

OFFLINE REINFORCEMENT LEARNING WITH VALUE-BASED EPISODIC MEMORY

Xiaoteng Ma^{1*}, Yiqin Yang^{1*}, Hao Hu^{2*}, Qihan Liu¹,
Jun Yang^{1†}, Chongjie Zhang^{2†}, Qianchuan Zhao¹, Bin Liang¹

¹Department of Automation, Tsinghua University

²Institute for Interdisciplinary Information Sciences, Tsinghua University

{ma-xt17, yangyiqi19, hu-h19, lqh20}@mails.tsinghua.edu.edu

{yangjun603, chongjie, zhaoqc, bliang}@tsinghua.edu.cn

ABSTRACT

Offline reinforcement learning (RL) shows promise of applying RL to real-world problems by effectively utilizing previously collected data. Most existing offline RL algorithms use regularization or constraints to suppress extrapolation error for actions outside the dataset. In this paper, we adopt a different framework, which learns the V -function instead of the Q -function to naturally keep the learning procedure within the offline dataset. To enable effective generalization while maintaining proper conservatism in offline learning, we propose Expectile V -Learning (EVL), which smoothly interpolates between the optimal value learning and behavior cloning. Further, we introduce implicit planning along offline trajectories to enhance learned V -values and accelerate convergence. Together, we present a new offline method called Value-based Episodic Memory (VEM). We provide theoretical analysis for the convergence properties of our proposed VEM method, and empirical results in the D4RL benchmark show that our method achieves superior performance in most tasks, particularly in sparse-reward tasks. Our code is public online at <https://github.com/YiqinYang/VEM>.

1 INTRODUCTION

Despite the great success of deep reinforcement learning (RL) in various domains, most current algorithms rely on interactions with the environment to learn through trial and error. In real-world problems, particularly in risky and safety-crucial scenarios, interactions with the environment can be expensive and unsafe, and only offline collected datasets are available, such as the expert demonstration or previously logged data. This growing demand has led to the emergence of *offline reinforcement learning* (offline RL) to conduct RL in a supervised manner.

The main challenge of offline RL comes from the actions out of the dataset’s support (Kumar et al., 2019; 2020). The evaluation of these actions that do not appear in the dataset relies on the generalization of the value network, which may exhibit *extrapolation error* (Fujimoto et al., 2019). This error can be magnified through bootstrapping, leading to severe estimation errors. A rapidly developing line of recent work (Fujimoto et al., 2019; Kumar et al., 2020; Ghasemipour et al., 2021; Yang et al., 2021) utilizes various methods to constrain optimistic estimation on unseen actions, such as restricting available actions with a learned behavior model (Fujimoto et al., 2019) or penalizing the unseen actions with additional regularization (Kumar et al., 2020). However, confining learning *within the distribution of the dataset* can be insufficient for reducing extrapolation errors.

Another line of methods, on the contrary, uses the returns of the behavior policy as the signal for policy learning, as adopted in Wang et al. (2018); Peng et al. (2019); Chen et al. (2020). By doing so, they keep the value learning procedure completely *within the dataset*. However, the behavior policy of the dataset can be imperfect and insufficient to guide policy learning. To achieve a trade-off between imitation learning and optimal value learning while confines learning *within the dataset*,

*Equal contribution. Listing order is random.

†Equal advising.

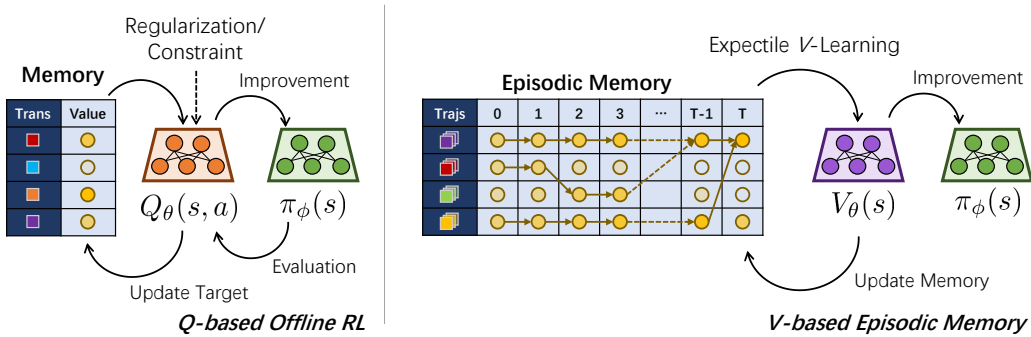


Figure 1: The diagram of algorithms. The left side denotes the general Q -based offline RL methods. The right side is the framework of our proposed approach (VEM). Q -based methods learn bootstrapped Q -values, but require additional constraint or penalty for actions out of the dataset. Our method, on the contrary, learns bootstrapped V -values while being completely confined within the dataset without any regularization.

we propose Expectile V -learning (EVL), which is based on a new expectile operator that smoothly interpolates between the Bellman expectation operator and optimality operator.

To better solve long-horizon and sparse-reward tasks, we further propose using value-based planning to improve the advantage estimation for policy learning. We adopt an implicit memory-based planning scheme that strictly plans within offline trajectories to compute the advantages effectively, as proposed in recent advances in episodic memory-based methods (Hu et al., 2021). Together, we present our novel framework for offline RL, Value-based Episodic Memory (VEM), which uses expectile V -learning to approximate the optimal value with offline data and conduct implicit memory-based planning to further enhance advantage estimation. With the properly learned advantage function, VEM trains the policy network in a simple regression manner. We demonstrate our algorithm in Figure 1, and a formal description of our algorithm is provided in Algorithm 1.

The contributions of this paper are threefold. First, we present a new offline V -learning method, EVL, and a novel offline RL framework, VEM. EVL learns the value function through the trade-offs between imitation learning and optimal value learning. VEM uses a memory-based planning scheme to enhance advantage estimation and conduct policy learning in a regression manner. Second, we theoretically analyze our proposed algorithm’s convergence properties and the trade-off between contraction rate, fixed-point bias, and variance. Specifically, we show that VEM is provably convergent and enjoys a low concentration rate with a small fixed-point bias. Finally, we evaluate our method in the offline RL benchmark D4RL (Fu et al., 2020). Comparing with other baselines, VEM achieves superior performance, especially in the sparse reward tasks like AntMaze and Adroit. The ablation study shows that VEM yields accurate value estimates and is robust to extrapolation errors.

2 BACKGROUND

Preliminaries. We consider a Markov Decision Process (MDP) M defined by a tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $P(\cdot | s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the transition distribution function, $r(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function and $\gamma \in [0, 1)$ is the discount factor. We say an environment is deterministic if $P(s' | s, a) = \delta(s' = f(s, a))$ for some deterministic transition function f , where $\delta(\cdot)$ is the Dirac function. The goal of an RL agent is to learn a policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, which maximizes the expectation of a discounted cumulative reward: $\mathcal{J}(\pi) = \mathbb{E}_{s_0 \sim \rho_0; a_t \sim (\cdot | s_t); s_{t+1} \sim P(\cdot | s_t; a_t)} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$, where ρ_0 is the distribution of the initial states.

Value-based Offline Reinforcement Learning Methods. Current offline RL methods can be roughly divided into two categories according to types of learned value function: Q -based and V -based methods. Q -based methods, such as BCQ (Fujimoto et al., 2019), learn Q -function for policy learning and avoid selecting unfamiliar actions via constraints or penalty. On the contrary, V -based methods (Peng et al., 2019; Siegel et al., 2020; Chen et al., 2020) learn the value of behavior policy $V(s)$ with the trajectories in the offline dataset \mathcal{D} and update policy as a regression problem. Based on the learned V -function, V -based methods like AWR (Peng et al., 2019) update the policy using advantage-weighted regression, where each state-action pair is weighted according

to the exponentiated advantage:

$$\max_{\mathcal{J}} (\phi) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} [\log \pi(a_t | s_t) \exp(R_t - V(s_t))]. \quad (1)$$

Episodic Memory-Based Methods. Inspired by psychobiology, episodic memory-based methods store experiences in a non-parametric table to fast retrieve past successful strategies when encountering similar states. Model-free episodic control (Blundell et al., 2016a) updates the memory table by taking the maximum return $R(s, a)$ among all rollouts starting from same state-action pair (s, a) . Hu et al. (2021) proposes *Generalizable Episodic Memory*, which extends this idea to the continuous domain, and proposes updating formula with a parametric memory Q^{EM} .

3 METHOD

In this section, we describe our novel offline method, value-based episodic memory, as depicted in Figure 1. VEM uses expectile V -learning (EVL) to learn V -functions while confines value learning *within the dataset* to reduce extrapolation error. EVL uses an expectile operator that interpolates between Bellman expectation operator and optimality operator to balance behavior cloning and optimal value learning. Further, VEM integrates memory-based planning to improve the advantage estimation and accelerate the convergence of EVL. Finally, generalized advantage-weighted learning is used for policy learning with enhanced advantage estimation. A formal description for the VEM algorithm is shown in Algorithm 1 in Appendix A.1.

3.1 EXPECTILE V-LEARNING

To achieve a balance between behavior cloning and optimal value learning, we consider the Bellman expectile operator defined as follows:

$$((\mathcal{T}^\tau)V)(s) := \arg \min_v \mathbb{E}_{a \sim \mu(\cdot|s)} \tau [\delta(s, a)]_+^2 + (1 - \tau) [\delta(s, a)]_-^2 \quad (2)$$

where μ is the behavior policy, $\delta(s, a) = \mathbb{E}_{s' \sim P(\cdot|s, a)} [r(s, a) + \gamma V(s') - v]$ is the expected one-step TD error, $[\cdot]_+ = \max(\cdot, 0)$ and $[\cdot]_- = \min(\cdot, 0)$. This operator resembles the *expectile* statistics (Newey & Powell, 1987; Rowland et al., 2019) and hence its name. We can see that when $\tau = 1/2$, this operator is reduced to Bellman expectation operator, while when $\tau \rightarrow 1$, this operator approaches Bellman optimality operator, as depicted in Lemma 3.

We use the following toy example to further illustrate the trade-offs achieved by EVL. Consider a random generated MDP. When the operator can be applied exactly, the Bellman optimality operator is sufficient to learn the optimal value V^* . However, applying operators with an offline dataset raises a noise on the actual operator due to the estimation error with finite and biased data. We simulate this effect by adding random Gaussian noise to the operator. Applying the optimality operator on offline datasets can lead to severe overestimation due to the maximization bias and bootstrapping. The value estimation learned by EVL, on the contrary, achieves a trade-off between learning optimal policy and behavior cloning and can be close to the optimal value with proper chosen τ , as depicted in Figure 2. The noise upon the operator largely depends on the size of the dataset. Estimation error can be significant with insufficient data. In this case, we need a small τ to be conservative and be close to behavior cloning. When the dataset is large and we are able to have an accurate estimation for the operator, we can use a larger τ to recover the optimal policy. By adjusting τ , the expectile operator can accommodate variant types of datasets. However, the expectile operator in Equation 2 does not have a

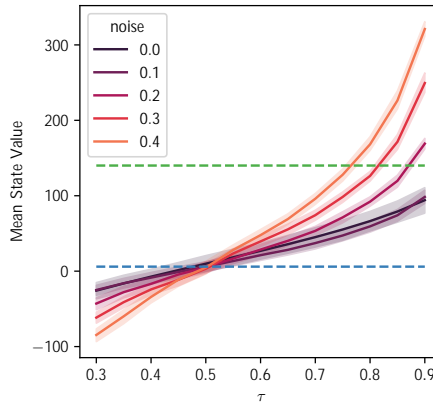


Figure 2: Trade-offs of EVL between generalization and conservatism in a random MDP. The green line shows the optimal value and the blue line shows the value of behavior policy. The curve is averaged over 20 MDPs.

closed-form solution. In practice, we consider the one-step gradient expectile operator

$$((T_g)V)(s) = V(s) + 2\epsilon E_{a \sim \pi(s)} [V(s; a)]_+ + (1 - \epsilon) [V(s; a)]_-; \quad (3)$$

where ϵ is the step-size. Please refer to Appendix B.1 for the detailed derivation. For notational convenience, we use \mathbb{E} to denote the one-step gradient expectile operator hereafter.

We consider the case where the dynamics are nearly-deterministic like robotic applications, and we remove the expectation over the next states in the operator. This leads to a practical algorithm, ExpectileV-Learning, where we train the value network to minimize the following loss:

$$J_V(\theta) = E_{(s; a; s^0) \sim D} \|\hat{V}(s) - V(s; a)\|^2; \quad (4)$$

$$\hat{V}(s) = V(s; a) + 2\epsilon [V(s; a)]_+ + (1 - \epsilon) [V(s; a)]_-;$$

where \hat{V} is the target value after applying one-step gradient expectile operator ($\hat{V}(s) = r(s; a) + \epsilon (V(s; a) - V(s; s^0)) + V(s; s^0)$). V -function and the target \hat{V} -function are parameterized by θ , respectively. EVL is guaranteed to converge with concentration rate $\epsilon^{-1} 2(1 - \epsilon) \max_{s, a} |V(s; a) - V(s; s^0)|$. Please refer to Section 4 for a detailed analysis.

3.2 IMPLICIT MEMORY-BASED PLANNING

Although EVL reduces the extrapolation error, it is still a challenging problem to bootstrap over long time horizons due to estimation errors with a fixed dataset. Therefore, we propose using value-based planning to conduct bootstrapping more efficiently. We adopt an implicit memory-based planning scheme that strictly plans within of the trajectories to avoid over-optimistic estimations in the planning phase. This is aligned with recent advances in episodic memory-based methods (Hu et al., 2021), but we conduct this planning on expectile values rather than Q -values. Specifically, we compare the best return so far along the trajectory with the value estimation. \hat{R}_t takes the maximum between them to get the augmented return

$$\hat{R}_t = \begin{cases} r_t + \max(\hat{R}_{t+1}, \hat{V}(s_{t+1})); & \text{if } t < T; \\ r_t; & \text{if } t = T; \end{cases} \quad (5)$$

where t denotes steps along the trajectory, T is the episode length, and \hat{V} is generalized from similar experiences. This procedure is conducted recursively from the last step to the first step along the trajectory, forming an implicit planning scheme within the dataset to aggregate experiences along and across trajectories. Further, the back-propagation process in Equation 5 can be unrolled and rewritten as follows:

$$\hat{R}_t = \max_{0 < n < n_{\max}} \hat{V}_{t,n}; \quad \hat{V}_{t,n} = \begin{cases} r_t + \hat{V}_{t+1,n-1} & \text{if } n > 0; \\ \hat{V}(s_t) & \text{if } n = 0; \end{cases} \quad (6)$$

where n denotes different length of rollout steps and $\hat{V}_{t,n} = 0$ for $n > T$.

3.3 GENERALIZED ADVANTAGE-WEIGHTED LEARNING

Based on \hat{R}_t calculated in Section 3.2, we can conduct policy learning in a regression form, as adopted in return-based of the RL methods (Nair et al., 2020; Siegel et al., 2020; Peng et al., 2019):

$$\max_{\theta} J(\theta) = E_{(s_t; a_t) \sim D} \log \sum_j \pi_j(s_t; a_t) f_j(\hat{A}(s_t; a_t)); \quad (7)$$

where $\hat{A}(s_t; a_t) = \hat{R}_t - \hat{V}(s_t)$ and f_j is an increasing, non-negative function. Please refer to Appendix C.1 for the detailed implementation of Equation 7. Note that \hat{R}_t is not the vanilla returns in the dataset, but the enhanced estimation calculated by implicit planning. It is compared with other return based methods. Please refer to Algorithm 1 and Section 4 for implementation details and theoretical analysis.

4 THEORETICAL ANALYSIS

In this section, we first derive the convergence property of expectile learning. Then, we demonstrate that memory-based planning accelerates the convergence of the EVL. Finally, we design a toy example to demonstrate these theoretical analyses empirically. Please refer to Appendix B for the detailed proofs of the following analysis.

4.1 CONVERGENCE PROPERTY OF THE EXPECTILE V-LEARNING

In this section, we assume the environment is deterministic. We derive the contraction property of T as the following statement:

Lemma 1. For any $\beta \in (0, 1)$, T is a β -contraction, where $\beta = 1 - \frac{2}{2 + \frac{1}{\min f; 1} - g}$.

Proof. We introduce two more operators to simplify the analysis:

$$(T_+ V)(s) = V(s) + E_a [r(s; a)]_+; (T_- V)(s) = V(s) + E_a [r(s; a)]_- \quad (8)$$

Next we show that both operators are non-expansion (i.e., $\|T_+ V_1 - T_+ V_2\|_k \leq \|V_1 - V_2\|_k$). Finally, we rewrite T based on T_+ and T_- and we prove that T is a β -contraction. Please refer to Appendix B.2 for the complete proof. \square

Based on Lemma 1, we give a discussion about the step-size and the fraction β :

About the step-size α . Generally, we always want a larger α . However, α must satisfy that $V(s) + 2\alpha \max_a r(s; a) \leq V(s) + g$ and $V(s) + 2\alpha \min_a r(s; a) \geq V(s) - g$, otherwise the V -value will be overestimated. Thus, we must have $\alpha \leq \frac{1}{2(\frac{1}{\min f; 1} - g)}$ and $\alpha \geq \frac{1}{2(\frac{1}{\max f; 1} + g)}$, which infers that $\alpha \in [\frac{1}{2(\frac{1}{\max f; 1} + g)}, \frac{1}{2(\frac{1}{\min f; 1} - g)}]$. When $\alpha = \frac{1}{2(\frac{1}{\max f; 1} + g)}$, we have $\beta = 1 - \frac{2}{2 + \frac{1}{\min f; 1} - g} = 1 - \frac{\frac{1}{\max f; 1} + g}{1 + \frac{1}{\max f; 1} + g}$.

About the fraction β . It is easy to verify that β approaches to 1 when $\frac{1}{\max f; 1} \rightarrow 0$ or $\frac{1}{\min f; 1} \rightarrow 1$, which means that with a larger $\frac{1}{\max f; 1}$ the contractive property is getting weaker. The choice of α makes a trade-off between the learning stability and the optimality of values. We further point out that when $\alpha = \frac{1}{2(\frac{1}{\max f; 1} + g)}$, the Expectile V -learning degrades as a special case of the generalized self-imitation learning (Tang, 2020), which loses the contractive property.

Next, we prove that T is monotonous improving with respect to

Lemma 2. For any $\beta \in (0, 1)$, if $\beta > 0$, we have $T^\beta V(s) \geq T V(s)$, $\forall s \in \mathcal{S}$.

Based on the Lemma 2, we derive that T is monotonous improving with respect to

Proposition 1. Let V^* denote the fixed point of T . For any $\beta \in (0, 1)$, if $\beta > 0$, we have $V^\beta(s) \leq V^*(s)$, $\forall s \in \mathcal{S}$.

Further, we derive that T gradually approaches V^* with respect to:

Lemma 3. Let V^* denote the fixed point of Bellman optimality operator \mathcal{B} . In the deterministic MDP, we have $\lim_{n \rightarrow \infty} T^n V = V^*$.

Based on the above analysis, we have the following conclusion:

Remark 1. By choosing a suitable α , we can achieve the trade-off between the contraction rate and the fixed point bias. Particularly, a larger α introduces a smaller fixed point bias between V^β and V^* , and produces a larger contraction rate simultaneously.

4.2 VALUE-BASED EPISODIC MEMORY

In this part, we demonstrate that the memory-based planning effectively accelerates the convergence of the EVL. We first define the VEM operator as:

$$(T_{\text{vem}} V)(s) = \max_{1 \leq n \leq n_{\max}} f(T)^n T V(s); \quad (9)$$

(a) The maximal rollout step n_{\max} . (b) The different behavior policies.

Figure 3: A toy example in the random MDP. In both figures, the color darkens with a larger n_{\max} (2; 3; 4; 5). The size of the spots is proportional to the relative scale of the third variable: (a) Change n_{\max} . From magenta to blue, n_{\max} is set as 1; 2; 3; 4 in order. (b) Change the behavior policies, where $\pi(s) = \text{softmax}(Q(s; \cdot))$. From light yellow to dark red, the is set as 0; 1; 2; 3 in order.

where n_{\max} is the maximal rollout step for memory control. Then, we derive that multi-step estimation operator T_{vem} does not change the fixed point and contraction property of

Lemma 4. Given $\gamma \in (0, 1)$ and $n_{\max} \in \mathbb{N}^+$, T_{vem} is a γ -contraction. If $\gamma > \frac{1}{2}$, T_{vem} has the same fixed point as T .

Next, we derive that the contraction rate of T_{vem} depends on the dataset quality. Further, we demonstrate that the convergence rate of T_{vem} is quicker than T even the behavior policy is random:

Lemma 5. When the current value estimates are much lower than the value of behavior policy, T_{vem} provides an optimistic update. Formally, we have

$$\|T_{\text{vem}} V(s) - V(s)\| \leq \gamma^{n(s)-1} (kV - V_n; k_1 + kV_n; V - k_1; \gamma \geq \frac{1}{2}); \quad (10)$$

where $n(s) = \arg \max_{0 \leq n \leq n_{\max}} f(T)^n V(s)g$, V_n is the fixed point of $(T)^n$ and it is the optimal rollout value starting from

This lemma demonstrates that T_{vem} can provide an optimistic update for pessimistic value estimates. Specifically, the scale of the update depends on the quality of the datasets. If the behavior policy is expert, which means π_n is close to V . Then, following the lemma, the contraction rate will be near to $\gamma^{n(s)-1}$. Moreover, if the initial value estimates are pessimistic (e.g., the initialized value function with zeros), we will have $n(s) = n_{\max}$, indicating that the value update will be extremely fast towards a lower bound of V . On the contrary, if π is random, we have $n(s) = 1$ and the value update will be slow towards V .

Remark 2. By choosing a suitable n_{\max} , we can achieve the trade-off between the contraction rate and the estimation variance, i.e., a large n_{\max} yields a fast update towards a lower bound of fixed point and tolerable variances empirically. Meanwhile, the choice of n_{\max} does not introduce additional bias, and the fixed point bias is totally controlled by

4.3 TOY EXAMPLE

We design a toy example in the random deterministic MDP to empirically demonstrate the above analysis. Following (Rowland et al., 2020), we adopt three indicators, including update variance, fixed-point bias, and contraction rate, which is shown in Figure 3. Specifically, the contraction rate is $\sup_{V \in \mathcal{V}} \|kT_{\text{vem}} V - T_{\text{vem}} V\| = kV - V_n$, the bias is $\|kT_{\text{vem}} V - V\| = k_1$ and the variance is $E \|kT_{\text{vem}} V - T_{\text{vem}} V\|_2^2 \leq \frac{1}{2}$, where \hat{T}_{vem} is the stochastic approximation of T_{vem} and V_{vem} is the fixed point of T_{vem} . First, the experimental results in Figure 3(a) demonstrate that the relationship of n -step estimation and. Formally, the contraction rate decreases as n_{\max} becomes larger, and the fixed-point bias increases as n_{\max} becomes smaller, which are consistent with Lemma 1 and Lemma 2. Figure 3(a) also shows that the variance is positively correlated with n_{\max} . Second, the experimental results in Figure 3(b) demonstrate that the relationship of dataset quality. The higher dataset quality corresponds to the lower contraction rate and variance, which is consistent with Lemma 5.

(a) Large (b) Medium (c) Umaze

Figure 4: Visualization of the value estimation in various AntMaze tasks. Darker colors correspond to the higher value estimation. Each map has several terminals (golden stars) and one of which is reached by the agent (the light red star). The red line is the trajectory of the ant.

5 RELATED WORK

Offline Reinforcement Learning. Offline RL methods (Kumar et al., 2019; Siegel et al., 2020; Argenson & Dulac-Arnold, 2020; Wu et al., 2021; Dadashi et al., 2021; Kostrikov et al., 2021; Jin et al., 2021; Rashidinejad et al., 2021) can be roughly divided into policy constraint, pessimistic value estimation, and model-based methods. Policy constraint methods aim to keep the policy to be close to the behavior under a probabilistic distance (Fujimoto et al., 2019; Peng et al., 2019; Nair et al., 2020). Pessimistic value estimation methods like CQL (Kumar et al., 2020) enforces a regularization constraint on the critic loss to penalize overgeneralization. Model-based methods attempt to learn a model from offline data, with minimal modification to the policy learning (Kidambi et al., 2020; Yu et al., 2020; Janner et al., 2019). However, these methods have to introduce additional behavioral policy models, dynamics models, or regularization terms (Zhang et al., 2020b;a; Lee et al., 2021). Another line of methods uses empirical return as the signal for policy learning, which confines learning within the dataset but leads to limited performance (Levine et al., 2020; Geist et al., 2019; Wang et al., 2021).

Episodic Control. Episodic control aims to store good past experiences in a non-parametric memory and rapidly latch into past successful policies when encountering similar states instead of waiting for many optimization steps (Blundell et al., 2016b). Pritzel et al. (2017) and Lin et al. (2018) introduce a parametric memory, which enables better generalization through neural networks. Our work is closely related to recent advances in Hu et al. (2021), which adopts an implicit planning scheme to enable episodic memory updates in continuous domains. Our method follows this implicit scheme, but conducts planning with expected values to avoid overgeneralization on actions out of dataset support.

6 EXPERIMENTS

In our experiments, we aim to answer the following questions: 1) How does our method perform compared to state-of-the-art offline RL algorithms on the D4RL benchmark dataset? 2) How does implicit planning affect the performance on sparse reward tasks? 3) Can experience learning effectively reduce the extrapolation error compared with other offline methods? 4) How does the critical parameter affect the performance of our method?

6.1 EVALUATION ENVIRONMENTS

We ran VEM on AntMaze, Adroit, and MuJoCo environments to evaluate its performance on various types of tasks. Precisely, the AntMaze navigation tasks control an 8-DoF quadruped robot to reach a specific or randomly sampled goal in three types of maps. The reward in the AntMaze domain is highly sparse. The Adroit domain involves controlling a 24-DoF simulated hand tasked with hammering a nail, opening a door, twirling a pen, or picking up and moving a ball. On the adroit tasks, these datasets are the following, “human”: transitions collected by a human operator,

Type	Env	VEM(Ours)	VEM(=0.5)	BAIL	BCQ	CQL	AWR
xed	umaze	87.5 1.1	85.0 1.5	62.5 2.3	78.9	74.0	56.0
play	medium	78.0 3.1	71.0 2.5	40.0 15.0	0.0	61.2	0.0
play	large	57.0 5.0	45.0 2.5	23.0 5.0	6.7	11.8	0.0
diverse	umaze	78.0 1.1	75.0 5.0	75.0 1.0	55.0	84.0	70.3
diverse	medium	77.0 2.2	60.0 5.0	50.0 10.0	0.0	53.7	0.0
diverse	large	58.0 2.1	48.0 2.7	30.0 5.0	2.2	14.9	0.0
human	door	11.2 4.2	6.9 1.1	0.0 0.1	-0.0	9.1	0.4
human	hammer	3.6 1.0	2.5 1.0	0.0 0.1	0.5	2.1	1.2
human	relocate	1.3 0.2	0.0 0.0	0.0 0.1	0.5	2.1	-0.0
human	pen	65.0 2.1	55.2 3.1	32.5 1.5	68.9	55.8	12.3
cloned	door	3.6 0.3	0.0 0.0	0.0 0.1	0.0	3.5	0.0
cloned	hammer	2.7 1.5	0.5 0.1	0.1 0.1	0.4	5.7	0.4
cloned	pen	48.7 3.2	27.8 2.2	46.5 3.5	44.0	40.3	28.0
expert	door	105.5 0.2	104.8 0.2	104.7 0.3	99.0	-	102.9
expert	hammer	128.3 1.1	102.3 5.6	123.5 3.1	114.9	-	39.0
expert	relocate	109.8 0.2	101.0 1.5	94.4 2.7	41.6	-	91.5
expert	pen	111.7 2.6	115.2 1.3	126.7 0.3	114.9	-	111.0
random	walker2d	6.2 4.7	6.2 4.7	3.9 2.5	4.9	7.0	1.5
random	hopper	11.1 1.0	10.8 1.2	9.8 0.1	10.6	10.8	10.2
random	halfcheetah	16.4 3.6	2.6 2.1	0.0 0.1	2.2	35.4	2.5
medium	walker2d	74.0 1.2	16.6 0.1	73.0 1.0	53.1	79.2	17.4
medium	hopper	56.6 2.3	56.6 2.3	58.2 1.0	54.5	58.0	35.9
medium	halfcheetah	47.4 0.2	45.3 0.2	42.6 1.2	40.7	44.4	37.4

Table 1: Performance of VEM with four of the RL baselines on the AntMaze, Adroit, and MuJoCo domains with the normalized score metric proposed by D4RL benchmark, averaged over three random seeds with standard deviation. Scores range from 0 to 100, where 0 corresponds to a random policy performance, and 100 indicates an expert. We use the results in Fu et al. (2020) for AWR and BCQ, and use the results in Kumar et al. (2020) for CQL. The results of BAIL come from our implementation according to the official code (<https://github.com/lanyavik/BAIL>).

“cloned”: transitions collected by a policy trained with behavioral cloning interacting in the environment + initial demonstrations, “expert”: transitions collected by a fine-tuned RL policy interacting in the environment. As for the MuJoCo tasks, the datasets are “random”: transitions collected by a random policy, “medium”: transitions collected by a policy with suboptimal performance. The complete implementation details are presented in Appendix C.

6.2 PERFORMANCE ON D4RL TASKS

As shown in Table 1, VEM achieves state-of-the-art performance on most AntMaze tasks and has a significant improvement over other methods on most Adroit tasks. VEM also achieves good performances in MuJoCo domains. We find that VEM has low value estimation errors in all tasks, which promotes its superior performance. However, as a similar training framework, BAIL only has reasonable performances on simple tasks, such as MuJoCo. Please refer to Appendix D.2 for the complete training curves and value estimation error on D4RL.

To further analyze the superior performance of VEM in the sparse reward tasks, we visualize the learned value estimation in AntMaze tasks, which is shown in Figure 4. Experimental results show that VEM has the higher value estimates on the critical place of the map (e.g., corners) since various trajectories in the datasets are connected. The accurate value estimation leads to its success on complex sparse reward tasks.

6.3 ANALYSIS OF VALUE ESTIMATION

As both Expectile Value Learning (EVL) and Batch Constrained Value Learning (BCQ) (Fujimoto et al., 2019) aim to avoid using the unseen state-action pairs to eliminate the extrapolation error, we replace EVL in VEM with BCQ (named BCQ-EM) to evaluate the effectiveness of the EVL module.

The experimental results in Figure 9 in Appendix D.1 indicate that the performance of BCQ-EM is mediocre, and BCQ reaches performance significantly below VEM. We observe a strong correlation between the training instability and the explosion of the value estimation. This result should not come as a surprise since the Adroit tasks have a larger action space compared with MuJoCo domains and narrow human demonstrations. Therefore, the generative model in BCQ cannot guarantee completely the unseen actions are avoided. In contrast, VEM avoids fundamentally unseen actions by keeping the learning procedure within the support of an of ine dataset, indicating the necessity of the EVL module. Please refer to Appendix C for the implementation details.

We evaluate $\beta \in \{0.1, 0.2, \dots, 0.9\}$ to investigate the effect of the critical hyper-parameter in EVL, which is shown in Figure 7 in Appendix D.1. The experimental results demonstrate that the estimated value increases with a larger β , which is consistent with the analysis in Section 4.1. Moreover, we observe that β is set at a low value in some complex high-dimensional robotic tasks or narrow human demonstrations, such as Adroit-cloned/human, to get the conservative value estimates. However, if β is set too high (e.g., $\beta = 0.9$ in the pen-human task), the estimated value will explode and poor performance. This is as expected since the over-large β leads to the overestimation error caused by neural networks. The experimental results demonstrate that we can balance behavior cloning and optimal value learning by choosing β in terms of different tasks.

6.4 ABLATIONS

Episodic Memory Module. Our first study aims to answer the impact of memory-based planning on performance. We replace the episodic memory module in VEM with standard value estimation (named VEM-1step or VEM-nstep). The experimental results in Figure 8 in Appendix D.1 indicate that implicit planning along of ine trajectories effectively accelerates the convergence of EVL.

Expectile Loss. In addition to the Expectile loss, we explored other forms of loss. Formally, we compare the Expectile loss and quantile loss, a popular form in Distributional RL algorithms (Dabney et al., 2018), which is shown in Figure 5 in Appendix D.1. The experimental results indicate that the Expectile loss is better since it is more stable when dealing with extreme values.

7 CONCLUSION

In this paper, we propose a novel of ine RL method, VEM, based on a value learning algorithm, EVL. EVL naturally avoids actions outside the dataset and provides a smooth tradeoff between generalization and conversation for of ine learning. Further, VEM enables effective implicit planning along of ine trajectories to accelerate the convergence of EVL and achieve better advantage estimation. Unlike most existing of ine RL methods, we keep the learning procedure totally within the dataset's support without any auxiliary modular, such as environment model or behavior policy. The experimental results demonstrate that VEM achieves superior performance in most D4RL tasks and learns the accurate values to guide policy learning, especially in sparse reward tasks. We hope that VEM will inspire more works on of ine RL and promote practical RL methods in the future.

8 REPRODUCIBILITY

To ensure our work is reproducible, we provide our code in the supplementary materials. In the future, we will publish all source code on Github. The detailed implementation of our algorithm is presented as follows. The value network is trained according to Equation 4. The actor-network is trained according to Equation 7. The hyper-parameters and network structure used in VEM are shown in Appendix C.3. All experiments are run on the standard of ine tasks, D4RL (<https://github.com/rail-berkeley/d4rl/tree/master/d4rl>).

REFERENCES

- Arthur Argenson and Gabriel Dulac-Arnold. Model-based of ine planning. arXiv preprint arXiv:2008.05556, 2020.
- Charles Blundell, Benigno Uria, Alexander Pritzel, Yazhe Li, Avraham Ruderman, Joel Z Leibo, Jack Rae, Daan Wierstra, and Demis Hassabis. Model-free episodic control. arXiv preprint arXiv:1606.04460, 2016a.
- Charles Blundell, Benigno Uria, Alexander Pritzel, Yazhe Li, Avraham Ruderman, Joel Z Leibo, Jack Rae, Daan Wierstra, and Demis Hassabis. Model-free episodic control. arXiv preprint arXiv:1606.04460, 2016b.
- Xinyue Chen, Zijian Zhou, Zheng Wang, Che Wang, Yanqiu Wu, and Keith Ross. BAIL: Best-action imitation learning for batch deep reinforcement learning. Advances in Neural Information Processing Systems, 33, 2020.
- Will Dabney, Mark Rowland, Marc G Bellemare, and Ori Munos. Distributional reinforcement learning with quantile regression. Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- Robert Dadashi, Shideh Rezaeifar, Nino Vieillard, Edward Hussenot, Olivier Pietquin, and Matthieu Geist. Of ine reinforcement learning with pseudometric learning. arXiv preprint arXiv:2103.01948, 2021.
- Justin Fu, Aviral Kumar, Ori Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for deep data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In International Conference on Machine Learning, pp. 2052–2062. PMLR, 2019.
- Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A theory of regularized markov decision processes. In International Conference on Machine Learning, pp. 2160–2169. PMLR, 2019.
- Seyed Kamyar Seyed Ghasemipour, Dale Schuurmans, and Shixiang Shane Gu. EMaQ: Expected-max Q-learning operator for simple yet effective of ine and online RL. International Conference on Machine Learning, pp. 3682–3691. PMLR, 2021.
- Hao Hu, Jianing Ye, Zhizhou Ren, Guangxiang Zhu, and Chongjie Zhang. Generalizable episodic memory for deep reinforcement learning. arXiv preprint arXiv:2103.06469, 2021.
- Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. Advances in Neural Information Processing Systems, 32, 12519–12530, 2019.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for of ine RL? In International Conference on Machine Learning, pp. 5084–5096. PMLR, 2021.
- Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based of ine reinforcement learning. arXiv preprint arXiv:2005.05951, 2020.
- Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ori Nachum. Of ine reinforcement learning with shier divergence critic regularization. In International Conference on Machine Learning, pp. 5774–5783. PMLR, 2021.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy Q-learning via bootstrapping error reduction. Advances in Neural Information Processing Systems, 32:11784–11794, 2019.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-learning for of ine reinforcement learning. arXiv preprint arXiv:2006.04779, 2020.
- Jongmin Lee, Wonseok Jeon, Byung-Jun Lee, Joelle Pineau, and Kee-Eung Kim. OptiDICE: Of ine policy optimization via stationary distribution correction estimation. arXiv preprint arXiv:2106.10783, 2021.

- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Of ine reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Zichuan Lin, Tianqi Zhao, Guangwen Yang, and Lintao Zhang. Episodic memory deep Q-networks. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pp. 2433–2439, 2018.
- Ashvin Nair, Murtaza Dalal, Abhishek Gupta, and Sergey Levine. Accelerating online reinforcement learning with of ine datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- Whitney K Newey and James L Powell. Asymmetric least squares estimation and testing-metrica: *Journal of the Econometric Society*, pp. 819–847, 1987.
- Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- Alexander Pritzel, Benigno Uria, Sriram Srinivasan, Adria Puigdomenech Badia, Oriol Vinyals, Demis Hassabis, Daan Wierstra, and Charles Blundell. Neural episodic control. *International Conference on Machine Learning*, pp. 2827–2836. PMLR, 2017.
- Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging of ine reinforcement learning and imitation learning: A tale of pessimism. *arXiv preprint arXiv:2103.12021*, 2021.
- Mark Rowland, Robert Dadashi, Saurabh Kumar, Oriol Munos, Marc G Bellemare, and Will Dabney. Statistics and samples in distributional reinforcement learning. *International Conference on Machine Learning*, pp. 5528–5536. PMLR, 2019.
- Mark Rowland, Will Dabney, and Oriol Munos. Adaptive trade-offs in off-policy learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 34–44. PMLR, 2020.
- Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael Neunert, Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Keep doing what worked: Behavioral modelling priors for of ine reinforcement learning. *arXiv preprint arXiv:2002.08396*, 2020.
- Yunhao Tang. Self-imitation learning via generalized lower bound Q-learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- Qing Wang, Jiechao Xiong, Lei Han, Peng Sun, Han Liu, and Tong Zhang. Exponentially weighted imitation learning for batched historical data. *Advances in Neural Information Processing Systems* 31:6288, 2018.
- Ruosong Wang, Yifan Wu, Ruslan Salakhutdinov, and Sham M Kakade. Instabilities of of ine rl with pre-trained neural representations. *arXiv preprint arXiv:2103.04947*, 2021.
- Yue Wu, Shuangfei Zhai, Nitish Srivastava, Joshua Susskind, Jian Zhang, Ruslan Salakhutdinov, and Hanlin Goh. Uncertainty weighted Actor-Critic for of ine reinforcement learning. *arXiv preprint arXiv:2105.08140*, 2021.
- Yiqin Yang, Xiaoteng Ma, Li Chenghao, Zewu Zheng, Qiyuan Zhang, Gao Huang, Jun Yang, and Qianchuan Zhao. Believe what you see: Implicit constraint approach for of ine multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* 34, 2021.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. MOPO: Model-based of ine policy optimization. *Advances in Neural Information Processing Systems* 33:14129–14142, 2020.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. GenDICE: Generalized of ine estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020a.
- Shangdong Zhang, Bo Liu, and Shimon Whiteson. GradientDICE: Rethinking generalized of ine estimation of stationary values. *International Conference on Machine Learning*, pp. 11194–11203. PMLR, 2020b.

A ALGORITHM

A.1 VALUE-BASED EPISODIC MEMORY CONTROL

Algorithm 1 Value-based Episodic Memory Control

```

Initialize critic networks  $V_1; V_2$  and actor network with random parameters  $\theta_1; \theta_2$ ;
Initialize target networks  $V_1^0; V_2^0$ 
Initialize episodic memory  $M$ 
for t = 1 to T do
  for i = 1; 2g do
    Sample  $N$  transitions  $s_t; a_t; r_t; s_{t+1}; R_t^{(i)}$  from  $M$ 
    Update  $V_i = \min_j \sum_{j=1}^N R_t^{(i)} + V_j(s_t)$ 
    Update  $\theta_i = \arg \max_j \sum_{j=1}^N r \log(\pi(a_t | s_t)) - \min_j R_t^{(i)} - \text{mean}_j V_j(s_t)$ 
  end for
  if t mod  $u$  then
     $V_i^0 = V_i + (1 - \alpha) V_i^0$ 
    Update Memory
  end if
end for

```

Algorithm 2 Update Memory

```

for trajectories in buffer  $M$  do
  for  $s_t; a_t; r_t; s_{t+1}$  in reverse( $\mathcal{D}$ ) do
    for i = 1; 2g do
      Compute  $R_t^{(i)}$  with Equation 6 and save into buffer  $M$ 
    end for
  end for
end for

```

A.2 AN APPROACH FOR AUTO-TUNING

When we have a good estimation V^* , for example, when there is some expert data in the dataset, we can auto-tune such that the value learned by EVL is close to the estimation V^* . This can be done by calculating the Monte-Carlo return estimates of each state and selecting good return values as the estimation of optimal value V^* . Based on this target, we develop a method for auto-tuning

By parameterizing $\pi = \text{sigmoid}(\cdot)$ with a differentiable parameter β , we can auto-tune by minimizing the following loss $J(\beta) = \sum_s (E\hat{V}(s) - V^*(s))^2$. If $(E\hat{V}(s) - V^*(s)) < 0$, the differentiable parameter β will become larger and the value estimation $E\hat{V}(s)$ will become larger accordingly. Similarly, β and $E\hat{V}(s)$ will become smaller if $(E\hat{V}(s) - V^*(s)) > 0$. The experimental results in Figure 10 in Appendix D.1 show that auto-tuning can lead to similar performance compared with manual selection.

B THEORETICAL ANALYSIS

B.1 COMPLETE DERIVATION.

The expectile regression loss (Rowland et al., 2019) is defined as

$$ER(q; \%) = E_Z \% [I(Z > q) + (1 - \%) I(Z \leq q)] (Z - q)^2; \quad (11)$$

where $\%$ is the target distribution and the minimiser of this loss is called the $\%$ expectile of $\%$ the corresponding loss in reinforcement learning is

$$\begin{aligned} J_V(\cdot) &= E [(r(s; a) + V_0(s^0) - V(s))^2 + (1 - \%) (r(s; a) + V_0(s^0) - V(s))^2] \\ &= E [(y - V(s))^2 + (1 - \%) (y - V(s))^2]; \end{aligned} \quad (12)$$

Then, taking the gradient of the value objective with respect to $V(s)$, we have

$$\begin{aligned} rJ_V(\cdot) &= \int_X (a_j s) [2(y - V(s))_+ I(y > V(s)) - 2(1 - \%) (y - V(s))_+ I(y \leq V(s))] \\ &= \int_X (a_j s) [2(y - V(s))_+ - 2(1 - \%) (y - V(s))] \\ &= \int_X (a_j s) [2(\cdot)_+ - 2(1 - \%) (\cdot)]; \end{aligned} \quad (13)$$

Therefore,

$$\begin{aligned} \hat{V}(s) &= V(s) - rJ_V(\cdot) \\ &= V(s) + 2 E_a [[(s; a)]_+ + (1 - \%) [(s; a)]]; \end{aligned} \quad (14)$$

B.2 PROOF OF LEMMA 1

Lemma 1. For any $\beta \in [0, 1)$, T is a β -contraction, where $\beta = 1 - 2(1 - \%) \min\{\beta_1, \beta_2\}$.

Proof. Note that $T_{1=2}$ is the standard policy evaluation Bellman operator for whose fixed point is V^* . We see that for any V_1, V_2 ,

$$\begin{aligned} T_{1=2} V_1(s) - T_{1=2} V_2(s) &= V_1(s) + E_a [\beta_1(s; a)] (V_2(s) + E_a [\beta_2(s; a)]) \\ &= (1 - \beta_1)(V_1(s) - V_2(s)) + E_a [\beta_1(s; a) + V_1(s^0) - \beta_2(s; a) - V_2(s^0)] \\ &= (1 - \beta_1)kV_1 - V_2k_1 + \beta_1kV_1 - V_2k_1 \\ &= (1 - \beta_1(1 - \beta_2))kV_1 - V_2k_1; \end{aligned} \quad (15)$$

We introduce two more operators to simplify the analysis:

$$\begin{aligned} T_+ V(s) &= V(s) + E_a [\beta_1(s; a)]_+; \\ T_- V(s) &= V(s) + E_a [\beta_2(s; a)]_-; \end{aligned} \quad (16)$$

Next we show that both operators are non-expansive (i.e. $kT_+ V_1 - T_+ V_2k_1 \leq kV_1 - V_2k_1$). For any V_1, V_2 , we have

$$\begin{aligned} T_+ V_1(s) - T_+ V_2(s) &= V_1(s) - V_2(s) + E_a [[\beta_1(s; a)]_+ - [\beta_2(s; a)]_+] \\ &= E_a [[\beta_1(s; a)]_+ + V_1(s) - ([\beta_2(s; a)]_+ + V_2(s))]; \end{aligned} \quad (17)$$

The relationship between $[\beta_1(s; a)]_+ + V_1(s)$ and $[\beta_2(s; a)]_+ + V_2(s)$ exists in four cases, which are

- $\beta_1 \geq 0; \beta_2 \geq 0$, then $[\beta_1(s; a)]_+ + V_1(s) - ([\beta_2(s; a)]_+ + V_2(s)) = (V_1(s^0) - V_2(s^0))$.
- $\beta_1 < 0; \beta_2 < 0$, then $[\beta_1(s; a)]_+ + V_1(s) - ([\beta_2(s; a)]_+ + V_2(s)) = V_1(s) - V_2(s)$.
- $\beta_1 \geq 0; \beta_2 < 0$, then

$$\begin{aligned} &[\beta_1(s; a)]_+ + V_1(s) - ([\beta_2(s; a)]_+ + V_2(s)) \\ &= (r(s; a) + V_1(s^0) - V_2(s)) \\ &< (r(s; a) + V_1(s^0) - (r(s; a) + V_2(s^0))) \\ &= (V_1(s^0) - V_2(s^0)); \end{aligned} \quad (18)$$

where the inequality comes from $(r(s; a) + V_1(s^0) - V_2(s)) < V_2(s)$.

- $\beta_1 < 0; \beta_2 \geq 0$, then

$$\begin{aligned} & [\beta_1(s; a)]_+ + V_1(s) - ([\beta_2(s; a)]_+ + V_2(s)) \\ &= V_1(s) - (r(s; a) + V_2(s^0)) \\ & \quad V_1(s) - V_2(s); \end{aligned} \tag{19}$$

where the inequality comes from $(r(s; a) + V_2(s^0)) \geq V_2(s)$.

Therefore, we have $T_+ V_1(s) \geq T_+ V_2(s) - k(V_1 - V_2)k_1$. With the $T_+; T_-$, we rewrite T as

$$\begin{aligned} T V(s) &= V(s) + 2 E_a [[\beta_1(s; a)]_+ + (1 - \beta_1) [\beta_2(s; a)]] \\ &= (1 - \beta_2) V(s) + 2 (V(s) + E_a [[\beta_1(s; a)]_+] + 2 (1 - \beta_1) (V(s) + E_a [[\beta_2(s; a)]]) \\ &= (1 - \beta_2) V(s) + 2 T_+ V(s) + 2 (1 - \beta_1) T V(s); \end{aligned} \tag{20}$$

And

$$\begin{aligned} T_{1=2} V(s) &= V(s) + E_a [[\beta_1(s; a)]] \\ &= V(s) + (T_+ V(s) + T V(s) - 2V(s)) \\ &= (1 - \beta_2) V(s) + (T_+ V(s) + T V(s)); \end{aligned} \tag{21}$$

We first focus on $\beta_1 < \frac{1}{2}$. For any $V_1; V_2$, we have

$$\begin{aligned} & T V_1(s) - T V_2(s) \\ &= (1 - \beta_2) (V_1(s) - V_2(s)) + 2 (T_+ V_1(s) - T_+ V_2(s)) + 2 (1 - \beta_1) (T V_1(s) - T V_2(s)) \\ &= (1 - \beta_2 - 2(1 - \beta_2)) (V_1(s) - V_2(s)) + 2 T_{1=2} V_1(s) - 2 T_{1=2} V_2(s) + \\ & \quad 2(1 - \beta_2) (T V_1(s) - T V_2(s)) \\ &= (1 - \beta_2 - 2(1 - \beta_2)) kV_1 - V_2k_1 + 2(1 - \beta_1) (1 - \beta_2) kV_1 - V_2k_1 + 2(1 - \beta_2) kV_1 - V_2k_1 \\ &= (1 - \beta_2) (1 - \beta_1) kV_1 - V_2k_1 \end{aligned} \tag{22}$$

Similarly, when $\beta_1 > \frac{1}{2}$, we have $T V_1(s) - T V_2(s) = (1 - \beta_2) (1 - \beta_1) kV_1 - V_2k_1$. \square

B.3 PROOF OF LEMMA 2

Lemma 2. For any $\beta \in (0, 1)$, if $\beta \geq \frac{1}{2}$, we have $T_0 \geq T$; $\forall s \in \mathcal{S}$.

Proof. Based on Equation 20, we have

$$\begin{aligned} & T_0 V(s) - T V(s) \\ &= (1 - \beta_2) V(s) + 2 \beta_1 T_+ V(s) + 2 (1 - \beta_1) T V(s) \\ & \quad - ((1 - \beta_2) V(s) + 2 T_+ V(s) + 2 (1 - \beta_1) T V(s)) \\ &= 2 (1 - \beta_1) (T_+ V(s) - T V(s)) \\ &= 2 (1 - \beta_1) E_a [[\beta_1(s; a)]_+ - [\beta_2(s; a)]] \geq 0; \end{aligned} \tag{23}$$

\square

B.4 PROOF OF LEMMA 3

Lemma 3. Let V^* denote the fixed point of Bellman optimality operator. In the deterministic MDP, we have $\lim_{t \rightarrow \infty} V^t = V^*$.

Proof. We first show that V^* is also a fixed point for T_+ . Based on the definition of T_+ , we have $V^*(s) = \max_a [r(s; a) + V^*(s^0)]$, which infers that $(s; a) \in \arg \max_{s \in \mathcal{S}; a \in \mathcal{A}}$. Thus, we have $T_+ V^*(s) = V^*(s) + E_a [[\beta_1(s; a)]_+] = V^*(s)$. By setting $(1 - \beta_1) = 0$, we eliminate the effect of T_- . Further by the contractive property of T_+ , we obtain the uniqueness of V^* . The proof is completed. \square

B.5 PROOF OF LEMMA 4

Lemma 4. Given $\beta \in (0, 1)$ and $T \in \mathcal{N}^+$, T_{vem} is a β -contraction. If $\beta > \frac{1}{2}$, T_{vem} has the same fixed point as T .

Proof. We prove the contraction first. For any V_1, V_2 , we have

$$\begin{aligned} |T_{\text{vem}} V_1(s) - T_{\text{vem}} V_2(s)| &= \max_{1 \leq n \leq n_{\max}} |f(T)^n \circ T V_1(s) - f(T)^n \circ T V_2(s)| \\ &= \max_{1 \leq n \leq n_{\max}} |j(T)^n \circ T V_1(s) - j(T)^n \circ T V_2(s)| \\ &= \max_{1 \leq n \leq n_{\max}} \beta^n |kV_1 - V_2k_1| \\ &= \beta |kV_1 - V_2k_1|. \end{aligned} \quad (24)$$

Next we show that V^* , the fixed point of T , is also the fixed point of T_{vem} when $\beta > \frac{1}{2}$. By definition, we have $V^* = T V^*$. Following Lemma 2, we have $V^* = T V^* \circ T_{1=2} V^* = T V^*$. Repeatedly applying T and using its monotonicity, we have $T V^* \circ (T)^n \circ V^* \circ 1 \leq \max_{1 \leq n \leq n_{\max}} T V^* \circ (T)^n \circ V^* \circ 1$. Thus, we have $T_{\text{vem}} V^*(s) = \max_{1 \leq n \leq n_{\max}} |f(T)^n \circ T V^*(s) - V^*(s)| = 0$. \square

B.6 PROOF OF LEMMA 5

Lemma 5. When the current value estimates $V_n(s)$ are much lower than the value of behavior policy, T_{vem} provides an optimistic update. Formally, we have

$$|jT_{\text{vem}} V(s) - V(s)| \leq \beta |kV - V_n; k_1 + kV_n; V - k_1|; \beta \leq \frac{1}{2}; \quad (25)$$

where $s^* = \arg \max_{1 \leq n \leq n_{\max}} |f(T)^n \circ T V(s) - V(s)|$ and $V_n; \circ$ is the fixed point of $(T)^n \circ (s^*) \circ T$.

Proof. The lemma is a direct result of the triangle inequality. We have

$$\begin{aligned} |T_{\text{vem}} V(s) - V(s)| &= |(T)^{n(s^*)} \circ T V(s) - V(s)| \\ &= |(T)^{n(s^*)} \circ T V(s) - (T)^{n(s^*)} \circ T V_n; (s^*) + V_n; (s^*) - V(s)| \\ &\leq \beta |kV - V_n; k_1 + kV_n; V - k_1| \\ &= \beta |kV - V_n; k_1 + kV_n; V - k_1|. \end{aligned} \quad (26)$$

B.7 PROOF OF PROPOSITION 1

Proposition 1. Let V^* denote the fixed point of T . For any $\beta \in (0, 1)$, if $\beta > \frac{1}{2}$, we have $|V_{\text{vem}}(s) - V^*(s)| \leq \beta |kV - V^*; k_1 + kV^*; V - k_1|$.

Proof. With the Lemma 2, we have $V_{\text{vem}} = T V_{\text{vem}} \circ T_{1=2} V_{\text{vem}}$. Since V^* is the fixed point of T , we have $T V^* = V^*$. Putting the results together, we obtain $V_{\text{vem}} = T V_{\text{vem}} \circ T_{1=2} V_{\text{vem}}$. Repeatedly applying $T_{1=2}$ and using its monotonicity, we have $V_{\text{vem}} = T_{1=2} V_{\text{vem}} \circ (T_{1=2})^n \circ V_{\text{vem}} = V_{\text{vem}}$. \square

C DETAILED IMPLEMENTATION

C.1 GENERALIZED ADVANTAGE-WEIGHTED LEARNING

In practice, we adopt Leaky-ReLU or Softmax functions.

Leaky-ReLU:

$$\begin{aligned} \max_{\theta} J(\theta) &= \mathbb{E}_{(s,a) \sim D} \sum_{h=1}^H \log \sum_{i=1}^I \pi_i(s; a) f_i(\hat{A}(s; a)); \\ \text{where } f(\hat{A}(s; a)) &= \begin{cases} \hat{A}(s; a) & \text{if } \hat{A}(s; a) > 0 \\ \beta \hat{A}(s; a) & \text{if } \hat{A}(s; a) \leq 0 \end{cases} \end{aligned} \quad (27)$$

Softmax:
$$\max_{\theta} J(\theta) = \mathbb{E}_{(s;a) \sim D} \log \sum_j \pi_j(s; a) \frac{\exp(\theta_j(s; a))}{\sum_{(s_i; a_i) \in \text{Batch}} \exp(\theta_j(s_i; a_i))} \quad (28)$$

C.2 BCQ-EM

The value network of BCQ-EM is trained by minimizing the following loss:

$$\min_{\theta} J_Q(\theta) = \mathbb{E}_{(s_t; a_t; s_{t+1}) \sim D} (R_t - Q(s_t; a_t))^2 \quad (29)$$

$$R_t = \max_{0 < n < n_{\max}} Q_{t;n}; \quad Q_{t;n} = \begin{cases} r_t + \gamma Q_{t+1;n} & \text{if } n > 0; \\ Q(s_t; a_t) & \text{if } n = 0; \end{cases} \quad (30)$$

where a_t corresponds to the perturbed actions, sampled from the generative model $G_{\phi}(\theta)$.

The perturbation network of BCQ-EM is trained by minimizing the following loss:

$$\min_{\psi} J(\psi) = \mathbb{E}_{s \sim D} [Q(s; a_i + \psi(s; a_i))] ; \quad \psi \sim G_{\psi}(s) \quad (31)$$

where $\psi(s; a_i)$ is a perturbation model, which outputs an adjustment to an action in the range $[-\epsilon; \epsilon]$. We adopt conditional variational auto-encoder to represent the generative model $G_{\phi}(\theta)$ and it is trained to match the state-action pairs sampled from D by minimizing the cross-entropy loss-function.

C.3 HYPER-PARAMETER AND NETWORK STRUCTURE

Table 2: Hyper-parameter Sheet

Hyper-Parameter	Value
Critic Learning Rate	1e-3
Actor Learning Rate	1e-3
Optimizer	Adam
Target Update Rate γ	0.005
Memory Update Period	100
Batch Size	128
Discount Factor	0.99
Gradient Steps per Update	200
Maximum Length	Episode Length

Table 3: Hyper-Parameter used in VEM across different tasks

	umaze	medium	large
AntMaze- fixed	0.4	0.3	0.3
AntMaze-diverse	0.3	0.4	0.1
Adroit-human	door 0.4	hammer 0.4	pen 0.4
Adroit-cloned	door 0.2	hammer 0.3	pen 0.1
Adroit-expert	door 0.3	hammer 0.3	pen 0.3
MuJoCo-medium	walker2d 0.3	halfcheetah 0.4	hopper 0.5
MuJoCo-random	walker2d 0.5	halfcheetah 0.6	hopper 0.7

We use a fully connected neural network as a function approximation with 256 hidden units and ReLU as an activation function. The structure of the actor network is [(state dim, 256); (256, 256); (256, action dim)]. The structure of the value network is [(state dim, 256); (256, 256); (256, 1)].

D ADDITIONAL EXPERIMENTS ON D4RL

D.1 ABLATION STUDY

(a) door-human (b) hammer-human (c) relocate-human (d) pen-human

Figure 5: Comparison results between expectile loss and quantile loss on Adroit tasks. We respectively name our algorithm with expectile loss and quantile loss as VEM and VEM (abs).

(a) pen-human (b) door-human (c) hammer-human

(d) pen-human (e) door-human (f) hammer-human

Figure 6: The results of VEM $\hat{\rho}$ with various ρ in Adroit tasks. The results in the upper row are the performance. The results in the bottom row are the estimation value.

(a) pen-human (b) door-human (c) hammer-human (d) relocate-human

(e) pen-human (f) door-human (g) hammer-human (h) relocate-human

Figure 7: Comparison results between VEM with TD3+BC. We adopt different hyper-parameters $\beta \in \{0.5, 2.5, 4.5\}$ in TD3+BC to test its performance. The upper row are the performance. The results in the bottom row are the estimation error (the unit is $\times 10^{-3}$).

(a) medium-play (b) medium-diverse (c) large-play (d) large-diverse

Figure 8: The comparison between episodic memory and step value estimation on AntMaze tasks.

(a) door-human (b) hammer-human (c) relocate-human (d) pen-human

(e) door-human (f) hammer-human (g) relocate-human (h) pen-human

Figure 9: The comparison between VEM, BCQ-EM and BCQ on Adroit-human tasks. The results in the upper row are the performance. The results in the bottom row are the estimation error, where the unit is 10^{13} .

D.2 COMPLETE TRAINING CURVES AND VALUE ESTIMATION ERROR

(a) Episode return (b) value

Figure 10: Comparison between $\text{xed}(\text{VEM})$ and auto-tuning ($\text{VEM}(\text{auto})$) in the door-human task.

(a) door-human (b) hammer-human (c) relocate-human (d) pen-human

Figure 11: Value estimation of $VEM_{h(\max)}$ in adroit-human tasks, where τ_{\max} is the maximal rollout step for memory control (see Equation 11). We set $\epsilon=0.5$ in all tasks.

(a) antmaze-umaze (b) antmaze-umaze-diverse (c) antmaze-medium-play (d) antmaze-medium-diverse

(e) antmaze-large-play (f) antmaze-large-diverse (g) door-human (h) hammer-human

(i) relocate-human (j) pen-human (k) door-cloned (l) hammer-cloned

(m) pen-cloned (n) hopper-medium (o) walker2d-medium (p) halfcheetah-medium

(q) hopper-random (r) walker2d-random (s) halfcheetah-random

Figure 12: The training curves of VEM and BAIL on D4RL tasks.

