

---

# Uncertainty-Based Offline Reinforcement Learning with Diversified Q-ensemble

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Offline reinforcement learning (offline RL), which aims to find an optimal policy  
2 from a previously collected static dataset, bears algorithmic difficulties due to  
3 function approximation errors from out-of-distribution (OOD) data points. To  
4 this end, offline RL algorithms adopt either a constraint or a penalty term that  
5 explicitly guides the policy to stay close to the given dataset. However, prior  
6 methods typically require accurate estimation of the behavior policy or sampling  
7 from OOD data points, which themselves can be a non-trivial problem. Moreover,  
8 these methods under-utilize the generalization ability of deep neural networks  
9 and often fall into suboptimal solutions too close to the given dataset. In this  
10 work, we propose an uncertainty-based offline RL method that takes into account  
11 the confidence of the Q-value prediction and does not require any estimation or  
12 sampling of the data distribution. We show that the clipped Q-learning, a technique  
13 widely used in online RL, can be leveraged to successfully penalize OOD data  
14 points with high prediction uncertainties. Surprisingly, we find that it is possible  
15 to substantially outperform existing offline RL methods on various tasks by simply  
16 increasing the number of Q-networks along with the clipped Q-learning. Based  
17 on this observation, we propose an ensemble-diversified actor-critic algorithm that  
18 reduces the number of required ensemble networks down to a tenth compared to  
19 the naive ensemble while achieving state-of-the-art performance on most of the  
20 D4RL benchmarks considered.

## 21 1 Introduction

22 Over the recent years, deep reinforcement learning (deep RL) has achieved considerable success  
23 in various domains such as robotics [20], recommendation systems [6], and strategy games [26].  
24 However, a major drawback of RL algorithms is that they adopt an active learning procedure, where  
25 training steps require active interactions with the environment. This trial-and-error procedure can be  
26 prohibitive when scaling RL to real-world applications such as autonomous driving and healthcare,  
27 as exploratory actions can cause critical damage to the agent or the environment [19]. *Offline RL*,  
28 also known as *batch RL*, aims to overcome this problem by learning policies using only previously  
29 collected data without further interactions with the environment [2, 11, 19].

30 Even though offline RL is a promising direction to lead a more *data-driven* way of solving RL  
31 problems, recent works show offline RL faces new algorithmic challenges [19]. Typically, if the  
32 coverage of the dataset is not sufficient, vanilla RL algorithms suffer severely from extrapolation  
33 error, overestimating the Q-values of out-of-distribution (OOD) state-action pairs [15]. To this end,  
34 most offline RL methods apply some constraints or penalty terms on top of the existing RL algorithms  
35 to enforce the learning process to be more conservative. For example, some prior works explicitly  
36 regularize the policy to be close to the behavior policy that was used to collect the data [11, 15]. A

37 more recent work instead penalizes the Q-values of OOD state-action pairs to enforce the Q-values to  
 38 be more pessimistic [16].

39 While these methods achieve significant performance gains over vanilla RL methods, they either  
 40 require an estimation of the behavior policy or explicit sampling from OOD data points, which  
 41 themselves can be non-trivial to solve. Furthermore, these methods do not utilize the generalization  
 42 ability of the Q-function networks and prohibit the agent from approaching any OOD state-actions  
 43 without any consideration on whether they are good or bad. However, if we can identify OOD data  
 44 points where we can predict their Q-values with high confidence, it is more effective not to restrain  
 45 the agent from choosing those data points.

46 From this intuition, we propose an uncertainty-based model-free offline RL method that effectively  
 47 quantifies the uncertainty of the Q-value estimates by an ensemble of Q-function networks and does  
 48 not require any estimation or sampling of the data distribution. To achieve this, we first show that a  
 49 well-known technique from online RL, the clipped Q-learning [10], can be successfully leveraged as  
 50 an uncertainty-based penalization term. Our experiments reveal that we can achieve state-of-the-art  
 51 performance on various offline RL tasks by solely using this technique with increased ensemble  
 52 size. To further improve the practical usability of the method, we develop an ensemble diversifying  
 53 objective that significantly reduces the number of required ensemble networks. We evaluate our  
 54 proposed method on D4RL benchmarks [9] and verify that the proposed method outperforms the  
 55 previous state-of-the-art by a large margin on various types of environments and datasets.

## 56 2 Preliminaries

57 We consider an environment formulated as a Markov Decision Process (MDP) defined by a tuple  
 58  $(\mathcal{S}, \mathcal{A}, T, r, \mu_0, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $T(\mathbf{s}' | \mathbf{s}, \mathbf{a})$  is the transition  
 59 probability distribution,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $\mu_0$  is the initial state distribution,  
 60 and  $\gamma \in (0, 1]$  is the discount factor. The goal of reinforcement learning is to find an optimal  
 61 policy  $\pi(\mathbf{a} | \mathbf{s})$  that maximizes the cumulative discounted reward  $\mathbb{E}_{\mathbf{s}_t, \mathbf{a}_t} [\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$ , where  
 62  $\mathbf{s}_0 \sim \mu_0(\cdot)$ ,  $\mathbf{a}_t \sim \pi(\cdot | \mathbf{s}_t)$ , and  $\mathbf{s}_{t+1} \sim T(\cdot | \mathbf{s}_t, \mathbf{a}_t)$ .

63 One of the major approaches for obtaining such a policy is Q-learning [12, 20] which learns a  
 64 state-action value function  $Q_\phi(\mathbf{s}, \mathbf{a})$  parameterized by a neural network that represents the expected  
 65 cumulative discounted reward when starting from state  $\mathbf{s}$  and action  $\mathbf{a}$ . Standard actor-critic approach  
 66 [14] learns this Q-function by minimizing the Bellman residual  $(Q_\phi(\mathbf{s}, \mathbf{a}) - \mathcal{B}^{\pi_\theta} Q_\phi(\mathbf{s}, \mathbf{a}))^2$ , where  
 67  $\mathcal{B}^{\pi_\theta} Q_\phi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\mathbf{s}' \sim T(\cdot | \mathbf{s}, \mathbf{a})} [r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta(\cdot | \mathbf{s}')} Q_\phi(\mathbf{s}', \mathbf{a}')] ]$  is the Bellman operator. In the  
 68 context of offline RL, where transitions are sampled from a static dataset  $\mathcal{D}$ , the objective for the  
 69 Q-network becomes minimizing

$$J_q(Q_\phi) := \mathbb{E}_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \sim \mathcal{D}} \left[ \left( Q_\phi(\mathbf{s}, \mathbf{a}) - (r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_\theta(\cdot | \mathbf{s}')} [Q_{\phi'}(\mathbf{s}', \mathbf{a}')] ) \right)^2 \right], \quad (1)$$

70 where  $Q_{\phi'}$  represents the target Q-network softly updated for algorithmic stability [20]. The policy,  
 71 which is also parameterized by a neural network, is updated in an alternating fashion to maximize the  
 72 expected Q-value:  $J_p(\pi_\theta) := \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi_\theta(\cdot | \mathbf{s})} [Q_\phi(\mathbf{s}, \mathbf{a})]$ .

73 However, as the policy is updated to maximize the Q-values, the actions  $\mathbf{a}'$  sampled from the current  
 74 policy in Equation (1) can be biased towards OOD actions with erroneously high Q-values. In the  
 75 offline RL setting, such errors cannot be corrected by feedback from the environment as in online  
 76 RL. To handle the error propagation from these OOD actions, most offline RL algorithms regularize  
 77 either the policy [11, 15] or the Q-function [16] to be biased towards the given dataset. However, the  
 78 policy regularization methods typically require an accurate estimation of the behavior policy, which  
 79 limits their scope of applications to simple environments. The previous state-of-the-art method CQL  
 80 [16] instead learns conservative Q-values without estimating the behavior policy by penalizing the  
 81 Q-values of OOD actions by

$$\min_{\phi} J_q(Q_\phi) + \alpha \left( \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \mu(\cdot | \mathbf{s})} [Q_\phi(\mathbf{s}, \mathbf{a})] - \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} [Q_\phi(\mathbf{s}, \mathbf{a})] \right),$$

82 where  $\mu$  is an approximation of the policy that maximizes the current Q-function. While CQL  
 83 does not need explicit behavior policy estimation, it requires sampling from an appropriate action  
 84 distribution  $\mu(\cdot | \mathbf{s})$ , which might be computationally ineffective in high-dimensional action space.

85 **3 Uncertainty penalization with Q-ensemble**

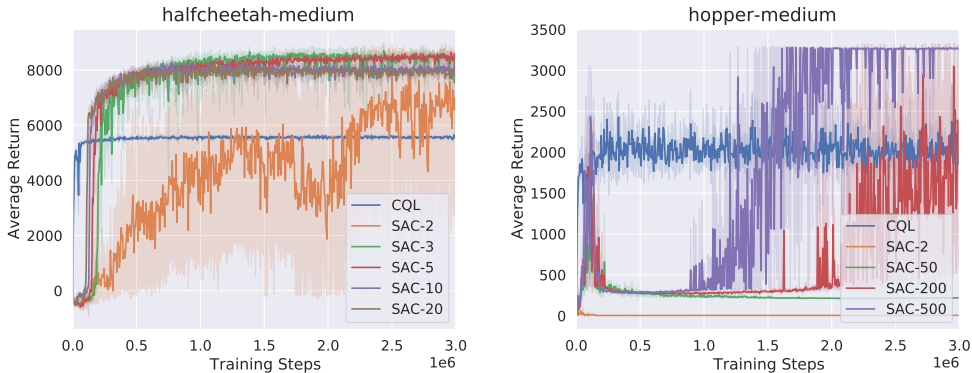


Figure 1: Performance of SAC- $N$  on halfcheetah-medium and hopper-medium datasets while varying  $N$ , compared to CQL. ‘Average Return’ denotes the undiscounted return of each policies on evaluation. Results averaged over 4 seeds.

86 In this section, we turn our attention to a conventional technique from online RL, Clipped Double  
 87 Q-learning [10], which uses the minimum value of two parallel Q-networks as the Bellman target:  $y =$   
 88  $r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_{\theta}(\cdot | \mathbf{s}')}$   $\left[ \min_{j=1,2} Q_{\phi'_j}(\mathbf{s}', \mathbf{a}') \right]$ . Although this technique was originally proposed in  
 89 online RL to mitigate the overestimation from general prediction errors, some offline RL algorithms  
 90 [11, 15, 28] also utilize this technique to enforce their Q-value estimates to be more pessimistic.  
 91 However, the isolated effect of the clipped Q-learning in offline RL was not fully analyzed in the  
 92 previous works, as they use the technique only as an auxiliary term that adds up to their core methods.  
 93 To examine the ability of clipped Q-learning to prevent the overestimation in offline RL on its own,  
 94 we modify SAC [12] by increasing the number of Q-ensembles from 2 to  $N$ :

$$\min_{\phi_i} \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}' \sim \mathcal{D}} \left[ \left( Q_{\phi_i}(\mathbf{s}, \mathbf{a}) - \left( r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi_{\theta}(\cdot | \mathbf{s}')} \left[ \min_{j=1, \dots, N} Q_{\phi'_j}(\mathbf{s}', \mathbf{a}') - \beta \log \pi_{\theta}(\mathbf{a}' | \mathbf{s}') \right] \right) \right)^2 \right]$$

$$\max_{\theta} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi_{\theta}(\cdot | \mathbf{s})} \left[ \min_{j=1, \dots, N} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) - \beta \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \right], \tag{2}$$

95 for  $i = 1, \dots, N$ . We denote this modified algorithm as SAC- $N$ .

96 Figure 1 shows the preliminary experiments on D4RL halfcheetah-medium and hopper-medium  
 97 datasets [9] while varying  $N$ . Note that these datasets are constructed from suboptimal behavior  
 98 policies. Surprisingly, as we gradually increase  $N$ , we can successfully find policies that outperform  
 99 the previous state-of-the-art method (CQL) by a large margin. In fact, as we will present in Section 5,  
 100 SAC- $N$  outperforms CQL on various types of environments and data-collection policies.

101 To understand why this simple technique works so well, we can first interpret the clipping procedure  
 102 (choosing the minimum value from the ensemble) as penalizing state-action pairs with high-variance  
 103 Q-value estimates, which encourages the policy to favor actions that appeared in the dataset [11].  
 104 Note that the dataset samples will naturally have lower variance compared to the OOD samples as the  
 105 Bellman residual term in Equation (2) explicitly aligns the Q-value predictions for the dataset samples.  
 106 More formally, we can regard this difference in variance as accounting for *epistemic uncertainty* [8]  
 107 which refers to the uncertainty stemming from limited data and knowledge.

108 Utilization of the clipped Q-value relates to methods that consider the confidence bound of the  
 109 Q-value estimates [24]. Online RL methods typically utilize the Q-ensemble to form an optimistic  
 110 estimate of the Q-value, by adding the standard deviation to the mean of the Q-ensembles [18]. This  
 111 optimistic Q-value, also known as the upper-confidence bound (UCB), can encourage the exploration  
 112 of unseen actions with high uncertainty. However, in offline RL, the dataset available during training  
 113 is fixed, and we have to focus on *exploiting* the given data. For this purpose, it is natural to utilize the  
 114 lower-confidence bound (LCB) of the Q-value estimates, for example by subtracting the standard  
 115 deviation from the mean, which allows us to avoid risky state-actions.

116 The clipped Q-learning algorithm, which chooses the worst-case Q-value instead to compute the  
 117 pessimistic estimate, can also be interpreted as utilizing the LCB of the Q-value predictions. Sup-  
 118 pose  $Q(\mathbf{s}, \mathbf{a})$  follows a Gaussian distribution with mean  $\mu(\mathbf{s}, \mathbf{a})$  and variance  $\sigma^2(\mathbf{s}, \mathbf{a})$ . Also, let  
 119  $\{Q_j(\mathbf{s}, \mathbf{a})\}_{j=1}^N$  be realizations of  $Q(\mathbf{s}, \mathbf{a})$ . Then, we can approximate the expected minimum of the  
 120 realizations following the work of Royston [23] as

$$\mathbb{E} \left[ \min_{j=1, \dots, N} Q_j(\mathbf{s}, \mathbf{a}) \right] \approx \mu(\mathbf{s}, \mathbf{a}) - \Phi^{-1} \left( \frac{N - \frac{\pi}{8}}{N - \frac{\pi}{4} + 1} \right) \sigma(\mathbf{s}, \mathbf{a}), \quad (3)$$

121 where  $\Phi$  is the CDF of the standard Gaussian distribution. This relation indicates that using the  
 122 clipped Q-value is similar to penalizing the ensemble mean of the Q-values with the standard deviation  
 123 scaled by a coefficient dependent on  $N$ .

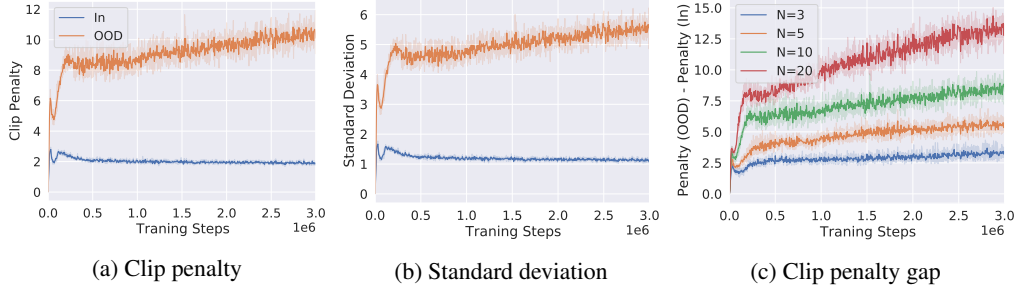


Figure 2: (a) and (b) each plots the size of the clip penalty and the standard deviation of the Q-value estimates for in-distribution (behavior) and OOD (random) actions while training SAC-10 on halfcheetah-medium dataset. (c) plots the gap of the clip penalty between the in-distribution and OOD actions while varying  $N$ . Results averaged over 4 seeds.

124 We now move on to the empirical analysis of the clipped Q-learning. Figure 2a compares the  
 125 strength of the uncertainty penalty on in-distribution and OOD actions. Specifically, we compare  
 126 actions sampled from two types of policies: (1) the behavior policy which was used to collect the  
 127 dataset, and (2) the random policy which samples actions uniformly from the action space. For each  
 128 policy, we measure the size of the penalty from the clipping as  $\mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi(\cdot | \mathbf{s})} \left[ \frac{1}{N} \sum_{j=1}^N Q_{\phi_j}(\mathbf{s}, \mathbf{a}) - \min_{j=1, \dots, N} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \right]$ . Figure 2a shows that the clipping term penalizes the random state-action  
 129 pairs much stronger than the in-distribution pairs throughout the training. For comparison, we also  
 130 measure the standard deviation of the Q-values for each policy. The results in Figure 2b show that as  
 132 we conjectured, the Q-value predictions for the OOD actions have a higher variance. We also find that  
 133 the size of the penalty and the standard deviation are highly correlated, as we noted in Equation (3).

134 Since we confirmed that OOD actions have higher variance on Q-value estimates, the effect of  
 135 increasing  $N$  becomes obvious: it strengthens the penalty applied to the OOD samples compared to  
 136 the dataset samples. To verify this, we measured the relative penalty applied to the OOD samples in  
 137 Figure 2c and found that indeed the OOD samples are penalized relatively further as  $N$  increases.

## 138 4 Ensemble gradient diversification

139 Even though SAC- $N$  outperforms existing methods on vari-  
 140 ous tasks, it sometimes requires an excessively large number  
 141 of ensembles to learn stably (e.g.,  $N = 500$  for hopper-  
 142 medium). While investigating its reason, we found that  
 143 the performance of SAC- $N$  is highly correlated with the  
 144 diversity of the Q-functions' input gradients  $\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a})$ ,  
 145 which increases with  $N$ . Figure 3 measures the minimum  
 146 cosine similarity between the (normalized) gradients of the  
 147 Q-functions  $\min_{i \neq j} \langle \nabla_{\mathbf{a}} Q_{\phi_i}(\mathbf{s}, \mathbf{a}), \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle$  to exam-  
 148 ine the alignment of the gradients while varying  $N$  on the  
 149 D4RL hopper-medium dataset. The results imply that the  
 150 performance of the learned policy degrades significantly  
 151 when the Q-functions share a similar local structure.

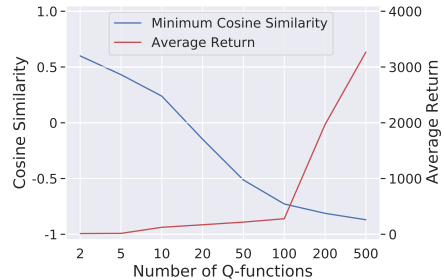


Figure 3: Plot of the minimum cosine similarity between the input gradients of Q-functions and the average return while varying the number of Q-functions.

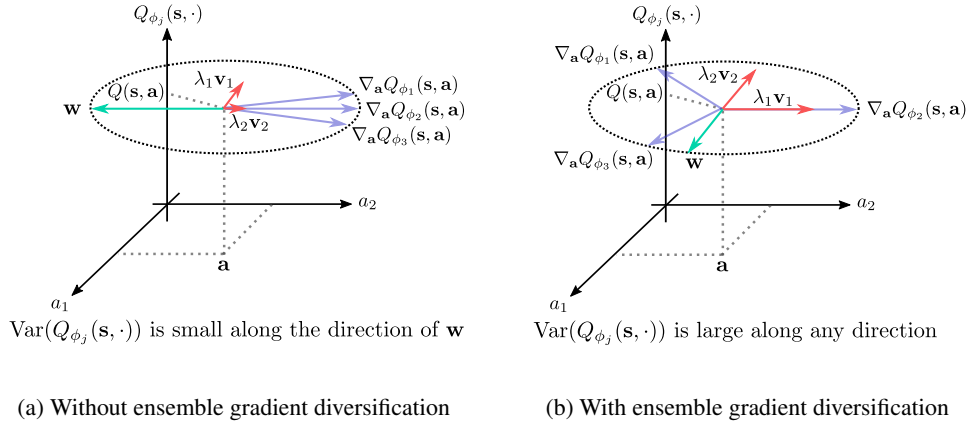


Figure 4: Illustration of the ensemble gradient diversification. The vector  $\lambda_i \mathbf{v}_i$  represents the normalized eigenvector  $\mathbf{v}_i$  of the variance matrix  $\text{Var}(\nabla_{\mathbf{a}} Q_j(\mathbf{s}, \mathbf{a}))$  multiplied by its eigenvalue  $\lambda_i$ .

152 We now show that the alignment of the input gradients can induce insufficient penalization of  
 153 near-distribution data points, which leads to requiring a large number of ensemble networks. Let  
 154  $\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a})$  be the gradient of the  $j$ -th Q-function with respect to the action  $\mathbf{a}$  and assume the  
 155 gradient is normalized for simplicity. If the gradients of the Q-functions are well-aligned as illustrated  
 156 in Figure 4a, then there exists a unit vector  $\mathbf{w}$  such that the Q-values for the OOD actions along the  
 157 direction of  $\mathbf{w}$  have a low variance. To show this, we first assume the Q-value predictions for the  
 158 in-distribution state-action pairs coincide, *i.e.*,  $Q_{\phi_j}(\mathbf{s}, \mathbf{a}) = Q(\mathbf{s}, \mathbf{a})$  for  $j = 1, \dots, N$ . Note that this  
 159 assumption can be easily satisfied by optimizing the Bellman error. Then, using the first-order Taylor  
 160 approximation, the sample variance of the Q-values at an OOD action along  $\mathbf{w}$  can be represented as

$$\begin{aligned}
 \text{Var}(Q_{\phi_j}(\mathbf{s}, \mathbf{a} + k\mathbf{w})) &\approx \text{Var}(Q_{\phi_j}(\mathbf{s}, \mathbf{a}) + k\langle \mathbf{w}, \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle) \\
 &= \text{Var}(Q(\mathbf{s}, \mathbf{a}) + k\langle \mathbf{w}, \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle) \\
 &= k^2 \text{Var}(\langle \mathbf{w}, \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle) \\
 &= k^2 \mathbf{w}^T \text{Var}(\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a})) \mathbf{w},
 \end{aligned} \tag{4}$$

161 where  $\langle \cdot, \cdot \rangle$  denotes an inner-product and  $k \in \mathbb{R}^+$ . The variance of the Q-values along  $\mathbf{w}$  is minimized  
 162 when  $\mathbf{w}$  is the normalized eigenvector corresponding to the smallest eigenvalue of  $\text{Var}(\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}))$ .  
 163 Let us denote this direction as  $\bar{\mathbf{w}}$ . From Proposition 1, if there exists  $\epsilon > 0$  such that for all  $i \neq j$ ,  
 164  $\langle \nabla_{\mathbf{a}} Q_{\phi_i}(\mathbf{s}, \mathbf{a}), \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle \geq 1 - \epsilon$ , then the variance of the Q-values along  $\bar{\mathbf{w}}$  is bounded by

$$\text{Var}(Q_{\phi_j}(\mathbf{s}, \mathbf{a} + k\bar{\mathbf{w}})) \leq \frac{1}{|\mathcal{A}|} \frac{N-1}{N} k^2 \epsilon, \tag{5}$$

165 where  $|\mathcal{A}|$  is the action space dimension.

166 **Lemma 1.** *The total variance of the matrix  $\text{Var}(\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}))$  is equal to  $1 - \|\bar{\mathbf{q}}\|_2^2$ , where  $\bar{\mathbf{q}} =$   
 167  $\frac{1}{N} \sum_{j=1}^N \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a})$ .*

168 **Proposition 1.** *If  $\exists \epsilon > 0$ , *s.t.*,  $\langle \nabla_{\mathbf{a}} Q_{\phi_i}(\mathbf{s}, \mathbf{a}), \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle \geq 1 - \epsilon$ ,  $\forall i \neq j$ , then the smallest  
 169 eigenvalue of the matrix  $\text{Var}(\nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}))$  is less than or equal to  $\frac{1}{|\mathcal{A}|} \frac{N-1}{N} \epsilon$ .*

170 We provide the proof in Supplementary material Appendix A. Equation (5) implies that if there exists  
 171 such  $\epsilon > 0$  that is small, which means the gradients of Q-function are well-aligned, then the variance  
 172 of the Q-values along  $\bar{\mathbf{w}}$  will also be small. This in turn degrades the ability of the ensembles to  
 173 penalize OOD actions, which ultimately leads to requiring a large number of ensemble networks.

174 To address this problem, we propose a regularizer that effectively increases the variance of the  
 175 Q-values for near-distribution OOD actions. One of the obvious ways to maximize this variance is to

176 compute the smallest eigenvalue of the variance matrix in Equation (4) and maximize it, which can  
 177 be formulated as

$$\underset{\phi}{\text{maximize}} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} \left[ \lambda_{\min} \left( \text{Var} \left( \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \right) \right) \right],$$

178 where  $\lambda_{\min}$  denotes the smallest eigenvalue and  $\phi$  denotes the collection of the parameters  $\{\phi_j\}_{j=1}^N$ .  
 179 There are several methods to compute the smallest eigenvalue, such as the power method or the  
 180 QR algorithm [27]. However, these iterative methods require construction of huge computation  
 181 graphs, which makes it computationally inefficient to optimize the eigenvalue using back-propagation.  
 182 Instead, we aim to maximize the sum of all eigenvalues, which is equal to the total variance. By  
 183 Lemma 1, it is equivalent to minimizing the norm of the average gradients:

$$\underset{\phi}{\text{minimize}} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} \left[ \left\langle \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{a}} Q_{\phi_i}(\mathbf{s}, \mathbf{a}), \frac{1}{N} \sum_{j=1}^N \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \right\rangle \right]. \quad (6)$$

184 With simple modification, we can reformulate Equation (6) as diversifying the gradients of each  
 185 Q-function network for in-distribution actions:

$$\underset{\phi}{\text{minimize}} J_{\text{ES}}(Q_{\phi}) := \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \mathcal{D}} \left[ \frac{1}{N-1} \sum_{1 \leq i \neq j \leq N} \underbrace{\langle \nabla_{\mathbf{a}} Q_{\phi_i}(\mathbf{s}, \mathbf{a}), \nabla_{\mathbf{a}} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) \rangle}_{\text{ES}_{\phi_i, \phi_j}(\mathbf{s}, \mathbf{a})} \right].$$

186 Concretely, our final objective formulates as measuring the pairwise alignment of the gradients  
 187 using cosine similarity, which we denote as the Ensemble Similarity (ES) metric  $\text{ES}_{\phi_i, \phi_j}(\mathbf{s}, \mathbf{a})$ , and  
 188 minimizing the ES values for every pair in the Q-ensemble with regard to the dataset state-actions.  
 189 The illustration of the ensemble gradient diversification is shown in Figure 4b.

---

### Algorithm 1 Ensemble-Diversified Actor Critic (EDAC)

---

- 1: Initialize policy parameters  $\theta$ , Q-function parameters  $\{\phi_j\}_{j=1}^N$ , target Q-function parameters  $\{\phi'_j\}_{j=1}^N$ , and offline data replay buffer  $\mathcal{D}$
- 2: **repeat**
- 3:   Sample a mini-batch  $B = \{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')\}$  from  $\mathcal{D}$
- 4:   Compute target Q-values (shared by all Q-functions):

$$y(r, \mathbf{s}') = r + \gamma \left( \min_{j=1, \dots, N} Q_{\phi'_j}(\mathbf{s}', \mathbf{a}') - \beta \log \pi_{\theta}(\mathbf{a}' | \mathbf{s}') \right), \quad \mathbf{a}' \sim \pi_{\theta}(\cdot | \mathbf{s}')$$

- 5:   Update each Q-function  $Q_{\phi_i}$  with gradient descent using

$$\nabla_{\phi_i} \frac{1}{|B|} \sum_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \in B} \left( \left( Q_{\phi_i}(\mathbf{s}, \mathbf{a}) - y(r, \mathbf{s}') \right)^2 + \frac{\eta}{N-1} \sum_{1 \leq i \neq j \leq N} \text{ES}_{\phi_i, \phi_j}(\mathbf{s}, \mathbf{a}) \right)$$

- 6:   Update policy with gradient ascent using

$$\nabla_{\theta} \frac{1}{|B|} \sum_{\mathbf{s} \in B} \left( \min_{j=1, \dots, N} Q_{\phi_j}(\mathbf{s}, \mathbf{a}) - \beta \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \right), \quad \mathbf{a} \sim \pi_{\theta}(\cdot | \mathbf{s})$$

- 7:   Update target networks with  $\phi'_i \leftarrow \rho \phi'_i + (1 - \rho) \phi_i$
- 

190 We name the resulting actor-critic algorithm as Ensemble-Diversified Actor Critic (EDAC) and  
 191 present the detailed procedure in Algorithm 1 (differences with the SAC algorithm marked in blue).  
 192 Note that Algorithm 1 reduces to SAC- $N$  when  $\eta = 0$ , and further reduces to vanilla SAC when also  
 193  $N = 2$ .

## 194 5 Experiments

195 We evaluate our proposed methods against the previous offline RL algorithms on the standard D4RL  
 196 benchmark [9]. Concretely, we perform our evaluation on MuJoCo Gym (Section 5.1) and Adroit

197 (Section 5.2) domains. We consider the following baselines: SAC, the backbone algorithm of our  
 198 method, CQL, the previous state-of-the-art on the D4RL benchmark, REM [2], an offline RL method  
 199 which utilized Q-network ensemble on discrete control environments, and BC, the behavior cloning  
 200 method. We evaluate each method under the normalized average return metric where the average  
 201 return is scaled such that 0 and 100 each equals the performance of a random policy and an online  
 202 expert policy. For the implementation details of our algorithm and the baselines, please refer to  
 203 Supplementary material Appendix B and C.

## 204 5.1 Evaluation on D4RL MuJoCo Gym tasks

205 We first evaluate each method on D4RL MuJoCo Gym tasks which consist of three environments,  
 206 halfcheetah, hopper, and walker2d, each with six datasets from different data-collecting policies.  
 207 In detail, the considered policies are *random*: a uniform random policy, *expert*: a fully trained  
 208 online expert, *medium*: a suboptimal policy with approximately 1/3 the performance of the expert,  
 209 *medium-expert*: a mixture of medium and expert policies, *medium-replay*: the replay buffer of a  
 210 policy trained up to the performance of the medium agent, and *full-replay*: the final replay buffer of  
 211 the expert policy. Each dataset consists of 1M transitions except for medium-expert which combines  
 212 the medium and expert datasets.

Table 1: Normalized average returns on D4RL Gym tasks, averaged over 4 random seeds. CQL (Paper) denotes the results reported in the original paper.

Task Name	BC	SAC	REM	CQL (Paper)	CQL (Reproduced)	SAC- <i>N</i> (Ours)	EDAC (Ours)
halfcheetah-random	2.2±0.0	29.7±1.4	-0.8±1.1	<b>35.4</b>	31.3±3.5	28.0±0.9	28.4±1.0
halfcheetah-medium	43.2±0.6	55.2±27.8	-0.8±1.3	44.4	46.9±0.4	<b>67.5±1.2</b>	<b>65.9±0.6</b>
halfcheetah-expert	91.8±1.5	-0.8±1.8	4.1±5.7	104.8	97.3±1.1	<b>105.2±2.6</b>	<b>106.8±3.4</b>
halfcheetah-medium-expert	44.0±1.6	28.4±19.4	0.7±3.7	62.4	95.0±1.4	<b>107.1±2.0</b>	<b>106.3±1.9</b>
halfcheetah-medium-replay	37.6±2.1	0.8±1.0	6.6±11.0	46.2	45.3±0.3	<b>63.9±0.8</b>	<b>61.3±1.9</b>
halfcheetah-full-replay	62.9±0.8	<b>86.8±1.0</b>	27.8±35.4	-	76.9±0.9	84.5±1.2	84.6±0.9
hopper-random	3.7±0.6	9.9±1.5	3.4±2.2	10.8	5.3±0.6	<b>31.3±0.0</b>	<b>25.3±10.4</b>
hopper-medium	54.1±3.8	0.8±0.0	0.7±0.0	86.6	61.9±6.4	<b>100.3±0.3</b>	<b>101.6±0.6</b>
hopper-expert	107.7±9.7	0.7±0.0	0.8±0.0	109.9	106.5±9.1	<b>110.3±0.3</b>	<b>110.1±0.1</b>
hopper-medium-expert	53.9±4.7	0.7±0.0	0.8±0.0	<b>111.0</b>	96.9±15.1	110.1±0.3	110.7±0.1
hopper-medium-replay	16.6±4.8	7.4±0.5	27.5±15.2	48.6	86.3±7.3	<b>101.8±0.5</b>	<b>101.0±0.5</b>
hopper-full-replay	19.9±12.9	41.1±17.9	19.7±24.6	-	101.9±0.6	<b>102.9±0.3</b>	<b>105.4±0.7</b>
walker2d-random	1.3±0.1	0.9±0.8	6.9±8.3	7.0	5.4±1.7	<b>21.7±0.0</b>	<b>16.6±7.0</b>
walker2d-medium	70.9±11.0	-0.3±0.2	0.2±0.7	74.5	79.5±3.2	<b>87.9±0.2</b>	<b>92.5±0.8</b>
walker2d-expert	108.7±0.2	0.7±0.3	1.0±2.3	<b>121.6</b>	109.3±0.1	107.4±2.4	115.1±1.9
walker2d-medium-expert	90.1±13.2	1.9±3.9	-0.1±0.0	98.7	109.1±0.2	<b>116.7±0.4</b>	<b>114.7±0.9</b>
walker2d-medium-replay	20.3±9.8	-0.4±0.3	12.5±6.2	32.6	76.8±10.0	<b>78.7±0.7</b>	<b>87.1±2.3</b>
walker2d-full-replay	68.8±17.7	27.9±47.3	-0.2±0.3	-	94.2±1.9	<b>94.6±0.5</b>	<b>99.8±0.7</b>
Average	49.9	16.2	6.2	-	73.7	<b>84.5</b>	<b>85.2</b>

213 The experiment results in Table 1 show EDAC and SAC-*N* both outperform or are competitive with  
 214 the previous state-of-the-art on all of the tasks considered. Notably, the performance gap is especially  
 215 high for random, medium, and medium-replay datasets, where the performances of the previous  
 216 works are relatively low. Both the proposed methods achieve average normalized scores over 80,  
 217 reducing the gap with the online expert by 40% compared to CQL. While the performance of EDAC  
 218 is marginally better than the performance of SAC-*N*, EDAC achieves this result with a much smaller  
 219 Q-ensemble size. As noted in Figure 5, on hopper tasks, SAC-*N* requires 200 to 500 Q-networks,  
 220 while EDAC requires less than 50.

221 Figure 6 compares the distance between the actions chosen by each method and the dataset actions.  
 222 Concretely, we measure  $\mathbb{E}_{(s, \mathbf{a}) \sim \mathcal{D}, \hat{\mathbf{a}} \sim \pi_{\theta}(\cdot | s)} [(\hat{\mathbf{a}} - \mathbf{a})^2]$  for EDAC, SAC-*N*, CQL, and a random policy  
 223 on \*-medium datasets. We find that our proposed methods choose from a more diverse range of  
 224 actions compared to CQL. This shows the advantage of the uncertainty-based penalization which  
 225 considers the prediction confidence other than penalizing all OOD actions.

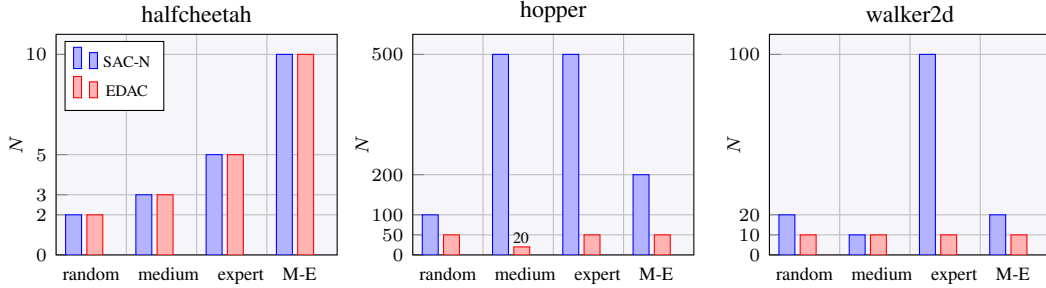


Figure 5: Minimum number of Q-ensembles ( $N$ ) required to achieve the performance reported in Table 1. M-E denotes medium-expert. We omit the results of medium-replay and full-replay as SAC- $N$  already works well with a small number of ensembles (less than or equal to 5). For more details of the experiment, please refer to Supplementary material Appendix C.

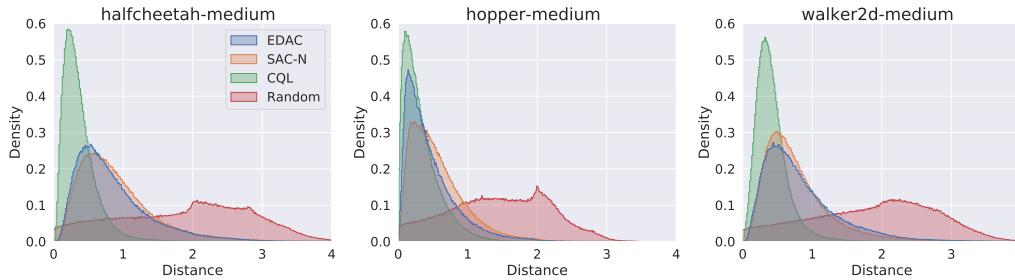


Figure 6: Histograms of the distances between the actions from each methods (EDAC, SAC- $N$ , CQL, and a random policy) and the actions from the dataset. For more details of the experiment, please refer to Supplementary material Appendix C.

## 226 5.2 Evaluation on D4RL Adroit tasks

227 We also experiment on the more complex D4RL Adroit tasks that require controlling a 24-DoF robotic  
 228 hand to perform tasks such as aligning a pen, hammering a nail, opening a door, or relocating a ball.  
 229 We use two types of datasets for each environment: *human*, containing 25 trajectories of human  
 230 demonstrations, and *cloned*, a 50-50 mixture between the demonstration data and the behavioral  
 231 cloned policy on the demonstrations. Note that for the Adroit tasks, we could not reproduce the CQL  
 232 results from the paper completely. For the detailed procedure of reproducing the results of CQL,  
 233 please refer to Supplementary material Appendix D.

Table 2: Normalized average returns on D4RL Adroit tasks, averaged over 4 random seeds.

Task Name	BC	SAC	REM	CQL (Paper)	CQL (Reproduced)	SAC- $N$ (Ours)	EDAC (Ours)
pen-human	25.8±8.8	4.3±3.8	5.4±4.3	<b>55.8</b>	35.2±6.6	9.5±1.1	52.1±8.6
hammer-human	3.1±3.2	0.2±0.0	0.3±0.0	2.1	0.6±0.5	0.3±0.0	0.8±0.4
door-human	2.8±0.7	-0.3±0.0	-0.3±0.0	9.1	1.2±1.8	-0.3±0.0	<b>10.7±6.8</b>
relocate-human	0.0±0.0	-0.3±0.0	-0.3±0.0	0.35	0.0±0.0	-0.1±0.1	0.1±0.1
pen-cloned	38.3±11.9	-0.8±3.2	-1.0±0.1	40.3	27.2±11.3	<b>64.1±8.7</b>	<b>68.2±7.3</b>
hammer-cloned	0.7±0.3	0.1±0.1	-0.3±0.0	5.7	1.4±2.1	0.2±0.2	0.3±0.0
door-cloned	0.0±0.0	-0.3±0.1	-0.3±0.0	3.5	2.4±2.4	-0.3±0.0	<b>9.6±8.3</b>
relocate-cloned	0.1±0.0	-0.1±0.1	-0.2±0.2	-0.1	0.0±0.0	0.0±0.0	0.0±0.0

234 The evaluation results are summarized in Table 2. For pen-\* tasks, where the considered algorithms  
 235 achieve meaningful performance, EDAC outperforms or matches with the previous state-of-the-art.  
 236 Especially, for pen-cloned, both EDAC and SAC- $N$  achieve 75% higher score compared to CQL.  
 237 Unlike the results from the Gym tasks, we find that SAC- $N$  falls behind in some datasets, for

238 example, pen-human, which could in part due to the size of the dataset being exceptionally small  
239 (5000 transitions). However, our method with ensemble diversification successfully overcomes this  
240 difficulty.

## 241 6 Related Works

242 **Model-free offline RL** A popular approach for offline RL is to regularize the learned policy to be  
243 close to the behavior policy where the offline dataset was collected. BCQ [11] uses a generative  
244 model to produce actions with high similarity to the dataset and trains a restricted policy to choose  
245 the best action from the neighborhood of the generated actions. Another line of work, such as BEAR  
246 [15] or BRAC [28], stabilizes policy learning by penalizing the divergence from the dataset measured  
247 by KL divergence or MMD. While these policy-constraint methods demonstrate high performance on  
248 datasets from expert behavior policies, they fail to find optimal policies from datasets with suboptimal  
249 policies due to the strict policy constraints [9]. Also, these methods require an accurate estimation of  
250 the behavior policy, which might be difficult in complex settings with multiple behavior sources or  
251 high-dimensional environments. To address these issues, CQL [16] directly regularizes Q-functions  
252 by introducing a term that minimizes the Q-values for out-of-distribution actions and maximizes the  
253 Q-values for in-distribution actions. Without such explicit regularizations, REM [2] proposes to use a  
254 random convex combination of Q-network ensembles on environments with discrete action spaces  
255 [4].

256 **Estimation bias in Q-learning** While Q-learning is one of the most popular algorithms in reinforcement  
257 learning, it suffers from overestimation bias due to the maximum operation  $\max_{\mathbf{a}' \in \mathcal{A}} Q(\mathbf{s}', \mathbf{a}')$   
258 used during Q-function updates [10, 25]. This overestimation bias, together with the bootstrapping,  
259 can lead to a catastrophic build-up of errors during the Q-learning process. To resolve this issue, TD3  
260 [10] introduces a clipped version of Double Q-learning [25] that takes the minimum value of two  
261 critics. Subsequently, Maxmin Q-learning [17] theoretically shows that the overestimation bias can  
262 be controlled by the number of ensembles in the clipped Q-learning. The overestimation problem in  
263 Q-learning can be exacerbated in the offline setting since the extrapolation error cannot be corrected  
264 with further interactions with the environment, and existing offline RL algorithms handle the bias by  
265 introducing constrained policy optimization [11, 15] or conservative Q-learning frameworks [16].

266 **Uncertainty measures in RL** Uncertainty estimates have been widely used in RL for various  
267 purposes including exploration, Q-learning, and planning. Bootstrapped DQN [21] leverages an  
268 ensemble of Q-functions to quantify the uncertainty of the Q-value, and utilizes it for efficient  
269 exploration. Following this work, the UCB exploration algorithm [5] constructs an upper confidence  
270 bound [3] of the Q-values using the empirical mean and standard deviation of Q-ensembles, which is  
271 used to promote efficient exploration by applying the principle of optimism in the face of uncertainty  
272 [7]. Osband et al. [22] proposes a randomly initialized Q-ensemble that reflects the concept of prior  
273 functions in Bayesian inference and Abbas et al. [1] introduces an uncertainty incorporated planning  
274 with imperfect models. The notion of uncertainty has also been considered in offline RL, mostly  
275 in the framework of model-based offline RL. Especially, MOPO [29] and MOREL [13] measure  
276 the uncertainty of the model’s prediction to formulate an uncertainty-penalized policy optimization  
277 problem in the offline RL setting. These methods introduce an ensemble of dynamics models  
278 for the quantification of the uncertainty, whereas our work adopts an ensemble of Q-functions for  
279 uncertainty-aware Q-learning.

## 280 7 Conclusion

281 We have shown that clipped Q-learning can be efficiently leveraged to construct an uncertainty-  
282 based offline RL method that outperforms previous methods on various datasets. Based on this  
283 observation, we proposed Ensemble-Diversifying Actor-Critic (EDAC) that effectively reduces the  
284 required number of ensemble networks for quantifying and penalizing the epistemic uncertainty. Our  
285 method does not require any explicit estimation of the data collecting policy or sampling from the  
286 out-of-distribution data and respects the epistemic uncertainty of each data point during penalization.  
287 EDAC, while requiring up to 90% less number of ensemble networks compared to the vanilla  
288 Q-ensemble, exhibits state-of-the-art performance on various datasets.

## 289 References

- 290 [1] Zaheer Abbas, Samuel Sokota, Erin Talvitie, and Martha White. Selective dyna-style planning  
291 under limited model capacity. In *ICML*, 2020.
- 292 [2] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on  
293 offline reinforcement learning. In *ICML*, 2020.
- 294 [3] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff  
295 using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):  
296 1876–1902, 2009.
- 297 [4] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning  
298 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*  
299 *Research*, 47:253–279, 2013.
- 300 [5] Richard Y Chen, Szymon Sidor, Pieter Abbeel, and John Schulman. Ucb exploration via  
301 q-ensembles. *arXiv preprint arXiv:1706.01502*, 2017.
- 302 [6] Xinshi Chen, Shuang Li, Hui Li, Shaohua Jiang, and Le Song. Generative adversarial user  
303 model for reinforcement learning based recommendation system. In *ICML*, 2019.
- 304 [7] Kamil Ciosek, Quan Vuong, Robert Loftin, and Katja Hofmann. Better exploration with  
305 optimistic actor-critic. In *NeurIPS*, 2019.
- 306 [8] William R Clements, Bastien Van Delft, Benoît-Marie Robaglia, Reda Bahi Slaoui, and  
307 Sébastien Toth. Estimating risk and uncertainty in deep reinforcement learning. *arXiv preprint*  
308 *arXiv:1905.09638*, 2019.
- 309 [9] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
310 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 311 [10] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in  
312 actor-critic methods. In *ICML*, 2018.
- 313 [11] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning  
314 without exploration. In *ICML*, 2019.
- 315 [12] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
316 maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, 2018.
- 317 [13] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:  
318 Model-based offline reinforcement learning. In *NeurIPS*, 2020.
- 319 [14] Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *NeurIPS*, 2000.
- 320 [15] Aviral Kumar, Justin Fu, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning  
321 via bootstrapping error reduction. In *NeurIPS*, 2019.
- 322 [16] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for  
323 offline reinforcement learning. In *NeurIPS*, 2020.
- 324 [17] Qingfeng Lan, Yangchen Pan, Alona Fyshe, and Martha White. Maxmin q-learning: Controlling  
325 the estimation bias of q-learning. In *ICLR*, 2020.
- 326 [18] Kimin Lee, Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Sunrise: A simple unified  
327 framework for ensemble learning in deep reinforcement learning. In *ICML*, 2021.
- 328 [19] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
329 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 330 [20] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
331 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
332 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

- 333 [21] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via  
334 bootstrapped dqn. In *NeurIPS*, 2016.
- 335 [22] Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep rein-  
336 forcement learning. In *NeurIPS*, 2018.
- 337 [23] JP Royston. Expected normal order statistics(exact and approximate). *Applied Statistics*, 31(2):  
338 161–5, 1982.
- 339 [24] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine  
340 learning algorithms. *arXiv preprint arXiv:1206.2944*, 2012.
- 341 [25] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double  
342 q-learning. In *AAAI*, 2016.
- 343 [26] Oriol Vinyals, Igor Babuschkin, Junyoung Chung, Michael Mathieu, Max Jaderberg, Wo-  
344 jciech M Czarnecki, Andrew Dudzik, Aja Huang, Petko Georgiev, Richard Powell, et al.  
345 Alphastar: Mastering the real-time strategy game starcraft ii. *DeepMind blog*, 2, 2019.
- 346 [27] David S Watkins. Understanding the qr algorithm. *SIAM review*, 24(4):427–440, 1982.
- 347 [28] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
348 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 349 [29] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea  
350 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *NeurIPS*, 2020.

## 351 Checklist

352 The checklist follows the references. Please read the checklist guidelines carefully for information on  
353 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
354 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
355 the appropriate section of your paper or providing a brief inline description. For example:

- 356 • Did you include the license to the code and datasets? **[Yes]** See Section xxx.
- 357 • Did you include the license to the code and datasets? **[No]** The code and the data are  
358 proprietary.
- 359 • Did you include the license to the code and datasets? **[N/A]**

360 Please do not modify the questions and only use the provided macros for your answers. Note that the  
361 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
362 block and only keep the Checklist section heading above along with the questions/answers below.

- 363 1. For all authors...
  - 364 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
365 contributions and scope? **[Yes]** We reflect the main claims through Section 2 to  
366 Section 5.
  - 367 (b) Did you describe the limitations of your work? **[Yes]** See Section 4.
  - 368 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
  - 369 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
370 them? **[Yes]** We have checked the ethics review guidelines.
- 371 2. If you are including theoretical results...
  - 372 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** The assump-  
373 tions are provided with the proposition in Section 4.
  - 374 (b) Did you include complete proofs of all theoretical results? **[Yes]** We provide the proofs  
375 in Supplementary material Appendix A.
- 376 3. If you ran experiments...

- 377 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
378 mental results (either in the supplemental material or as a URL)? [Yes] We include the  
379 code and the instructions along with the Supplementary material.
- 380 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
381 were chosen)? [Yes] The implementation details are listed in Supplementary material  
382 Appendix B and C.
- 383 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
384 ments multiple times)? [Yes] We report the standard deviation of the evaluation metric  
385 on multiple runs.
- 386 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
387 type of GPUs, internal cluster, or cloud provider)? [Yes] We describe the computation  
388 resource in Supplementary material Appendix B.
- 389 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 390 (a) If your work uses existing assets, did you cite the creators? [Yes] Our work cites the  
391 dataset used for our experiments.
- 392 (b) Did you mention the license of the assets? [N/A]
- 393 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 394
- 395 (d) Did you discuss whether and how consent was obtained from people whose data you're  
396 using/curating? [N/A]
- 397 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
398 information or offensive content? [N/A]
- 399 5. If you used crowdsourcing or conducted research with human subjects...
- 400 (a) Did you include the full text of instructions given to participants and screenshots, if  
401 applicable? [N/A]
- 402 (b) Did you describe any potential participant risks, with links to Institutional Review  
403 Board (IRB) approvals, if applicable? [N/A]
- 404 (c) Did you include the estimated hourly wage paid to participants and the total amount  
405 spent on participant compensation? [N/A]