occupancy measure $v_\pi(s,a)$ that satisfies the Bellman flow
697 constraint [77], and hence π and $v_\pi(s,a)$ share a bijective relationship [77, 26],

$$\pi(a | s) = \frac{v_\pi(s,a)}{u_\pi(s)} \quad (20)$$

698 with

$$u_\pi(s) = \sum_{a'} v_\pi(s,a') = p_0(s) + \gamma \sum_{s',a'} T(s | s',a') v_\pi(s',a') \quad (21)$$

699 By rearranging Eq. (20) to

$$v_\pi(s,a) = \pi(a | s) u_\pi(s) \quad (22)$$

700 and substituting Eq. (22) into Eq. (21), we can rewrite Eq. (21) as (defining $\mathcal{P}^\pi \triangleq \sum_a T(s |$
701 $s',a) \pi(a | s')$),

$$\begin{aligned} p_0(s) &= u_\pi(s) - \gamma \sum_{s',a} T(s | s',a) \pi(a | s') u_\pi(s') \\ &\triangleq u_\pi(s) - \gamma \sum_{s'} \mathcal{P}^\pi(s | s') u_\pi(s') \end{aligned} \quad (23)$$

702 which in matrix form is

$$\begin{aligned} \mathbf{p}_0 &= \mathbf{u}_\pi - \gamma \mathbf{P}^\pi \mathbf{u}_\pi \\ &= (\mathbb{I} - \gamma \mathbf{P}^\pi) \mathbf{u}_\pi, \end{aligned} \quad (24)$$

703 where $\mathbf{p}_0, \mathbf{u}_\pi \in \mathbb{R}^{|\mathcal{S}|}$ are column vectors and $\mathbf{P}^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ are matrices. Solving for \mathbf{u}_π , we get

$$\mathbf{u}_\pi = (\mathbb{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{p}_0 \quad (25)$$

704 The inverse matrix $(\mathbb{I} - \gamma \mathbf{P}^\pi)^{-1}$ exists because for $\gamma < 1$, $(\mathbb{I} - \gamma \mathbf{P}^\pi)$ is a strictly diagonally
705 dominant matrix [77]. Thus, $(\mathbb{I} - \gamma \mathbf{P}^\pi)^{-1} = \sum_{t=0}^{\infty} (\gamma \mathbf{P}^\pi)^t$, where $\sum_{t=0}^{\infty} (\gamma \mathbf{P}^\pi)^t$ forms a valid
706 *Neumann series* [81]. We let $\mathbf{A}^\pi = \sum_{t=0}^{\infty} (\gamma \mathbf{P}^\pi)^t$, so Eq. (25) can be written as $\mathbf{u}_\pi = \mathbf{A}^\pi \mathbf{p}_0$. We
707 can therefore express Eq. (22) in matrix form as:

$$\begin{aligned} \mathbf{v}_\pi &= \mathbf{\Pi} \odot (\mathbf{u}_\pi^T \otimes \mathbf{1})^T \\ &= \mathbf{\Pi} \odot (\mathbf{p}_0^T (\mathbf{A}^\pi)^T \otimes \mathbf{1})^T, \end{aligned} \quad (26)$$

708 where $\mathbf{\Pi}, \mathbf{v}_\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, $\mathbf{1} \in \mathbb{R}^{|\mathcal{A}|}$ is a column vector of ones, \otimes presents the Kronecker product,
709 and \odot denotes the Hadamard product.

710 If we consider the case of a parameterized policy $\mathbf{\Pi}(\theta)$, then the derivative of \mathbf{v}_π with respect to θ is

$$\begin{aligned} \nabla_\theta \mathbf{v}_\pi &= \nabla_\theta \left[\mathbf{\Pi} \odot (\mathbf{p}_0^T (\mathbf{A}^\pi)^T \otimes \mathbf{1})^T \right] \\ &= \nabla_\theta \mathbf{\Pi} \odot (\mathbf{p}_0^T (\mathbf{A}^\pi)^T \otimes \mathbf{1})^T + \mathbf{\Pi} \odot \nabla_\theta (\mathbf{p}_0^T (\mathbf{A}^\pi)^T \otimes \mathbf{1})^T \\ &= \nabla_\theta \mathbf{\Pi} \odot (\mathbf{p}_0^T (\mathbf{A}^\pi)^T \otimes \mathbf{1})^T + \mathbf{\Pi} \odot (\mathbf{p}_0^T (\nabla_\theta \mathbf{A}^\pi)^T \otimes \mathbf{1})^T \end{aligned} \quad (27)$$

711 The first term in Eq. (27) is differentiable since the policy is parameterized by θ . We expand $\nabla_\theta \mathbf{A}^\pi$
 712 as follows:

$$\begin{aligned}
 \nabla_\theta \mathbf{A}^\pi &= \nabla_\theta \left(\sum_{t=0}^{\infty} (\gamma \mathbf{P}^\pi)^t \right) \\
 &= \nabla_\theta (\mathbb{I} + (\gamma \mathbf{P}^\pi) + (\gamma \mathbf{P}^\pi)^2 + \dots + (\gamma \mathbf{P}^\pi)^t + \dots) \\
 &= 0 + \nabla_\theta (\gamma \mathbf{P}^\pi) + 2(\gamma \mathbf{P}^\pi) \nabla_\theta (\gamma \mathbf{P}^\pi) + 3(\gamma \mathbf{P}^\pi)^2 \nabla_\theta (\gamma \mathbf{P}^\pi) + \dots + t(\gamma \mathbf{P}^\pi)^{t-1} \nabla_\theta (\gamma \mathbf{P}^\pi) + \dots \\
 &= (\mathbb{I} + 2(\gamma \mathbf{P}^\pi) + 3(\gamma \mathbf{P}^\pi)^2 + \dots + t(\gamma \mathbf{P}^\pi)^{t-1} + \dots) \nabla_\theta (\gamma \mathbf{P}^\pi) \\
 &= \left(\sum_{t=0}^{\infty} (t+1)(\gamma \mathbf{P}^\pi)^t \right) \nabla_\theta (\gamma \mathbf{P}^\pi) \\
 &= \sum_{t=0}^{\infty} (t+1)(\gamma \mathbf{P}^\pi)^t \gamma \nabla_\theta \mathbf{P}^\pi \\
 &\equiv \sum_{t=0}^{\infty} (t+1)(\gamma \mathbf{P}^\pi)^t \gamma \nabla_\theta \left[\sum_{s',a} T(s|s',a) \pi(a|s') \right] \\
 &= \sum_{t=0}^{\infty} (t+1)(\gamma \mathbf{P}^\pi)^t \gamma \left[\sum_{s',a} T(s|s',a) \nabla_\theta \pi(a|s') \right].
 \end{aligned} \tag{28}$$

713 Since Eq. (28) shows that $\nabla_\theta \mathbf{A}^\pi$ is differentiable, $\nabla_\theta \mathbf{v}_\pi$ is also differentiable based on Eq. (27).
 714 Proceeding similarly, given the same conditions, we see that all higher derivatives of \mathbf{v}_π also exist
 715 with respect to θ . Thus, the space of parametrized occupancy measures v_π forms a differentiable
 716 manifold.

717 A.4 Proof of proposition 2

718 Regret is a common metric for evaluating agents, that measures the total loss an agent incurs over
 719 policy updates by using its policy in lieu of the optimal one, defined as [61],

$$\text{Regret} = \mathbb{E}_{s \sim \mu} \left[\sum_k (V^*(s) - V_{\pi_k}(s)) \right] \tag{29}$$

720 where $V^* = V_{\pi^*}$ is the value function of the optimal policy π^* while $V_{\pi_k}(s)$ is the value function of
 721 policy π_k , and μ is the initial state distribution.

722 Since $J_\mu^\pi = \mathbb{E}_{s \sim \mu} [V_\pi(s)]$, we can conclude from Eq. (29) that

$$\begin{aligned}
 \text{Regret} &= \mathbb{E}_{s \sim \mu} \left[\sum_k (V^*(s) - V_{\pi_k}(s)) \right] \\
 &= \sum_k [\mathbb{E}_{s \sim \mu} (V^*(s) - V_{\pi_k}(s))] \\
 &= \sum_k (J_\mu^{\pi^*} - J_\mu^{\pi_k}) \\
 &= \sum_k |J_\mu^{\pi^*} - J_\mu^{\pi_k}| \\
 &\leq \sum_k \frac{L_{\mathcal{R}}}{\rho} \mathcal{W}_1(v_{\pi^*}, v_{\pi_k})
 \end{aligned} \tag{30}$$

723 The last inequality is due to Eq. (18).

724 A.5 Proof of proposition 3

725 Let us begin the proof by defining the visitation probability at any step $h \in [H]$ in an episode,
 726 following policy $\pi(a|s)$. Specifically,

$$q_\pi^h(s, a) \triangleq \mathbb{P}(s_h = s, a_h = a) \quad \forall h \in [H] \quad \text{and} \quad q_\pi^h(s, a) \triangleq 0 \quad \forall h \in \mathbb{N} \wedge h > H. \quad (31)$$

727 Thus, we rewrite Eq. (7) as $v_\pi^H(s, a) = \frac{1}{H} \sum_{h=1}^H q_\pi^h(s, a)$.

728 Then, following [37], we can write the Linear Programming formulation for solving episodic MDP
 729 \mathbb{M}^H as

$$\begin{aligned} & \max_{\{q_\pi^h\}_{h=1}^H} \sum_{h,s,a} r(s, a) q_\pi^h(s, a) \\ & \text{subject to} \sum_a q_\pi^h(s, a) = \sum_{s',a} T(s | s', a) q_\pi^{h-1}(s', a) \quad \forall h \in [H] \wedge h > 1, \\ & q_\pi^1(s, a) = \pi(a|s) \mu(s), \\ & q_\pi^h(s, a) \geq 0 \quad \forall h \in [H], (s, a) \in \mathcal{S} \times \mathcal{A}, \end{aligned} \quad (32)$$

730 where $\mu(s)$ is the initial state distribution and $T(s | s', a)$ is the transition probability. The constraints
 731 of this optimization problem are often referred to as *Bellman Flow Constraints*.

732 This implies that

$$\begin{aligned} & \sum_{h=2}^{H+1} \sum_a q_\pi^h(s, a) = \sum_{h=2}^{H+1} \sum_{s',a} T(s | s', a) q_\pi^{h-1}(s', a) \\ \implies & \sum_a q_\pi^1(s, a) + \sum_{h=2}^{H+1} \sum_a q_\pi^h(s, a) = \sum_{h=2}^{H+1} \sum_{s',a} T(s | s', a) q_\pi^{h-1}(s', a) + \sum_a q_\pi^1(s, a) \\ \implies & \sum_a \sum_{h=1}^{H+1} q_\pi^h(s, a) = \sum_{h=2}^{H+1} \sum_{s',a} T(s | s', a) q_\pi^{h-1}(s', a) + \sum_a q_\pi^1(s, a) \\ \implies & H \sum_a v_\pi^H(s, a) = \sum_{s',a} T(s | s', a) \left(\sum_{h=2}^{H+1} q_\pi^{h-1}(s', a) \right) + \mu(s) \\ \implies & H \sum_a v_\pi^H(s, a) = H \sum_{s',a} T(s | s', a) v_\pi^H(s', a) + \mu(s) \\ \implies & \sum_a v_\pi^H(s, a) = \sum_{s',a} T(s | s', a) v_\pi^H(s', a) + \frac{1}{H} \mu(s) \\ \implies & u_\pi^H(s) \triangleq \sum_a v_\pi^H(s, a) = \sum_{s',a} T(s | s', a) \pi(a|s') u_\pi^H(s') + \frac{1}{H} \mu(s). \end{aligned} \quad (33)$$

733 Now, we denote \mathbf{u}_π^H and $\bar{\mu}$ as corresponding column vectors and the transition matrix $\mathbf{P}^\pi \triangleq$
 734 $\left[\sum_{s',a} T(s | s', a) \pi(a|s') \right]$. Thus, we obtain

$$(\mathbb{I} - \mathbf{P}^\pi) \mathbf{u}_\pi^H = \frac{1}{H} \bar{\mu}. \quad (34)$$

735 We can therefore express the finite horizon occupancy measure in matrix form as

$$\mathbf{v}_\pi^H = \mathbf{\Pi} \odot \left((\mathbf{u}_\pi^H)^T \otimes \mathbf{1} \right)^T \quad (35)$$

736 where $\mathbf{\Pi}, \mathbf{v}_\pi \in \mathbb{R}^{|S| \times |A|}$, $\mathbf{1} \in \mathbb{R}^{|A|}$ is a column vector of ones, \otimes presents the Kronecker product, \odot
 737 denotes the Hadamard product.

738 If we consider the case of a parameterized policy $\Pi(\theta)$, the derivative of \mathbf{v}_π^H with respect to θ is

$$\begin{aligned}\nabla_\theta \mathbf{v}_\pi^H &= \nabla_\theta \left[\Pi \odot ((\mathbf{u}_\pi^H)^T \otimes \mathbf{1})^T \right] \\ &= \nabla_\theta \Pi \odot ((\mathbf{u}_\pi^H)^T \otimes \mathbf{1})^T + \Pi \odot \nabla_\theta ((\mathbf{u}_\pi^H)^T \otimes \mathbf{1})^T \\ &= \nabla_\theta \Pi \odot ((\mathbf{u}_\pi^H)^T \otimes \mathbf{1})^T + \Pi \odot ((\nabla_\theta \mathbf{u}_\pi^H)^T \otimes \mathbf{1})^T\end{aligned}\quad (36)$$

739 The first term in Eq. (36) is differentiable since the policy is parameterized by θ . We show that $\nabla_\theta \mathbf{u}_\pi^H$
740 exists using Eq. (34) as follows:

$$\begin{aligned}\nabla_\theta [(\mathbb{I} - \mathbf{P}^\pi) \mathbf{u}_\pi^H] &= \frac{1}{H} \nabla_\theta \bar{\mu} \\ \nabla_\theta (\mathbb{I} - \mathbf{P}^\pi) \mathbf{u}_\pi^H + (\mathbb{I} - \mathbf{P}^\pi) \nabla_\theta \mathbf{u}_\pi^H &= 0 \\ -\nabla_\theta \mathbf{P}^\pi \mathbf{u}_\pi^H + (\mathbb{I} - \mathbf{P}^\pi) \nabla_\theta \mathbf{u}_\pi^H &= 0 \\ (\mathbb{I} - \mathbf{P}^\pi) \nabla_\theta \mathbf{u}_\pi^H &= \nabla_\theta \mathbf{P}^\pi \mathbf{u}_\pi^H\end{aligned}\quad (37)$$

741 In Eq. (37), we observe that $(\mathbb{I} - \mathbf{P}^\pi)^{-1}$ may or may not exist. To address this, we use the group
742 inverse, a generalized matrix inverse that extends the concept of inversion to both singular and
743 invertible matrices [19]. For a square matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$, its group inverse $\mathbf{C}^\#$ satisfies the conditions
744 $\mathbf{C}\mathbf{C}^\#\mathbf{C} = \mathbf{C}$, $\mathbf{C}^\#\mathbf{C}\mathbf{C}^\# = \mathbf{C}^\#$ and $\mathbf{C}\mathbf{C}^\# = \mathbf{C}^\#\mathbf{C}$.

745 The group inverse is useful in analysing Markov chains [56] and coincides with the standard inverse
746 when \mathbf{C} is invertible. Using group inverse $(\mathbb{I} - \mathbf{P}^\pi)^\#$, Eq. (37) can be expressed as,

$$\nabla_\theta \mathbf{u}_\pi^H = (\mathbb{I} - \mathbf{P}^\pi)^\# \nabla_\theta \mathbf{P}^\pi \mathbf{u}_\pi^H + \mathbf{w} \quad (38)$$

747 where \mathbf{w} is a vector in the null space of $(\mathbb{I} - \mathbf{P}^\pi)$. Note that $\mathbf{w} = \mathbf{0}$ for ergodic Markov chains [56].
748 An episodic MDP induces an ergodic Markov chain that admits a unique stationary distribution [11]
749 (see also Appendix B.3). Thus, Eq. (38) simplifies to

$$\nabla_\theta \mathbf{u}_\pi^H = (\mathbb{I} - \mathbf{P}^\pi)^\# \nabla_\theta \mathbf{P}^\pi \mathbf{u}_\pi^H \quad (39)$$

$$= (\mathbb{I} - \mathbf{P}^\pi)^\# \left[\sum_{s', a} T(s|s', a) \nabla_\theta \pi(a|s') \right] \mathbf{u}_\pi^H \quad (40)$$

750 which shows that $\nabla_\theta \mathbf{u}_\pi^H$ is differentiable and so is $\nabla_\theta \mathbf{v}_\pi^H$. Proceeding similarly, given the same
751 conditions, we see that all higher derivatives of \mathbf{v}_π^H also exist with respect to θ . Thus, the space of
752 parametrized finite-horizon occupancy measures \mathbf{v}_π^H forms a differentiable manifold \mathcal{M}^H .

753 A.6 Optimal Transport Dataset Distance (OTDD)

754 Suppose we have two datasets, each consisting of feature-label pairs, $\mathcal{D}_A = \{(t_A^i, u_A^i)\}_{i=1}^m \sim$
755 $P_A(t, u)$ and $\mathcal{D}_B = \{(t_B^i, u_B^i)\}_{i=1}^n \sim P_B(t, u)$ with $t_A, t_B \in \mathcal{T}$ and $u_A, u_B \in \mathcal{U}_A, \mathcal{U}_B$. These
756 datasets can be used to create empirical distributions $\hat{P}_A(t, u)$ and $\hat{P}_B(t, u)$. OTDD is the p-
757 Wasserstein distance between the datasets \mathcal{D}_A and \mathcal{D}_B - which is essentially the distance between
758 their empirical distributions \hat{P}_A and \hat{P}_B - with the cost function defined as the metric of the joint
759 space $\mathcal{T} \times \mathcal{U}$ [5].

760 Naturally, the metric on this joint space can be defined as $d_{\mathcal{T}\mathcal{U}}((t, u), (t', u')) =$
761 $(d_{\mathcal{T}}(t, t')^p + d_{\mathcal{U}}(u, u')^p)^{1/p}$, for $p \geq 1$. However, in most applications $d_{\mathcal{T}}$ is readily available,
762 while $d_{\mathcal{U}}$ might be scarce, especially in supervised learning (SL) between labels from unrelated label
763 sets [5]. Further, we want $d_{\mathcal{T}}$ and $d_{\mathcal{U}}$ to have the same units to be addable. To overcome these issues,
764 $d_{\mathcal{U}}$ is expressed in terms of $d_{\mathcal{T}}$ by mapping labels u to distributions over the feature space $\mathcal{P}(\mathcal{T})$
765 as $u \rightarrow \alpha_u(T) \triangleq P(T | U = u) \in \mathcal{P}(\mathcal{T})$. Therefore, the distance between the labels u and u' is
766 defined as the p-Wasserstein distance between $\alpha_u(T)$ and $\alpha_{u'}(T)$,

$$\begin{aligned}d_{\mathcal{U}}(u, u') &= \mathcal{W}_p(\alpha_u(T), \alpha_{u'}(T)) \\ &= \left(\min_{\pi \in \Pi(\alpha_u, \alpha_{u'})} \int_{\mathcal{T} \times \mathcal{T}} (d_{\mathcal{T}}(t, t'))^p d\pi(t, t') \right)^{1/p}\end{aligned}\quad (41)$$

767 The metric on the joint space becomes,

$$d_{\mathcal{T}\mathcal{U}}((t, u), (t', u')) = (d_{\mathcal{T}}(t, t')^p + \mathcal{W}_p^p(\alpha_u(T), \alpha_{u'}(T)))^{1/p} \quad (42)$$

768 Let $\mathcal{Z} = \mathcal{T} \times \mathcal{U}$, then the p-Wasserstein distance between $\hat{P}_A(t, u)$ and $\hat{P}_B(t, u)$ is a "nested"
769 Wasserstein distance:

$$\begin{aligned} \mathcal{W}_p^p(\hat{P}_A, \hat{P}_B) &= \min_{\pi \in \Pi(P_A, P_B)} \int_{\mathcal{Z} \times \mathcal{Z}} (d_{\mathcal{Z}}(z, z'))^p d\pi \\ &= \min_{\pi \in \Pi(P_A, P_B)} \int_{\mathcal{T}\mathcal{U} \times \mathcal{T}\mathcal{U}} (d_{\mathcal{T}}(t, t')^p + \mathcal{W}_p^p(\alpha_u, \alpha_{u'})) d\pi \end{aligned} \quad (43)$$

770 $\mathcal{W}_p^p(\hat{P}_A, \hat{P}_B)$ is the OTDD between datasets \mathcal{D}_A and \mathcal{D}_B , often expressed as $d_{OT}(\mathcal{D}_A, \mathcal{D}_B)$. This is
771 used in transfer learning to determine the distance (or similarity) between datasets. Figure 6 illustrates
772 OTDD when applied to RL using datasets of state-action pairs.

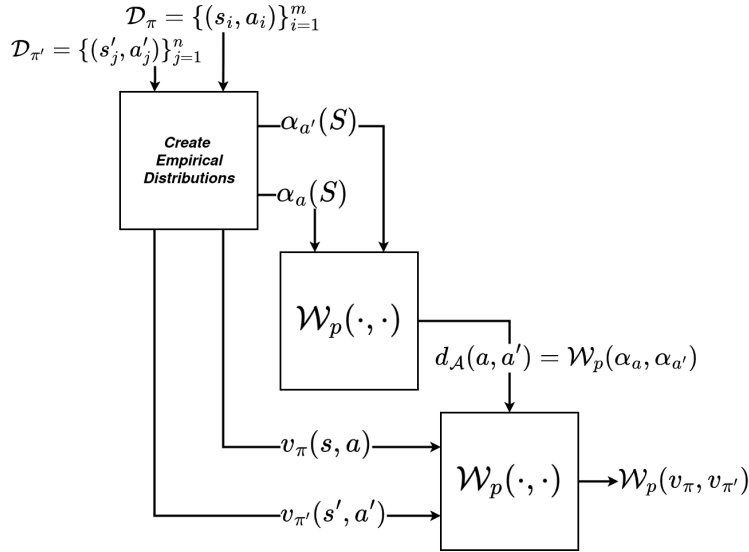


Figure 6: Illustration of OTDD when applied to RL.

773 A.7 Proof of Proposition 4

774 We compute the error in occupancy measure for both the infinite and finite horizon cases. In infinite
775 horizon MDPs, the occupancy measure is defined as the expected discounted number of visits of a
776 state-action pair (s, a) in a trajectory [43]: $\mu = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mu_t$, where $\mu_t = P(s_t, a_t \mid \pi, \eta)$ is
777 the state-action probability distribution at time step t with the initial state distribution η following
778 the policy π . In finite horizon MDPs, the occupancy measure is the expected number of visits of a
779 state-action pair (s, a) in an episode of length H [4]: $\mu = \frac{1}{H} \sum_{t=1}^H \mu_t$.

780 First, we derive error bounds for the infinite horizon MDP in which $\gamma < 1$ and the occupancy measure
781 is approximated using a finite number of samples collected up to a finite number of time steps T .
782 Later, we derive error bounds for the finite horizon MDP.

783 A.7.1 Infinite Horizon MDPs

784 **Estimated Occupancy Measure.** For convenience, we express the occupancy measure as $\mu =$
785 $(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mu_t$, where $\mu_t = P(s_t, a_t \mid \pi, \eta)$ is the state-action probability distribution at time
786 step t with the initial state distribution η following the policy π . To compute μ , we roll out N episodes
787 (each of multiple time steps) using π , and take N number of samples at t to approximate μ_t . Thus,
788 the empirical occupancy measure $\hat{\mu}$ is given by $\hat{\mu} = \rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^N$, where $\rho = \frac{1}{\sum_{t=0}^T \gamma^t}$. Note that
789 the total number of samples in the policy dataset \mathcal{D}_π is $|\mathcal{D}_\pi| = N(T + 1)$.

790 **Occupancy Measure Estimation Error.** Consider two occupancy measures $\mu = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mu_t$
 791 and $\nu = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \nu_t$ (with estimates $\hat{\mu} = \rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}$ and $\hat{\nu} = \rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}$). For
 792 independent sets $\{\mu_t\}_{t \geq 0}$ and $\{\nu_t\}_{t \geq 0}$, the Wasserstein distance has the following additive property
 793 [63],

$$\mathcal{W}_p(\sum_t \mu_t, \sum_t \nu_t) \leq \sum_t \mathcal{W}_p(\mu_t, \nu_t) \quad (44)$$

794 While for $a \in \mathbb{R}$ [63],

$$\mathcal{W}_p(a\mu, av) = |a| \mathcal{W}_p(\mu, v) \quad (45)$$

795 Therefore, for our scenario where $p = 1$, the Wasserstein distance between μ and ν is given by:

$$\begin{aligned} \mathcal{W}_1(\mu, \nu) &= \mathcal{W}_1((1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mu_t, (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \nu_t) \\ &\leq (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathcal{W}_1(\mu_t, \nu_t) \end{aligned} \quad (46)$$

796 while for $\hat{\mu}$ and $\hat{\nu}$,

$$\mathcal{W}_1(\hat{\mu}, \hat{\nu}) \leq \rho \sum_{t=0}^T \gamma^t \mathcal{W}_1(\hat{\mu}_t^{N_\mu}, \hat{\nu}_t^{N_\nu}) \quad (47)$$

797 In the RL problems we consider, the state-action space $\mathcal{Z} = \mathcal{S} \times \mathcal{A}$ is commonly defined as the
 798 subset of the Euclidean space $\mathcal{Z} \in \mathbb{R}^B$, where usually $B \geq 2$. Theorems 1 and 3 in [74] establish the
 799 following error bounds between the true and empirical probability distributions,

$$\begin{aligned} \mathbb{E}[\mathcal{W}_1(\hat{\mu}_t^{N_\mu}, \mu_t)] &\leq \mathcal{E}_2 N_\mu^{-\frac{1}{2}} \\ \mathbb{E}[\mathcal{W}_1(\hat{\nu}_t^{N_\nu}, \nu_t)] &\leq \mathcal{E}_2 N_\nu^{-\frac{1}{2}} \end{aligned} \quad (48)$$

800 where

$$\mathcal{E}_2 \leq 4B^{1/2} \text{diam}(\mathcal{Z}) \cdot \begin{cases} 2 + (1/2)\log_2 |\mathcal{Z}| & \text{if } B = 2 \\ |\mathcal{Z}|^{1/2-1/B} [2 + 1/(2^{B/2-1} - 1)] & \text{if } B > 2 \end{cases}$$

801 Note that $|\mathcal{Z}|$ and $\text{diam}(\mathcal{Z})$ denote the cardinality and diameter of \mathcal{Z} , respectively.

802 Suppose $a = \mathcal{W}_1(\hat{\mu}, \hat{\nu})$, $b = \mathcal{W}_1(\hat{\mu}, \mu)$, $c = \mathcal{W}_1(\hat{\nu}, \mu)$, $d = \mathcal{W}_1(\mu, \nu)$, and $e = \mathcal{W}_1(\hat{\nu}, \nu)$. Then by
 803 performing two reverse triangle inequalities,

$$\begin{aligned} |a - c| &\leq b \quad \text{and} \quad |c - d| \leq e \\ \implies |a - d| &\leq b + e \end{aligned} \quad (49)$$

804 Eq. (49) implies that,

$$\begin{aligned} \mathbb{E}[|\mathcal{W}_1(\hat{\mu}, \hat{\nu}) - \mathcal{W}_1(\mu, \nu)|] &\leq \mathbb{E}[\mathcal{W}_1(\hat{\mu}, \mu) + \mathcal{W}_1(\hat{\nu}, \nu)] \\ &= \mathbb{E}[\mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}, \mu) + \mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}, \nu)] \\ &= \mathbb{E}[\mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}, \mu)] + \mathbb{E}[\mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}, \nu)] \\ &\quad + \mathbb{E}[\mathcal{W}_1((1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, \mu) - \mathcal{W}_1((1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, \mu)] \\ &\quad + \mathbb{E}[\mathcal{W}_1((1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, \nu) - \mathcal{W}_1((1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, \nu)] \end{aligned} \quad (50)$$

805 By virtue of triangle inequalities, we get

$$\begin{aligned} \mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}) &\geq \mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}, \mu) - \mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, \mu) \\ \mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}) &\geq \mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}, \nu) - \mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, \nu) \end{aligned} \quad (51)$$

806 Therefore, the right-hand-side (R.H.S) of Eq. (50) can be further simplified as

$$\begin{aligned} \text{R.H.S} &\leq \mathbb{E}[\mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu})] + \mathbb{E}[\mathcal{W}_1(\rho \sum_{t=0}^T \gamma^t \hat{\nu}_t^{N_\nu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu})] \\ &\quad + \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, \mu)] + \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, \nu)] \end{aligned} \quad (52)$$

807 For simplicity, we denote $\hat{\mu}_\infty = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}$ (similarly $\hat{\nu}_\infty$) and $\hat{\mu}_T = \rho \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu}$
 808 (similarly $\hat{\nu}_T$), where $\rho = \frac{1}{\sum_{t=0}^T \gamma^t} = \frac{1-\gamma}{1-\gamma^{T+1}}$. Using Theorem 4 in [25], the 1-Wasserstein metric
 809 \mathcal{W}_1 and the total variation distance d_{TV} satisfy the following,

$$\begin{aligned} \mathcal{W}_1(\hat{\mu}_\infty, \hat{\mu}_T) &\leq \text{diam}(\mathcal{Z}) \cdot d_{TV}(\hat{\mu}_\infty, \hat{\mu}_T) \\ &= \text{diam}(\mathcal{Z}) \cdot \frac{1}{2} \sum_{z \in \mathcal{Z}} |\hat{\mu}_\infty(z) - \hat{\mu}_T(z)| \end{aligned} \quad (53)$$

810 However,

$$\begin{aligned} \hat{\mu}_\infty - \hat{\mu}_T &= (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} - \frac{1-\gamma}{1-\gamma^{T+1}} \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \\ &= (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} - \frac{1-\gamma}{1-\gamma^{T+1}} \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \\ &\quad + (1-\gamma) \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} - (1-\gamma) \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \\ &= (1-\gamma) \left(\sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} - \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \right) + \left((1-\gamma) - \frac{1-\gamma}{1-\gamma^{T+1}} \right) \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \\ &= (1-\gamma) \sum_{t=T+1}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} - \gamma^{T+1} \frac{1-\gamma}{1-\gamma^{T+1}} \sum_{t=0}^T \gamma^t \hat{\mu}_t^{N_\mu} \\ &\leq (1-\gamma) \sum_{t=T+1}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} \\ &= \gamma^{T+1} \frac{1-\gamma}{\gamma^{T+1}} \sum_{t=T+1}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu} \\ &= \gamma^{T+1} \hat{\mu}_{T+1, \infty} \end{aligned} \quad (54)$$

811 where $\frac{1-\gamma}{\gamma^{T+1}}$ normalizes $\sum_{t=T+1}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}$. We utilize Eq. (54) in Eq. (53) as,

$$\begin{aligned} \mathcal{W}_1(\hat{\mu}_\infty, \hat{\mu}_T) &\leq \text{diam}(\mathcal{Z}) \cdot \frac{1}{2} \sum_{z \in \mathcal{Z}} |\hat{\mu}_\infty(z) - \hat{\mu}_T(z)| \\ &\leq \text{diam}(\mathcal{Z}) \cdot \frac{1}{2} \sum_{z \in \mathcal{Z}} |\gamma^{T+1} \hat{\mu}_{T+1, \infty}(z)| \\ &= \frac{\gamma^{T+1}}{2} \text{diam}(\mathcal{Z}) \end{aligned} \quad (55)$$

Eq. (55) also applies for $\mathcal{W}_1(\hat{\nu}_\infty, \hat{\nu}_T)$, therefore by substituting these into Eq. (52),

$$\begin{aligned}
\text{R.H.S} &\leq \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, \mu)] + \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, \nu)] + \gamma^{T+1} \text{diam}(\mathcal{Z}) \\
&= \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\mu}_t^{N_\mu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \mu_t)] \\
&\quad + \mathbb{E}[\mathcal{W}_1((1-\gamma) \sum_{t=0}^{\infty} \gamma^t \hat{\nu}_t^{N_\nu}, (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \nu_t)] + \gamma^{T+1} \text{diam}(\mathcal{Z}) \\
&\leq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \left(\mathbb{E}[\mathcal{W}_1(\hat{\mu}_t^{N_\mu}, \mu_t)] + \mathbb{E}[\mathcal{W}_1(\hat{\nu}_t^{N_\nu}, \nu_t)] \right) + \gamma^{T+1} \text{diam}(\mathcal{Z}).
\end{aligned} \tag{56}$$

By substituting Eq. (48) into Eq. (56)

$$\begin{aligned}
\text{R.H.S} &\leq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \left(\mathcal{E}_2 N_\mu^{-\frac{1}{2}} + \mathcal{E}_2 N_\nu^{-\frac{1}{2}} \right) + \gamma^{T+1} \text{diam}(\mathcal{Z}) \\
&= \mathcal{E}_2 \left(N_\mu^{-\frac{1}{2}} + N_\nu^{-\frac{1}{2}} \right) + \gamma^{T+1} \text{diam}(\mathcal{Z})
\end{aligned} \tag{57}$$

Therefore, Eq. (50) becomes:

$$\mathbb{E}[|\mathcal{W}_1(\hat{\mu}, \hat{\nu}) - \mathcal{W}_1(\mu, \nu)|] \leq \mathcal{E}_2 \left(N_\mu^{-\frac{1}{2}} + N_\nu^{-\frac{1}{2}} \right) + \gamma^{T+1} \text{diam}(\mathcal{Z}) \tag{58}$$

Over the full trajectory in the occupancy measure space. The true distance between consecutive policies π_i and π_{i+1} after an update is $\mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}})$, which is induced by the i^{th} policy update. We estimate this distance using datasets of the policies, i.e. approximated distributions, using $\mathcal{W}_1(\hat{v}_{\pi_i}, \hat{v}_{\pi_{i+1}})$.

For M roll out episodes of each π_i , we use Eq. (58), with $N_\mu = N_\nu = M$, to derive the following error bounds,

$$\mathbb{E}[|\mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) - \mathcal{W}_1(\hat{v}_{\pi_i}, \hat{v}_{\pi_{i+1}})|] \leq 2\mathcal{E}_2 M^{-\frac{1}{2}} + \gamma^{T+1} \text{diam}(\mathcal{Z}) \tag{59}$$

which is consistent with learning from \mathcal{D}_{π_i} and then $\mathcal{D}_{\pi_{i+1}}$. By summing sequentially through policies encountered during RL training, we compute the total distance over a path of N segments obtained via policy updates:

$$\sum_{i=0}^{N-1} \mathbb{E}[|\mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) - \mathcal{W}_1(\hat{v}_{\pi_i}, \hat{v}_{\pi_{i+1}})|] \leq 2N\mathcal{E}_2 M^{-\frac{1}{2}} + N\gamma^{T+1} \text{diam}(\mathcal{Z}) \tag{60}$$

Since $|\sum_t x_t| \leq \sum_t |x_t|$ then,

$$\mathbb{E} \left[\left| \sum_{i=0}^{N-1} \mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) - \sum_{i=0}^{N-1} \mathcal{W}_1(\hat{v}_{\pi_i}, \hat{v}_{\pi_{i+1}}) \right| \right] \leq \frac{2N\mathcal{E}_2}{\sqrt{M}} + N\gamma^{T+1} \text{diam}(\mathcal{Z}) \tag{61}$$

A.7.2 Finite Horizon MDPs

Occupancy Measure Estimated Error. Consider two occupancy measures $\mu = \frac{1}{H} \sum_{t=1}^H \mu_t$ and $\nu = \frac{1}{H} \sum_{t=1}^H \nu_t$ with estimates $\hat{\mu} = \frac{1}{H} \sum_{t=1}^H \hat{\mu}_t^{N_\mu}$ and $\hat{\nu} = \frac{1}{H} \sum_{t=1}^H \hat{\nu}_t^{N_\nu}$. From Eq. (49), we have

$$\begin{aligned}
&\mathbb{E}[|\mathcal{W}_1(\hat{\mu}, \hat{\nu}) - \mathcal{W}_1(\mu, \nu)|] \\
&\leq \mathbb{E}[\mathcal{W}_1(\hat{\mu}, \mu) + \mathcal{W}_1(\hat{\nu}, \nu)] \\
&= \mathbb{E}[\mathcal{W}_1(\frac{1}{H} \sum_{t=1}^H \hat{\mu}_t^{N_\mu}, \frac{1}{H} \sum_{t=1}^H \mu_t) + \mathcal{W}_1(\frac{1}{H} \sum_{t=1}^H \hat{\nu}_t^{N_\nu}, \frac{1}{H} \sum_{t=1}^H \nu_t)] \\
&\leq \frac{1}{H} \sum_{t=1}^H \mathbb{E}[\mathcal{W}_1(\hat{\mu}_t^{N_\mu}, \mu_t)] + \frac{1}{H} \sum_{t=1}^H \mathbb{E}[\mathcal{W}_1(\hat{\nu}_t^{N_\nu}, \nu_t)] \\
&\leq \mathcal{E}_2 \left(N_\mu^{-\frac{1}{2}} + N_\nu^{-\frac{1}{2}} \right)
\end{aligned} \tag{62}$$

Therefore **for the total path in the occupancy measure space** with M roll out episodes of each π_i , the error bound is

$$\mathbb{E} \left[\left| \sum_{i=0}^{N-1} \mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) - \sum_{i=0}^{N-1} \mathcal{W}_1(\hat{v}_{\pi_i}, \hat{v}_{\pi_{i+1}}) \right| \right] \leq \frac{2N\mathcal{E}_2}{\sqrt{M}} \quad (63)$$

by assigning $N_\mu = N_\nu = M$ in Eq. (62), which concludes the proof.

A.8 Proof of Proposition 5

By definition of η_{sub} , we get

$$\begin{aligned} \eta_{sub} &= \frac{\sum_{i=0}^{N-2} \mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) + \mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi_N})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \\ &= \frac{\sum_{i=0}^{N-2} \mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) + \mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi_N})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})} \times \frac{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \\ &\geq \frac{\sum_{i=0}^{N-2} \mathcal{W}_1(v_{\pi_i}, v_{\pi_{i+1}}) + \mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi^*}) - \mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})} \times \frac{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \\ &= \left(\eta - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right) \frac{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})}. \end{aligned} \quad (64)$$

The inequality above is true due to the triangle inequality $\mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi_N}) + \mathcal{W}_1(v_{\pi_N}, v_{\pi^*}) \geq \mathcal{W}_1(v_{\pi_{N-1}}, v_{\pi^*})$.

By applying triangle inequality, we also get

$$\mathcal{W}_1(v_{\pi_0}, v_{\pi^*}) + \mathcal{W}_1(v_{\pi_N}, v_{\pi^*}) \geq \mathcal{W}_1(v_{\pi_0}, v_{\pi_N}).$$

This implies that

$$\frac{\mathcal{W}_1(v_{\pi_0}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \geq 1 - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})}. \quad (65)$$

Eq. (64) and Eq. (65) together yield

$$\begin{aligned} \eta_{sub} &\geq \left(\eta - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right) \left(1 - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right) \\ &= \eta - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} - \eta \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} + \left(\frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right)^2 \\ &\geq \eta \left(1 - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right) - \frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \\ &\geq \eta \left(1 - \frac{2\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right). \end{aligned}$$

The second last inequality is due to non-negativity of $\left(\frac{\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})} \right)^2$. The last inequality is due to the fact that $\eta \geq 1$.

Thus, we conclude that

$$\frac{\eta - \eta_{sub}}{\eta} \leq \frac{2\mathcal{W}_1(v_{\pi_N}, v_{\pi^*})}{\mathcal{W}_1(v_{\pi_0}, v_{\pi_N})}.$$

A.9 Wasserstein Spaces as Geodesic Spaces

Given probability measures $\mu, \nu \in \mathcal{P}(\mathcal{X})$ on a metric space $\mathcal{X} \subset \mathbb{R}^B$ with metric $d_{\mathcal{X}}(x, x')$, the Wasserstein distance $\mathcal{W}_p(\mu, \nu)$ is the minimal transport cost for $c(x, x') = (d_{\mathcal{X}}(x, x'))^p$ with $p \geq 1$ [79]. The Wasserstein distance $\mathcal{W}_p(\mu, \nu)$ takes a distance on \mathcal{X} and creates out of it a distance on $\mathcal{P}(\mathcal{X})$ [64]. Proposition 5.1 in [70] asserts that \mathcal{W}_p is a distance over $\mathcal{P}(\mathcal{X})$.

846 **Definition A.9** (Wasserstein Space). [70] *Given a Polish space \mathcal{X} , for each $p \in [1, \infty)$, the space*
847 *$\mathcal{P}(\mathcal{X})$ endowed with the distance \mathcal{W}_p is a Wasserstein space \mathbb{W}_p of order p .*

848 Theorem 5.27 in [70] states that if \mathcal{X} is a convex space, then the space \mathbb{W}_p is a geodesic space (length
849 space). Thus, the geodesic (shortest path distance) between $\mu, \nu \in \mathcal{P}(\mathcal{X})$ is given by $\mathcal{W}_p(\mu, \nu)$ [41].
850 It was mentioned in Appendix A.7.1 that the RL problems we consider consist of the state-action
851 space $\mathcal{Z} = \mathcal{S} \times \mathcal{A} \in \mathbb{R}^B : B \geq 2$ (subsets of the Euclidean space). Given that Euclidean spaces are
852 convex spaces [12], our space of occupancy measures \mathcal{M} is a Wasserstein space $\mathbb{W}_1 = (\mathcal{M}, \mathcal{W}_1)$ and
853 thus a geodesic space. Therefore, $\mathcal{W}_1(\mu, \nu)$ measures the shortest path on the surface of the manifold
854 \mathcal{M} between probability distributions μ and ν .

B Additional Experimental Analysis and Results

B.1 Environment Description

2D-Gridworld environment of size 5×5 with actions: {up, right, down, left}. The start and goal states are always located at top-left and bottom-right, respectively. In the gridworld, we perform experiments on three settings namely:- A) *Deterministic, dense rewards setting*: State transitions are deterministic. The reward function is given by $\|s_t - s_g\|_1$, where s_t is the agent state at timestep t and s_g is the goal state. B) *Deterministic, sparse rewards setting*: State transitions are deterministic and all states issue a reward of -0.04 except the goal state with reward of 1. C) *Stochastic, dense rewards setting*: Actions have a probability of 0.8 in the instructed direction and 0.1 in each adjacent direction. Reward function is as defined in setting A.

2D-Gridworlds (Task Difficulty). Figure 7 depicts the configurations of the 5 tasks that were used to assess ESL with respect to task hardness. They are all deterministic with actions: {up, right, down, left}, and mostly have the start-state at the top-left and the goal-state at the bottom-right - with only one task that has the goal-state at the center. In the order of appearance: a) $[5 \times 5]$ *dense*: has size 5×5 and dense rewards, b) $[5 \times 5]$ *sparse (hard)*: has size 5×5 and sparse rewards, c) $[5 \times 5]$ *sparse (easy)*: has size 5×5 , sparse rewards, and goal-state at the center, d) $[15 \times 15]$ *dense*: has size 15×15 and dense rewards, and e) $[15 \times 15]$ *sparse*: has size 15×15 and sparse rewards. The reward functions for both dense and sparse rewards are as previously described above for **2D-Gridworld**.

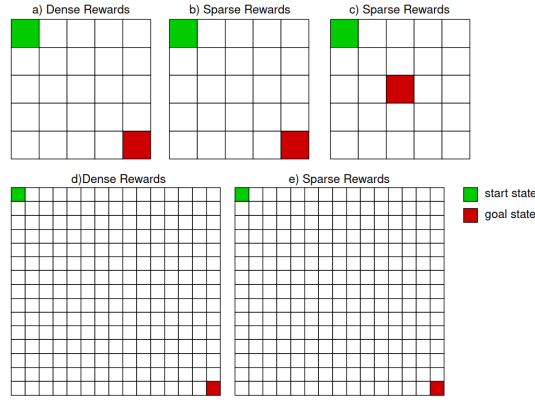


Figure 7: Five gridworld tasks with the same action space, but different rewards, state space and location of the goal state.

B.2 OMR(k): OMR over number of updates

OMR is defined for the entire policy trajectory by Equation 6 as,

$$\kappa \triangleq \frac{\sum_{k \in K^+} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}{\sum_{k=0}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}.$$

To observe how it changes with respect to updates, we compute OMR from update i onwards till the end of the learning trajectory, i.e. over subsequences with a decreasing number of policy updates with increasing i , using:

$$\kappa(i) \triangleq \frac{\sum_{k \in K^+, k \geq i} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}{\sum_{k=i}^{N-1} \mathcal{W}_1(v_{\pi_k}, v_{\pi_{k+1}})}, \text{ such that } i \in [0, N - T] \quad (66)$$

where $T \approx 0.9N$ to ensure that the last subsequence of policy updates have at least 10% of the total updates in the trajectory.

880 B.3 Computation of Occupancy Measures

881 The finite-horizon occupancy measure is defined as [4],

$$v_{\pi}^H(s, a) = \frac{1}{H} \sum_{t=1}^H \mathbb{P}(s_t = s, a_t = a \mid \pi, \mu)$$

882 for which the probability of the state-action pair selected is time-dependent. If we restrict our analysis
883 to stationary policies where $\pi(a_t | s_t) = \pi(a | s)$, then the probability of the state-action pair becomes
884 time-independent and thus,

$$v_{\pi}^H(s, a) = \mathbb{P}(s, a \mid \pi, \mu)$$

885 This implies that the use of stationary policies in finite-horizon MDPs, as observed in practice with
886 many episodic MDPs [54, 3, 52], induces stationary occupancy measures - where the expected
887 number of state-action pair visits are independent of the time-step. Work by [11] provides extensive
888 details about the existence of stationarity in episodic MDPs and shows (in Theorem 4) that,

$$\mathbb{E}_{(s,a) \sim v_{\pi}^H} [\bar{\mathcal{R}}(s, a)] = \frac{\mathbb{E}_{\zeta \sim M_{\pi}} \left[\sum_{t=1}^{H(\zeta)} R(s_t, a_t) \right]}{\mathbb{E}_{\zeta \sim M_{\pi}} [H(\zeta)]} \quad (67)$$

889 where ζ is the episodic state-action pair trajectory, $H(\zeta)$ is the episode length corresponding to ζ ,
890 and M_{π} is the markov chain induced by policy π . We verified the correctness of our v_{π}^H computation
891 by calculating the relative error derived from Equation 67 to check its validity. The relative error is
892 given as

$$\text{Rel. Error \%} = 100 * \frac{\mathbb{E}_{(s,a) \sim v_{\pi}^H} [\bar{\mathcal{R}}(s, a)] \mathbb{E}_{\zeta \sim M_{\pi}} [H(\zeta)] - \mathbb{E}_{\zeta \sim M_{\pi}} \left[\sum_{t=1}^{H(\zeta)} R(s_t, a_t) \right]}{\mathbb{E}_{\zeta \sim M_{\pi}} \left[\sum_{t=1}^{H(\zeta)} R(s_t, a_t) \right]} \quad (68)$$

893 Figure 8 depicts Rel. Error% vs number of updates in the stochastic 2D-Gridworld environment with
894 dense rewards. We observe that increasing the number of rollouts M reduces the estimation error
895 of v_{π}^H . For $M = 10$, the absolute relative error can be as high as 50% with the mean less than 10%.
While for $M = 500$, the maximum absolute relative error is 4%.

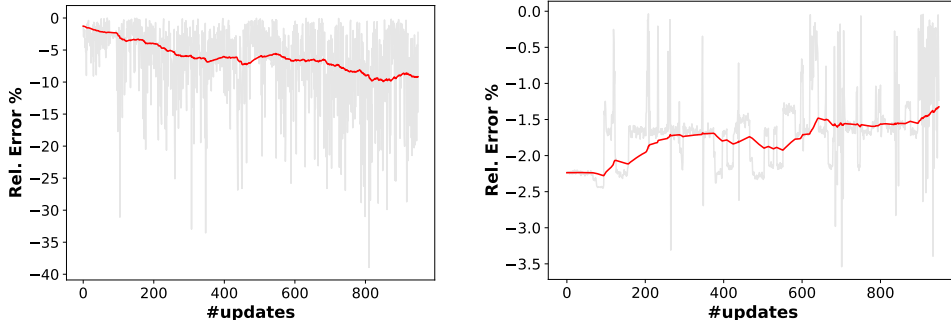


Figure 8: Rel. Error% vs number of updates plots in the 2D-Gridworld environment where v_{π}^H is estimated using $M = 10$ rollouts (left) and $M = 500$ rollouts (right).

897 B.4 Effects of the number of rollouts - SAC

898 The policy dataset D_{π_i} in a deterministic environment is made up of (s,a) pairs generated from a
899 single episode of the policy π_i . In a deterministic environment, this sequence remains the same across
900 repeats of episodes, for each policy π_i (deterministic) at update step i . Therefore, a single rollout is
901 sufficient to estimate the occupancy measure v_{π_i} . In a stochastic environment, rollouts are impacted
902 by the environment's stochasticity. Thus, multiple rollouts are needed to estimate the occupancy
903 measure accurately. As the number of rollouts increases, the occupancy measure should converge
904 and become less noisy.

Table 5 shows that, in a stochastic setting, the ESL values converge as the number of rollouts increases. OMR appears to be invariant across various the number of rollouts, and the mean number of updates appear to be consistent around 2900 (with exception for #rollouts = 1). The results indicate that from about 6 rollouts, the estimated occupancy measures become less noisy. This aligns with Eq. (61), which shows that increasing the number of rollouts reduces estimation error.

| #rollouts | ESL | OMR | UC |
|-----------|-------------|-------------|------------|
| 1 | 849.1±468.5 | 0.500±0.004 | 1849±742.2 |
| 3 | 618.6±257.3 | 0.501±0.005 | 2413±1397 |
| 6 | 445.4±245.8 | 0.501±0.042 | 2462±2043 |
| 9 | 428.1±234.4 | 0.503±0.004 | 2281±1743 |

Table 5: Evaluation of SAC algorithm in the **stochastic, dense-rewards setting** for 5x5 gridworld with **40 maximum steps per episode** across various number of rollouts. The effects of #rollouts on the Effort of Sequential Learning (ESL), Optimal Movement Ratio (OMR), and number of updates to convergence (UC) are observed.

B.5 η_{sub} can be a reasonable proxy for η , when optimal policy is not fully reached

We compare ESL when the optimal policy was reached, denoted η , versus when it was not, denoted η_{sub} , in Tables 6 and 7. First, we observe that the number of rollouts impacts the metric values. Second, η_{sub} values are always greater than η values. Note that UCRL2 and PSRL update their policies only at the end of each episode, whereas SAC and DQN update theirs after each time step. Hence, $UC_{sub} = 499$ for both UCRL2 and PSRL.

The ESL values (both η and η_{sub}) in Table 7 are lower than those in Table 6, as expected since more data samples reduce estimation error. The distance from the initial policies to the final policies are not so different. Using both Tables 6 and 7, we notice that comparing algorithms with η_{sub} yields the same efficiency ranking (e.g. PSRL, UCRL, SAC and DQN) as η . This indicates that η_{sub} reliably predicts results provided by η for comparing algorithms.

The results presented in Table 2 for stochastic dense-rewards setting are consistent with those in Table 7 because the number of rollouts used was $Nr = 6$.

| Algo. | η | η_{sub} | d | c | UC | UC_{sub} |
|-------|-----------|--------------|-----------|-----------|-----------|------------|
| SAC | 849±468 | 3623±4166 | 5.63±1.50 | 5.26±2.10 | 1850±742 | 7451±3535 |
| | 230±155 | 613±999 | 5.65±0.93 | 5.45±2.15 | 284±180 | 499±0.0 |
| UCRL2 | 86.2±44.4 | 389±102 | 4.96±1.26 | 5.29±1.49 | 97.2±52.5 | 499±0.0 |
| | 564±478 | 3911±1710 | 5.52±1.39 | 6.54±2.05 | 1213±1061 | 9097±1904 |
| PSRL | 86.2±44.4 | 389±102 | 4.96±1.26 | 5.29±1.49 | 97.2±52.5 | 499±0.0 |
| | 564±478 | 3911±1710 | 5.52±1.39 | 6.54±2.05 | 1213±1061 | 9097±1904 |
| DQN | 86.2±44.4 | 389±102 | 4.96±1.26 | 5.29±1.49 | 97.2±52.5 | 499±0.0 |
| | 564±478 | 3911±1710 | 5.52±1.39 | 6.54±2.05 | 1213±1061 | 9097±1904 |

Table 6: Evaluation of algorithms in the **stochastic, dense-rewards setting** for 5x5 gridworld with **40 maximum steps per episode** with the number of rollouts $Nr = 1$. The total number of training episodes is 500. When the algorithm converged at optimality, η is the *Effort of Sequential Learning*, $d = \mathcal{W}_1(\pi_0, \pi^*)$ is distance from initial policy to the optimal policy, and UC is the number of updates to convergence. When the algorithm did not converge at the optimal policy, rather a non-optimal π_N , we use η_{sub} , $c = \mathcal{W}_1(\pi_0, \pi_N)$, and UC_{sub} to denote the aforementioned quantities. 40 training trials were used.

B.6 Effects of Hyperparameters - UCRL2

Table 8 illustrates the effects of hyperparameter values in the UCRL2 algorithm. The environment is deterministic dense-rewards setting with 200 training episodes. We observe that high exploration rates ($\delta \rightarrow 0$) appear to align with high ESL and UC, while high exploitation rates ($\delta \rightarrow 1$) appear to align with low ESL and UC. OMR appears to be invariant across various δ values.

| Algo. | η | η_{sub} | d | c | UC | UC _{sub} |
|-------|------------|--------------|------------|------------|------------|-------------------|
| SAC | 445 \pm | 853 \pm | 5.63 \pm | 7.26 \pm | 2463 \pm | 6293 \pm |
| | 246 | 127 | 1.23 | 1.45 | 2043 | 441 |
| UCRL2 | 198 \pm | 510 \pm | 5.36 \pm | 4.58 \pm | 268 \pm | 499 \pm |
| | 121 | 274 | 0.84 | 1.90 | 155 | 0.0 |
| PSRL | 55.4 \pm | 361 \pm | 4.97 \pm | 3.91 \pm | 76.1 \pm | 499 \pm |
| | 33.6 | 43.6 | 1.34 | 0.48 | 50.6 | 0.0 |
| DQN | 458 \pm | 1971 \pm | 4.88 \pm | 6.52 \pm | 1586 \pm | 13713 \pm |
| | 311 | 250 | 1.06 | 0.31 | 1077 | 6907 |

Table 7: Evaluation of algorithms in the **stochastic, dense-rewards setting** for 5x5 gridworld with **40 maximum steps per episode** with the number of rollouts $N_r = 6$. The total number of training episode is 500. When the algorithm converged at optimality, η is the *Effort of Sequential Learning*, $d = \mathcal{W}_1(\pi_0, \pi^*)$ is distance from initial policy to the optimal policy, and UC is the number of updates to convergence. When the algorithm did not converge at the optimal policy however some π_N , we use η_{sub} , $c = \mathcal{W}_1(\pi_0, \pi_N)$, and UC_{sub} to denote the aforementioned quantities. 40 training trials were used.

| δ | ESL | OMR | UC | SR% |
|----------|-------------------|-------------------|-------------------|-----|
| 0.1 | 47.76 \pm 7.768 | 0.512 \pm 0.033 | 62.26 \pm 9.977 | 100 |
| 0.3 | 39.29 \pm 5.860 | 0.515 \pm 0.034 | 58.08 \pm 7.746 | 100 |
| 0.5 | 38.26 \pm 6.747 | 0.511 \pm 0.036 | 56.92 \pm 9.111 | 100 |
| 0.7 | 37.48 \pm 5.094 | 0.507 \pm 0.029 | 56.68 \pm 7.460 | 100 |
| 0.9 | 36.40 \pm 5.301 | 0.510 \pm 0.036 | 54.86 \pm 7.326 | 100 |

Table 8: Evaluation of UCRL2 algorithm in the **deterministic, dense-rewards setting** for 5x5 gridworld with **15 maximum steps per episode**. Different confidence parameter $\delta \in (0, 1)$ were evaluated to see their effects on Effort of Sequential Learning (ESL), Optimal Movement Ratio (OMR), number of updates to convergence (UC), and success rate (SR). Note that as $\delta \rightarrow 0$, the agent approaches absolute exploration, and with $\delta \rightarrow 1$ absolute exploitation.

928 B.7 Extended Discussion of Usefulness of ESL and OMR

929 The quantities like regret and number of updates (UC) are outcomes of the exploratory processes,
930 and thus reflect only a partial view of the underlying exploration mechanisms. We propose ESL and
931 OMR to complement regret and number of updates as metrics but not to replace them.

932 1. Complementarity of ESL and OMR with respect to UC:

933 a. **Case 1.** Let us consider two RL algorithms that reach optimality with the same number of updates,
934 i.e. they have the same UC. *How would one be able to distinguish the exploratory processes of*
935 *these algorithms?* ESL and OMR are the summary metrics of the policy trajectory during learning.
936 These can reveal which algorithm’s exploratory process is more direct versus meandering, smooth
937 versus noisy, or has large versus small coverage area in the policy space (Figures 3 and 4, top rows).
938 Therefore, ESL and OMR quantify with granularity the characteristics of the exploratory process of
939 an RL algorithm for any given environment.

940 b. **Case 2.** Let us consider the case when optimality is not reached but the maximum number of
941 updates is attained by two RL algorithms. *How would one be able to evaluate the exploratory*
942 *processes of these algorithms and systematically uncover which exploratory process demonstrates*
943 *desired characteristics?* Looking into the training trajectories of RL algorithms in an environment
944 and corresponding higher/lower ESLs (η_{sub} , Section 4.2), we can make a knowledgeable choice of
945 an RL algorithm exhibiting desired characteristics (e.g. high coverage, smooth exploration). We have
946 shown in Section 4.2 and results in Section 5.3 (also Appendix B.5) that ranking based on suboptimal
947 ESL is aligned with true ESL, and additionally, the visualization of the training trajectories (Figures 3
948 and 4) can indicate the characteristics of corresponding RL algorithms even when optimal policy is
949 not reached.

950 c. **Experimental Evidence.** UCRL2 is known to be provably regret-optimal and is designed to
951 continuously explore. SAC does not have such rigorous theoretical guarantees but is known to be
952 practically efficient. In Table 1, by UC, we observe that SAC is significantly suboptimal than UCRL2.
953 But SAC has lower ESL than UCRL2 as its exploration is smoother. Additionally, OMR for SAC is

954 higher than that of UCRL2. They together indicate that SAC takes smoother but larger number of
955 policy transitions aligned to optimal direction for exploration, while UCRL2 exhibits bigger policy
956 changes and in diverse manner trying to cover the environment faster.

957 **2. Complementarity of ESL and OMR with respect to Regret:**

958 UCRL2 and PSRL have the same order of regret bound [61]. But PSRL leads to smoother policy
959 transitions that are much more orientated towards optimality (as shown in Figure 3), while UCRL2
960 leads to less smooth policy transitions that do not taper as it approaches optimality. This information
961 is not evident from regret but from corresponding ESLs and OMRs (Table 1).

962 **3. Insights for Algorithm Design:**

963 Knowing ESL (or suboptimal ESL) and OMR can assist with developing algorithms that emphasize
964 certain exploratory characteristics. We can develop algorithms with grades of coverage or directness,
965 while also being able to visualize this. Ultimately, depending on the environment, we can choose
966 which characteristics of exploratory process are well suited. In contrast, looking only at the final
967 outcomes of RL algorithms like regret and number of updates does not include these nuances.

C Specifications of the RL Algorithms under Study

C.1 Methods for simulation results (Discrete MDP)

Model parameter initialisation. We initialised model parameters for deep learning RL algorithms like DQN and SAC by uniformly sampling weight values between $-3 \cdot 10^{-4}$ and $3 \cdot 10^{-4}$ and the biases at 0. For tabular Q-learning algorithms, we randomly initialized the Q-values between -1.0 and 1.0 . For UCRL and PSRL, the policy model was randomly initialized. Note that all Wasserstein distances were computed using a python package POT [22]. Additionally, L1 norm was used in our Wasserstein metric cost function as the ground metric for the 2D gridworld environment.

Results in Figure 3. The problem setting was deterministic with dense-rewards and 15 maximum number of steps per episode. The total number of episodes was 200. The convergence criterion was satisfied when maximum returns were produced by an algorithm over 5 consecutive updates. The results showcase a single representative run of each algorithm. The confidence parameter $\delta = 0.1$ was utilized for UCRL2. The α parameter for SAC was autotuned using the approach in [29] along with hyperparameters described in Table 9. While DQN began with $\epsilon = 1.0$ and the value decayed as $\epsilon[t + 1] = \max\{0.9999 \times \epsilon[t], 0.0001\}$. Table 10 shows hyperparameters for DQN. Note that the ADAM [40] optimizer was used in all the neural network models.

Table 9: SAC Hyperparameters.

| Parameter | Value |
|---|-------------------|
| learning rate | $5 \cdot 10^{-4}$ |
| discount(γ) | 0.99 |
| replay buffer size | 10^4 |
| number of hidden layers (all networks) | 1 |
| number of hidden units per layer | 32 |
| number of samples per minibatch | 64 |
| nonlinearity | ReLU |
| entropy target | -4 |
| target smoothing coefficient (τ) | 0.01 |
| target update interval | 1 |
| gradient steps | 1 |
| initial exploration steps | |
| before model starts updating | 500 |

Table 10: DQN Hyperparameters.

| Parameter | Value |
|---|-------------------|
| learning rate | $5 \cdot 10^{-2}$ |
| discount(γ) | 0.99 |
| replay buffer size | 10^4 |
| number of hidden layers (all networks) | 1 |
| number of hidden units per layer | 32 |
| number of samples per minibatch | 64 |
| nonlinearity | ReLU |
| target smoothing coefficient (τ) | 0.001 |
| target update interval | 1 |
| gradient steps | 1 |
| initial exploration steps | |
| before ϵ decays | 500 |

Results in Tables 1 and 2. The problem settings had 40 maximum number of steps per episode, and the convergence criterion was satisfied when maximum returns were produced by an algorithm over 5 consecutive updates. The means and standard deviations for each algorithm were computed over 50 runs. The total number of episodes was 200 for results in Table 1, and 500 in Table 2. For results in Figure 5, the Q-learning with decaying ϵ -greedy where $\epsilon = 0.9$ was employed in the gridworld tasks described in Appendix B.1. A convergence criterion of 50 consecutive model updates with maximum

990 returns was utilized. We aggregated the result over 40 training trials and the maximum number of
 991 steps per episode was 60.

992 C.2 Methods for simulation results (Continuous MDP)

993 **Model parameter initialisation.** We initialised model parameters for the deep learning SAC
 994 algorithm by uniformly sampling weight values between $-3 \cdot 10^{-4}$ and $3 \cdot 10^{-4}$ and the biases at 0. For
 995 the DDPG algorithm, the output layer weight values were initialised using Xavier Initialization [27],
 996 while the rest were uniformly sampled between $-3 \cdot 10^{-3}$ and $3 \cdot 10^{-3}$. This was done on both the
 997 actor and critic networks. The ADAM [40] optimizer was used in all the neural network models.
 998 In both algorithms, 1) a discount factor $\gamma = 0.99$ was used, 2) 500 initial steps were taken before
 999 updating model weights, and 3) replay buffer size was 10^6 . Tables 11 and 12 display hyperparameters
 1000 for DDPG and SAC, respectively.

1001 **Results in Figure 4.** The problem setting was Mountain Car continuous [58] with 999 maximum
 1002 number of steps per episode [13]. The total number of training episodes was 100. The convergence
 1003 criterion was satisfied when maximum returns were produced by an algorithm over 10 consecutive
 1004 updates. The results showcase a single representative run of each algorithm. For results in **Table**
 1005 **3**, the mean and standard deviations for each algorithm were computed over 5 runs. While RL
 1006 training was conducted in a continuous state-action space, we discretized it for Wasserstein distance
 1007 calculations between occupancy measures, using 4 bins for actions and 10 bins for states. Note that all
 1008 Wasserstein distances were computed using a python package POT [22]. Additionally, L2 norm was
 1009 used in our Wasserstein metric cost function as the ground metric for the Mountain Car environment.

Table 11: DDPG Hyperparameters.

| Parameter | Value |
|--|-------------------|
| number of samples per minibatch | 128 |
| nonlinearity | ReLU |
| target smoothing coefficients (τ) | 0.001 |
| target update interval | 1 |
| gradient steps | 1 |
| number of hidden layers (all networks) | 2 |
| number of hidden units per layer | 64 |
| Actor learning rate | $5 \cdot 10^{-4}$ |
| Critic learning rate | $5 \cdot 10^{-3}$ |

Table 12: SAC Hyperparameters.

| Parameter | Value |
|---|-------------------|
| learning rate | $3 \cdot 10^{-3}$ |
| number of hidden layers (all networks) | 2 |
| number of hidden units per layer | 64 |
| number of samples per minibatch | 128 |
| nonlinearity | ReLU |
| target smoothing coefficient (τ) | 0.001 |
| target update interval | 1 |
| gradient steps | 1 |

1010 D Supplementary Results

1011 In this section we present enlarged versions of results in Figure 3 (see Section D.1) and additional
 1012 plots that support the results in the main paper (see Section D.3).

1013 D.1 Enlarged Visualisation of the Occupancy Measure Trajectories

1014 Figures 9 - 11 are enlarged versions of enlarged versions of Figure 3. For each algorithm, there is a
 1015 visualisation of the policy trajectory and visualisation of the state visitation below it.

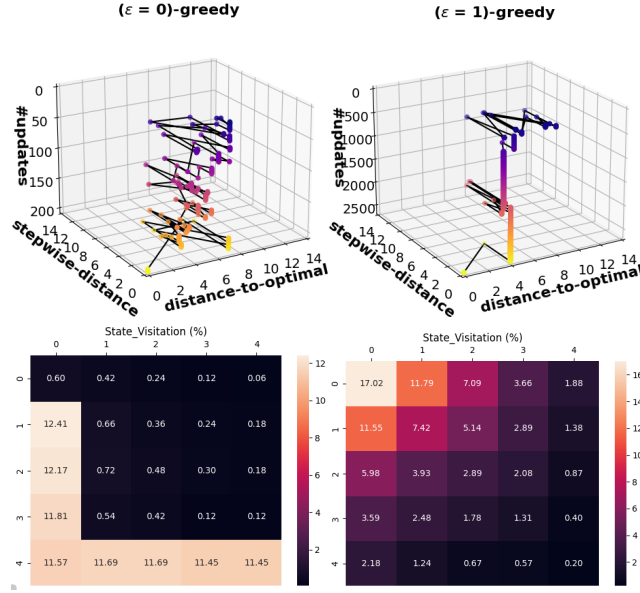


Figure 9: Top row: Scatter plots of *distance-to-optimal* and *stepwise-distance* over updates for $\epsilon(=0)$ -greedy and $\epsilon(=1)$ -greedy Q-learning. Bottom row: State visitations.

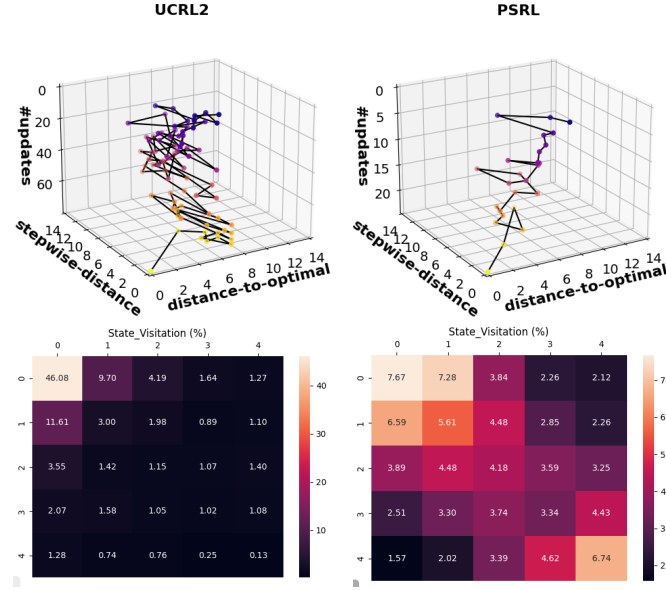


Figure 10: Top row: Scatter plots of *distance-to-optimal* and *stepwise-distance* over updates for UCRL2 and PSRL. Bottom row: State visitations.

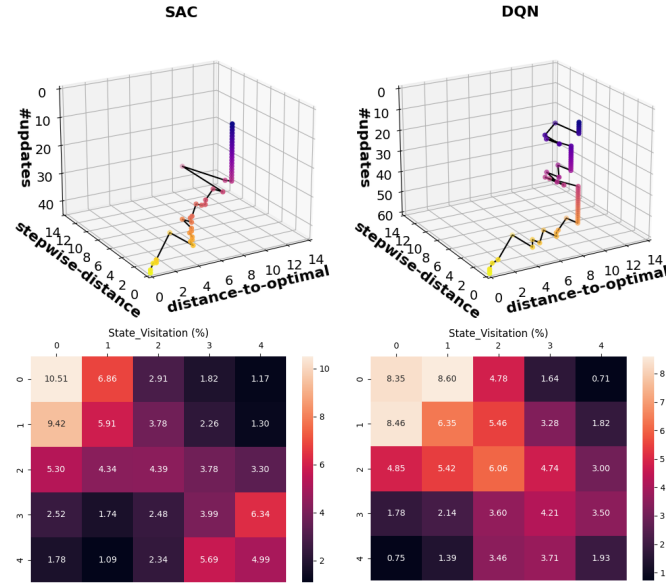


Figure 11: Top row: Scatter plots of *distance-to-optimal* and *stepwise-distance* over updates for SAC and DQN. Bottom row: State visitations.

1016 D.2 Performance Plots

1017 This section contains Return plots of the algorithms. This allows us to assess the learning of algorithms
 1018 from the performance perspective. Figures 12 and 13 depict performance evolution that corresponds
 1019 to settings in Figures 3 and 4, respectively. Note that while all algorithms find the optimal policy,
 1020 UCRL2 and $\epsilon(=1)$ -greedy Q-learning fail to remain there if training continues without truncation. As
 1021 a result, their performance does not improve over time compared to those that stabilize at the optimal
 1022 policy.

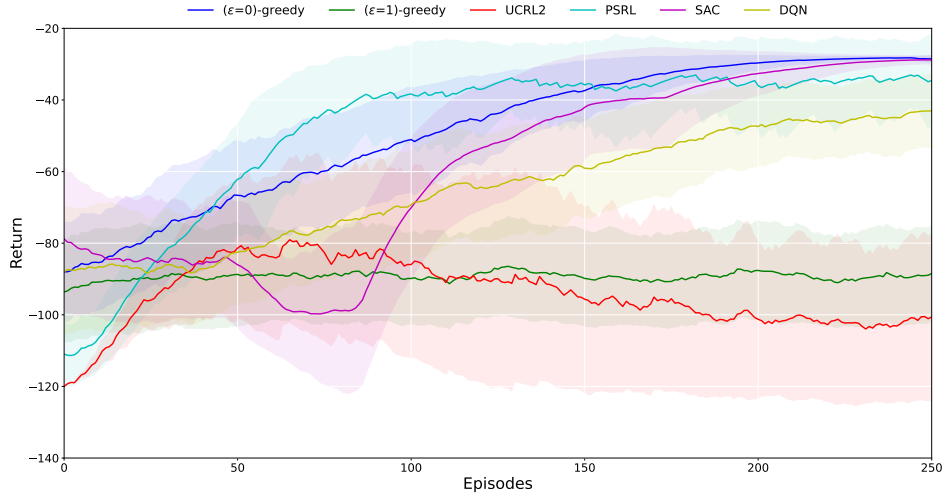


Figure 12: Return plots of algorithms: $\epsilon(=0)$ -greedy and $\epsilon(=1)$ -greedy Q-learning, UCRL2, PSRL, SAC, and DQN averaged over 5 runs in the deterministic 5×5 Gridworld with dense rewards.

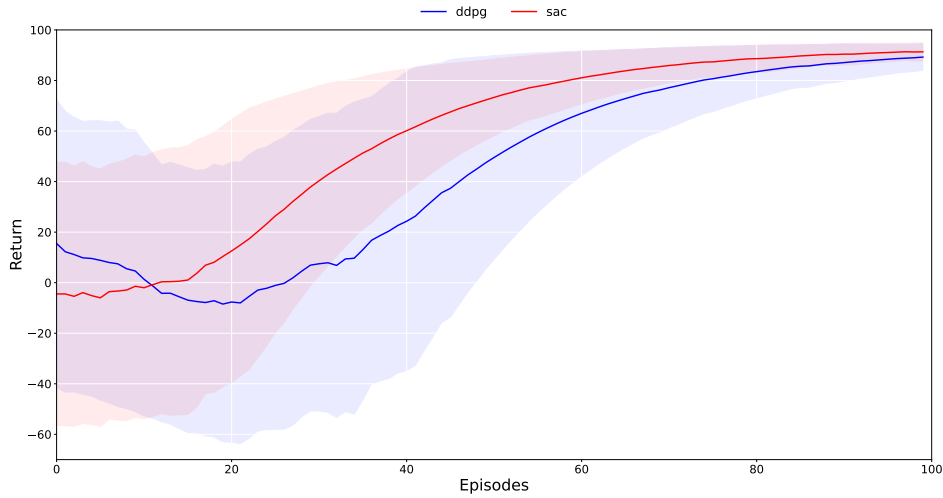


Figure 13: Return plots of algorithms: DDPG and SAC averaged over 5 runs in the continuous Mountain Car problem.

1023 D.3 Evolution of *stepwise-distance*, *distance-to-optimal*, and $\text{OMR}(k)$

1024 In this section we present 2 dimensional versions of the policy trajectories in Figures 3 and 4,
 1025 along with corresponding OMR evolution plots. These are *stepwise-distance* vs. updates, *distance-*
 1026 *to-optimal* vs. updates, and $\text{OMR}(k)$ plots for the algorithms. Figure 14 presents plots for the
 1027 continuous environment Mountain Car, while Figure 15) presents plots for the discrete environment
 1028 2D Gridworld.

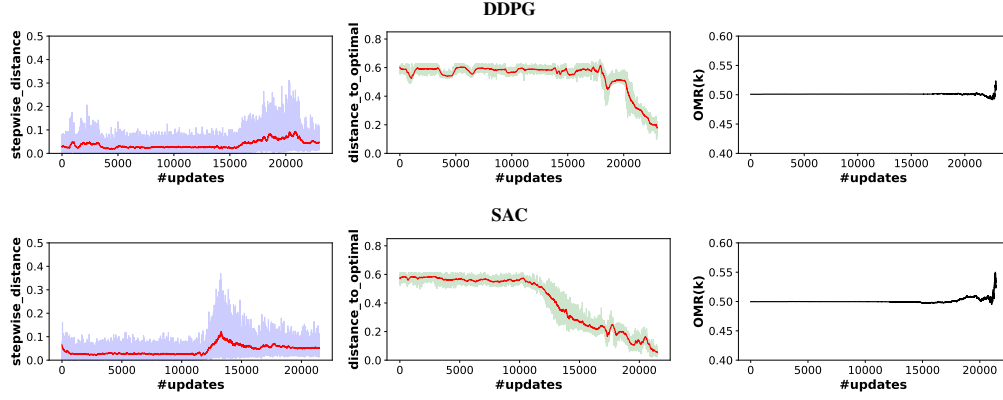


Figure 14: Plots in the first column are *stepwise-distance* vs. number of updates, second column *distance-to-optimal* vs. number of updates, and third $\text{OMR}(k)$ vs. number of updates. Top row plots belong to DDPG algorithm, while bottom row plots belong to SAC.

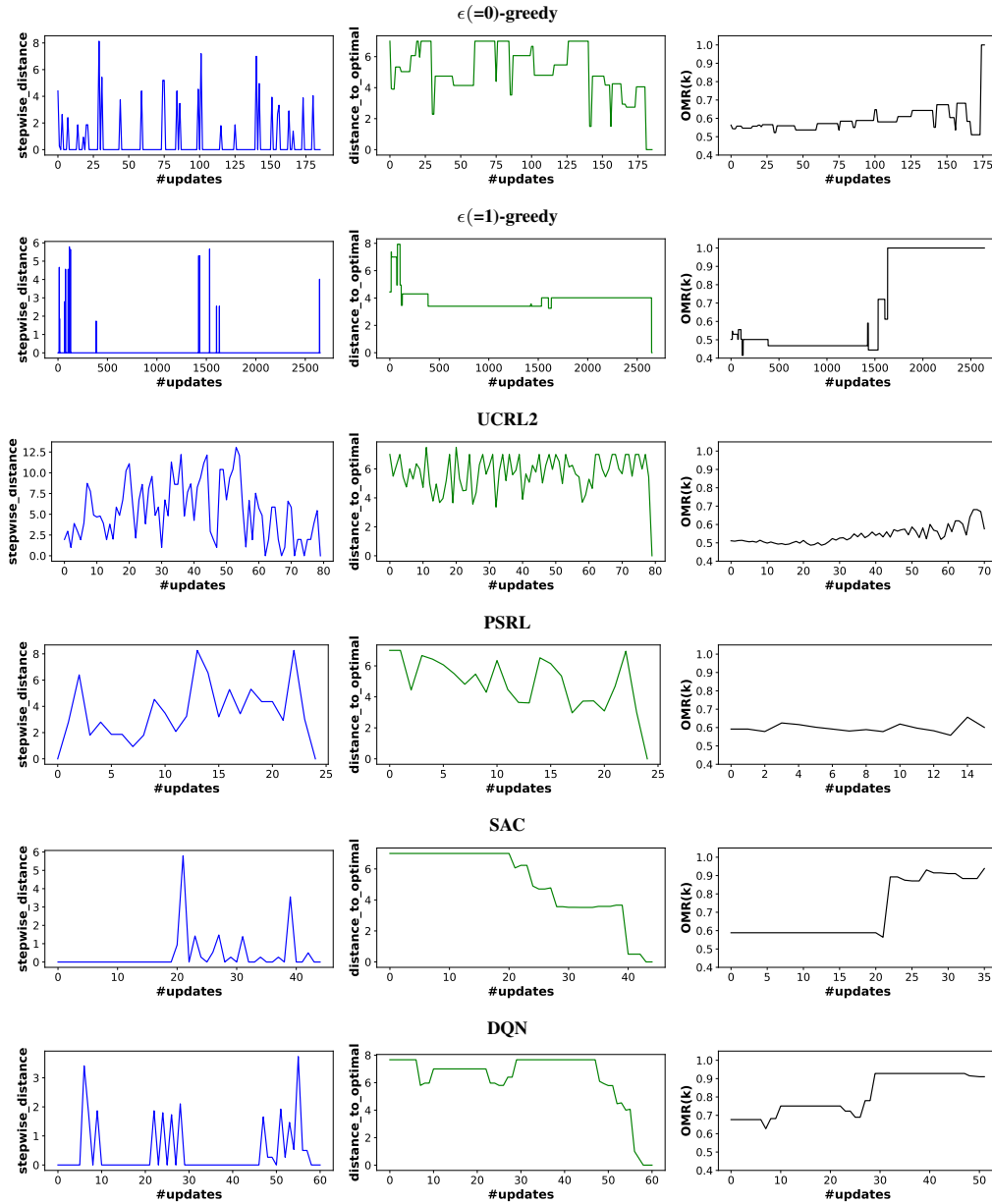


Figure 15: Plots in the first column are *stepwise-distance* vs. number of updates, second column *distance-to-optimal* vs. number of updates, and third $OMR(k)$ vs. number of updates. The plots in the row belong to algorithms in the following order from top to bottom: $\epsilon(=0)$ -greedy, $\epsilon(=1)$ -greedy, UCRL2, PSRL, SAC, and DQN.

1029 E Technical Appendices and Supplementary Material

1030 Technical appendices with additional results, figures, graphs and proofs may be submitted with
 1031 the paper submission before the full submission deadline (see above), or as a separate PDF in the
 1032 ZIP file before the supplementary material deadline. There is no page limit for the technical
 1033 appendices.