

SHOULD I RUN OFFLINE REINFORCEMENT LEARNING OR BEHAVIORAL CLONING?

Anonymous authors

Paper under double-blind review

ABSTRACT

Offline reinforcement learning (RL) algorithms can acquire effective policies by utilizing only previously collected experience, without any online interaction. While it is widely understood that offline RL is able to extract good policies even from highly suboptimal data, in practice offline RL is often used with data that resembles demonstrations. In this case, one can also use behavioral cloning (BC) algorithms, which mimic a subset of the dataset via supervised learning. It seems natural to ask: *When should we prefer offline RL over BC?* In this paper, our goal is to characterize environments and dataset compositions where offline RL leads to better performance than BC. In particular, we characterize the properties of environments that allow offline RL methods to perform better than BC methods even when only provided with expert data. Additionally, we show that policies trained on suboptimal data that is sufficiently noisy can attain better performance than even BC algorithms with expert data, especially on long-horizon problems. We validate our theoretical results via extensive experiments on both diagnostic and high-dimensional domains including robot manipulation, maze navigation and Atari games, when learning from a variety of data sources. We observe that modern offline RL methods trained on suboptimal, noisy data in sparse reward domains outperform cloning the expert data in several practical problems.

1 INTRODUCTION

Offline reinforcement learning (RL) algorithms aim to leverage large, existing datasets of previously collected data to produce effective policies that generalize across a wide range of scenarios, without the need for costly active data collection. Many recent offline RL algorithms (Fujimoto et al., 2018; Kumar et al., 2019; Wu et al., 2019; Kumar et al., 2020; Yu et al., 2020; Sinha et al., 2021; Kostrikov et al., 2021) can work well even highly suboptimal data. With recent advances, the performance of offline RL algorithms has improved significantly, and a number of these approaches have been studied theoretically (Wang et al., 2021; Zanette, 2020; Rashidinejad et al., 2021). While it is clear that offline RL algorithms are a good choice when the available data is either random or highly suboptimal, such methods are also often used with datasets that come from demonstrations, or other near-optimal data sources. In these cases, imitation learning algorithms, such as behavioral cloning (BC), can be used to train policies via supervised learning. It then seems natural to ask: *When should we prefer to use offline RL over imitation learning?*

A rigorous theoretical and empirical characterization of when offline RL perform better than imitation learning, or algorithms that blend in imitation learning with suboptimal data (e.g., filtered behavior cloning) is still absent. This makes it quite confusing for practitioners to understand what class of methods to apply to a given problem. Existing empirical studies comparing offline RL to IL have been mixed, and are likely heavily confounded by modeling choices and algorithm hyperparameters. Some works show BC approaches significantly outperform offline RL on demonstration data (Mandlekar et al., 2021). Others show that offline RL methods appear to lead to greatly improved performance over imitation learning, especially in environments that require “stitching” parts of suboptimal trajectories (Fu et al., 2020). Recent theoretical results (Rashidinejad et al., 2021; Jin et al., 2020) show that pessimistic offline RL algorithms are, in general, minimax optimal in contextual bandits with minimal assumptions on the dataset, but do not compare offline RL algorithms and imitation learning. We therefore wish to understand under which environment or dataset conditions offline RL algorithms will outperform imitation learning.

Our contribution in this paper is to characterize when and how offline RL can outperform BC. We do this using both theory and empirical results. Theoretically, we compare performance guarantees for example pessimistic offline RL algorithms to those of BC and general imitation learning methods. First, we derive general guarantees for offline RL and BC that scale with the suboptimality of the behavior policy that collected the dataset. When the data is optimal (*i.e.*, from expert demonstrations), we note that both RL and BC achieve the same worst-case bounds. However, when the MDP satisfies certain structural assumptions, the error incurred by offline RL algorithms can scale significantly more favorably with the horizon. Such structure includes horizon-independent returns (*i.e.*, sparse rewards), or a low volume of states where it is “critical” to take the same action as the expert. Meanwhile, when the data is suboptimal, we show that can be preferable to use RL, not only over BC with the same dataset, but even over BC with an optimal dataset of the same size in long horizon tasks. This occurs when we have access to data generated by a sufficiently noisy behavior policy, which is often much easier to obtain than demonstrations. Finally, we consider filtered behavior cloning methods that use the reward to inform learning, and characterize conditions when offline RL can perform better.

Empirically, we validate our theoretical conclusions on diagnostic gridworld domains (Fu et al., 2019) and large-scale benchmark problems in robotic manipulation and navigation and Atari games, using human data (Fu et al., 2020), scripted data (Singh et al., 2020) and data generated from RL policies (Agarwal et al., 2020b). We verify that in multiple relatively long-horizon problem domains where the conditions we consider are likely to be satisfied, practical deep offline RL methods do outperform behavioral cloning, and related methods, especially when allowed to use noisy data.

2 RELATED WORK

Offline RL (Lange et al., 2012; Levine et al., 2020) has shown promise in domains such as robotic manipulation (Kalashnikov et al., 2018b; Mandlekar et al., 2020; Singh et al., 2020; Kalashnikov et al., 2021), NLP (Jaques et al., 2020) and healthcare (Shortreed et al., 2011; Wang et al., 2018). The major challenge in offline RL is distribution shift (Fujimoto et al., 2018; Kumar et al., 2019), where the learned policy might execute out-of-distribution actions. Prior offline RL methods can broadly be characterized into two categories: (1) *policy-constraint* methods that regularize the learned policy to be “close” to the behavior policy either explicitly (Fujimoto et al., 2018; Kumar et al., 2019; Liu et al., 2020b; Wu et al., 2019; Fujimoto & Gu, 2021) or implicitly (Siegel et al., 2020; Peng et al., 2019; Nair et al., 2020), or via importance sampling (Liu et al., 2019; Swaminathan & Joachims, 2015; Nachum et al., 2019), and (2) *conservative* methods that learn a lower-bound, or conservative, estimate of return and optimize the policy against it (Kumar et al., 2020; Kostrikov et al., 2021; Kidambi et al., 2020; Yu et al., 2020; 2021). Our goal is not to devise a new offline RL algorithm, but rather to understand when existing offline RL methods from each category can outperform BC.

When do offline RL methods outperform BC? Rashidinejad et al. (2021) derive a conservative offline RL algorithm based on lower-confidence bounds that provably outperforms BC in the simpler contextual bandits (CB) setting, but do not extend it to MDPs. While this CB result signals the possibility that offline RL can outperform BC in theory, this generalization is not trivial, as RL suffers from compounding errors (Munos, 2003; 2005; Wang et al., 2021). Laroché et al. (2019); Nadjahi et al. (2019); Kumar et al. (2020); Liu et al. (2020b); Xie et al. (2021) present safe policy improvement bounds expressed as improvements over the behavior policy, which imitation aims to recover, but these bounds do not clearly indicate when offline RL is better or worse. Empirically, Fu et al. (2020) show that offline RL considerably outperforms BC for tasks that require “stitching” trajectory segments to devise an optimal policy. In contrast, Mandlekar et al. (2021); Brandfonbrener et al. (2021); Chen et al. (2021) suggests that simple BC or filtered BC that only uses the top fraction of the data is the better alternative on other tasks. We provide a theoretical characterization of problems where we would expect offline RL to be better than BC, and empirical results that verify that offline RL performs well on such problems, which span domains as robotics, navigation and games (Fu et al., 2020; Singh et al., 2020; Bellemare et al., 2013).

Our theoretical analysis combines tools from a number of prior works. We analyze the total error incurred by RL via an error propagation analysis (Munos, 2003; 2005; Farahmand et al., 2010; Chen & Jiang, 2019; Xie & Jiang, 2020; Liu et al., 2020b), which gives rise to bounds with *concentration coefficients* that bound the total distributional shift between the learned policy and the data distribution (Xie & Jiang, 2020; Liu et al., 2020b). We use tools from Ren et al. (2021), which provide horizon-free bounds for standard (non-conservative) offline Q-learning but under strict cov-

erage assumptions. While we use a LCB-style algorithm (Rashidinejad et al., 2021) for analysis, our conservative offline RL algorithm uses tighter Bernstein bonuses (Zhang et al., 2021; Agarwal et al., 2020a) compared to the standard Hoeffding bonus used by the prior work, which makes our instantiation of the LCB paradigm enjoy stronger suboptimality guarantees.

3 PROBLEM SETUP AND PRELIMINARIES

The goal in reinforcement learning is to learn a policy $\pi(\cdot|s)$ that maximizes the expected cumulative discounted reward in a Markov decision process (MDP), which is defined by a tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$. \mathcal{S}, \mathcal{A} represent state and action spaces, $P(s'|s, a)$ and $r(s, a)$ represent the dynamics and mean reward function, and $\gamma \in (0, 1)$ represents the discount factor. The effective horizon of the MDP is given by $H = 1/(1 - \gamma)$. The Q-function, $Q^\pi(s, a)$ for a given policy π is equal to the discounted long-term reward attained by executing a at the state s and then following policy π thereafter. Q^π satisfies the recursion: $\forall s, a \in \mathcal{S} \times \mathcal{A}, Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a), a' \sim \pi(\cdot|s')} [Q(s', a')]$. The value function V^π considers the expectation of the Q-function over the policy $V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)]$. Meanwhile, the Q-function of the optimal policy, Q^* , satisfies the recursion: $Q^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} [\max_{a'} Q^*(s', a')]$, and the optimal value function is given by $V^*(s) = \max_a Q^*(s, a)$. Finally, the expected cumulative discounted reward is given by $J(\pi) = \mathbb{E}_{s_0 \sim \rho} [V^\pi(s_0)]$.

In offline RL, we are provided with a dataset \mathcal{D} of transitions, $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$ of size $|\mathcal{D}| = N$. We assume that the dataset \mathcal{D} is generated i.i.d. from a distribution $\mu(s, a)$ that specifies the effective behavior policy $\pi_\beta(a|s) := \mu(s, a) / \sum_a \mu(s, a)$. Note that the data might itself be generated by running a non-Markovian policy, but the marginal (s, a) distribution of this data would always correspond to an effective behavior policy π_β (Puterman, 1994). Let $n(s, a)$ be the number of times (s, a) appear in \mathcal{D} , and $\hat{P}(\cdot|s, a)$ and $\hat{r}(s, a)$ denote the empirical dynamics and reward distributions in \mathcal{D} , which may be different from P and r due to stochasticity. Following Rashidinejad et al. (2021), the goal of policy learning is to maximize performance of the learned policy $\hat{\pi}$ averaged over the randomness in \mathcal{D} :

$$\text{SubOpt}(\hat{\pi}) = \mathbb{E}_{\mathcal{D} \sim \mu} [J(\pi^*) - J(\hat{\pi})] = \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{s_0 \sim \rho} [V^*(s_0) - V^{\hat{\pi}}(s_0)] \right].$$

We will evaluate both BC and offline RL performance using this suboptimality metric (Equation 3).

Dataset and MDP conditions. Here we introduce some conditions on the offline dataset and MDP structure that we make for our analysis. The first characterizes the distribution shift between the data distribution $\mu(s, a)$ and the normalized state-action marginal of π^* , given by $d^*(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi^*)$, via a *concentrability coefficient* C^* .

Condition 3.1 (Rashidinejad et al. (2021), Concentrability of the data distribution). *Define C^* to be the smallest, finite constant that satisfies: $d^*(s, a) / \mu(s, a) \leq C^* \forall s \in \mathcal{S}, a \in \mathcal{A}$.*

Intuitively, the coefficient C^* scales with how suboptimal the data $\mu(s, a)$ is relative to the optimal π^* , where $C^* = 1$ corresponds to data from π^* .

The next condition we consider is that the discounted return for any trajectory in the MDP is bounded by a constant, which w.l.o.g., we assume to be 1.

Condition 3.2 (Ren et al. (2021), the value of any trajectory is bounded by 1). *The infinite-horizon discounted return for any trajectory $\tau = (s_0, a_0, r_0, s_1, \dots)$ is bounded as $\sum_{t=0}^{\infty} \gamma^t r_t \leq 1$.*

This condition holds in sparse-reward environments, particularly those where an agent succeeds or fails at its task once per episode. This is common in domains such as robotics (Singh et al., 2020; Kalashnikov et al., 2018b) and games (i.e., winning or losing) (Bellemare et al., 2013). However, it is violated in dense reward tasks, such as MuJoCo locomotion benchmarks (Fu et al., 2020), where a reward is assigned at every time step. It is important to note that this assumption is only used in analysis to derive bounds that are sub-linear in horizon H ; our results can be generalized to environments where the returns are bounded by a function of horizon.

Notation. Let $n \wedge 1 = \max\{n, 1\}$ be shorthand. Denote $\iota = \text{polylog}(|\mathcal{S}|, H, N)$. For simplicity of analysis, we let ι change with context as done in Ren et al. (2021), so ι is a different polylogarithmic quantity each time it appears. For d -dimensional vectors \mathbf{x}, \mathbf{y} , we use $\mathbf{x}(i)$ to denote its i -th entry, and define $\mathbb{V}(\mathbf{x}, \mathbf{y}) = \sum_i \mathbf{x}(i) \mathbf{y}(i)^2 - (\sum_i \mathbf{x}(i) \mathbf{y}(i))^2$. If \mathbf{x} is a probability vector, i.e. $\mathbf{x}_i \geq 0$ and $\sum_i \mathbf{x}(i) = 1$, then we can write $\mathbb{V}(\mathbf{x}, \mathbf{y}) = \sum_i \mathbf{x}(i) (\mathbf{y}(i) - (\sum_{i'} \mathbf{x}(i') \mathbf{y}(i'))^2$.

4 THEORETICAL COMPARISON OF BC AND OFFLINE RL

In this section, we present performance guarantees for BC and offline RL, and characterize scenarios where offline RL algorithms will outperform BC. We first present general upper bounds for both algorithms in Section 4.1 under the conditions discussed in Section 3. Then, we compare the performance of BC and RL when provided with the same data generated by an expert in Section 4.2 and when RL is given noisy, suboptimal data in Section 4.3. Our goal is to characterize the conditions on the environment and offline dataset where RL can outperform BC.

4.1 GENERAL PERFORMANCE GUARANTEES OF BC AND OFFLINE RL

In this section, we analyze a generic algorithm for BC and offline RL, respectively. Each algorithm attains optimal scaling in suboptimality, and can be treated as theoretical representatives for their respective family of methods. For brevity, we consider a conservative offline RL algorithm (as defined in Section 2) in the main paper and defer analysis of a policy-constraint one to the Appendix. Both algorithms are detailed in Algorithms 1 and 2, respectively.

Guarantees for BC. For analysis purposes, we consider a BC algorithm that matches the empirical behavior policy on states in the offline dataset, and takes uniform random actions outside the support of the dataset. This BC algorithm was also used in prior work (Rajaraman et al., 2020), and is no worse than other schemes for acting at out-of-support states, in general. Denote the learned BC policy as $\hat{\pi}_\beta$, then $\forall \mathbf{s} \in \mathcal{D}, \hat{\pi}_\beta(\mathbf{a}|\mathbf{s}) \leftarrow n(\mathbf{s}, \mathbf{a})/n(\mathbf{s})$, and $\forall \mathbf{s} \notin \mathcal{D}, \hat{\pi}_\beta(\mathbf{a}|\mathbf{s}) \leftarrow 1/|\mathcal{A}|$. We adapt the results presented in Rajaraman et al. (2020) to the setting with Conditions 3.1 and 3.2. BC can only incur a non-zero asymptotic suboptimality (i.e., does not decrease to 0 as $N \rightarrow \infty$) in scenarios where $C^* = 1$, as it aims to match the data distribution $\mu(\mathbf{s}, \mathbf{a})$, and a non-expert dataset will inhibit the cloned policy from matching the expert π^* . The performance for BC is given in Theorem 4.1.

Theorem 4.1 (Performance of BC). *Under Conditions 3.1 and 3.2, the suboptimality of BC satisfies*

$$\text{SubOpt}(\hat{\pi}_\beta) \lesssim \frac{(C^* - 1)H}{2} + \frac{|\mathcal{S}|H\epsilon}{N}.$$

A proof of Theorem 4.1 is presented in Appendix B.1. The first term is the additional suboptimality incurred due to discrepancy between the behavior and optimal policies. The second term in this bound is derived by bounding the expected visitation frequency of the learned policy $\hat{\pi}_\beta$ onto states not observed in the dataset. The analysis is similar to that for existing bounds for imitation learning (Ross & Bagnell, 2010; Rajaraman et al., 2020). We achieve $\tilde{O}(H)$ suboptimality rather than $\tilde{O}(H^2)$ due to Condition 3.2, since the worst-case suboptimality of any trajectory is 1 rather than H .

Guarantees for conservative offline RL. We consider guarantees for a class of offline RL algorithms that maintain conservative value estimator such that the estimated value lower-bounds the true one, i.e., $\hat{V}^\pi \leq V^\pi$ for policy π . Existing offline RL algorithms achieve this by subtracting a penalty from reward either explicitly (Yu et al., 2020; Kidambi et al., 2020) or implicitly (Kumar et al., 2020). We only analyze one such algorithm that does the former, but we believe the algorithm can serve as a theoretical model for general conservative offline RL algorithms, where analyzing similar algorithms can be accomplished using the same outline. The algorithm we consider is similar in spirit to VI-LCB proposed by Rashidinejad et al. (2021) that subtracts penalty $b(\mathbf{s}, \mathbf{a})$ from the reward during value iteration; we consider a different penalty that results in a tighter bound. The estimated Q-values are obtained by iteratively solving the following Bellman backup: $\hat{Q}(\mathbf{s}, \mathbf{a}) \leftarrow \hat{r}(\mathbf{s}, \mathbf{a}) - b(\mathbf{s}, \mathbf{a}) + \sum_{\mathbf{s}'} \hat{P}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) \max_{\mathbf{a}'} \hat{Q}(\mathbf{s}', \mathbf{a}')$. The learned policy is then given by $\hat{\pi}^*(\mathbf{s}) \leftarrow \arg \max_{\mathbf{a}} \hat{Q}(\mathbf{s}, \mathbf{a})$. Building on work in online RL (Zhang et al., 2021), our specific $b(\mathbf{s}, \mathbf{a})$ is derived using Bernstein’s inequality, and is shown below:

$$b(\mathbf{s}, \mathbf{a}) \leftarrow \sqrt{\frac{\mathbb{V}(\hat{P}(\cdot|\mathbf{s}, \mathbf{a}), \hat{V})\epsilon}{(n(\mathbf{s}, \mathbf{a}) \wedge 1)}} + \sqrt{\frac{\hat{r}(\mathbf{s}, \mathbf{a})\epsilon}{(n(\mathbf{s}, \mathbf{a}) \wedge 1)}} + \frac{\epsilon}{(n(\mathbf{s}, \mathbf{a}) \wedge 1)}.$$

The performance of the learned policy $\hat{\pi}^*$ can then be bounded as:

Theorem 4.2 (Performance of our generic conservative offline RL algorithm). *Under Conditions 3.1 and 3.2, the policy $\hat{\pi}^*$ found by the generic conservative offline RL algorithm satisfies*

$$\text{SubOpt}(\hat{\pi}^*) \lesssim \sqrt{\frac{C^*|\mathcal{S}|H\epsilon}{N}} + \frac{C^*|\mathcal{S}|H\epsilon}{N}.$$

We defer a proof for Theorem 4.2 to Appendix B.2. On a high level, we first show that our algorithm is always conservative, *i.e.*, $\forall s, \hat{V}(s) \leq V^{\pi^*}(s)$, and then bound the total suboptimality incurred as a result of being conservative. Our bound in Theorem 4.2 improves on existing bounds in two ways: (1) by considering pessimistic value estimates, we are able to remove the strict coverage assumptions used by Ren et al. (2021), and (2) we eliminate an additional $|\mathcal{S}|$ factor by introducing s -absorbing MDPs for each state as done in Agarwal et al. (2020a). In addition, compared to the pessimistic algorithm in Rashidinejad et al. (2021), our bound enjoys better scaling in H .

4.2 COMPARISON UNDER EXPERT DATA

We first compare the performance bounds from Section 4.1 when the offline dataset is generated from expert demonstrations. In relation to Condition 3.1, this corresponds to small C^* . Specifically, we consider $C^* \in [1, 1 + \tilde{O}(1/N)]$ so that the suboptimality of BC in Theorem 4.1 scales as $\tilde{O}(|\mathcal{S}|H/N)$. In this regime, we perform a nuanced comparison by analyzing specific scenarios where RL may outperform BC. We consider the case of $C^* = 1$ and $C^* = 1 + \tilde{O}(1/N)$ separately.

What happens when $C^* = 1$? In this case, we derive an information-theoretic lower-bound of $|\mathcal{S}|H/N$ for any offline algorithm. Our result in Theorem 4.3 utilizes the analysis of Rajaraman et al. (2020) but additionally factoring in Condition 3.2.

Theorem 4.3 (Information-theoretic lower-bound for offline learning with $C^* = 1$). *For any learner $\hat{\pi}$, there exists an MDP \mathcal{M} satisfying Assumption 3.2, and a deterministic expert π^* , such that the expected suboptimality of the learner is lower-bounded:*

$$\sup_{\mathcal{M}, \pi^*} \text{SubOpt}(\hat{\pi}) \gtrsim \frac{|\mathcal{S}|H}{N}$$

The proof of Theorem 4.3 uses the same hard instance from Theorem 6.1 of Rajaraman et al. (2020), except that one factor of H is dropped due to Condition 3.2. The other factor of H arises from the performance difference lemma and is retained. In this case, where BC achieves the lower bound up to logarithmic factors, we argue that we cannot improve over BC. This is because the suboptimality of BC is entirely due to encountering states that do not appear in the dataset; without additional assumptions on the ability to generalize to unseen states, offline RL must incur the same suboptimality, as both methods would choose actions uniformly at random. Hence, Theorem 4.3 shows the negative result that no algorithm can outperform BC when $C^* = 1$ exactly.

However, we argue that even with expert demonstrations as data, $C^* = 1$ is an unrealistic assumption. Naively, it seems plausible that the expert who collected the dataset did not perform optimally at every transition; this is often true for humans, or stochastic experts as ϵ -greedy or maximum-entropy policies. In addition, a scenario where $C^* > 1$ when the expert always behaves optimally is under distribution shift of the environment. One practical example of this is when the initial state distribution changes between dataset collection and evaluation (*e.g.*, in robotics (Singh et al., 2020), or self-driving (Bojarski et al., 2016)). Since the normalized state-action marginals $d^*(s, a), \mu(s, a)$ are impacted by $\rho(s)$, this would lead to $C^* > 1$ even when the expert behaves exactly as π^* .

What happens when $C^* = 1 + \tilde{O}(1/N)$? Here C^* is small enough that BC still achieves the same optimal $\tilde{O}(|\mathcal{S}|H/N)$ performance guarantee. However, there is suboptimality incurred by BC for even states that appear in the dataset due to distribution shift, which allows us to argue about structural properties of MDPs that allow offline RL to perform better across those states, particularly for problems with large effective horizon H . We motivate one such structure below.

In several practical RL problems, only a small fraction of states that appear in a trajectory require precise action selection. This means that the return of any trajectory can mostly be explained by the actions taken in those states, which we call *critical states*. This can occur when there exist a large proportion of states where it is not costly to recover after deviating from the optimal trajectory, or when there exist a large volume of optimal trajectories. For instance, in a robotic grasping task, if the robot is not close to the object, the robot can take many possible actions and still pick up the object in the end (Kalashnikov et al., 2018a). Similarly, for navigation problems where there exist wide, unobstructed areas, many trajectories exist to traverse those areas (Savva et al., 2019). In contrast, in environments such as cliffwalk, every state is critical, as an incorrect action at a given state will cause the agent to fall off the cliff (Schaal et al., 2015). Formally, critical states are defined as:

Definition 4.1 (Critical states). A state \mathbf{s} is said to be non-critical, i.e., $\mathbf{s} \in \mathcal{S} \setminus \mathcal{C}$ if the advantage of any action $\mathbf{a} \in \mathcal{A}$ is close to 0 under the optimal policy, i.e., $|\max_{\mathbf{a}'} Q^*(\mathbf{s}, \mathbf{a}') - Q^*(\mathbf{s}, \mathbf{a})| \leq \varepsilon/H$.

Condition 4.1 (Volume of critical states is small.). $\exists p_c \in (0, 1)$ that satisfies: $|\mathcal{C}| \leq p_c |\mathcal{S}|$.

We can show that, if the MDP satisfies having a small fraction or volume of critical states, then conservative offline RL enjoys stronger guarantees than BC.

Corollary 4.1 (Performance of conservative offline RL with critical states). *Under Conditions 4.1, 3.1 and 3.2, the policy $\hat{\pi}^*$ found by the conservative offline RL satisfies*

$$\text{SubOpt}(\hat{\pi}^*) \leq \sqrt{\frac{p_c C^* |\mathcal{S}| H \iota}{N}} + \frac{p_c C^* |\mathcal{S}| H \iota}{N} + \varepsilon.$$

For $\varepsilon = \mathcal{O}(\sqrt{H})$, if the environment satisfies $p_c = \mathcal{O}(1/\sqrt{H})$, meaning we encounter $\mathcal{O}(\sqrt{H})$ critical states on average in any trajectory, then we achieve better scaling in H with offline RL than BC. Note that BC does not enjoy the benefit of few critical states because it is agnostic to the reward of the environment and therefore limited by the $\mathcal{O}(H)$ suboptimality of the behavior policy.

4.3 COMPARISON UNDER NOISY DATA

In practice, it is often much more tractable to obtain suboptimal demonstrations rather than expert ones. From Theorem 4.1, we see that for $C^* = 1 + \Omega(1/\sqrt{N})$, BC will incur suboptimality that is worse asymptotically than offline RL. In contrast, from Theorem 4.2, we note that offline RL does not scale nearly as poorly with increasing C^* . Since offline RL is not as reliant on the performance of the behavior policy due to explicitly modeling the rewards, we hypothesize that RL can actually benefit from suboptimal but well-explored data. In this section, we aim to answer the following question: *Can offline RL with suboptimal data outperform BC with an equal amount of expert data?*

We show in Corollary 4.2 that if the suboptimal dataset \mathcal{D} satisfies an additional coverage condition, then running conservative offline RL will attain a $\tilde{\mathcal{O}}(\sqrt{H})$ suboptimality in the horizon. This implies, perhaps surprisingly, that for long horizon tasks offline RL using noisy \mathcal{D} can outperform even BC on expert data for this task. This conclusion has important consequences in practice when learning from demonstrations – it indicates that in tasks with large horizon H , and limited capability to collect demonstrations, it is advisable to run offline RL on noisy, suboptimal data. Suboptimal data is generally cheap to collect in domains such as robotics (Kalashnikov et al., 2018b; Singh et al., 2020; Mandlekar et al., 2021) by running noisy scripted policies, or computer systems (Liu et al., 2020a; Trofin et al., 2021) by running heuristics; such data also be automatically generated via data-augmentation, as was done in self-driving (Bojarski et al., 2016) and robotics (Rao et al., 2020).

Our coverage condition is that the data distribution sufficiently covers the optimal policy distribution.

Condition 4.2 (Coverage of the optimal policy). $\exists b \in [\log H/N, 1)$ such that μ satisfies: $\forall (\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ where $d^*(\mathbf{s}, \mathbf{a}) \geq b/H$, we have $\mu(\mathbf{s}, \mathbf{a}) \geq b$.

Intuitively, this means that the data distribution puts sufficient mass on states that have non-negligible but low density in the optimal policy distribution. Note that this is a weaker condition than prior works that require (1) full coverage of the state-action space, and (2) enforce a constraint on the empirical state-action visitations $\hat{\mu}(\mathbf{s}, \mathbf{a})$ instead of the data distribution $\mu(\mathbf{s}, \mathbf{a})$ (Ren et al., 2021; Zhang et al., 2021). This condition is reasonable when the dataset is large and is collected by a behavior policy that performs some exploration on the state space, e.g., ϵ -greedy or a maximum-entropy expert policies.

Corollary 4.2 (Performance of conservative offline RL with noisy data). *If μ satisfies Condition 4.2, and under Conditions 3.1 and 3.2, the policy $\hat{\pi}^*$ found by conservative offline RL satisfies:*

$$\text{SubOpt}(\hat{\pi}^*) \lesssim \sqrt{\frac{H \iota}{bN}} + \frac{H \iota}{bN} + \sqrt{b \iota} + \frac{C^* |\mathcal{S}| \iota}{N}.$$

If $b = \mathcal{O}(H/N)$, then the bound in Corollary 4.2 has $\tilde{\mathcal{O}}(\sqrt{H})$ scaling, rather than $\tilde{\mathcal{O}}(H)$ for BC from the lower-bound in Theorem 4.3. Thus, when the data satisfies mild coverage conditions, offline RL performs better in long-horizon tasks compared to even BC with the same amount of expert data.

4.4 COMPARISON OF GENERALIZED BC METHODS AND OFFLINE RL

So far we have studied scenarios where offline RL can outperform naive BC. One might now wonder how offline RL methods perform relative to generalized BC methods that additionally use reward information to inform learning. We study two such approaches: (1) filtered BC (Chen et al., 2021), which only fits to the top k -percentage of trajectories in \mathcal{D} , measured by the total reward, and (2) BC with one-step policy improvement (Brandfonbrener et al., 2021), which fits a Q-function for the behavior policy, then uses the values to perform one-step of policy improvement over the behavior policy. In this section, we aim to answer how these methods perform relative to RL.

Filtered BC. In expectation, this algorithm uses αN samples of the offline dataset \mathcal{D} for $\alpha \in [0, 1]$ to perform BC on. This means that the upper bound (Theorem 4.1) will have worse scaling in N . For $C^* = 1$, this leads to a strictly worse bound than regular BC. However, for suboptimal data, the filtering step could decrease C^* by filtering out suboptimal trajectories, allowing filtered BC to outperform traditional BC. Nevertheless, from our analysis in Section 4.3, offline RL is still preferred to filtered BC because RL can leverage the noisy data and potentially achieve sub-linear suboptimality $O(\sqrt{H})$ in the horizon, whereas even filtered BC would always incur $O(H)$ suboptimality.

BC with policy improvement. This algorithm utilizes the entire dataset to estimate the Q-value of the behavior policy, $\hat{Q}^{\hat{\pi}^\beta}$, and performs one step of policy improvement using the estimated Q-function, typically via an advantage-weighted update: $\hat{\pi}^1(\mathbf{a}|\mathbf{s}) = \hat{\pi}^\beta(\mathbf{a}|\mathbf{s}) \exp(\eta H \hat{A}^{\hat{\pi}^\beta}(\mathbf{s}, \mathbf{a})) / \mathbb{Z}_1(\mathbf{s})$. *When would this algorithm perform poorly compared to offline RL?* Intuitively, this would happen when multiple steps of policy improvement are needed to effectively discover high-advantage actions under the behavior policy. This is the case when the behavior policy puts low density on high-advantage transitions. In Theorem 4.4, we show that more than one step of policy improvement can improve the policy under Condition 4.2 for the softmax policy parameterization (Agarwal et al., 2021a).

Theorem 4.4 (One-step is worse than k -step policy improvement). *Assume that the learned policies are represented via a softmax parameterization (Equation 3, Agarwal et al. (2021a)). Let $\hat{\pi}^k$ denote the policy obtained after k -steps of policy improvement using exponentiated advantage weights. Then, under Condition 4.2, the performance difference between $\hat{\pi}^k$ and $\hat{\pi}^1$ is lower-bounded by:*

$$J(\hat{\pi}^k) - J(\hat{\pi}^1) \gtrsim \frac{k}{H\eta} \mathbb{E}_{\mathbf{s} \sim \mu} \left[\frac{1}{k} \sum_{t=1}^k \log \mathbb{Z}_t(\mathbf{s}) \right] - \sqrt{\frac{C^* H \epsilon}{N}}.$$

A proof of Theorem 4.4 is provided in Appendix B.5. This result implies that when the average exponentiated empirical advantage $1/k \sum_{i=1}^k \log \mathbb{Z}_t(\mathbf{s})$ is large enough (i.e., $\geq c_0$ for some universal constant), which is usually the case when the behavior policy is highly suboptimal, then for $k = O(H)$, multiple steps of policy improvement will improve performance, i.e., $J(\hat{\pi}^k) - J(\hat{\pi}^1) = \hat{O}(H - \sqrt{H/N})$, where the gap increases with a longer horizon. This is typically the case when the structure of the MDP allow for stitching parts of poor-performing trajectories. One example is in navigation, where trajectories that fail may still contain segments of a successful trajectory.

5 EMPIRICAL EVALUATION OF BC AND OFFLINE RL

Having characterized scenarios where offline RL methods can outperform BC in theory, we now validate our results empirically. Concretely, we aim to answer the following questions: (1) Does offline RL trained on expert data outperform BC on expert data in practice? (2) Does offline RL trained on suboptimal, noisy data outperform BC on expert data?, and (3) How does full offline RL compare to the generalized BC methods studied in Section 4.4? We will first validate our findings on a tabular gridworld domain, and then on several high-dimensional offline RL problems.

Diagnostic experiments in gridworld. We first evaluate tabular versions of the BC and offline RL algorithms analyzed in Section 4.1 on sparse-reward 10×10 gridworlds environments (Fu et al., 2019). Complete details about the setup can be found in Appendix D.1. On a high-level, we consider three different environments, each with varying number of critical states, from “Single Critical” with exactly one, to “Cliffwalk” where every state is critical and veering off yields zero reward. The baselines we consider are: naive BC (BC), conservative RL (RL-C), policy-constraint RL (RL-PC), and generalized BC with one-step and k -step policy improvement (BC-PI, BC-kPI).

In the left plot of Figure 1, we show the return (normalized by return of the optimal policy) across all the different environments for optimal data ($C^* = 1$) and data generated from the optimal policy

but with a different initial state distribution ($C^* > 1$ but $\pi_\beta(\cdot|s) = \pi^*(\cdot|s)$). As expected from our discussion in Section 4.2, BC performs best under $C^* = 1$, but RL performs much better when $C^* > 1$; also BC with one-step policy improvement outperforms naive BC for $C^* > 1$, but does not beat RL. In Figure 1 (right), we vary C^* by interpolating the dataset with one generated by a random policy, where α is the proportion of random data. RL performs much better over all BC methods, when the data supporting our analysis in Section 4.3. Finally, BC with multiple policy improvement performs better than one step when the data is noisy, which validates Theorem 4.4.

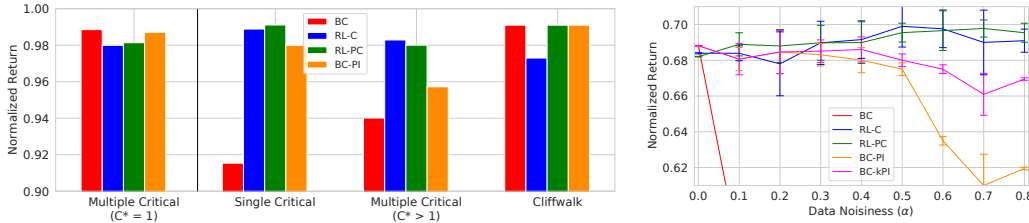


Figure 1: Offline RL vs BC on gridworld domains. *Left*: We compare offline RL and BC algorithms on three different gridworlds with varying number of critical points for expert and near-expert data. *Right*: Taking the “Multiple Critical” domain, we examine the effect of increasing the noisiness of the dataset by interpolating it with one generated by a random policy, and show that RL improves drastically with increased noise over BC.

Evaluation in high-dimensional tasks. Next, we turn to deep offline RL. We consider a diverse set of domains (shown on the right) and behavior policies that are representative of scenarios where we would decide between offline RL and BC: multi-stage robotic manipulation tasks from state (Adroit domains from Fu et al. (2020)) and image observations (Singh et al., 2020), antmaze navigation (Fu et al., 2020), and 7 Atari games (Agarwal et al., 2020b). We use the scripted expert provided by Fu et al. (2020) for antmaze and the one provided by Singh et al. (2020) for manipulation, an RL-trained expert for Atari, and human expert for Adroit (Rajeswaran et al., 2018). In scenarios where suboptimal data is used to train offline RL, we use failed attempts to solve the task from a noisy expert policy (*i.e.*, previous policies in the replay buffer for Atari, and noisy scripted experts for antmaze and manipulation). All these tasks utilize sparse rewards such that the return of any trajectory is bounded by a constant much smaller than the horizon. We use CQL (Kumar et al., 2020) as the base offline RL method, and utilize the approach by Brandfonbrener et al. (2021) as a representative BC-PI method.



Tuning offline RL and BC. Naïvely running offline RL can lead to poor performance, as noted by prior works (Mandlekar et al., 2021; Florence et al., 2021). This is also true for BC, but, some solutions such as early stopping based on validation losses, can help improve performance. We claim that a similar tuning strategy is also crucial for offline RL. In our experiments we utilize the offline workflow proposed by Kumar et al. (2021c) to perform policy selection, and address overfitting and underfitting, purely offline. When the Q-values learned by CQL are extremely negative (typically on the Adroit domains), we utilize dropout with probability 0.4 on the layers of the Q-function to combat overfitting. On the other hand, when the Q-values exhibit a relatively stable trend (*e.g.*, in Antmaze or Atari), we utilize the DR3 regularizer (Kumar et al., 2021a) to increase capacity. Consistent with prior work, we find that naïve offline RL generally performs worse than BC without offline tuning, but we find that *offline-tuned* offline RL generally outperforms BC. To make a stronger comparison, we tuned BC using the online rollouts. We applied regularizers such as dropout on the BC policy to prevent overfitting in Adroit, and utilized a larger ResNet (He et al., 2016) architecture for the robotic manipulation tasks and Atari domains. For BC, we report the performance of the *best* checkpoint found during training, giving BC an unfair advantage, but we still find that *offline-tuned* offline RL performs better. More details about tuning can be found in Appendix E.

Answers to questions (1) to (3). For (1), we run CQL and BC on expert data in each task, and present the comparison in Table 1 and Figure 2. Observe that while naïve CQL performs comparable or worse than BC in this case, after tuning on offline data following the procedure proposed by Kumar et al. (2021c), CQL outperforms BC. This tuning does not require any additional online rollouts, and as discussed before, and corrects for underfitting and overfitting completely offline. Note that while BC performs better or comparable to RL for antmaze (large) with expert data, it performs worse than RL when the data admits a more diverse initial state distribution such that $C^* \neq 1$, even though the behavior policy matches the expert.

Domain / Behavior Policy	Task/Data Quality	BC	Naïve CQL	Tuned CQL
AntMaze (scripted)	Medium, Expert	53.2%±8.7%	20.8% ± 1.0%	55.9% ± 3.2%
	Large, Expert	4.83% ±0.8%	0.0% ± 0.0%	0.0% ± 0.0%
	Medium, Expert w/ diverse initial	55.2%±6.7%	19.0% ± 5.2%	67.0% ± 7.3%
	Large, Expert w/ diverse initial	1.3%±0.5%	0.0% ± 0.0	5.1% ± 6.9%
Manipulation (scripted)	pick-place-open-grasp, Expert	14.5%±1.8%	12.3%±5.3%	23.5% ±6.0%
	close-open-grasp, Expert	17.4%±3.1%	20.0%±6.0%	49.7% ±5.4%
	open-grasp, Expert	33.2%±8.1%	22.8%±5.3%	51.9% ±6.8%
Adroit (Human)	hammer-human-v1, Expert	71.0% ± 9.3%	62.5% ± 39.0%	78.1% ± 6.7%
	door-human-v1, Expert	86.3% ± 6.5%	70.3% ± 27.2%	79.1% ± 4.7%
	pen-human-v1, Expert	73.0% ± 9.1%	64.0% ± 6.9%	74.1% ± 6.1%
	relocate-human-v1, Expert	0.0% ± 0.0%	0.0% ± 0.0%	0.0% ± 0.0%

Table 1: Offline RL vs. BC with expert dataset compositions averaged over 3 seeds. While naïve offline RL often performs comparable or worse than BC, the performance of offline RL improves drastically after offline tuning. Also note that offline RL can improve when provided with diverse initial states in the Antmaze domain. Additionally, note that offline-tuned offline RL outperforms BC significantly in the manipulation domains.

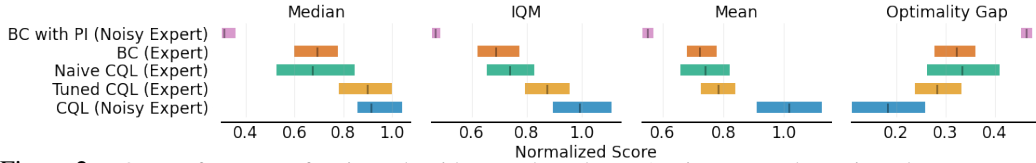


Figure 2: IQM performance of various algorithms evaluated on 7 Atari games under various dataset compositions (per game scores in Table 3). Note that offline-tuned CQL with expert data (“Tuned CQL”) outperforms cloning the expert data (“BC (Expert)”), even though naïve CQL is comparable to BC in this setting. When CQL is provided with noisy-expert data, it significantly outperforms cloning the expert policy.

For (2), we compare offline RL trained on noisy-expert data with BC trained on an equal amount of expert data, on domains where noisy-expert data is easy to generate: specifically (a) manipulation domains (Table 2) and (b) Atari games (Figure 2). Observe that CQL attains better performance compared to BC and also improves over only using expert data. The performance difference also increases with H , *i.e.*, open-grasp ($H = 40$) vs pick-place-open-grasp ($H = 80$) vs Atari domains ($H = 27000$). This validates that offline RL with noisy-expert data can outperform BC with expert data, particularly on long-horizon tasks.

Task	BC (Expert)	CQL (Noisy Expert)
pick-place-open-grasp	14.5% ± 1.8%	85.7% ± 3.1%
close-open-grasp	17.4% ± 3.1%	90.3% ± 2.3%
open-grasp	33.2% ± 8.1%	92.4% ± 4.9%

Table 2: CQL with noisy-expert data vs BC with expert data with equal dataset size on manipulation tasks. CQL outperforms BC as well as CQL with only expert data.

Finally, for (3), we compare CQL to a representative BC-PI method (Brandfonbrener et al., 2021) trained using noisy-expert data on Atari domains, which present multiple stitching opportunities. The BC-PI method estimates the Q-function of the behavior policy using SARSA and then performs one-step of policy improvement. The results in Figure 2 support what is predicted by our theoretical results, *i.e.*, BC-PI still performs significantly worse than CQL with noisy-expert data, even though we utilized online rollouts for tuning BC-PI and report the best hyperparameters found.

6 DISCUSSION

In this paper, we sought to understand when offline RL methods are preferable over BC ones, when both are provided with optimal or near-optimal demonstration data. While in the worst case, both methods attain similar performance on expert data, additional structural assumptions on the environment can provide offline RL with a significant advantage. Perhaps surprisingly, we also show that running RL on noisy-expert, suboptimal data attains more favorable guarantees compared to running BC on expert data for the same task, using equal amounts of data. Empirically, we observe that tuned offline RL algorithms can outperform BC on various practical problems domains, with different expert policy distributions. While our theoretical analysis identifies several cases where offline RL can perform better than BC, and our empirical results support the conclusions obtained, there is still plenty of room for further investigation. Our theoretical analysis can be improved to handle function approximation, which would allow us to analyze modern deep offline RL methods. Our empirical results suggest that better tuning strategies for offline RL will be crucial for good performance, which is also a promising avenue for future work.

REPRODUCIBILITY STATEMENT

For theoretical results, we provide explanations of all the assumptions and a complete proof of the claims in Appendix B. We provide complete experimental details regarding the tasks and tuning strategy in the Appendix E. Additionally, we follow the recommendations of Agarwal et al. (2021b) for reliable evaluation in deep RL and report the statistical uncertainty in reported results including aggregate performance metrics (Figure 2). We provide individual scores in Table 3 and we will open-source our code.

REFERENCES

- Alekh Agarwal, Sham Kakade, and Lin F Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pp. 67–83. PMLR, 2020a.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021a.
- Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2020b.
- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron Courville, and Marc G Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Information Processing Systems*, 2021b.
- Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Int. Res.*, 47(1):253–279, May 2013. ISSN 1076-9757.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- David Brandfonbrener, William F Whitney, Rajesh Ranganath, and Joan Bruna. Offline rl without off-policy evaluation. *arXiv preprint arXiv:2106.08909*, 2021.
- Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. *ICML*, 2019.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *arXiv preprint arXiv:2106.01345*, 2021.
- Lasse Espeholt, Hubert Soyer, Rémi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA: scalable distributed deep-rl with importance weighted actor-learner architectures. In *International Conference on Machine Learning*, 2018.
- Amir-massoud Farahmand, Csaba Szepesvári, and Rémi Munos. Error propagation for approximate policy and value iteration. In *Advances in Neural Information Processing Systems (NIPS)*, 2010.
- Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In *5th Annual Conference on Robot Learning*, 2021. URL <https://openreview.net/forum?id=rif3a5NAXU6>.
- Justin Fu, Aviral Kumar, Matthew Soh, and Sergey Levine. Diagnosing bottlenecks in deep Q-learning algorithms. *arXiv preprint arXiv:1902.10250*, 2019.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning, 2020.

- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *arXiv preprint arXiv:2106.06860*, 2021.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. *arXiv preprint arXiv:1812.02900*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Natasha Jaques, Judy Hanwen Shen, Asma Ghandeharioun, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Shane Gu, and Rosalind Picard. Human-centric dialog training via offline reinforcement learning. *arXiv preprint arXiv:2010.05848*, 2020.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? *arXiv preprint arXiv:2012.15085*, 2020.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pp. 651–673, 2018a.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pp. 651–673. PMLR, 2018b.
- Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski, Chelsea Finn, Sergey Levine, and Karol Hausman. Mt-opt: Continuous multi-task robotic reinforcement learning at scale. *arXiv preprint arXiv:2104.08212*, 2021.
- Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- Ilya Kostrikov, Jonathan Tompson, Rob Fergus, and Ofir Nachum. Offline reinforcement learning with fisher divergence critic regularization. *arXiv preprint arXiv:2103.08050*, 2021.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing Systems*, pp. 11761–11771, 2019.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *arXiv preprint arXiv:2006.04779*, 2020.
- Aviral Kumar, Rishabh Agarwal, Aaron Courville, Tengyu Ma, George Tucker, and Sergey Levine. Value-based deep reinforcement learning requires explicit regularization. In *RL for Real Life Workshop & Overparameterization: Pitfalls and Opportunities Workshop, ICML*, 2021a. URL https://drive.google.com/file/d/1Fg43H5oagQp-ksjpWBf_aDYEzAFMVJm6/view.
- Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-parameterization inhibits data-efficient deep reinforcement learning. In *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=O9bnihsFfxU>.
- Aviral Kumar, Anikait Singh, Stephen Tian, Chelsea Finn, and Sergey Levine. A workflow for offline model-free robotic reinforcement learning. In *5th Annual Conference on Robot Learning*, 2021c. URL <https://openreview.net/forum?id=fy4ZBWxYbIo>.
- Sascha Lange, Thomas Gabel, and Martin A. Riedmiller. Batch reinforcement learning. In *Reinforcement Learning*, volume 12. Springer, 2012.
- Romain Laroché, Paul Trichelair, and Remi Tachet Des Combes. Safe policy improvement with baseline bootstrapping. In *International Conference on Machine Learning*, pp. 3652–3661. PMLR, 2019.

- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Evan Liu, Milad Hashemi, Kevin Swersky, Parthasarathy Ranganathan, and Junwhan Ahn. An imitation learning approach for cache replacement. In *International Conference on Machine Learning*, pp. 6237–6247. PMLR, 2020a.
- Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient with state distribution correction. *CoRR*, abs/1904.08473, 2019.
- Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch reinforcement learning without great exploration. *arXiv preprint arXiv:2007.08202*, 2020b.
- Ajay Mandlekar, Fabio Ramos, Byron Boots, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Dieter Fox. Iris: Implicit reinforcement without interaction at scale for learning control from offline robot manipulation data. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4414–4420. IEEE, 2020.
- Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Fei-Fei Li, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *5th Annual Conference on Robot Learning*, 2021. URL <https://openreview.net/forum?id=JrsfBJtDFdI>.
- Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization. *arxiv preprint arxiv:0907.3740*, 2009.
- Rémi Munos. Error bounds for approximate policy iteration. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ICML’03, pp. 560–567. AAAI Press, 2003. ISBN 1577351894.
- Rémi Munos. Error bounds for approximate value iteration. In *AAAI Conference on Artificial intelligence (AAAI)*, pp. 1006–1011. AAAI Press, 2005.
- Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- Kimia Nadjahi, Romain Laroche, and Rémi Tachet des Combes. Safe policy improvement with soft baseline bootstrapping. *arXiv preprint arXiv:1907.05079*, 2019.
- Ashvin Nair, Murtaza Dalal, Abhishek Gupta, and Sergey Levine. Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.
- Nived Rajaraman, Lin F Yang, Jiantao Jiao, and Kannan Ramachandran. Toward the fundamental limits of imitation learning. *arXiv preprint arXiv:2009.05990*, 2020.
- Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In *Robotics: Science and Systems*, 2018.
- Kanishka Rao, Chris Harris, Alex Irpan, Sergey Levine, Julian Ibarz, and Mohi Khansari. RL-cyclelegan: Reinforcement learning aware simulation-to-real. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11157–11166, 2020.
- Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *arXiv preprint arXiv:2103.12021*, 2021.
- Tongzheng Ren, Jialian Li, Bo Dai, Simon S Du, and Sujay Sanghavi. Nearly horizon-free offline reinforcement learning. *arXiv preprint arXiv:2103.14077*, 2021.

- Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 661–668, 2010.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A Platform for Embodied AI Research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *International Conference on Learning Representations (ICLR)*, 2015.
- Susan M Shortreed, Eric Laber, Daniel J Lizotte, T Scott Stroup, Joelle Pineau, and Susan A Murphy. Informing sequential clinical decision-making through reinforcement learning: an empirical study. *Machine learning*, 84(1-2):109–136, 2011.
- Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael Neunert, Thomas Lampe, Roland Hafner, and Martin Riedmiller. Keep doing what worked: Behavioral modelling priors for offline reinforcement learning. *arXiv preprint arXiv:2002.08396*, 2020.
- Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. Cog: Connecting new skills to past experience with offline reinforcement learning. *arXiv preprint arXiv:2010.14500*, 2020.
- Samarth Sinha, Ajay Mandlekar, and Animesh Garg. S4RL: Surprisingly simple self-supervision for offline reinforcement learning in robotics. In *5th Annual Conference on Robot Learning*, 2021. URL https://openreview.net/forum?id=8xC5NNej-1_.
- Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback through counterfactual risk minimization. *J. Mach. Learn. Res.*, 16:1731–1755, 2015.
- Mircea Trofin, Yundi Qian, Eugene Brevdo, Zinan Lin, Krzysztof Choromanski, and David Li. Mlgo: a machine learning guided compiler optimizations framework. *arXiv preprint arXiv:2101.04808*, 2021.
- L. Wang, Wei Zhang, Xiaofeng He, and H. Zha. Supervised reinforcement learning with recurrent neural network for dynamic treatment recommendation. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018.
- Ruosong Wang, Dean Foster, and Sham M. Kakade. What are the statistical limits of offline $\{rl\}$ with linear function approximation? In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=30EvkP2aQLD>.
- Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Tengyang Xie and Nan Jiang. Q^* approximation schemes for batch reinforcement learning: A eoretical comparison. 2020.
- Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. *arXiv preprint arXiv:2106.06926*, 2021.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *arXiv preprint arXiv:2005.13239*, 2020.
- Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. Combo: Conservative offline model-based policy optimization. *arXiv preprint arXiv:2102.08363*, 2021.
- Andrea Zanette. Exponential lower bounds for batch reinforcement learning: Batch rl can be exponentially harder than online rl. *arXiv preprint arXiv:2012.08005*, 2020.
- Yufeng Zhang, Qi Cai, Zhuoran Yang, Yongxin Chen, and Zhaoran Wang. Can temporal-difference and q-learning learn representation? a mean-field theory. *arXiv preprint arXiv:2006.04761*, 2020.

Zihan Zhang, Xiangyang Ji, and Simon Du. Is reinforcement learning more difficult than bandits? a near-optimal algorithm escaping the curse of horizon. In *Conference on Learning Theory*, pp. 4528–4531. PMLR, 2021.