

---

# Effective Offline RL Needs Going Beyond Pessimism: Representations and Distributional Shift

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Standard off-policy reinforcement learning (RL) methods based on temporal dif-  
2 ference (TD) learning generally fail to learn good policies when applied to static  
3 offline datasets. Conventionally, this is attributed to distribution shift, where the  
4 Bellman backup queries high-value out-of-distribution (OOD) actions for the next  
5 time step, which then leads to systematic overestimation. However, this expla-  
6 nation is incomplete, as conservative offline RL methods that directly address  
7 overestimation still suffer from stability problems in practice. This suggests that  
8 although OOD actions may account for part of the challenge, the difficulties with  
9 TD learning in the offline setting are also deeply connected to other aspects such  
10 as the quality of representations of learned function approximators. In this work,  
11 we demonstrate that merely imposing pessimism is not sufficient for good per-  
12 formance, and demonstrate empirically that regularizing representations actually  
13 accounts for a large part of the improvement observed in modern offline RL meth-  
14 ods. Building on this insight, we identify concrete metrics that enable effective  
15 diagnosis of the quality of the learned representation, and are able to adequately  
16 predict performance of the underlying method. Finally, we show that a simple  
17 approach for handling representations, without any changing any other aspect of  
18 conservative offline RL algorithms, can lead to better performance in several offline  
19 RL problems.

## 20 1 Introduction

21 Offline reinforcement learning (RL), combined with powerful deep net function approximators,  
22 has the potential for solving decision-making tasks where online interaction is either expensive  
23 or unsafe, circumventing a major barrier to the deployment of RL in the real-world. Temporal  
24 difference (TD) learning methods, such as Q-learning, provide a natural framework for building  
25 offline RL algorithms [30], fitting a parametric value function by sequentially regressing to targets  
26 generated from its own previous snapshot using only offline data. However, directly applying TD  
27 to a static offline dataset often fails to learn effective policies, as the maximization in the target  
28 value computation will find erroneously high-valued out-of-distribution (OOD) actions, resulting  
29 in systematic overestimation. A variety of offline RL methods, such as those that apply value  
30 conservatism [26, 58] or behavioral constraint [14, 24, 53, 13, 18, 23, 22], have been proposed to  
31 address this issue with OOD actions in TD learning by inducing some form of pessimism. While  
32 all these methods lead to promising improvement in performance on offline RL tasks, determining  
33 why one method for addressing the OOD actions issue is better than another has proven challenging,  
34 which in turn makes it difficult to develop insights and guidelines for designing better offline RL  
35 algorithms. In fact, in theory, majority of these approaches essentially optimize the very same RL  
36 objective subject to a divergence constraint against the behavior policy that generates the data, and

37 would, behave identically in a tabular problem setting. Hence, a natural question to ask is: does the  
38 improvement observed from these methods really stem from their ability to induce pessimism?

39 In this paper, we will show that a significant part of the benefit of offline RL approaches that aim to  
40 address OOD actions actually stems from the effect they have on the learned representations, rather  
41 than merely from their ability to avoiding overestimation. We first show that the even if we can prevent  
42 the value of the learned Q-function at OOD actions from being overestimated, training Q-functions  
43 against Bellman targets computed using OOD actions still induces Q-function representations that  
44 give rise to poor policy performance, which indicates that overestimation is not sufficient to explain  
45 poor performance in offline RL. Second, we empirically demonstrate that an offline RL method  
46 that does not apply any pessimism, but only regularizes the representation learned for the dataset  
47 and OOD actions to be different using adversarial training, can actually perform quite well. The  
48 method we develop resembles the conservative Q-learning (CQL) [26] approach, but crucially only  
49 regularizes the representations and not the final Q-values. Our analysis shows that this approach  
50 recovers 68% of the performance of CQL, indicating that the performance of CQL, in large part,  
51 comes from the implicit regularization obtained by penalizing OOD actions.

52 Based on this analysis, we propose a metric that evaluates the quality of the representation learned by  
53 offline RL methods based on the ability to accurately reconstruct the dataset actions from the learned  
54 representation. We demonstrate that comparing this reconstruction error to a dynamic programming  
55 approach that does not utilize OOD actions gives us a good measure of representational quality, that  
56 is predictive of performance. Finally, we discover that good representations can actually be obtained  
57 by a surprisingly simple method: interpolating between TD and supervised learning via an ensemble  
58 of N-step returns, similar to  $TD(\lambda)$ . We not only find that utilizing an ensemble of N-step returns  
59 approach attains better performance, but, more interestingly, we argue that this *cannot* be attributed  
60 to standard explanations of a better bias-variance tradeoff.

61 Our main contributions are to demonstrate, via an extensive empirical study, that merely addressing  
62 the OOD action issue in offline RL via pessimism is not sufficient for TD-based offline RL methods,  
63 and that the quality of learned representation is crucial for good performance. Our analysis provides  
64 guidance on how to measure representational quality, and shows how simple methods such as an  
65 ensemble of N-step returns already attain better performance on benchmark tasks from D4RL [12] as  
66 a result of improved representational quality. We hope that our analysis provides concrete takeaways  
67 for researchers in offline RL and highlights a largely overlooked line of challenges beyond behavior  
68 regularization that is crucial in devising more effective and reliable offline RL methods.

## 69 2 Related Work

70 Modern offline RL methods based on Q-learning typically utilize dynamic programming to train  
71 a value function together with a mechanism to prevent backing up out-of-distribution (OOD) ac-  
72 tions [30]. This can be done by applying an explicit constraint that forces the learned policy to be  
73 “close” to the behavior policy under a variety of divergence measures [18, 54, 37, 42, 54, 24, 23,  
74 22, 50, 13], or by directly learning a conservative value function, either via a pessimistic training  
75 objective [26, 56, 36, 58] or by utilizing pessimistic bonuses [57, 39, 19, 54] in the backup. Other  
76 offline methods include model-based methods [20, 57, 2, 45, 38, 29, 58] that also utilize rollouts  
77 under a learned dynamics model to train the value function while also avoiding out-of-distribution  
78 actions. While most of these methods differ from each other in implementation details and empirical  
79 performance, in theory and in tabular problem settings, most of these methods can be traced back to  
80 the same objective that attempts to constrain the policy from choosing OOD actions. It is not entirely  
81 clear why one method should work better than another, or how one should go about designing better  
82 offline RL methods. In this paper, we show that, to a large extent, the benefits of offline RL methods  
83 comes from better representational quality, and how improving representational quality alone can  
84 lead to reasonable performance without any form of pessimism.

85 Prior works have sought to analyze several aspects of the representations induced by TD-based  
86 methods with function approximation largely in the standard online RL setting [1, 5, 25, 48, 31, 32]  
87 and in the offline RL setting [28, 27]. In the linear setting, [15, 55], study which representations  
88 can induce stable convergence of TD and [44, 33] have tried to devise convergent TD methods  
89 for arbitrary representations, but these prior works do not attempt to study the effect of pessimism  
90 on representations, or how OOD actions affect representations. Recent work [27, 28] study the  
91 learning dynamics of Q-learning in an overparameterized setting and observes excessively low-rank  
92 and aliased feature representations at the fixed points found by TD-learning. These prior works

93 propose some metrics to evaluate representational quality, and we do evaluate these in our analyses  
 94 in Section 5, but find that these metrics generally behave well, even though performance can be  
 95 improved with simple representational regularization. As we show, the metric we propose is more  
 96 predictive of algorithm performance. Moreover, these prior works do not quite study the interplay  
 97 between pessimism and representations that we do.

98 Finally, we note that our proposed approach of utilizing an ensemble of  $N$ -step returns is not new.  
 99 Most notably, it is related to TD( $\lambda$ ) which has been instantiated in various forms [41, 21, 51, 9]. Prior  
 100 works have also used  $N$ -step returns for a fixed value of  $N$  in methods that perform off-policy TD  
 101 learning [49, 17, 10]. Besides the fact that most of these works are based in an online RL setting, the  
 102 crucial distinction behind these prior works and our paper is that our work goes beyond the standard  
 103 explanation of bias-variance tradeoff for  $N$ -step returns [40], and analyzes  $N$ -step returns from a  
 104 different perspective: improving the quality of learned representations. We emphasize that our goal  
 105 is not to produce a novel algorithm, but rather to understand the efficacy of different components  
 106 towards the representations learned by the Q-function.

### 107 3 Preliminaries

108 The RL problem is formally defined by a Markov decision processes (MDPs) defined as  $\mathcal{M} =$   
 109  $(\mathcal{S}, \mathcal{A}, T, r, \mu_0, \gamma)$ , where  $\mathcal{S}, \mathcal{A}$  denote the state and action spaces, and  $T(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ ,  $r(\mathbf{s}, \mathbf{a})$  represent  
 110 the dynamics and reward function respectively.  $\mu_0(s)$  denotes the initial state distribution, and  
 111  $\gamma \in (0, 1)$  denotes the discount factor. The objective of RL is to learn a policy that maximizes the  
 112 return (discounted sum of rewards):  $\max_{\pi} J(\pi) := \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \pi} [\sum_t \gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$ . In offline RL, we are  
 113 provided with an offline dataset,  $\mathcal{D} = \{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')\}$ , of transitions collected using a behavior policy  
 114  $\pi_{\beta}$ , and our goal is to find the best possible policy only using the given dataset.

115 Directing training a  $Q$ -value function from the offline dataset often suffers from OOD actions [14, 24,  
 116 30], and therefore effective offline RL algorithms must enforce some constraint to prevent querying  
 117 the target Q-function on unseen actions. This constraint could be a behavior constraint, where the  
 118 learned policy  $\pi$  is constrained to be close to the behavior policy  $\pi_{\beta}$ . In this work, we build our  
 119 analysis on top of conservative Q-learning (CQL) [26], which applies a regularizer  $\mathcal{R}(\theta)$  to prevent  
 120 overestimation of Q-values for OOD actions.  $\mathcal{R}(\theta)$  minimizes the Q-values under the policy  $\pi(\mathbf{a}|\mathbf{s})$ ,  
 121 and counterbalances this term by maximizing the values of the actions in  $\mathcal{D}$ . Formally:

$$\min_{\theta} \alpha \left( \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi} [Q_{\theta}(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} [Q_{\theta}(\mathbf{s}, \mathbf{a})] \right) + \frac{1}{2} \mathbb{E}_{\substack{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D} \\ \mathbf{a}' \sim \pi}} \left[ (Q_{\theta}(\mathbf{s}, \mathbf{a}) - r - \gamma \bar{Q}(\mathbf{s}', \mathbf{a}'))^2 \right], \quad (1)$$

122 where  $\bar{Q}$  denotes the target  $Q$ -function. On the other hand, training a  $Q$ -value function for the  
 123 behavior policy, that only relies on action samples from the offline dataset is fairly easy and does  
 124 not suffer from the problem of OOD actions. A standard approach of learning such a  $Q$ -function is  
 125 what we refer to as “offline SARSA” [43], which only queries the action observed in the dataset at  
 126 the subsequent timestep to compute the Bellman target for training the Q-function. The objective for  
 127 SARSA can be written as:

$$\min_{\theta} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}', \mathbf{a}' \sim \mathcal{D}} \left[ (Q_{\theta}(\mathbf{s}, \mathbf{a}) - r - \gamma \bar{Q}(\mathbf{s}', \mathbf{a}'))^2 \right]. \quad (2)$$

128 Since the next step  $Q$ -values are computed using dataset actions, it eliminates the need to query  
 129  $Q$ -function for the values of any OOD actions. In effect, this procedure only relies on supervision  
 130 observed in the dataset (i.e., actions, the corresponding rewards and the next states) to learn repre-  
 131 sentations. Prior works [28] have argued that avoiding out-of-distribution actions altogether enables  
 132 SARSA to enjoy benefits of implicit regularization [52, 3] that otherwise may hurt TD learning.

133 In order to understand representational quality, we focus our analysis on the last layer feature  
 134 representation  $\phi(\mathbf{s}, \mathbf{a})$  learned by the neural network, following the conventions in prior work [8, 28,  
 135 27, 31, 32]. These prior works have also attempted to show that certain characteristics of the learned  
 136 representations  $\phi(\mathbf{s}, \mathbf{a})$  of a value network can explain certain pathologies with Q-learning.

### 137 4 To What Extent Do OOD Actions Explain the Instability in Offline RL?

138 Most prior works in offline RL focus on addressing the action distribution shift problem, proposing  
 139 a wide variety of methods in preventing the policies from taking OOD actions during the training  
 140 process. However, it remains unclear why different methods for mitigating OOD actions seem to

141 attain significantly different performance, and whether being *better* at preventing OOD actions is  
 142 actually the key to better results. It therefore seems natural to ask: to what degree is good (or bad)  
 143 performance of offline RL approaches really dependent on their ability to be pessimistic? In this  
 144 section, we study this question by performing a controlled empirical study. We perform experiments  
 145 to investigate both the sufficiency and necessity of being pessimistic and present them next.

#### 146 4.1 Is Pessimism Sufficient for Good Performance?

147 While several recent offline RL methods that correct for OOD actions by adding some form of  
 148 pessimism work well, in most of these approaches, the pessimism-inducing penalty (e.g., value  
 149 conservatism penalty like in CQL) or constraint (e.g., behavioral constraints) also affects the rep-  
 150 resentation learned by the internal layers of the Q-function (or the policy). In this section, we  
 151 argue via an empirical study on top of the CQL algorithm that, to a large extent, the benefits of this  
 152 pessimism-inducing mechanism stem from its impact on the learned representation and not so much  
 153 from its ability to combat overestimation.

154 **Empirical results showing insufficiency of pessimism.** To decouple the effects of pessimism in  
 155 handling overestimation and representational quality, we train a CQL [26] agent on the hopper-  
 156 medium-replay-v2 environment from the D4RL [11] suite, and make the following modification: we  
 157 let the last layer representation  $\phi(s, a)$  of the Q-network be updated by the TD-error (second term in  
 158 Equation 1) and the conservatism regularizer ( $\mathcal{R}(\theta)$ ) is **not** allowed to affect this representation. That  
 159 said, this regularizer  $\mathcal{R}(\theta)$  is allowed to affect the final layer weights of the Q-function. As a result,  
 160 while the CQL regularizer can still curb overestimation by manipulating the last layer Q-values, it is  
 161 unable to affect the representations, thereby inhibiting pessimism from providing any representational  
 162 benefits. For comparison, we also train a regular CQL agent on the same environments. For both  
 163 runs, we apply the same weight on the conservatism penalty.

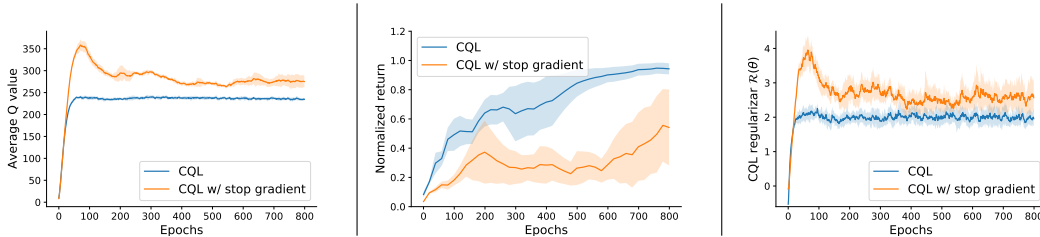


Figure 1: **CQL w/ stop gradient vs CQL** in hopper-medium-replay task. **Left:** CQL w/ stop gradient is able to prevent overestimation and results in non-divergent Q-values. **Middle:** the performance of CQL w/ stop gradient is significantly lower than regular CQL. **Right:** Values of the CQL regularizer are quite comparable between CQL and CQL w/ stop gradient, even though the observed performance is quite different.

164 As shown in Figure 1, once we prevent the CQL conservatism penalty from affecting the represen-  
 165 tation, the performance decreases significantly. In the left part of the figure, we see that when the  
 166 CQL regularizer is not allowed to affect the learned Q-function representations (denoted “CQL w/  
 167 stop gradient”), we are still able to attain stable and non-divergent Q-values, thereby avoiding the  
 168 issues typically observed with standard TD methods. However, CQL w/ stop gradient performs  
 169 significantly worse than base CQL (Figure 1, middle). As shown in Figure 1 (right), the value of the  
 170 CQL regularizer (i.e., the amount of pessimism) is still quite comparable in both cases, differing only  
 171 by about 0.5, which is quite small relative to the average magnitude of the learned Q-values ( $\sim 300$ ),  
 172 however there is a significant performance difference. This difference indicates that while pessimism  
 173 might be beneficial in lowering the value of OOD actions, it also contributes significantly to other  
 174 factors such as representation learning, and this representation learning benefit accounts for much of  
 175 the improvement from CQL, since without it the method performs much worse.

**Takeaway 4.1.** Besides preventing OOD actions, pessimism-inducing mechanisms in offline RL algorithms can also contribute to representation learning, and simply ensuring pessimism, without affecting representations might not be sufficient for good performance.

#### 177 4.2 How Much Performance Improvement Do Good Representations Account for?

178 While the above results suggest that pessimism alone does not account for full performance of  
 179 offline RL methods, and the quality of the learned representation has a crucial role to play in  
 180 determining the performance of value-based offline RL, it is not quite clear how much performance

181 do good representations account for, how much performance is accounted to by other factors and  
 182 what even a good representation even means. In this section, we attempt to answer this question  
 183 by construction: we perform an empirical study that completely removes any sort of pessimism,  
 184 but applies a representational regularizer. We show that it is still possible to obtain reasonable  
 185 performance if the learned representation is carefully regularized, despite the fact that the method we  
 186 test has no explicit mechanism for ensuring pessimistic estimates for OOD actions or constraining  
 187 the policy to remain in-distribution.

188 **Experiment setup.** As shown in Equation 1, the  
 189 CQL regularizer ( $\mathcal{R}(\theta)$  in Equation 1) pushes down  
 190 the Q-value at OOD actions and pushes up the Q-  
 191 value for in-distribution dataset actions. If this kind  
 192 of a pessimism penalty truly induces beneficial repre-  
 193 sentational regularization, a nature conjecture is that  
 194 representations that trained to minimize just the CQL  
 195 regularizer independently of the TD error must also  
 196 be useful, and must contain enough information to  
 197 distinguish dataset actions from OOD actions. On its  
 198 own, the CQL regularizer (Equation 1) resembles the  
 199 objective of the discriminator in generative adversar-  
 200 ial networks (GAN) [16] which serves a similar function of distinguishing dataset examples from  
 201 generated examples. Based on this intuition, in the next experiment, we construct an offline RL  
 202 method that utilizes a GAN objective, but only to train a *separate* linear output head on top of the  
 203 Q-function network, whereas the Q-values are simply trained to minimize TD error with no form of  
 204 pessimism whatsoever. A schematic illustration of this approach is shown in Figure 3. Specifically,  
 205 we adopt the least square GAN [34] objective due to its simplicity and stability. Concretely, let us  
 206 denote the linear discriminator weight as  $w_d$ , then given the Q-network representation  $\phi_\theta(s, a)$ , our  
 207 explicit regularization objective can be written as

$$\min_{\theta, w_d} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi} [(\phi_\theta(s, a)^\top w_d + 1)^2] + \mathbb{E}_{s, a \sim \mathcal{D}} [(\phi_\theta(s, a)^\top w_d - 1)^2]. \quad (3)$$

208 We apply this regularization on top of standard off-policy SAC [47], without any form of pessimism,  
 209 and evaluate the algorithm in the same environment as Section 4.1. For comparison, we also train an  
 210 naïve SAC agent with identical hyperparameters but without this second head.

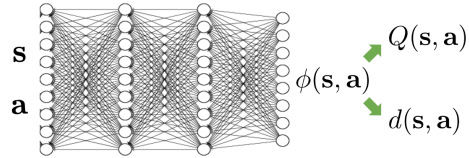


Figure 2: A schematic illustration of our approach for representational regularization that trains a Q-function with an auxiliary discriminator head for distinguishing potentially out-of-distribution policy actions from in-distribution dataset actions.

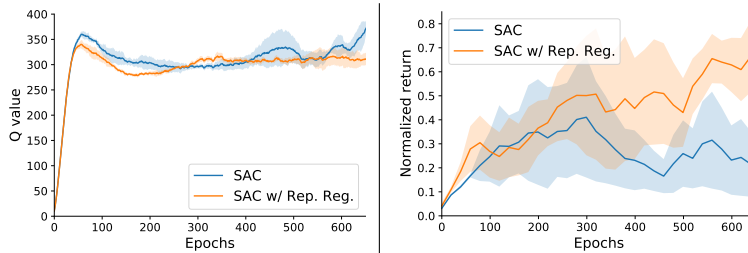


Figure 3: **SAC with representation regularization vs regular SAC** on hopper-medium-replay-v2 task. **Left:** SAC with representation regularization learns similar Q-values to regular SAC. **Right:** Representation regularization significantly improves the performance even without pessimism.

211 As shown in Figure 3, this modified algorithm can attain reasonable performance, significantly  
 212 outperforming naïve SAC, despite having no explicit mechanism to ensure pessimism, conservatism,  
 213 or policy constraints. Since the additional GAN term only influences the last layer representation,  
 214 its benefits can be attributed entirely to learning better representations. While the method is not as  
 215 effective as dedicated offline RL approaches such as CQL, this result, together with the experiment  
 216 from Section 4.1 strongly suggests that representation learning is not only important for offline RL,  
 217 but it also explains a large fraction of the performance gains for methods such as CQL. This in turn  
 218 implies that, in designing better offline RL methods, we should put particular emphasis on their effect  
 219 on representation learning, rather than simply on enforcing pessimism.

**Takeaway 4.2.** *The ability to learn good representations can explain a large fraction of the performance gains for practical offline RL methods. Explicit regularization techniques that gives good representations can be effective in offline RL, even in the absence of pessimism.*

221 **5 What Constitutes a Good Representation for Offline RL?**

222 Our empirical analysis from the previous section suggests that pessimistic offline RL methods do  
 223 affect the representations learned by offline RL algorithms such as CQL, and utilizing only the TD  
 224 error can give rise to representations that fail to adequately distinguish the dataset action from actions  
 225 from the learned policy. This distinction is crucial: since an offline RL algorithm only observes  
 226 ground truth supervision only in the form of instantaneous rewards and the subsequent environment  
 227 state, for dataset actions, the ability to successfully associate the right (long-term) reward with the  
 228 right dataset action is critical for attaining good performance. Can we formalize this intuition into a  
 229 diagnostic metric for measuring the “goodness” of the learned representation?

230 The most natural choice of such a metric inspired by our experimental analysis in Section 4.2 is  
 231 the accuracy of the separate discriminator head trained to distinguish dataset actions from policy  
 232 actions. In our preliminary experiments, we find that while a discriminator accuracy near 50% is  
 233 clearly indicative of poor performance, a reasonable discriminator accuracy (say  $\geq 60 - 70\%$ ) does  
 234 not necessarily indicate the absence of any representational issues. This is because while even a  
 235 somewhat correct representation can attain high accuracy, the representation may still not be rich  
 236 enough to match the fidelity needed for Q-value estimation. Therefore, we propose to utilize a  
 237 more complete metric for tracking the extent of action information in the learned representation: we  
 238 propose to train a non-linear model to reconstruct both the dataset and policy actions from the learned  
 239 representation  $\phi(\mathbf{s}, \mathbf{a})$ , and suggest tracking the reconstruction error of this model in aggregate over  
 240 dataset actions. This metric can be formalized as:

**Metric 5.1.** Train a parametric model,  $\Delta : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathcal{A}$  on the dataset:  $\mathcal{D}_\Delta := \mathcal{D}_\Delta^\pi \cup \mathcal{D}_\Delta^{\pi_\beta}$ ,  
 where  $\mathcal{D}_\Delta^{\pi_\beta} := \{(\mathbf{s}_i, \phi(\mathbf{s}_i, \mathbf{a}_i)), \mathbf{a}_i\}_{i=1}^N$  and  $\mathcal{D}_\Delta^\pi := \{(\mathbf{s}_i, \phi(\mathbf{s}_i, \pi(\mathbf{s}_i)), \pi(\mathbf{s}_i))\}_{i=1}^N$ . Then, track  
 the error metric:

$$\mathcal{L}_{recons}(\Phi) := \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{s}_i, \mathbf{a}_i) \in \mathcal{D}} \|\mathbf{a}_i - \Delta(\mathbf{s}_i, \phi(\mathbf{s}_i, \mathbf{a}_i))\|_2^2. \quad (4)$$

241  
 242 Since the reconstruction error,  $\mathcal{L}_{recons}(\Phi)$ , can take on a range of values, how should we choose values  
 243 to decide whether a representation is good enough or not? Specifically, what is a baseline value of this  
 244 quantity that can be considered a “gold standard” for comparison? To identify a good value of this  
 245 good standard, we seek to intuitively understand how OOD actions would impact the representations  
 246 learned by a value-based offline RL algorithm. We can do so by utilizing the following informal model  
 247 of the behavior of neural networks that is implied by several theories of deep learning [3, 4, 46, 7]:  
 248 sufficiently expressive and overparameterized neural networks are believed to learn the “simplest”  
 249 function that can fit the training data (i.e., match the actual label on the training datapoints). That  
 250 is to say that the learned function retains only information about the training data that is absolutely  
 251 critical for making predictions, and attempts to lose any unnecessary information.

252 When instantiated in the context of TD-learning, this intuitive model implies that the simplicity of the  
 253 function approximator would depend on its ability to fit the Bellman constraints on the training data.  
 254 If several of the actions used to compute Bellman targets are out-of-distribution, in principle, a simpler  
 255 function approximator can be learned by assigning arbitrary values to them, as Q-values at such  
 256 actions are hallucinated by the function approximator itself. On the other hand, if all the actions used  
 257 to produce Bellman targets also appear in the dataset (i.e., these actions also appear on the left hand  
 258 side of some Bellman constraint), the resulting function approximator is the most constrained, and  
 259 likely least simple. This implies that a good baseline that can serve as a gold standard for comparing  
 260  $\mathcal{L}_{recons}$  is the reconstruction error attained by offline SARSA (Equation 2). This means that closer the  
 261 value of  $\mathcal{L}_{recons}(\Phi_{\text{offline RL}})$  to  $\mathcal{L}_{recons}(\Phi_{\text{SARSA}})$ , the more desirable the learned representation.

262 **Empirical results.** To empirically validate the efficacy of our reconstruction error metric, we compute  
 263 the values of  $\mathcal{L}_{recons}$  for a variety of D4RL [12] tasks and compare them to the values attained by  
 264 SARSA. Observe in Figure 4 that while in some cases (e.g. kitchen), the reconstruction error for  
 265 naïve CQL is much larger than SARSA, indicating excessive loss of information about the dataset, in  
 266 other cases (antmaze and antmaze-heterogeneous), the reconstruction error for naïve CQL is smaller,  
 267 indicating that CQL hallucinates information about the dataset action. As an additional point of  
 268 reference, we also plot this metric for an approach that utilizes an  $N$ -step Bellman backup with CQL,  
 269 and observe that this approach attains a value of  $\mathcal{L}_{recons}$  closer to that of SARSA. Furthermore, even  
 270 though the policies produced by naïve SARSA don’t perform well (as confirmed by prior works [6]),

271 the value of  $\mathcal{L}_{\text{recons}}$  to that of SARSA, the better the performance of the resulting method. This  
 272 empirically corroborates our intuition about the efficacy of this metric.

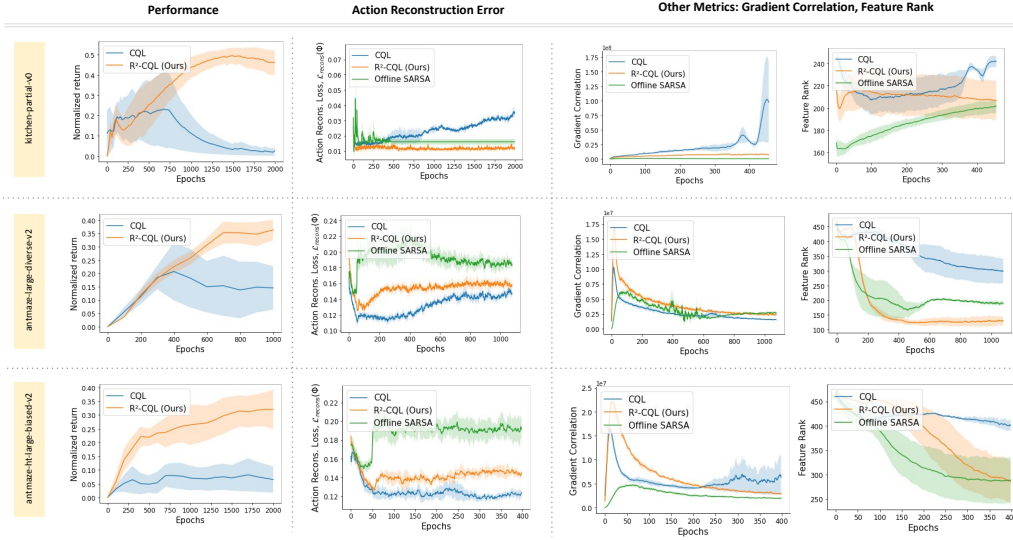


Figure 4: **Performance and metrics of  $R^2$ -CQL vs regular CQL, in comparison with SARSA.** Observe that measuring of closeness of the reconstruction error on the dataset actions (Metric 5.1) to the corresponding value for SARSA is able to accurately predict the performance trends, while other prior metrics may not.

273 Additionally, we also measure the predictive power of existing metrics from prior works such as  
 274 feature rank penalty [27] and feature dot products [28] in predicting the performance difference  
 275 between CQL and our approach. While these prior works used extreme values of these metrics (e.g.,  
 276 extremely low rank or extremely large dot products) to diagnose pathologies in TD, our analysis  
 277 shows that representational issues can still arise when these metrics behave relatively stably (see  
 278 Figure 4).

**Takeaway 5.1.** *The closer the value of the reconstruction error metric of an offline RL method based on TD-learning method that utilizes out-of-distribution actions, to that of SARSA, the better we would expect the performance of the learned policy to be.*

## 280 6 $R^2$ -CQL: A Simple Approach for Improving Representations For CQL

281 How can we improve the representations learned by offline RL algorithms? Our analysis above  
 282 suggests that this would involve constraining the learned representation to be closer to that learned  
 283 via offline SARSA, which only utilizes dataset actions for which ground truth supervision is available.  
 284 That is, we wish to devise an approach that can introduce a form of *representational regularization*,  
 285 which makes the representations closer to that of offline SARSA.

286 A simple approach that meets these requirements, and imposes a form of representational regu-  
 287 larization, is one that utilizes a Bellman backup operator which interpolates between complete  
 288 bootstrapping and estimating the value for SARSA. To this end, we propose to utilize an ensemble of  
 289  $n$ -step return estimators in conjunction with offline RL methods, similar to  $TD(\lambda)$  [43]. Concretely,  
 290 for a given choice of values of  $n = \{n_0, n_1, \dots, n_k\}$ , we utilize the following Bellman operator to  
 291 generate regression targets for TD:

$$\tilde{\mathcal{B}}^\pi Q(s_0, \mathbf{a}_0) := \frac{1}{k} \sum_{j=1}^k \left( \sum_{l=0}^{n_j-1} \gamma^l r(s_l, \mathbf{a}_l) + \gamma^{n_j} Q(s_{n_j}, \mathbf{a}_{n_j}) \right). \quad (5)$$

292 We will now discuss how we can convert this approach into a practical method for offline RL.

293 **Practical instantiation.** Our practical algorithm only modifies the CQL training objective (Equa-  
 294 tion 1) to now use the Bellman backup operator shown in Equation 5, with no other changes. We  
 295 inherit the value of  $\alpha$  directly from CQL, without tuning it, and do not modify any other hyperpa-  
 296 rameters. We utilize values of  $n = \{1, 3, 5\}$  across all domains. Note that unlike prior methods

297 based on explicit regularization such as the feature rank [28] or dot products [27], our approach does  
 298 not require any specific hyperparameter to be tuned per domain, highlighting the simplicity of this  
 299 approach.

300 **Empirical results.** We empirically validate our n-step approach by evaluating both the value of  
 301  $\mathcal{L}_{\text{recons}}$  and the performance across a wide range of offline RL tasks from D4RL [12]. Following  
 302 the protocol in [28], we present two sets of performance numbers in Table 1: the final performance  
 303 attained by the algorithm after a fixed number of gradient steps (denoted “Final Performance”) and  
 304 the average performance attained over the course of training (denoted “Average Performance”), which  
 305 is a measure of the stability of the offline RL algorithm over the course of training. We additionally  
 306 already presented the value of the reconstruction error on a subset of domains in Figure 4.

307 Observe that on all the tasks, our approach,  $R^2$ -CQL attains a better or comparable performance both  
 308 measured by the final performance of the algorithm and the average performance across iterations,  
 309 which demonstrates the stability of training. The gap between naïve CQL and the n-step approach is  
 310 larger under the average performance metric, indicating that the latter is much more stable. Finally,  
 311 perhaps unsurprisingly, the representational metrics do indicate that utilizing the mixture of n-step  
 312 Bellman targets does lead to reconstruction error values closer to that of offline SARSA.

313 While this simple approach does lead to improvements in performance, perhaps the more important  
 314 question is *why* does it actually improve performance. Traditionally, in the on-policy setting, the  
 315 utility of an ensemble of  $N$ -step returns via approaches such as TD( $\lambda$ ) [43] or GAE [41] primarily  
 316 emerges from an ability to better manage a bias-variance tradeoff: by controlling an algorithmic  
 317 hyperparameter, the bias induced in learning a parametric Q-function can be effectively traded against  
 318 the variance of a Monte-Carlo return estimator. However, in this case, we utilize  $N$ -step returns in an  
 319 offline setting, with an already pessimistic algorithm (CQL). Since CQL already aims to underestimate  
 320 the return of the learned policy, we would expect  $N$ -step Bellman targets to only be *more* conservative,  
 321 since they bias the Q-function towards the values of the behavior policy and therefore be more biased  
 322 than CQL. Typically, this bias issue is solved by utilizing importance corrections [9, 35], but we do  
 323 not use any such correction. Therefore, not only does  $R^2$ -CQL use a high variance Bellman target,  
 324 but also a more biased one, and yet it outperforms CQL. This again indicates that the representation  
 325 learning benefits of this approach are likely much more useful towards improving performance despite  
 326 the bias.

Task	Final Performance		Average Performance	
	CQL	$R^2$ -CQL	CQL	$R^2$ -CQL
kitchen-mixed	0.000 $\pm$ 0.000	0.362 $\pm$ 0.013	0.085 $\pm$ 0.114	0.330 $\pm$ 0.098
kitchen-partial	0.138 $\pm$ 0.138	0.475 $\pm$ 0.075	0.089 $\pm$ 0.111	0.414 $\pm$ 0.139
kitchen-complete	0.000 $\pm$ 0.000	0.025 $\pm$ 0.025	0.163 $\pm$ 0.143	0.100 $\pm$ 0.106
antmaze-medium-play	0.435 $\pm$ 0.315	0.670 $\pm$ 0.090	0.569 $\pm$ 0.200	0.602 $\pm$ 0.216
antmaze-medium-diverse	0.680 $\pm$ 0.070	0.645 $\pm$ 0.045	0.511 $\pm$ 0.214	0.538 $\pm$ 0.212
antmaze-large-play	0.005 $\pm$ 0.005	0.320 $\pm$ 0.000	0.098 $\pm$ 0.105	0.265 $\pm$ 0.104
antmaze-large-diverse	0.095 $\pm$ 0.035	0.420 $\pm$ 0.010	0.162 $\pm$ 0.083	0.303 $\pm$ 0.145
antmaze-ht-large	0.090 $\pm$ 0.090	0.380 $\pm$ 0.160	0.082 $\pm$ 0.057	0.283 $\pm$ 0.125
antmaze-ht-large-biased	0.000 $\pm$ 0.000	0.310 $\pm$ 0.190	0.067 $\pm$ 0.057	0.302 $\pm$ 0.098
antmaze-ht-medium	0.000 $\pm$ 0.000	0.320 $\pm$ 0.140	0.155 $\pm$ 0.118	0.290 $\pm$ 0.121
antmaze-ht-medium-biased	0.000 $\pm$ 0.000	0.220 $\pm$ 0.040	0.126 $\pm$ 0.192	0.234 $\pm$ 0.083

Table 1: Final and average performance for  $R^2$ -CQL and CQL across 7 D4RL tasks and 4 heterogeneous antmaze tasks. All performances are evaluated with 2 random seeds for 1000 epochs. We see that  $R^2$ -CQL improves the final and average performance over naïve CQL significantly.

## 327 7 Discussion and Conclusion

328 In this paper, we demonstrate that while addressing the overestimation due to OOD actions is  
 329 important for offline RL, a crucial, but largely overlooked, factor for obtaining good performance  
 330 in value-based offline RL algorithms is good representation quality. We show through extensive  
 331 empirical results that, perhaps surprisingly, pessimism in practical offline RL algorithms such as CQL  
 332 contributes to the performance not only as a way to prevent overestimation, but more significantly



333 as a way to induce good representations. We also show that pessimism is not the only way to attain  
334 good representations and methods that attain good representations can still work well. Based on  
335 this experimental analysis, we propose a practical metric that quantitatively tracks the quality of  
336 learned representation, and show that simply utilizing an ensemble of  $N$ -step returns to compute  
337 Bellman targets can provide a strong representation regularization and thus significantly improve  
338 the performance of conservative offline RL algorithm. We hope that our discovery can highlight the  
339 importance of representation learning in offline RL, and thus open up new opportunities to devise  
340 stronger offline RL methods.

341 While we provide a practical method  $R^2$ -CQL to regularize representations, by no means we claim  
342 that it is an optimal method. Therefore a natural step for future work direction is to seek for better  
343 ways to understand and improve the quality of learned representations. We believe that such future  
344 search has the potential of bringing deep insights and profound influences to the field of offline RL  
345 and hope that our analysis sheds light on some of these questions.

## 346 References

- 347 [1] Joshua Achiam, Ethan Knight, and Pieter Abbeel. Towards characterizing divergence in deep  
348 q-learning. *arXiv preprint arXiv:1903.08894*, 2019.
- 349 [2] Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. *arXiv preprint*  
350 *arXiv:2008.05556*, 2020.
- 351 [3] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit  
352 acceleration by overparameterization. *arXiv preprint arXiv:1802.06509*, 2018.
- 353 [4] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix  
354 factorization. In *Advances in Neural Information Processing Systems*, pages 7413–7424, 2019.
- 355 [5] Emmanuel Bengio, Joelle Pineau, and Doina Precup. Interference and generalization in temporal  
356 difference learning. *arXiv preprint arXiv:2003.06350*, 2020.
- 357 [6] David Brandfonbrener, William F Whitney, Rajesh Ranganath, and Joan Bruna. Offline rl  
358 without off-policy evaluation. *arXiv preprint arXiv:2106.08909*, 2021.
- 359 [7] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural  
360 networks trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338.  
361 PMLR, 2020.
- 362 [8] Will Dabney, André Barreto, Mark Rowland, Robert Dadashi, John Quan, Marc G Bellemare,  
363 and David Silver. The value-improvement path: Towards better representations for reinforcement  
364 learning. *arXiv preprint arXiv:2006.02243*, 2020.
- 365 [9] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymir Mnih, Tom Ward,  
366 Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl  
367 with importance weighted actor-learner architectures. *arXiv preprint arXiv:1802.01561*, 2018.
- 368 [10] William Fedus, Prajit Ramachandran, Rishabh Agarwal, Yoshua Bengio, Hugo Larochelle,  
369 Mark Rowland, and Will Dabney. Revisiting fundamentals of experience replay. *arXiv preprint*  
370 *arXiv:2007.06700*, 2020.
- 371 [11] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven  
372 reinforcement learning. In *arXiv*, 2020. URL <https://arxiv.org/pdf/2004.07219>.
- 373 [12] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
374 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 375 [13] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
376 *arXiv preprint arXiv:2106.06860*, 2021.
- 377 [14] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning  
378 without exploration. *arXiv preprint arXiv:1812.02900*, 2018.
- 379 [15] Dibya Ghosh and Marc G Bellemare. Representations for stable off-policy reinforcement  
380 learning. *arXiv preprint arXiv:2007.05520*, 2020.
- 381 [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
382 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural*  
383 *information processing systems*, 27, 2014.

- 384 [17] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dab-  
385 ney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining  
386 improvements in deep reinforcement learning. In *Thirty-Second AAAI Conference on Artificial  
387 Intelligence*, 2018.
- 388 [18] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza,  
389 Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement  
390 learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 391 [19] Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl?  
392 *arXiv preprint arXiv:2012.15085*, 2020.
- 393 [20] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:  
394 Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- 395 [21] Hajime Kimura, Shigenobu Kobayashi, et al. An analysis of actor-critic algorithms using  
396 eligibility traces: reinforcement learning with imperfect value functions. *Journal of Japanese  
397 Society for Artificial Intelligence*, 15(2):267–275, 2000.
- 398 [22] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
399 q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- 400 [23] Ilya Kostrikov, Jonathan Tompson, Rob Fergus, and Ofir Nachum. Offline reinforcement  
401 learning with fisher divergence critic regularization. *arXiv preprint arXiv:2103.08050*, 2021.
- 402 [24] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy  
403 q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing  
404 Systems*, pages 11761–11771, 2019.
- 405 [25] Aviral Kumar, Abhishek Gupta, and Sergey Levine. Discor: Corrective feedback in reinforce-  
406 ment learning via distribution correction. *arXiv preprint arXiv:2003.07305*, 2020.
- 407 [26] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for  
408 offline reinforcement learning. *arXiv preprint arXiv:2006.04779*, 2020.
- 409 [27] Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-  
410 parameterization inhibits data-efficient deep reinforcement learning. In *International Con-  
411 ference on Learning Representations*, 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=O9bnihSFFXU)  
412 [O9bnihSFFXU](https://openreview.net/forum?id=O9bnihSFFXU).
- 413 [28] Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron Courville, George Tucker, and Sergey  
414 Levine. DR3: Value-Based Deep Reinforcement Learning Requires Explicit Regularization.  
415 *arXiv preprint arXiv:2112.04716*, 2021.
- 416 [29] Byung-Jun Lee, Jongmin Lee, and Kee-Eung Kim. Representation balancing offline model-  
417 based reinforcement learning. In *International Conference on Learning Representations*, 2021.  
418 URL [https://openreview.net/forum?id=QpNz8r\\_Ri2Y](https://openreview.net/forum?id=QpNz8r_Ri2Y).
- 419 [30] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
420 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 421 [31] Clare Lyle, Mark Rowland, Georg Ostrovski, and Will Dabney. On the effect of auxiliary tasks  
422 on representation dynamics. In *International Conference on Artificial Intelligence and Statistics*,  
423 pages 1–9. PMLR, 2021.
- 424 [32] Clare Lyle, Mark Rowland, and Will Dabney. Understanding and preventing capacity loss in  
425 reinforcement learning. *arXiv preprint arXiv:2204.09560*, 2022.
- 426 [33] Hamid R. Maei, Csaba Szepesvári, Shalabh Bhatnagar, Doina Precup, David Silver, and  
427 Richard S. Sutton. Convergent temporal-difference learning with arbitrary smooth function  
428 approximation. In *Proceedings of the 22nd International Conference on Neural Information  
429 Processing Systems*, 2009.
- 430 [34] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley.  
431 Least squares generative adversarial networks. In *Proceedings of the IEEE international  
432 conference on computer vision*, pages 2794–2802, 2017.
- 433 [35] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient  
434 off-policy reinforcement learning. In *Advances in Neural Information Processing Systems  
435 (NeurIPS)*, pages 1054–1062, 2016.

- 436 [36] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Al-  
437 gaedice: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- 438 [37] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:  
439 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 440 [38] Rafael Rafailov, Tianhe Yu, A. Rajeswaran, and Chelsea Finn. Offline reinforcement learning  
441 from images with latent space models. *Learning for Decision Making and Control (LADC)*,  
442 2021.
- 443 [39] Shideh Rezaeifar, Robert Dadashi, Nino Vieillard, Léonard Hussenot, Olivier Bachem, Olivier  
444 Pietquin, and Matthieu Geist. Offline reinforcement learning as anti-exploration. *arXiv preprint*  
445 *arXiv:2106.06431*, 2021.
- 446 [40] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous  
447 control using generalized advantage estimation. In *International Conference on Learning*  
448 *Representations (ICLR)*, 2016.
- 449 [41] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-  
450 dimensional continuous control using generalized advantage estimation. *arXiv preprint*  
451 *arXiv:1506.02438*, 2015.
- 452 [42] Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael  
453 Neunert, Thomas Lampe, Roland Hafner, and Martin Riedmiller. Keep doing what worked:  
454 Behavioral modelling priors for offline reinforcement learning. *arXiv preprint arXiv:2002.08396*,  
455 2020.
- 456 [43] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. Second  
457 edition, 2018.
- 458 [44] Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba  
459 Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning  
460 with linear function approximation. In *Proceedings of the 26th Annual International Conference*  
461 *on Machine Learning*, pages 993–1000, 2009.
- 462 [45] Phillip Swazinna, Steffen Udluft, and Thomas Runkler. Overcoming model bias for robust  
463 offline deep reinforcement learning. *arXiv preprint arXiv:2008.05533*, 2020.
- 464 [46] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In  
465 *2015 IEEE information theory workshop (itw)*, pages 1–5. IEEE, 2015.
- 466 [47] Kristian Hartikainen George Tucker Sehoon Ha Jie Tan Vikash Kumar Henry Zhu Abhishek  
467 Gupta Pieter Abbeel Tuomas Haarnoja, Aurick Zhou and Sergey Levine. Soft actor-critic  
468 algorithms and applications. Technical report, 2018.
- 469 [48] Hado Van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, and Joseph  
470 Modayil. Deep reinforcement learning and the deadly triad. *arXiv preprint arXiv:1812.02648*,  
471 2018.
- 472 [49] Hado van Hasselt, Matteo Hessel, and John Aslanides. When to use parametric models in  
473 reinforcement learning? *arXiv preprint arXiv:1906.05243*, 2019.
- 474 [50] Ziyu Wang, Alexander Novikov, Konrad Żoła, Jost Tobias Springenberg, Scott Reed, Bobak  
475 Shahriari, Noah Siegel, Josh Merel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized  
476 regression. *arXiv preprint arXiv:2006.15134*, 2020.
- 477 [51] Paweł Wawrzyński. Real-time reinforcement learning by sequential actor–critics and experience  
478 replay. *Neural networks*, 22(10):1484–1497, 2009.
- 479 [52] Colin Wei, Jason Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and  
480 optimization of neural nets vs their induced kernel. 2019.
- 481 [53] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
482 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 483 [54] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
484 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 485 [55] Chenjun Xiao, Bo Dai, Jincheng Mei, Oscar A Ramirez, Ramki Gummadi, Chris Harris, and  
486 Dale Schuurmans. Understanding and leveraging overparameterization in recursive value  
487 estimation. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=shbAgEsk3qM>.  
488

- 489 [56] Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-  
490 consistent pessimism for offline reinforcement learning. *Advances in neural information*  
491 *processing systems*, 34, 2021.
- 492 [57] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea  
493 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *arXiv preprint*  
494 *arXiv:2005.13239*, 2020.
- 495 [58] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea  
496 Finn. Combo: Conservative offline model-based policy optimization. *arXiv preprint*  
497 *arXiv:2102.08363*, 2021.