

# Unsupervised Data Generation for Offline Reinforcement Learning: A Perspective from Model

Shuncheng He<sup>1</sup>, Hongchang Zhang<sup>1</sup>, Jianzhun Shao<sup>1</sup>, Yuhang Jiang<sup>1</sup>, Xiangyang Ji<sup>1</sup>

<sup>1</sup>Tsinghua University  
hesc16@mails.tsinghua.edu.cn

## Abstract

Offline reinforcement learning (RL) recently gains growing interests from RL researchers. However, the performance of offline RL suffers from the out-of-distribution problem, which can be corrected by feedback in online RL. Previous offline RL research focuses on restricting the offline algorithm in in-distribution even in-sample action sampling. In contrast, fewer work pays attention to the influence of the batch data. In this paper, we first build a bridge over the batch data and the performance of offline RL algorithms theoretically, from the perspective of model-based offline RL optimization. We draw a conclusion that, with mild assumptions, the distance between the state-action pair distribution generated by the behavioural policy and the distribution generated by the optimal policy, accounts for the performance gap between the policy learned by model-based offline RL and the optimal policy. Secondly, we reveal that in task-agnostic settings, a series of policies trained by unsupervised RL can minimize the worst-case regret in the performance gap. Inspired by the theoretical conclusions, UDG (Unsupervised Data Generation) is devised to generate data and select proper data for offline training under task-agnostic settings. Empirical results demonstrate that UDG can outperform supervised data generation on solving unknown tasks.

## 1 Introduction

Reinforcement learning (RL) recently gains significant advances in sequential decision making problems, with applications ranging from the game of Go (Silver et al. 2016, 2017), video games (Mnih et al. 2015; Hessel et al. 2018), to autonomous driving (Kiran et al. 2021) and robotic control (Zhao, Queralta, and Westerlund 2020). However, the costly online trial-and-error process requires numerous samples of interactions with the environment which restricts RL from real world deployment. In the scenarios where online interaction is expensive or unsafe, we have to resort to offline experience (Levine et al. 2020). However, transplanting RL to offline can provoke disastrous error by falsely overestimating the out-of-distribution samples without correction from environment feedback. Despite recent advances on mitigating bootstrapped error by constraining the policy in data distribution or even in data samples (Fujimoto and Gu 2021), offline RL is still limited since they can barely generalize to out-of-distribution areas (Yu et al. 2020c). The inability of generalization of offline RL will be a serious issue when the

batch data deviates from the optimal policy especially under the settings of multi-task, task transfer or task-agnostic. As plenty of research on online RL succeeds (Sodhani, Zhang, and Pineau 2021; Yu et al. 2020b,a; Laskin et al. 2021; Eysenbach et al. 2019; Sharma et al. 2020), we hope offline RL can cope with task-agnostic problems either. To this end, how the batch data distributes becomes the primal concern.

Recent research empirically shows that diversity in offline data improves performance on task transfer and solving multiple tasks (Lambert et al. 2022; Yarats et al. 2022). The diverse dataset is obtained from unsupervised RL by competitively training diverse policies (Eysenbach et al. 2019), or exploration emphasized pre-training (Liu and Abbeel 2021a), and all of the generated data is fed to offline algorithms. However, these studies barely address the connection between the batch data and the performance of offline RL theoretically. How the diversity of data contributes to solving task-agnostic problems remains unclear.

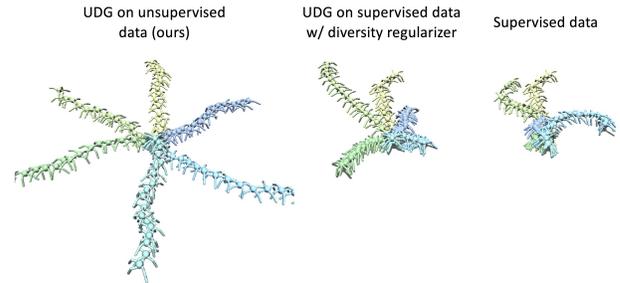


Figure 1: Rendered trajectories of offline trained policies in 6 Ant-Angle tasks. These tasks require the ant to move along 6 different directions. With diverse data buffers generated by unsupervisedly trained policies, our method UDG can solve all the tasks by offline reinforcement learning.

Our study addresses the connection between batch data and performance from a perspective of model based offline optimization MOPO (Yu et al. 2020c). MOPO establishes a lower bound of the expected return of offline trained policy and reveals that the model prediction error on the optimal data distribution mainly contributes to the performance gap. In this paper, we examine the model prediction error and find the connection between the performance gap and the

Wasserstein distance of the batch data distribution from the optimal distribution. We conclude that the offline trained policy will have higher return close to the optimal policy if the behavioural distribution is closer to the optimal distribution. We discover that, in task-agnostic scenarios, unsupervised RL methods which propel the policies far away from each other, approximately optimize the minimal regret to the optimal policy. Based on these theoretical analysis, we propose a framework named unsupervised data generation (UDG) as illustrated in Figure 2. In UDG, a series of policies are trained with diversity rewards. They are used to generate batch data stored in different buffers. Before the offline training stage, the buffers are relabeled with given reward function corresponding to the task, and the buffer with highest return is sent to the offline algorithm.

The contributions in this work are three-fold. First, to our best knowledge, we are the first to establish a theoretical bond between the behavioural batch data and the performance of offline RL algorithms on Lipschitz continuous environments. Second, we present an objective of minimal worst-case regret for data generation on task-agnostic problems. Third, we propose a new framework UDG for unsupervised offline RL and evaluate UDG on locomotive environments. Empirical results on locomotive tasks like Ant-Angle and Cheetah-Jump show that UDG outperforms conventional offline RL with random or supervised data. Further experiments validate the soundness of our theoretical findings.

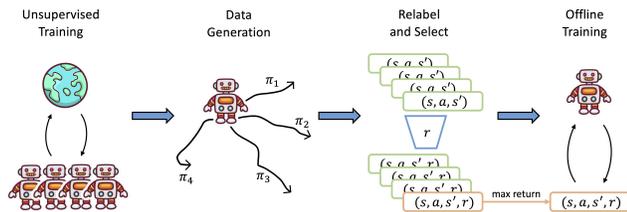


Figure 2: The framework of UDG. First a series of  $K$  policies are trained simultaneously with diversity rewards. Second, collect rollout experience  $(s, a, s')$  from each policy and construct a corresponding data buffer. Third, relabel the reward in the batch data with a designated reward function, and select the data buffer with the maximal average return. Finally train the agent on the chosen data by offline RL approaches.

## 2 Related work

**Offline RL.** Online reinforcement requires enormous samples and makes RL less feasible in many real-world tasks (Levine et al. 2020). Therefore, learning from batch data is a key route to overcome high sample complexity. However, vanilla online RL algorithms face the out-of-distribution problem that the Q function may output falsely high values on samples not in data distribution. To mitigate this issue, model-free offline RL methods cope with the out-of-distribution problem from two aspects, constraining the learned policy on the support of batch data (Fujimoto, Meger, and Precup 2019; Kumar et al. 2019; Wu, Tucker, and Nachum 2019;

Peng et al. 2019; Siegel et al. 2020; Cheng et al. 2022; Rezaeifar et al. 2022; Zhang et al. 2022; Fujimoto and Gu 2021), and suppressing Q values on out-of-distribution area (Agarwal, Schuurmans, and Norouzi 2019; Kumar et al. 2020). To generalize beyond the batch data, model-based offline RL methods employ transition models to produce extra samples for offline learning (Chen et al. 2021; Kidambi et al. 2020; Yu et al. 2020c; Matsushima et al. 2020). These methods can naturally generalize to areas where the model is accurate (Janner et al. 2019). Our theoretical work originates from MOPO (Yu et al. 2020c), a model-based offline algorithm by adding an uncertainty penalty to avoid unexpected exploitation when model is inaccurate. MOPO derives a performance lower bound w.r.t. model prediction error. In contrast, we investigate the lower bound w.r.t. data to show the connection between data distribution and performance. On previous theoretical analysis on offline RL, the data coverage condition is crucial for stable, convergent offline learning (Wang, Foster, and Kakade 2020; Chen and Jiang 2019). In this paper, we focus on Lipschitz continuous environments which are common in locomotive tasks such as HalfCheetah, Ant in MuJoCo (Todorov, Erez, and MuJoCo 2012). By investigating into the transition with Lipschitz geometry, our analysis makes no assumptions on data coverage.

**Unsupervised RL.** Reinforcement learning heavily relies on the reward feedback of the environment for a specific tasks. However, recent research on unsupervised RL demonstrate training policies without extrinsic rewards enables the agent to adapt to general tasks (Laskin et al. 2021; Eysenbach et al. 2019). Without supervision from extrinsic rewards, unsupervised RL methods can either be driven by curiosity/novelty (Pathak et al. 2017; Pathak, Gandhi, and Gupta 2019; Burda et al. 2018), maximum coverage of the state space (Liu and Abbeel 2021b; Campos et al. 2020; Yarats et al. 2021), and diversity of a series of policies (Florensa, Duan, and Abbeel 2017; Lee et al. 2019; Eysenbach et al. 2019; Sharma et al. 2020; He et al. 2022; Liu and Abbeel 2021a; Strouse et al. 2021; Kim, Park, and Kim 2021). These methods all provide a pseudo reward derived from their own criteria. Our work employs a diverse series of policies to generate batch data for offline learning under task-agnostic settings. Therefore we utilize an unsupervised training paradigm in alignment with DIAYN (Eysenbach et al. 2019), DADS (Sharma et al. 2020), WURL (He et al. 2022), and choose WURL as base algorithm in accordance with our theoretical analysis.

**Offline dataset.** D4RL (Fu et al. 2020) and RL Unplugged (Gulcehre et al. 2020) are most commonly used offline RL benchmarks. The datasets in these benchmarks consist of replay buffers during training, rollout samples generated by a policy of a specific level, or samples mixed from different policies. Apart from benchmark datasets, exploratory data gains growing interest (Wang et al. 2022). Explore2Offline (Lambert et al. 2022) and ExORL (Yarats et al. 2022) both investigate into the role of batch data and construct a more diverse dataset with unsupervised RL algorithms for task generalization. In addition, extensive experiments in ExORL empirically show exploratory data can improve the performance of offline RL algorithms and even TD3 (Fujimoto, Hoof, and Meger 2018). However, neither of these methods

have theoretical analysis on the connection between data diversity and offline performance. Our work completes the picture of how diverse data improves offline RL performance. And we point out that data selection before offline training has considerable influence on performance, which is not addressed in previous work.

### 3 Preliminaries

A Markov decision process (MDP) is formalized as  $M = (\mathcal{S}, \mathcal{A}, T, r, p_0, \gamma)$ , where  $\mathcal{S}$  denotes the state space while  $\mathcal{A}$  denotes the action space,  $T(s'|s, a)$  the transition dynamics,  $r(s, a)$  the reward function,  $p_0$  the initial state distribution and  $\gamma \in [0, 1)$  the discounted factor. RL aims to solve the MDP by finding a policy  $\pi(a|s)$  maximizing the expected accumulated discounted return  $\eta_M(\pi) := \mathbb{E}_{\pi, T, p_0} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ . The state value function  $V_M^\pi(s_t) := \mathbb{E}_{\pi, T} [\sum_{k=0}^{\infty} \gamma^k r(s_{t+k}, a_{t+k})]$  provides an expectation of discounted future return from  $s_t$  under policy  $\pi$  and MDP  $M$ . Let  $\mathbb{P}_{T,t}^\pi(s)$  be the probability of being in state  $s$  at step  $t$  when acting with policy  $\pi$ . The discounted occupancy measure of  $\pi$  under dynamics  $T$  is denoted by  $\rho_T^\pi(s, a) := \frac{1}{c} \pi(a|s) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{T,t}^\pi(s)$  where  $c = 1/(1 - \gamma)$  is the normalization constant. Likewise,  $\rho_T^\pi(s) := \frac{1}{c} \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_{T,t}^\pi(s)$  is the state occupancy distribution. The expected accumulated discounted return can be rewritten as  $\eta_M^\pi = c \mathbb{E}_{\rho_T^\pi} [r(s, a)]$ .  $\hat{\rho}_T^\pi = \frac{1}{K} \sum_{i=1}^K \delta(s_i, a_i)$  denotes the empirical distribution of  $\rho_T^\pi$  based on  $K$  samples sampled from  $\pi$  under transition dynamics  $T$ .  $\delta(\cdot)$  denotes Dirac distribution.

Model-based RL approaches learn an estimated model  $\hat{T}$  from interaction experience, which defines a model MDP  $\hat{M} = (\mathcal{S}, \mathcal{A}, \hat{T}, r, p_0, \gamma)$ . Similarly we have the expected return under the learned dynamics  $\eta_{\hat{M}}(\pi) = c \mathbb{E}_{\hat{\rho}_T^\pi} [r(s, a)]$ .

In offline settings, the RL algorithm optimizes the policy solely on a fixed dataset  $\mathcal{D}_\beta = (s, a, r, s')$  generated by the behavioural policy  $\pi^\beta$ .  $\pi^\beta$  can be one policy or a mixture of policies. Note that offline RL algorithms cannot interact with the environment or produce extra samples. In model-based offline RL, the algorithm first learn a transition model  $\hat{T}$  from the batch data  $\mathcal{D}_\beta$ . At the training stage, the algorithm executes  $k$ -step rollout using the estimated model from the state sample from  $\mathcal{D}_\beta$ . The generated data are added to another buffer  $\mathcal{D}_m$ . Both data buffers are used in offline policy optimization.

### 4 UDG: Unsupervised Data Generation

In order to find the connection between performance and data, we first review key propositions in MOPO (Yu et al. 2020c). As discussed before, offline RL faces a dilemma of out-of-distribution samples and lack of exploration. Model-based RL like MBPO (Janner et al. 2019) can naturally extend to the regions where the model predicts as well as the true dynamics. However, when the model is inaccurate, the algorithm may exploit the falsely high return regions, resulting in inferior test performance in true dynamics. MOPO first derives the performance lower bound represented by the model error, and then addresses the risk-return trade-off by incorpo-

rating the penalty represented by the error of the estimated dynamics into the reward of offline policy optimization.

We briefly summarize the derivation of the performance lower bound. First we introduce the telescoping lemma:

**Lemma 4.1.** *Let  $M$  and  $\hat{M}$  be two MDPs with the same reward function  $r$ , but different dynamics  $T$  and  $\hat{T}$  respectively. Denote  $G_M^\pi(s, a) := \mathbb{E}_{s' \sim \hat{T}(s, a)} [V_M^\pi(s')] - \mathbb{E}_{s' \sim T(s, a)} [V_M^\pi(s')]$ . Then*

$$\eta_{\hat{M}}(\pi) - \eta_M(\pi) = c \gamma \mathbb{E}_{(s, a) \sim \rho_T^\pi} [G_M^\pi(s, a)]. \quad (1)$$

If we have mild constraints on the value function  $V_M^\pi \in \mathcal{F}$  where  $\mathcal{F}$  is a bounded function class under a specific metric, then we can bound the gap  $G_M^\pi(s, a)$  with model error measured by corresponding integral probability measure (IPM)  $d_{\mathcal{F}}$  (Müller 1997),

$$\begin{aligned} |G_M^\pi(s, a)| &\leq \sup_{f \in \mathcal{F}} \left| \mathbb{E}_{s' \sim \hat{T}(s, a)} [f(s')] - \mathbb{E}_{s' \sim T(s, a)} [f(s')] \right| \\ &= d_{\mathcal{F}}(\hat{T}(s, a), T(s, a)). \end{aligned} \quad (2)$$

Since we cannot access the true model  $T$  in most cases, MOPO adopts an admissible error estimator  $u : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  for  $\hat{T}$ , and have an assumption that for all  $s \in \mathcal{S}, a \in \mathcal{A}$ ,  $d_{\mathcal{F}}(\hat{T}(s, a), T(s, a)) \leq u(s, a)$ . An uncertainty-penalized MDP  $\tilde{M} = (\mathcal{S}, \mathcal{A}, \hat{T}, \tilde{r}, p_0, \gamma)$  is defined given the error estimator  $u$ , with the reward  $\tilde{r}(s, a) := r(s, a) - \gamma u(s, a)$  penalized by model error.

By optimizing the policy in the uncertainty-penalized MDP  $\tilde{M}$ , MOPO has a following performance lower bound,

**Theorem 4.2 (MOPO).** *Given  $\hat{\pi} = \arg \max_{\pi} \eta_{\tilde{M}}(\pi)$  and  $\epsilon_u(\pi) := c \mathbb{E}_{(s, a) \sim \rho_T^\pi} [u(s, a)]$ , the expected discounted return of  $\hat{\pi}$  satisfies*

$$\eta_M(\hat{\pi}) \geq \sup_{\pi} \{ \eta_M(\pi) - 2\gamma \epsilon_u(\pi) \}. \quad (3)$$

The theorem 4.2 reveals the optimality gap between  $\pi^*$  and  $\hat{\pi}$ . Immediately we have  $\eta_M(\hat{\pi}) \geq \eta_M(\pi^*) - 2\gamma \epsilon_u(\pi^*)$ . This corollary indicates if the model error is small on the  $(s, a)$  occupancy distribution under the optimal policy  $\pi^*$  and dynamics  $\hat{T}$ , the optimality gap will be small. In order to find the deep connection between the batch data and the performance gap of model-based offline RL algorithms, in the following section, we directly analyze the model prediction deviation  $d_{\mathcal{F}}(\hat{T}(s, a), T(s, a))$  instead of the error estimator  $u(s, a)$ .

### The connection between batch data and offline RL performance

Before presenting a lower bound of  $\eta_M(\hat{\pi})$ , here we make a few assumptions to simplify the proof. Some of these assumptions can be loosen and do not change the conclusion of the main result. The generalization of the theoretical analysis is discussed in Appendix A.

**Definition 4.3.** *Given a bounded subset  $\mathcal{K}$  in the corresponding  $d$ -dimension Euclidean space  $\mathbb{R}^d$ . The diameter  $B_{\mathcal{K}}$  of set  $\mathcal{K}$  is defined as the minimum value of  $B$  such that there exists  $k_0 \in \mathbb{R}^d$ , for all  $k \in \mathcal{K}$ ,  $\|k_0 - k\| \leq B$ .*

**Assumption 4.4.** *The state space  $\mathcal{S}$ , the action space  $\mathcal{A}$  are both bounded subsets of corresponding Euclidean spaces, with diameter  $B_A \ll B_S$ . The state transition function  $T(s'|s, a)$  is deterministic and continuous.*

**Assumption 4.5.** *The state transition function  $T(s'|s, a)$  is  $L_T$ -Lipschitz. For any  $\pi$  the value function  $V^\pi(s)$  is  $L_r$ -Lipschitz.*

As a consequence, given two state-action pairs  $(s_1, a_1), (s_2, a_2)$ , the next-state deviation under transition  $T$  is upper bounded by

$$\|T(s_1, a_1) - T(s_2, a_2)\| \leq L_T \|(s_1, a_1) - (s_2, a_2)\|. \quad (4)$$

**Assumption 4.6.** *The prediction model  $\hat{T}(s'|s, a)$  is a non-parametric transition model, which means the model outputs the next state prediction by searching the nearest entry.*

Formally speaking,  $\hat{T}$  has an episodic memory storing all input experience  $\mathcal{D}_{\text{memory}} = \{(s_i, a_i, s'_i, r_i)\}_{i=1}^K$  (Pritzel et al. 2017). When feeding  $\hat{T}$  a query  $(s, a)$ , the model returns  $\hat{T}(s, a) = s'_k$  where  $k = \arg \min_i \|(s, a) - (s_i, a_i)\|$ . Assumption 4.6 implies  $\hat{T}(s, a)$  is a deterministic function. Therefore combined with these two assumptions 4.4, 4.5, the gap  $G_M^\pi(s, a)$  defined in Lemma 4.1 is then bounded by

$$\begin{aligned} |G_M^\pi(s, a)| &\leq L_r W_1(\hat{T}(s, a), T(s, a)) \\ &= L_r \|\hat{T}(s, a) - T(s, a)\|, \end{aligned} \quad (5)$$

where  $W_1$  is the 1-Wasserstein distance w.r.t. the Euclidean metric.

**Assumption 4.7.**  $\rho_T^{\pi^\beta}$  have a bounded support. The diameter of the support of distribution  $\rho_T^{\pi^\beta}$  is denoted as  $B_{\pi^\beta}$ .

Since the dataset size is out of our concern, we suppose the batch data is sufficient, such that  $\hat{\rho}_T^{\pi^\beta} \approx \rho_T^{\pi^\beta} \approx \rho_T^{\pi^*}$ . For conciseness, we use  $\rho_T^{\pi^*}$  in the following statements.

**Theorem 4.8.** *Given  $\hat{\pi} = \arg \max_\pi \eta_M(\pi)$ , the expected discounted return of  $\hat{\pi}$  satisfies*

$$\begin{aligned} \eta_M(\hat{\pi}) &\geq \eta_M(\pi^*) - C(W_1(\rho_T^{\pi^*}, \rho_T^{\pi^*}) + W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*})) \\ &\geq \eta_M(\pi^*) - 2C(W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*}) + B_{\pi^\beta} + B_A), \end{aligned} \quad (6)$$

where  $C = 2c\gamma L_r L_T$  is the constant related to the assumptions. If the batch data is collected from  $N$  different policies  $\pi_1^\beta, \dots, \pi_N^\beta$ , a tighter bound is obtained, where  $\rho_T^{\pi^\beta}$  denotes the mixture of distribution  $\rho_T^{\pi_1^\beta}, \dots, \rho_T^{\pi_N^\beta}$ ,

$$\begin{aligned} \eta_M(\hat{\pi}) &\geq \eta_M(\pi^*) - C(W_1(\rho_T^{\pi^*}, \rho_T^{\pi^*}) \\ &\quad + \min_i W_1(\rho_T^{\pi_i^\beta}, \rho_T^{\pi^*}) + 2B_{\pi^\beta} + 2B_A). \end{aligned} \quad (7)$$

The distance term  $D_1 := W_1(\rho_T^{\pi^*}(s, a), \rho_T^{\pi^\beta}(s, a))$  in the first line in Equation 6 is quite hard to estimate. However, we notice that the prediction model only outputs states in the episodic memory  $\mathcal{D}_{\text{memory}}$  which implies  $\text{supp}(\rho_T^{\pi^*}(s)) \subseteq$

$\text{supp}(\rho_T^{\pi^\beta}(s))$ . We can naturally suppose that  $\rho_T^{\pi^*}(s, a)$  will not be too distinct from  $\hat{\rho}_T^{\pi^\beta}(s, a)$ . Therefore we can assume that  $D_1 \approx D_2 := W_1(\rho_T^{\pi^\beta}(s, a), \rho_T^{\pi^*}(s, a))$ , leading to an approximate lower bound free of  $B_{\pi^\beta}$  and  $B_A$

$$\eta_M(\hat{\pi}) \geq \eta_M(\pi^*) - 2C(W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*})). \quad (8)$$

For the data compounded by a mixture of policies, the approximate lower bound is

$$\eta_M(\hat{\pi}) \geq \eta_M(\pi^*) - C(W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*}) + \min_i W_1(\rho_T^{\pi_i^\beta}, \rho_T^{\pi^*})). \quad (9)$$

The detailed proof of Theorem 4.8 is presented in Appendix A. The main idea is to show the performance gap  $\mathbb{E}_{(s,a) \sim \rho_T^{\pi^*}} |G_M^{\pi^*}(s, a)|$  can be bounded by the distance between  $\rho_T^{\pi^*}$  and  $\rho_T^{\pi^\beta}$ . The remaining part of proof utilizes triangle inequality to split the distance into two terms and then applies the assumptions to yield Equation 6.

**Interpretation:** Theorem 4.8 and Equation 8 suggest that the gap relies on  $\pi^\beta$  and  $\pi^*$ . We denote the gap as  $\mathcal{L}(\pi^\beta, \pi^*)$  such that  $\eta_M(\hat{\pi}) \geq \eta_M(\pi^*) - \mathcal{L}(\pi^\beta, \pi^*)$ . When the occupancy distribution of  $\pi^\beta$  is closer to the occupancy distribution of the optimal policy  $\pi^*$ , the return of the policy optimized by MOPO will be closer to the optimal. Especially when  $\pi^\beta = \pi^*$ , MOPO can reach the optimal return. Theorem 4.8 concentrates on the gap between  $\pi^\beta$  and  $\pi^*$ . However, since the derivation does not involve the optimality of  $\pi^*$ , the inequality 6 holds true for any other policy  $\pi$  instead of the optimal policy  $\pi^*$ . By substituting  $\pi^*$  with  $\pi^\beta$ , we will obtain  $\eta_M(\hat{\pi}) \geq \eta_M(\pi^\beta)$ , which means the performance of the learned policy will perform no worse than the behavioral policy. This conclusion is consistent with the theoretical analysis in MOPO.

The second line of Equation 6 indicates a wider range of  $\rho_T^{\pi^\beta}$  may enlarge the performance gap. This issue is mainly determined by the relation between  $\rho_T^{\pi^*}(s, a)$  and  $\rho_T^{\pi^\beta}(s, a)$ . In general cases,  $\pi^*(a|s)$  will output actions that lead the next states closer to the optimal occupancy distribution. As a consequence,  $\rho_T^{\pi^*}(s, a)$  may be closer to  $\rho_T^{\pi^*}(s, a)$  than  $\rho_T^{\pi^\beta}(s, a)$ . Therefore  $D_1$  will be smaller than  $D_2$ . Nevertheless,  $D_1 > D_2$  is still possible under some non-smooth dynamics or multi-modal situations. As a result, a broader distribution of  $\rho_T^{\pi^\beta}(s, a)$  may impair MOPO performance.

## The minimal worst-case regret approach

There are many cases where the optimal policy and the corresponding experience data is inaccessible for offline learning, e.g., (1) the reward function is unknown or partly unknown at the stage of data generation; (2) the batch data is prepared for multiple tasks with various reward functions; (3) training to optimal is expensive at the stage of data generation. Previous work in online RL suggests the diversity of policies plays the crucial role (Eysenbach et al. 2019). Especially in offline RL, where exploration is not feasible during training, the diversity of batch data should not be ignored.

Suppose we can train a series of  $N$  policies simultaneously without any external reward. Our goal is to improve the diversity of the experience collected by policies  $\{\pi_i\}_{i=1}^N$ , such that there is at least one subset of the experience will be close enough to the optimal policy determined by the lately designated reward function at the offline training stage. Combined with Theorem 4.8, this objective can be formulated by

$$\min_{\pi_1, \dots, \pi_N \in \Pi} \max_{\pi^* \in \Pi} \min_i \mathcal{L}(\pi^i, \pi^*). \quad (10)$$

The inner  $\min$  term  $\text{REGRET}(\{\pi_i\}_{i=1}^N, \pi^*) := \min_i \mathcal{L}(\pi_i, \pi^*)$  represents the regret of the series of policies confronting the true reward function and its associate optimal policy. Since we are able to choose one policy in the series after informed of the optimal policy, we take the minimum as the regret. The  $\max$  operator in the middle depicts the worst-case regret if any policy in the feasible policy set  $\Pi$  has the possibility to be the optimal one. The outer  $\min$  means the goal of optimizing  $\{\pi_i\}_{i=1}^N$  is to minimize the worst-case regret. If the approximate lower bound is considered in Equation 8, the objective is equivalent to

$$\min_{\pi_1, \dots, \pi_N \in \Pi} \max_{\pi^* \in \Pi} \min_i W_1(\rho_T^{\pi_i}, \rho_T^{\pi^*}). \quad (11)$$

Directly optimizing a series of policies according to the min-max objective in Equation 11 inevitably requires adversarial training. Previous practice suggests an adversarial policy should be introduced to maximize  $\text{REGRET}(\{\pi_i\}_{i=1}^N, \pi^*)$ , playing the role of the unknown optimal policy. The adversarial manner of training brings us two main concerns. (1) Adversarial training may incur instability and require much more steps to converge (Arjovsky and Bottou 2017; Arjovsky, Chintala, and Bottou 2017); (2) The regret only provides supervision signals to the policy nearest to  $\pi^*$ , which leads to low efficiency in optimization.

Eysenbach et al. (Eysenbach, Salakhutdinov, and Levine 2021) proposed similar objective regarding unsupervised reinforcement learning. Under the assumptions of finite and discrete state space and, the quantity of policies  $N$  should cover the number of distinct states in the state space  $|\mathcal{S}|$ , the minimal worst-case regret objective is equivalent to the maximal discriminability objective. Likewise, we propose a surrogate objective

$$\max_{\pi_1, \dots, \pi_N \in \Pi} \min_{i \neq j} W_1(\rho_T^{\pi_i}, \rho_T^{\pi_j}). \quad (12)$$

The surrogate objective shares the same spirit with WURL (He et al. 2022). Both of them encourage diversity of a series of policies w.r.t. Wasserstein distance in the probability space of state occupancy. Although the optimal solution of  $\{\pi_i\}_{i=1}^N$  does not match the optimal solution in Equation 11 in general situations, both of them represent a kind of diversity. The relation between two objectives equals to the relation between finite covering and finite packing problems, which are notoriously difficult to analyze even in low-dimension, convex settings (Böröczky Jr, Böröczky et al. 2004; Toth, O’Rourke, and Goodman 2017). Nevertheless, we assume the gap will be small and the surrogate objective will be a satisfactory proxy of Equation 11 as previous literature does in the application of computational graphics (Schlömer, Heck,

and Deussen 2011; Chen and Xu 2004). Refer to Appendix B for more details.

## Practical implementation

To achieve diversity, practical algorithms assign a pseudo reward  $\tilde{r}_i$  to policy  $\pi_i$ . The pseudo reward usually indicates the “novelty” of the current policy w.r.t. all other policies. Similar to WURL, we adopt pseudo reward  $\tilde{r}_i := \min_{j \neq i} W_1(\rho_T^{\pi_j}, \rho_T^{\pi_i})$  which is the minimum distance from all other policies. We compute the Wasserstein distance using amortized primal form estimation in consistent with WURL (He et al. 2022).

In semi-supervised cases, only part of reward function is known. For example, in Mujoco simulation environments in OpenAI Gym (Brockman et al. 2016), the complete reward function is composed of a reward related to the task, and general rewards related to agent’s health, control cost, safety constraint, etc. We can train the series of policies with a partial reward and a pseudo reward simultaneously by reweighting two rewards with a hyperparameter  $\lambda$ . Moreover, the complete reward and the diversity-induced pseudo reward can be combined to train a diverse series of policies for generalization purposes.

The policies are trained with Soft Actor-Critic method (Haarnoja et al. 2018). The network model of the actors are stored to generate experience  $\mathcal{D}_1, \dots, \mathcal{D}_K$  for offline RL. When a different reward function is used at the offline training stage. The reward will be relabeled with  $r(s, a)$ . At the offline learning stage, we choose the best buffer and feed it to MOPO. The overall algorithm is illustrated in Figure 2 and formally described in Algorithm 1.

---

Algorithm 1: Unsupervised data generation for offline RL in task-agnostic settings

---

**Require:**  $K$  policies  $\pi_1, \dots, \pi_K$ .  $K$  empty buffers  $\mathcal{D}_1, \dots, \mathcal{D}_K = \{\}$ . Maximum buffer size  $N$ .

- 1: Train  $\pi_i, i = 1, \dots, K$  with SAC w.r.t. diversity rewards  $\tilde{r}_i := \min_{j \neq i} W_1(\rho_T^{\pi_j}, \rho_T^{\pi_i})$ .
  - 2: Let each  $\pi_i$  interacts with environment for  $N$  steps and fill  $\mathcal{D}_i$  with transitions  $(s, a, s')$ .
  - 3: Acquire the task and relabel all transitions in  $\mathcal{D}_1, \dots, \mathcal{D}_K$  with given  $r(s, a)$ .
  - 4: Evaluate each buffer and calculate the average return  $\bar{G}_i, i = 1, \dots, K$ .
  - 5: Select the buffer  $\mathcal{D}_k$  where  $k = \arg \max_i \bar{G}_i$
  - 6: Train the policy by MOPO with  $\mathcal{D}_k$ .
- 

## 5 Experiments

Based on our framework of UDG in Figure 2, we conduct experiments on two locomotive environments requiring the agent to solve a series of tasks with different reward functions at the offline stage. Both of the tasks are re-designed Mujoco environments. Ant-Angle is a task modified from Ant environment. In Ant-Angle, the agent should actuate the ant to move from the initial position to a specific direction on the x-y plane. The agent is rewarded by the inner product of the moving direction and the desired direction. The

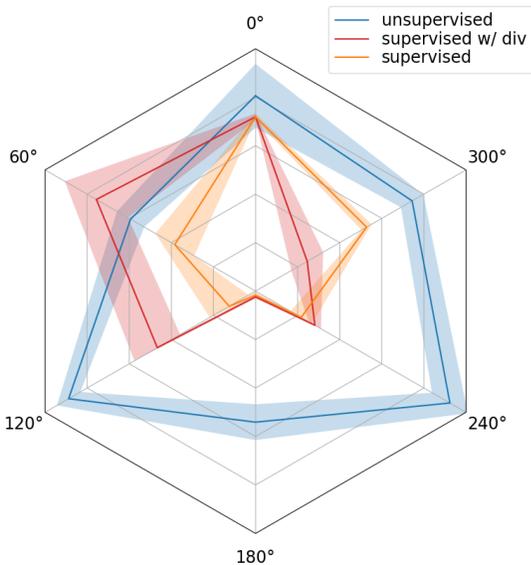


Figure 3: Results on Ant-Angle tasks. The data buffers of all three methods are evaluated by MOPO with 3 random seeds. The darker lines represent the average evaluation return. The lighter areas depict the standard deviation of the return.

goal is to construct a dataset while the desired direction is unknown until the offline stage. Cheetah-Jump is another task for evaluation, modified from HalfCheetah environment. The reward in Cheetah-Jump consists of three parts, control cost, velocity reward, and jumping reward. At the data generation stage, the agent can only have access to the control cost and the velocity reward for reducing energy cost of actuators and moving the cheetah forward. The jumping reward is added in offline training, by calculating the positive offset of the cheetah on the z axis. Likewise, a crawling reward can be added to encourage the cheetah to lower the body while moving forward.

Our experiments mainly focus on two aspects: (1) How does UDG framework perform on the two challenging tasks, Ant-Angle and Cheetah-Jump? (2) Can experimental results match the findings in the theoretical analysis?

### Evaluation on task-agnostic settings

To answer question (1), we construct three types of data buffer. The first is generated by unsupervisedly trained policies with the objective in Equation 12. The second is created by one supervisedly trained policy that maximizes a specific task reward. The third is the combination of two, which means the policies are trained with both task reward and diversity reward, reweighted by  $\lambda$ . We call the combination as supervised training with diversity regularizer. Note that the supervised method contains one data buffer. Another two methods have a series of 10 buffers w.r.t. 10 policies and only one buffer is selected during offline training.

We evaluate three kinds of data buffers on Ant-Angle. The supervised policy and the supervised policies with diversity regularizer are provided with reward to move in direction

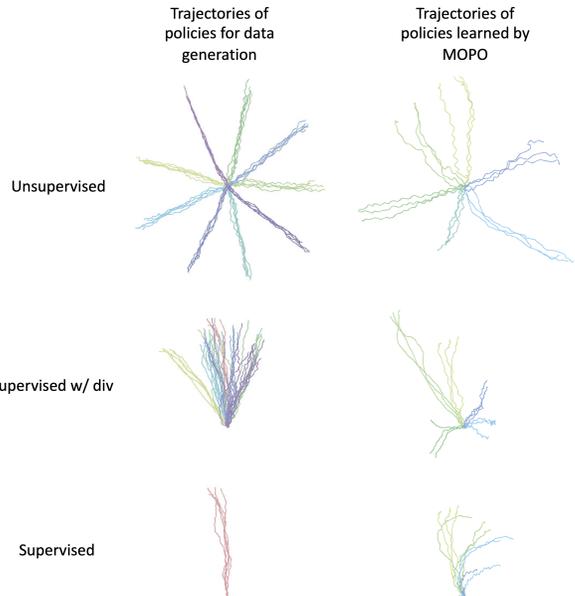


Figure 4: The trajectories of ant on the x-y plane. The lower left figure demonstrates UDG can solve all 6 offline tasks on account of the diversity of unsupervised learned policies. The trajectories in the lower right figure shows the policy learned by MOPO cannot generalize to the directions that largely deviate from  $0^\circ$ .

$0^\circ$ , the upper direction in Figure 4. The diversity reward is calculated on ant position instead of the whole state space, considering that the dimension of state space is extremely high. We evaluate three approaches on 6 offline tasks of moving along the directions of  $0^\circ$ ,  $60^\circ$ ,  $120^\circ$ ,  $180^\circ$ ,  $240^\circ$  and  $300^\circ$ . As Figure 4 shows, the policies trained with unsupervised RL are evenly distributed over the x-y plane. Therefore in the downstream offline tasks, no matter what direction the task needs, there exists at least one policy that could obtain relatively high reward. The trajectories of policies trained by MOPO confirm that UDG can handle all 6 tasks. Meanwhile the policy trained to move in the direction  $0^\circ$  generates narrow data and MOPO cannot perform well on other directions. The policies trained with combined reward have wider range of data distribution. Especially, an policy deviates to circa  $30^\circ$  and consequently the policy trained by MOPO acquires high reward in the  $60^\circ$  task.

Task	$c_z$	random	diverse
Cheetah-Jump	15	$1152.98 \pm 120$	<b><math>1721.25 \pm 56</math></b>
Cheetah-Crawl	-15	$1239.00 \pm 57$	<b><math>1348.19 \pm 274</math></b>

Table 1: Returns on two offline tasks Cheetah-Jump and Cheetah-Crawl. Two datasets consist of 5 policies trained with base rewards and base rewards plus diversity rewards respectively.

In Cheetah-Jump tasks, we relabel the data by adding an extra reward  $c_z(z - z_0)$  where  $z_0$  is the initial position on the

z axis, and  $c_z$  is the coefficient of the extra reward.  $c_z$  can either be positive or negative. For positive values of  $c_z$ , the cheetah is encouraged to jump while running forward. For negative values, crawling on the floor receives higher reward. We train 5 policies each for base rewards and base rewards plus diversity rewards, denoted by “random” and “diverse” respectively.

### Effects of the range of data distribution

Angle	top 1	top 2 mixed	all mixed
0°	1236.26±247	<b>1437.24±31</b>	989.13±65
60°	910.70±121	<b>1285.31±66</b>	593.88±434
120°	<b>1362.46±104</b>	917.10±218	281.88±301
180°	829.65±139	<b>1034.41±224</b>	717.68±120
240°	<b>1416.80±160</b>	1373.72±73	850.82±62
300°	<b>1141.68±100</b>	1087.26±137	817.77±37

Table 2: Returns on Ant-Angle tasks with different angles trained on different datasets. Top 1 dataset is the data buffer with highest return. Top 2 mixed dataset is a mixture of two highest-rewarded buffers. All mixed dataset is a mixture of all data buffers generated by unsupervisedly trained policies.

With the help of Ant-Angle environment and the policies learned by unsupervised RL, we conduct several experiments to verify the conclusions from theoretical derivations. Apart from the data buffer with maximum return, we build a data buffer denoted by “all mixed”, by mixing data generated by all 10 policies. We also mix the data from top 2 policies to create a “top 2 mixed” buffer.

Referring to the upper left figure in Figure 4, the “top 2 mixed” data buffer includes two policies lying on the left and the right side near direction 0°. the top two distributions have similar distance from the optimal distribution. Therefore the mixed distribution  $\rho_T^{\pi^\beta}$  has similar distance to the optimal compared with the nearest distribution  $W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*}) \approx \min_i W_1(\rho_T^{\pi^i}, \rho_T^{\pi^*})$ . However, when all policies are mixed, it is obvious  $W_1(\rho_T^{\pi^\beta}, \rho_T^{\pi^*}) > \min_i W_1(\rho_T^{\pi^i}, \rho_T^{\pi^*})$ . According to Equation 9, the top 2 mixed dataset will get higher return than all mixed dataset. Table 2 and results in Appendix D have verified this claim. From another aspect, the top 2 mixed data buffer has a wider distribution than the top 1 buffer. Therefore the top 2 mixed buffer has a larger radius  $B_{\pi^\beta}$  which may worsen performance according to Equation 6. Surprisingly, the top 2 mixed buffer makes higher return than top 1 single buffer. We can conjecture that  $B_{\pi^\beta}$  plays an insignificant role in the lower bound and the approximation in Equation 8 and 9 is proper. In addition, the wide spread of the mixed data may improve the generalization ability of the transition model in MOPO, which contributes to the higher return than top 1 data buffer.

## 6 Discussion

### Limitations

Our derivation is based on the continuous state space with assumption that the transition function and the value functions are Lipschitz. There are some tasks may break the assumptions, e.g., pixel based tasks like Atari, non-smooth reward functions in goal reaching tasks (Tassa et al. 2018). Therefore, it is necessary to verify the feasibility of UDG on these tasks in future deployment. We also adopt a non-parametric transition model in derivation. In practical model-based offline RL approaches, neural models have greater generalization ability than the non-parametric model. The influence of data distribution on neural models is not addressed by this work. In addition, whether can UDG be generalized with model-free offline RL algorithms remains unclear. Another limitation is at the unsupervised training stage, the diversity reward is calculated on the low dimensional space where the reward function is defined, e.g., the x-y plane in Ant-Angle. This requires prior knowledge of how the reward is computed. Nevertheless, the limitations mentioned above indicate interesting further research.

### Societal impact

The UDG framework contains a stage of unsupervised RL. At this stage, the agent is not provided with any reward for solving any task. During the process of training, the agent may unexpectedly exploit the states in unsafe regions. Especially when deployed in realistic environments, it could incur damage of the environment or the robotic agent itself, or cause injury in robot-human interactions. Any deployment of UDG in the real world should be carefully designed to avoid safety incidents.

### Conclusion

In this study we propose a framework UDG addressing data generation issues in offline reinforcement learning. In order to solve unknown tasks at the offline training stage, UDG first employs unsupervised RL and obtains a series of diverse policies for data generation. The experience generated by each policy is relabeled according the reward function adopted before the offline training stage. The final step of UDG is to select the data buffer with highest average return and to feed the data to model-based offline RL algorithms like MOPO. We provide theoretical analysis on the performance gap between the offline learned policy and the optimal policy w.r.t the distribution of the batch data. We also reveal that UDG is an approximate minimal worst-case regret approach under the task-agnostic setting. Our experiments evaluate UDG on two locomotive tasks, Ant-Angle and Cheetah-Jump. Empirical results on multiple offline tasks demonstrate UDG is overall better than data generated by a policy dedicated to solve a specific task. Additional experiments show that the range of data distribution has minor effects on performance and the distance from the optimal policy is the most important factor. It is also confirmed that choosing the data buffer with highest return is necessary for better performance.

## References

- Agarwal, R.; Schuurmans, D.; and Norouzi, M. 2019. Striving for simplicity in off-policy deep reinforcement learning. *arXiv preprint arXiv:1907.04543*.
- Arjovsky, M.; and Bottou, L. 2017. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*.
- Arjovsky, M.; Chintala, S.; and Bottou, L. 2017. Wasserstein generative adversarial networks. In *International conference on machine learning*, 214–223. PMLR.
- Böröczky Jr, K.; Böröczky, K.; et al. 2004. *Finite packing and covering*, volume 154. Cambridge University Press.
- Brockman, G.; Cheung, V.; Pettersson, L.; Schneider, J.; Schulman, J.; Tang, J.; and Zaremba, W. 2016. Openai gym. *arXiv preprint arXiv:1606.01540*.
- Burda, Y.; Edwards, H.; Storkey, A.; and Klimov, O. 2018. Exploration by random network distillation. *arXiv preprint arXiv:1810.12894*.
- Campos, V.; Trott, A.; Xiong, C.; Socher, R.; Giro-i Nieto, X.; and Torres, J. 2020. Explore, discover and learn: Unsupervised discovery of state-covering skills. In *International Conference on Machine Learning*, 1317–1327. PMLR.
- Chen, J.; and Jiang, N. 2019. Information-theoretic considerations in batch reinforcement learning. In *International Conference on Machine Learning*, 1042–1051. PMLR.
- Chen, L.; and Xu, J.-c. 2004. Optimal delaunay triangulations. *Journal of Computational Mathematics*, 299–308.
- Chen, X.-H.; Yu, Y.; Li, Q.; Luo, F.-M.; Qin, Z.; Shang, W.; and Ye, J. 2021. Offline model-based adaptable policy learning. *Advances in Neural Information Processing Systems*, 34: 8432–8443.
- Cheng, C.-A.; Xie, T.; Jiang, N.; and Agarwal, A. 2022. Adversarially trained actor critic for offline reinforcement learning. In *International Conference on Machine Learning*, 3852–3878. PMLR.
- Eysenbach, B.; Gupta, A.; Ibarz, J.; and Levine, S. 2019. Diversity is All You Need: Learning Skills without a Reward Function. In *International Conference on Learning Representations*.
- Eysenbach, B.; Salakhutdinov, R.; and Levine, S. 2021. The information geometry of unsupervised reinforcement learning. *arXiv preprint arXiv:2110.02719*.
- Florensa, C.; Duan, Y.; and Abbeel, P. 2017. Stochastic neural networks for hierarchical reinforcement learning. *arXiv preprint arXiv:1704.03012*.
- Fu, J.; Kumar, A.; Nachum, O.; Tucker, G.; and Levine, S. 2020. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*.
- Fujimoto, S.; and Gu, S. S. 2021. A minimalist approach to offline reinforcement learning. *Advances in neural information processing systems*, 34: 20132–20145.
- Fujimoto, S.; Hoof, H.; and Meger, D. 2018. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, 1587–1596. PMLR.
- Fujimoto, S.; Meger, D.; and Precup, D. 2019. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, 2052–2062. PMLR.
- Gulcehre, C.; Wang, Z.; Novikov, A.; Paine, T.; Gómez, S.; Zolna, K.; Agarwal, R.; Merel, J. S.; Mankowitz, D. J.; Paduraru, C.; et al. 2020. RL unplugged: A suite of benchmarks for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33: 7248–7259.
- Haarnoja, T.; Zhou, A.; Abbeel, P.; and Levine, S. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, 1861–1870. PMLR.
- He, S.; Jiang, Y.; Zhang, H.; Shao, J.; and Ji, X. 2022. Wasserstein unsupervised reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 6884–6892.
- Hessel, M.; Modayil, J.; Van Hasselt, H.; Schaul, T.; Ostrovski, G.; Dabney, W.; Horgan, D.; Piot, B.; Azar, M.; and Silver, D. 2018. Rainbow: Combining improvements in deep reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Janner, M.; Fu, J.; Zhang, M.; and Levine, S. 2019. When to trust your model: Model-based policy optimization. *Advances in neural information processing systems*, 32.
- Kidambi, R.; Rajeswaran, A.; Netrapalli, P.; and Joachims, T. 2020. Morel: Model-based offline reinforcement learning. *Advances in neural information processing systems*, 33: 21810–21823.
- Kim, J.; Park, S.; and Kim, G. 2021. Unsupervised skill discovery with bottleneck option learning. *arXiv preprint arXiv:2106.14305*.
- Kiran, B. R.; Sobh, I.; Talpaert, V.; Mannion, P.; Al Salhab, A. A.; Yogamani, S.; and Pérez, P. 2021. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6): 4909–4926.
- Kumar, A.; Fu, J.; Soh, M.; Tucker, G.; and Levine, S. 2019. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*, 32.
- Kumar, A.; Zhou, A.; Tucker, G.; and Levine, S. 2020. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33: 1179–1191.
- Lambert, N.; Wulfmeier, M.; Whitney, W.; Byravan, A.; Bloesch, M.; Dasagi, V.; Hertweck, T.; and Riedmiller, M. 2022. The challenges of exploration for offline reinforcement learning. *arXiv preprint arXiv:2201.11861*.
- Laskin, M.; Yarats, D.; Liu, H.; Lee, K.; Zhan, A.; Lu, K.; Cang, C.; Pinto, L.; and Abbeel, P. 2021. URLB: Unsupervised reinforcement learning benchmark. *arXiv preprint arXiv:2110.15191*.
- Lee, L.; Eysenbach, B.; Parisotto, E.; Xing, E.; Levine, S.; and Salakhutdinov, R. 2019. Efficient exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*.

- Levine, S.; Kumar, A.; Tucker, G.; and Fu, J. 2020. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*.
- Liu, H.; and Abbeel, P. 2021a. Aps: Active pretraining with successor features. In *International Conference on Machine Learning*, 6736–6747. PMLR.
- Liu, H.; and Abbeel, P. 2021b. Behavior from the void: Unsupervised active pre-training. *Advances in Neural Information Processing Systems*, 34: 18459–18473.
- Matsushima, T.; Furuta, H.; Matsuo, Y.; Nachum, O.; and Gu, S. 2020. Deployment-efficient reinforcement learning via model-based offline optimization. *arXiv preprint arXiv:2006.03647*.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *nature*, 518(7540): 529–533.
- Müller, A. 1997. Integral probability metrics and their generating classes of functions. *Advances in applied probability*, 29(2): 429–443.
- Pathak, D.; Agrawal, P.; Efros, A. A.; and Darrell, T. 2017. Curiosity-driven exploration by self-supervised prediction. In *International conference on machine learning*, 2778–2787. PMLR.
- Pathak, D.; Gandhi, D.; and Gupta, A. 2019. Self-supervised exploration via disagreement. In *International conference on machine learning*, 5062–5071. PMLR.
- Peng, X. B.; Kumar, A.; Zhang, G.; and Levine, S. 2019. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*.
- Pritzel, A.; Uria, B.; Srinivasan, S.; Badia, A. P.; Vinyals, O.; Hassabis, D.; Wierstra, D.; and Blundell, C. 2017. Neural episodic control. In *International conference on machine learning*, 2827–2836. PMLR.
- Rezaeifar, S.; Dadashi, R.; Vieillard, N.; Hussenot, L.; Bachem, O.; Pietquin, O.; and Geist, M. 2022. Offline reinforcement learning as anti-exploration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 8106–8114.
- Schlömer, T.; Heck, D.; and Deussen, O. 2011. Farthest-point optimized point sets with maximized minimum distance. In *Proceedings of the ACM SIGGRAPH Symposium on High Performance Graphics*, 135–142.
- Sharma, A.; Gu, S.; Levine, S.; Kumar, V.; and Hausman, K. 2020. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations (ICLR)*.
- Siegel, N. Y.; Springenberg, J. T.; Berkenkamp, F.; Abdolmaleki, A.; Neunert, M.; Lampe, T.; Hafner, R.; Heess, N.; and Riedmiller, M. 2020. Keep doing what worked: Behavioral modelling priors for offline reinforcement learning. *arXiv preprint arXiv:2002.08396*.
- Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; Van Den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587): 484–489.
- Silver, D.; Schrittwieser, J.; Simonyan, K.; Antonoglou, I.; Huang, A.; Guez, A.; Hubert, T.; Baker, L.; Lai, M.; Bolton, A.; et al. 2017. Mastering the game of go without human knowledge. *nature*, 550(7676): 354–359.
- Sodhani, S.; Zhang, A.; and Pineau, J. 2021. Multi-task reinforcement learning with context-based representations. In *International Conference on Machine Learning*, 9767–9779. PMLR.
- Strouse, D.; Baumli, K.; Warde-Farley, D.; Mnih, V.; and Hansen, S. 2021. Learning more skills through optimistic exploration. *arXiv preprint arXiv:2107.14226*.
- Tassa, Y.; Doron, Y.; Muldal, A.; Erez, T.; Li, Y.; Casas, D. d. L.; Budden, D.; Abdolmaleki, A.; Merel, J.; Lefrancq, A.; et al. 2018. Deepmind control suite. *arXiv preprint arXiv:1801.00690*.
- Todorov, E.; Erez, T.; and MuJoCo, Y. 2012. A physics engine for model-based control. In *Proceedings of the 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 5026–5033.
- Toth, C. D.; O’Rourke, J.; and Goodman, J. E. 2017. *Handbook of discrete and computational geometry*. CRC press.
- Wang, H.; Feng, D.; Ding, B.; and Li, W. 2022. Offline Imitation Learning Using Reward-free Exploratory Data. In *Proceedings of the 2022 5th International Conference on Algorithms, Computing and Artificial Intelligence*, 1–9.
- Wang, R.; Foster, D. P.; and Kakade, S. M. 2020. What are the statistical limits of offline RL with linear function approximation? *arXiv preprint arXiv:2010.11895*.
- Wu, Y.; Tucker, G.; and Nachum, O. 2019. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*.
- Yarats, D.; Brandfonbrener, D.; Liu, H.; Laskin, M.; Abbeel, P.; Lazaric, A.; and Pinto, L. 2022. Don’t change the algorithm, change the data: Exploratory data for offline reinforcement learning. *arXiv preprint arXiv:2201.13425*.
- Yarats, D.; Fergus, R.; Lazaric, A.; and Pinto, L. 2021. Reinforcement learning with prototypical representations. In *International Conference on Machine Learning*, 11920–11931. PMLR.
- Yu, T.; Kumar, S.; Gupta, A.; Levine, S.; Hausman, K.; and Finn, C. 2020a. Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*, 33: 5824–5836.
- Yu, T.; Quillen, D.; He, Z.; Julian, R.; Hausman, K.; Finn, C.; and Levine, S. 2020b. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, 1094–1100. PMLR.
- Yu, T.; Thomas, G.; Yu, L.; Ermon, S.; Zou, J. Y.; Levine, S.; Finn, C.; and Ma, T. 2020c. Mopo: Model-based offline policy optimization. *Advances in Neural Information Processing Systems*, 33: 14129–14142.
- Zhang, H.; Shao, J.; Jiang, Y.; He, S.; Zhang, G.; and Ji, X. 2022. State Deviation Correction for Offline Reinforcement Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 9022–9030.

Zhao, W.; Queraltà, J. P.; and Westerlund, T. 2020. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. In *2020 IEEE symposium series on computational intelligence (SSCI)*, 737–744. IEEE.