

A ALGORITHM

A.1 VALUE-BASED EPISODIC MEMORY CONTROL

Algorithm 1 Value-based Episodic Memory Control

```

Initialize critic networks  $V_{\theta_1}, V_{\theta_2}$  and actor network  $\pi_\phi$  with random parameters  $\theta_1, \theta_2, \phi$ 
Initialize target networks  $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2$ 
Initialize episodic memory  $\mathcal{M}$ 
for  $t = 1$  to  $T$  do
  for  $i \in \{1, 2\}$  do
    Sample  $N$  transitions  $(s_t, a_t, r_t, s_t, R_t^{(i)})$  from  $\mathcal{M}$ 

    Update  $\theta_i \leftarrow \min_{\theta_i} N^{-1} \sum (R_t^{(i)} - V_{\theta_i}(s_t))^2$ 

    Update  $\phi \leftarrow \max_{\phi} N^{-1} \sum \nabla \log \pi_\phi(a_t|s_t) \cdot f(\min_i R_t^{(i)} - \text{mean}_i V_{\theta_i}(s_t))$ 
  end for
  if  $t \bmod u$  then
     $\theta'_i \leftarrow \kappa \theta_i + (1 - \kappa) \theta'_i$ 
    Update Memory
  end if
end for

```

Algorithm 2 Update Memory

```

for trajectories  $\tau$  in buffer  $\mathcal{M}$  do
  for  $s_t, a_t, r_t, s_{t+1}$  in reversed( $\tau$ ) do
    for  $i \in \{1, 2\}$  do
      Compute  $R_t^{(i)}$  with Equation 8 and save into buffer  $\mathcal{M}$ 
    end for
  end for
end for

```

A.2 AN APPROACH FOR AUTO-TUNING τ

When we have a good estimation of V^* , for example, when there is some expert data in the dataset, we can auto-tune τ such that the value learned by EVL is close to the estimation of 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 \tilde{V}^* . Based on this target, we develop a method for auto-tuning τ .

By parameterizing $\tau = \text{sigmoid}(\xi)$ with a differentiable parameter $\xi \in \mathbb{R}$, we can auto-tune τ by minimizing the following loss $\mathcal{J}(\xi) = \xi(\mathbb{E}\hat{V}(s) - \tilde{V}^*)$. If $(\mathbb{E}\hat{V}(s) - \tilde{V}^*) < 0$, the differentiable parameter ξ will become larger and the value estimation $\mathbb{E}\hat{V}(s)$ will become larger accordingly. Similarly, ξ and $\mathbb{E}\hat{V}(s)$ will become smaller if $(\mathbb{E}\hat{V}(s) - \tilde{V}^*) > 0$. 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

$$\text{ER}(q; \varrho, \tau) = \mathbb{E}_{Z \sim \varrho} [\tau \mathbb{I}_{Z > q} + (1 - \tau) \mathbb{I}_{Z \leq q}] (Z - q)^2, \quad (13)$$

where ϱ is the target distribution and the minimiser of this loss is called the τ -expectile of ϱ . the corresponding loss in reinforcement learning is

$$\begin{aligned} \mathcal{J}_V(\theta) &= \mathbb{E}_\mu [\tau(r(s, a) + \gamma V_{\theta'}(s') - V_\theta(s))_+^2 + (1 - \tau)(r(s, a) + \gamma V_{\theta'}(s') - V_\theta(s))_-^2] \\ &= \mathbb{E}_\mu [\tau(y - V_\theta(s))_+^2 + (1 - \tau)(y - V_\theta(s))_-^2]. \end{aligned} \quad (14)$$

Then, taking the gradient of the value objective:

$$\begin{aligned} \nabla \mathcal{J}_V(\theta) &= \sum \mu(a | s) [-2\tau(y - V_\theta)_+ \mathbb{I}(y - V_\theta) - 2(1 - \tau)(y - V_\theta)_+ \mathbb{I}(y - V_\theta)] \\ &= \sum \mu(a | s) [-2\tau(y - V_\theta)_+ - 2(1 - \tau)(y - V_\theta)_+] \\ &= \sum \mu(a | s) [-2\tau(\delta)_+ - 2(1 - \tau)(\delta)_+]. \end{aligned} \quad (15)$$

Therefore,

$$\begin{aligned} \hat{V}(s) &= V_\theta(s) - \alpha \nabla \mathcal{J}_V(\theta) \\ &= V_\theta(s) + 2\alpha \mathbb{E}_{a \sim \mu} [\tau \delta(s, a)]_+ + (1 - \tau) [\delta(s, a)]_- \end{aligned} \quad (16)$$

B.2 PROOF OF LEMMA 1

Lemma 1. For any $\tau \in [0, 1]$, \mathcal{T}_τ^μ is a γ_τ -contraction, where $\gamma_\tau = 1 - 2\alpha(1 - \gamma) \min\{\tau, 1 - \tau\}$.

Proof. Note that $\mathcal{T}_{1/2}^\mu$ is the standard policy evaluation Bellman operator for μ , whose fixed point is V^μ . We see that for any V_1, V_2 ,

$$\begin{aligned} \mathcal{T}_{1/2}^\mu V_1(s) - \mathcal{T}_{1/2}^\mu V_2(s) &= V_1(s) + \alpha \mathbb{E}_{a \sim \mu} [\delta_1(s, a)] - (V_2(s) + \alpha \mathbb{E}_{a \sim \mu} [\delta_2(s, a)]) \\ &= (1 - \alpha)(V_1(s) - V_2(s)) + \alpha \mathbb{E}_{a \sim \mu} [r(s, a) + \gamma V_1(s') - r(s, a) - \gamma V_2(s')] \\ &\leq (1 - \alpha) \|V_1 - V_2\|_\infty + \alpha \gamma \|V_1 - V_2\|_\infty \\ &= (1 - \alpha(1 - \gamma)) \|V_1 - V_2\|_\infty. \end{aligned} \quad (17)$$

We introduce two more operators to simplify the analysis:

$$\begin{aligned} \mathcal{T}_+^\mu V(s) &= V(s) + \mathbb{E}_{a \sim \mu} [\delta(s, a)]_+, \\ \mathcal{T}_-^\mu V(s) &= V(s) + \mathbb{E}_{a \sim \mu} [\delta(s, a)]_-. \end{aligned} \quad (18)$$

Next we show that both operators are non-expansion (e.g., $\|\mathcal{T}_+^\mu V_1 - \mathcal{T}_+^\mu V_2\|_\infty \leq \|V_1 - V_2\|_\infty$). For any V_1, V_2 , we have

$$\begin{aligned} \mathcal{T}_+^\mu V_1(s) - \mathcal{T}_+^\mu V_2(s) &= V_1(s) - V_2(s) + \mathbb{E}_{a \sim \mu} [[\delta_1(s, a)]_+ - [\delta_2(s, a)]_+] \\ &= \mathbb{E}_{a \sim \mu} [[\delta_1(s, a)]_+ + V_1(s) - ([\delta_2(s, a)]_+ + V_2(s))]. \end{aligned} \quad (19)$$

The relationship between $[\delta_1(s, a)]_+ + V_1(s)$ and $[\delta_2(s, a)]_+ + V_2(s)$ exists in four cases, which are

- $\delta_1 \geq 0, \delta_2 \geq 0$, then $[\delta_1(s, a)]_+ + V_1(s) - ([\delta_2(s, a)]_+ + V_2(s)) = \gamma(V_1(s') - V_2(s'))$.
- $\delta_1 < 0, \delta_2 < 0$, then $[\delta_1(s, a)]_+ + V_1(s) - ([\delta_2(s, a)]_+ + V_2(s)) = V_1(s) - V_2(s)$.
- $\delta_1 \geq 0, \delta_2 < 0$, then

$$\begin{aligned} &[\delta_1(s, a)]_+ + V_1(s) - ([\delta_2(s, a)]_+ + V_2(s)) \\ &= (r(s, a) + \gamma V_1(s')) - V_2(s) \\ &< (r(s, a) + \gamma V_1(s')) - (r(s, a) + \gamma V_2(s')) \\ &= \gamma(V_1(s') - V_2(s')), \end{aligned} \quad (20)$$

where the inequality comes from $r(s, a) + \gamma V_2(s') < V_2(s)$.

- $\delta_1 < 0, \delta_2 \geq 0$, then

$$\begin{aligned} & [\delta_1(s, a)]_+ + V_1(s) - ([\delta_2(s, a)]_+ + V_2(s)) \\ &= V_1(s) - (r(s, a) + \gamma V_2(s')) \\ &\leq V_1(s) - V_2(s), \end{aligned} \quad (21)$$

where the inequality comes from $r(s, a) + \gamma V_2(s') \geq V_2(s)$.

Therefore, we have $\mathcal{T}_+^\mu V_1(s) - \mathcal{T}_+^\mu V_2(s) \leq \|V_1 - V_2\|_\infty$. With the $\mathcal{T}_+^\mu, \mathcal{T}_-^\mu$, we rewrite \mathcal{T}_τ^μ as

$$\begin{aligned} \mathcal{T}_\tau^\mu V(s) &= V(s) + 2\alpha \mathbb{E}_{a \sim \mu} [\tau [\delta(s, a)]_+ + (1 - \tau) [\delta(s, a)]_-] \\ &= (1 - 2\alpha) V(s) + 2\alpha \tau (V(s) + \mathbb{E}_{a \sim \mu} [\delta(s, a)]_+) + 2\alpha (1 - \tau) (V(s) + \mathbb{E}_{a \sim \mu} [\delta(s, a)]_-) \\ &= (1 - 2\alpha) V(s) + 2\alpha \tau \mathcal{T}_+^\mu V(s) + 2\alpha (1 - \tau) \mathcal{T}_-^\mu V(s). \end{aligned} \quad (22)$$

And

$$\begin{aligned} \mathcal{T}_{1/2}^\mu V(s) &= V(s) + \alpha \mathbb{E}_{a \sim \mu} [\delta(s, a)] \\ &= V(s) + \alpha (\mathcal{T}_+^\mu V(s) + \mathcal{T}_-^\mu V(s) - 2V(s)) \\ &= (1 - 2\alpha) V(s) + \alpha (\mathcal{T}_+^\mu V(s) + \mathcal{T}_-^\mu V(s)). \end{aligned} \quad (23)$$

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

$$\begin{aligned} & \mathcal{T}_\tau^\mu V_1(s) - \mathcal{T}_\tau^\mu V_2(s) \\ &= (1 - 2\alpha) (V_1(s) - V_2(s)) + 2\alpha \tau (\mathcal{T}_+^\mu V_1(s) - \mathcal{T}_+^\mu V_2(s)) + 2\alpha (1 - \tau) (\mathcal{T}_-^\mu V_1(s) - \mathcal{T}_-^\mu V_2(s)) \\ &= (1 - 2\alpha - 2\tau(1 - 2\alpha)) (V_1(s) - V_2(s)) + 2\tau (\mathcal{T}_{1/2}^\mu V_1(s) - \mathcal{T}_{1/2}^\mu V_2(s)) + \\ & \quad 2\alpha (1 - 2\tau) (\mathcal{T}_-^\mu V_1(s) - \mathcal{T}_-^\mu V_2(s)) \\ &\leq (1 - 2\alpha - 2\tau(1 - 2\alpha)) \|V_1 - V_2\|_\infty + 2\tau (1 - \alpha(1 - \gamma)) \|V_1 - V_2\|_\infty + 2\alpha (1 - 2\tau) \|V_1 - V_2\|_\infty \\ &= (1 - 2\alpha\tau(1 - \gamma)) \|V_1 - V_2\|_\infty \end{aligned} \quad (24)$$

Similarly, when $\tau > 1/2$, we have $\mathcal{T}_\tau^\mu V_1(s) - \mathcal{T}_\tau^\mu V_2(s) \leq (1 - 2\alpha(1 - \tau)(1 - \gamma)) \|V_1 - V_2\|_\infty$. \square

B.3 PROOF OF LEMMA 2

Lemma 2. For any $\tau, \tau' \in (0, 1)$, if $\tau' \geq \tau$, we have $\mathcal{T}_{\tau'}^\mu \geq \mathcal{T}_\tau^\mu, \forall s \in S$.

Proof. Based on Equation 22, we have

$$\begin{aligned} & \mathcal{T}_{\tau'}^\mu V(s) - \mathcal{T}_\tau^\mu V(s) \\ &= (1 - 2\alpha) V(s) + 2\alpha \tau' \mathcal{T}_+^\mu V(s) + 2\alpha (1 - \tau') \mathcal{T}_-^\mu V(s) \\ & \quad - ((1 - 2\alpha) V(s) + 2\alpha \tau \mathcal{T}_+^\mu V(s) + 2\alpha (1 - \tau) \mathcal{T}_-^\mu V(s)) \\ &= 2\alpha (\tau' - \tau) (\mathcal{T}_+^\mu V(s) - \mathcal{T}_-^\mu V(s)) \\ &= 2\alpha (\tau' - \tau) \mathbb{E}_{a \sim \mu} [[\delta(s, a)]_+ - [\delta(s, a)]_-] \geq 0. \end{aligned} \quad (25)$$

\square

B.4 PROOF OF LEMMA 3

Lemma 3. Let V^* denote the fixed point of Bellman optimality operator \mathcal{T}^* . In the deterministic MDP, we have $\lim_{\tau \rightarrow 1} V_\tau^* = V^*$.

Proof. We first show that V^* is also a fixed point for \mathcal{T}_+^μ . Based on the definition of \mathcal{T}^* , we have $V^*(s) = \max_a [r(s, a) + \gamma V^*(s')]$, which infers that $\delta(s, a) \leq 0, \forall s \in S, a \in A$. Thus, we have $\mathcal{T}_+^\mu V^*(s) = V^*(s) + \mathbb{E}_{a \sim \mu} [\delta(s, a)]_+ = V^*(s)$. By setting $(1 - \tau) \rightarrow 0$, we eliminate the effect of \mathcal{T}_-^μ . Further by the contractive property of \mathcal{T}_τ^μ , we obtain the uniqueness of V_τ^* . The proof is completed. \square

B.5 PROOF OF LEMMA 4

Lemma 4. *Given $\tau \in (0, 1)$ and $T \in \mathbb{N}^+$, \mathcal{T}_{vem} is a γ_τ -contraction. If $\tau > \frac{1}{2}$, \mathcal{T}_{vem} has the same fixed point as \mathcal{T}_τ^μ .*

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

$$\begin{aligned} \mathcal{T}_{\text{vem}} V_1(s) - \mathcal{T}_{\text{vem}} V_2(s) &= \max_{1 \leq n \leq n_{\max}} \{(\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V_1(s)\} - \max_{1 \leq n \leq T} \{(\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V_2(s)\} \\ &\leq \max_{1 \leq n \leq n_{\max}} |(\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V_1(s) - (\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V_2(s)| \\ &\leq \max_{1 \leq n \leq n_{\max}} \gamma^{n-1} \gamma_\tau \|V_1 - V_2\|_\infty \\ &\leq \gamma_\tau \|V_1 - V_2\|_\infty. \end{aligned} \quad (26)$$

Next we show that V_τ^* , the fixed point of \mathcal{T}_τ^μ , is also the fixed point of \mathcal{T}_{vem} when $\tau > \frac{1}{2}$. By definition, we have $V_\tau^* = \mathcal{T}_\tau^\mu V_\tau^*$. Following Lemma 2, we have $V_\tau^* = \mathcal{T}_\tau^\mu V_\tau^* \geq \mathcal{T}_{1/2}^\mu V_\tau^* = \mathcal{T}^\mu V_\tau^*$. Repeatedly applying \mathcal{T}^μ and using its monotonicity, we have $\mathcal{T}^\mu V_\tau^* \geq (\mathcal{T}^\mu)^{n-1} V_\tau^*, 1 \leq n \leq n_{\max}$. Thus, we have $\mathcal{T}_{\text{vem}} V_\tau^*(s) = \max_{1 \leq n \leq T} \{(\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V_\tau^*(s)\} = V_\tau^*(s)$. \square

B.6 PROOF OF LEMMA 5

Lemma 5. *When the current value estimates $V(s)$ are much lower than the value of behavior policy, \mathcal{T}_{vem} provides an optimistic update. Formally, we have*

$$|\mathcal{T}_{\text{vem}} V(s) - V_\tau^*(s)| \leq \gamma^{n^*(s)-1} \gamma_\tau \|V - V_{n^*, \tau}^\mu\|_\infty + \|V_{n^*, \tau}^\mu - V_\tau^*\|_\infty, \forall s \in S, \quad (27)$$

where $n^*(s) = \arg \max_{1 \leq n \leq T} \{(\mathcal{T}^\mu)^{n-1} \mathcal{T}_\tau^\mu V(s)\}$ and $V_{n^*, \tau}^\mu$ is the fixed point of $(\mathcal{T}^\mu)^{n^*(s)-1} \mathcal{T}_\tau^\mu$.

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

$$\begin{aligned} \mathcal{T}_{\text{vem}} V(s) - V_\tau^*(s) &= (\mathcal{T}^\mu)^{n^*(s)-1} \mathcal{T}_\tau^\mu V(s) - V_\tau^*(s) \\ &= (\mathcal{T}^\mu)^{n^*(s)-1} \mathcal{T}_\tau^\mu V(s) - (\mathcal{T}^\mu)^{n^*(s)-1} \mathcal{T}_\tau^\mu V_{n^*, \tau}^\mu(s) + V_{n^*, \tau}^\mu(s) - V_\tau^*(s) \\ &\leq \gamma^{n^*(s)-1} \gamma_\tau \|V - V_{n^*, \tau}^\mu\|_\infty + \|V_{n^*, \tau}^\mu - V_\tau^*\|. \end{aligned} \quad (28)$$

\square

B.7 PROOF OF PROPOSITION 1

Proposition 1. *Let V_τ^* denote the fixed point of \mathcal{T}_τ^μ . For any $\tau, \tau' \in (0, 1)$, if $\tau' \geq \tau$, we have $V_{\tau'}^*(s) \geq V_\tau^*(s), \forall s \in S$.*

Proof. With the Lemma 2, we have $\mathcal{T}_{\tau'}^\mu V_\tau^* \geq \mathcal{T}_\tau^\mu V_\tau^*$. Since V_τ^* is the fixed point of \mathcal{T}_τ^μ , we have $\mathcal{T}_\tau^\mu V_\tau^* = V_\tau^*$. Putting the results together, we obtain $V_\tau^* = \mathcal{T}_\tau^\mu V_\tau^* \leq \mathcal{T}_{\tau'}^\mu V_\tau^*$. Repeatedly applying $\mathcal{T}_{\tau'}^\mu$ and using its monotonicity, we have $V_\tau^* \leq \mathcal{T}_{\tau'}^\mu V_\tau^* \leq (\mathcal{T}_{\tau'}^\mu)^\infty V_\tau^* = V_{\tau'}^*$. \square

C DETAILED IMPLEMENTATION

C.1 GENERALIZED ADVANTAGE-WEIGHTED LEARNING

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

Leaky-ReLU:

$$\begin{aligned} \max_{\phi} J_\pi(\phi) &= \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\log \pi_\phi(a | s) \cdot f(\hat{A}(s, a)) \right], \\ \text{where } f(\hat{A}(s, a)) &= \begin{cases} \hat{A}(s, a) & \text{if } \hat{A}(s, a) > 0 \\ \frac{\hat{A}(s, a)}{\alpha} & \text{if } \hat{A}(s, a) \leq 0 \end{cases} \end{aligned} \quad (29)$$

Softmax:

$$\max_{\phi} J_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\log \pi_{\phi}(a | s) \cdot \frac{\exp(\frac{1}{\alpha} \hat{A}(s, a))}{\sum_{(s_i, a_i) \sim \text{Batch}} \exp(\frac{1}{\alpha} \hat{A}(s_i, a_i))} \right]. \quad (30)$$

C.2 BCQ-EM

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

$$\min_{\theta} \mathcal{J}_Q(\theta) = \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left[(R_t - Q_{\theta}(s_t, a_t))^2 \right] \quad (31)$$

$$R_t = \max_{0 < n \leq n_{\max}} Q_{t,n}, \quad Q_{t,n} = \begin{cases} r_t + \gamma Q_{t+1,n-1}(s_{t+1}, \hat{a}_{t+1}) & \text{if } n > 0, \\ Q(s_t, \hat{a}_t) & \text{if } n = 0, \end{cases} \quad (32)$$

where \hat{a}_t corresponds to the perturbed actions, sampled from the generative model $G_w(s_t)$.

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

$$\min_{\phi} \mathcal{J}_{\xi}(\phi) = -\mathbb{E}_{s \sim \mathcal{D}} [Q_{\theta}(s, a_i + \xi_{\phi}(s, a_i, \Phi))], \quad \{a_i \sim G_w(s)\}_{i=1}^n, \quad (33)$$

where $\xi_{\phi}(s, a_i, \Phi)$ is a perturbation model, which outputs an adjustment to an action a in the range $[-\Phi, \Phi]$. We adopt conditional variational auto-encoder to represent the generative model $G_w(s)$ and it is trained to match the state-action pairs sampled from \mathcal{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 (κ)	0.005
Memory Update Period	100
Batch Size	128
Discount Factor	0.99
Gradient Steps per Update	200
Maximum Length d	Episode Length T

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

	umaze	medium	large
AntMaze-fixed	0.4	0.3	0.3
AntMaze-diverse	umaze 0.3	medium 0.4	large 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 $[(\text{state dim}, 256), (256, 256), (256, \text{action dim})]$. The structure of the value network is $[(\text{state dim}, 256), (256, 256), (256, 1)]$.

D ADDITIONAL EXPERIMENTS ON D4RL

D.1 ABLATION STUDY

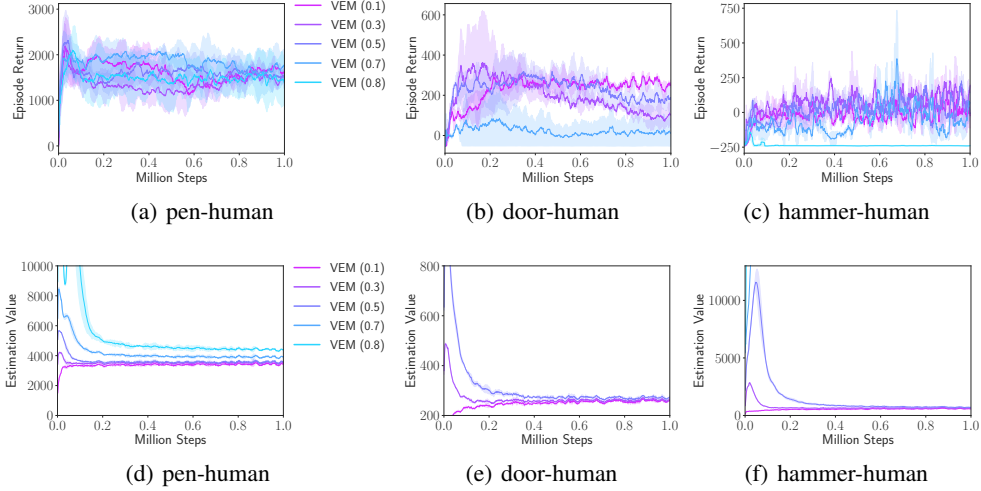


Figure 5: The results of VEM (τ) with various τ in Adroit tasks. The results in the upper row are the performance. The results in the bottom row are the estimation value.

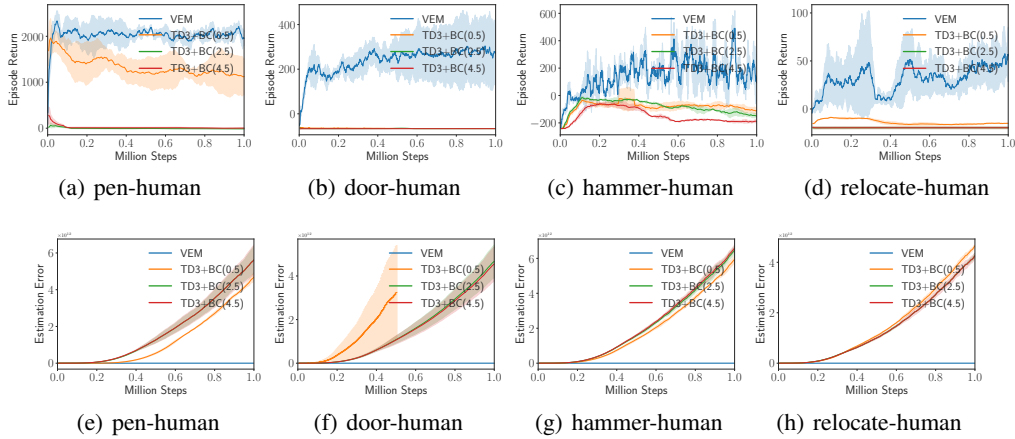


Figure 6: Comparison results between VEM with TD3+BC. We adopt different hyper-parameters $\alpha \in \{0.5, 2.5, 4.5\}$ in TD3+BC to test its performance. The upper row are the performance. The results in the bottom row are the estimation error (the unit is 10^{12}).

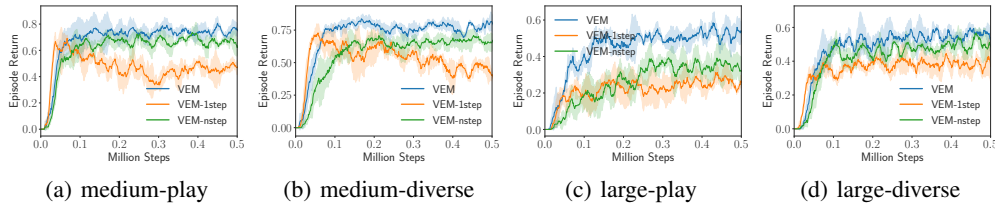


Figure 7: The comparison between episodic memory and n -step value estimation on AntMaze tasks.

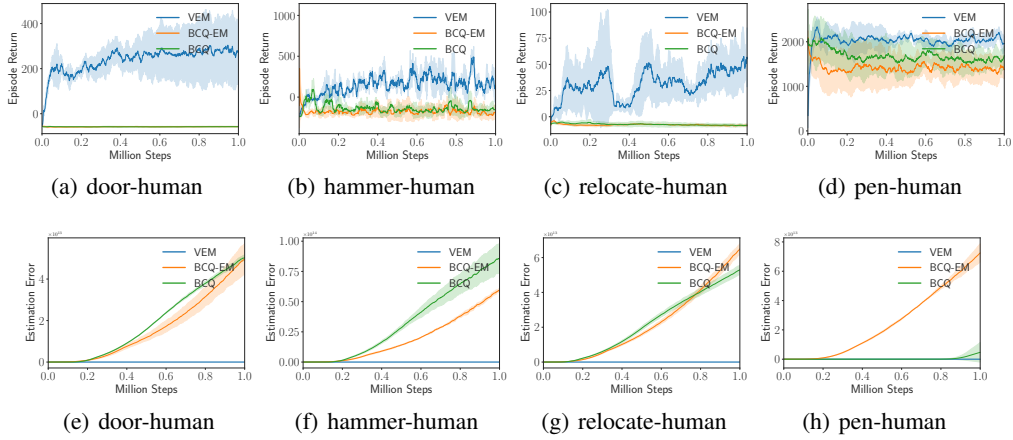


Figure 8: The comparison between VEM, BCQ-EM and BCQ on Adroit-human tasks. The results in the upper row are the performance. The results in the bottom row are the estimation error, where the unit is 10^{13} .

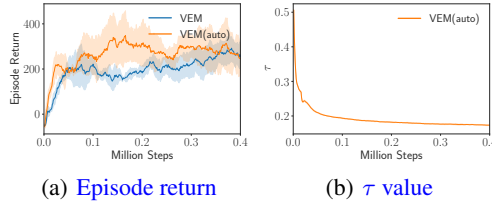


Figure 9: Comparison between fixed τ (VEM) and auto-tuning τ (VEM(auto)) in the door-human task.

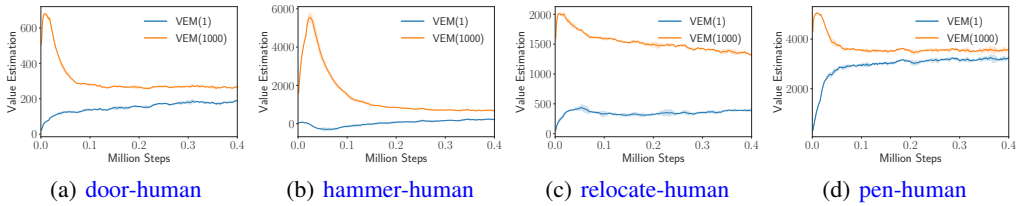


Figure 10: Value estimation of VEM (n_{\max}) in adroit-human tasks, where n_{\max} is the maximal rollout step for memory control (see Equation 11). We set $\tau = 0.5$ in all tasks.

D.2 COMPLETE TRAINING CURVES AND VALUE ESTIMATION ERROR

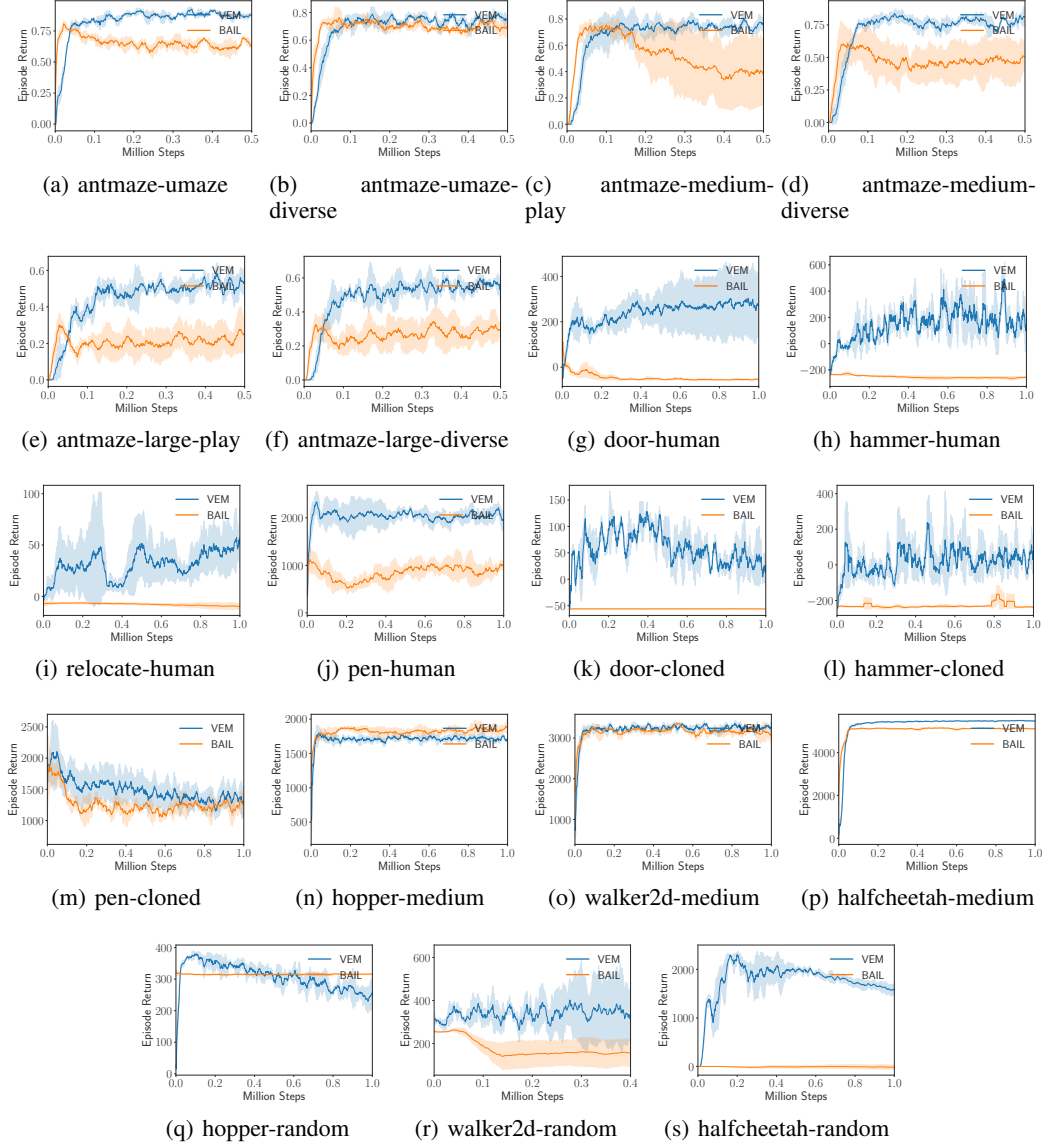


Figure 11: The training curves of VEM and BAIL on D4RL tasks.

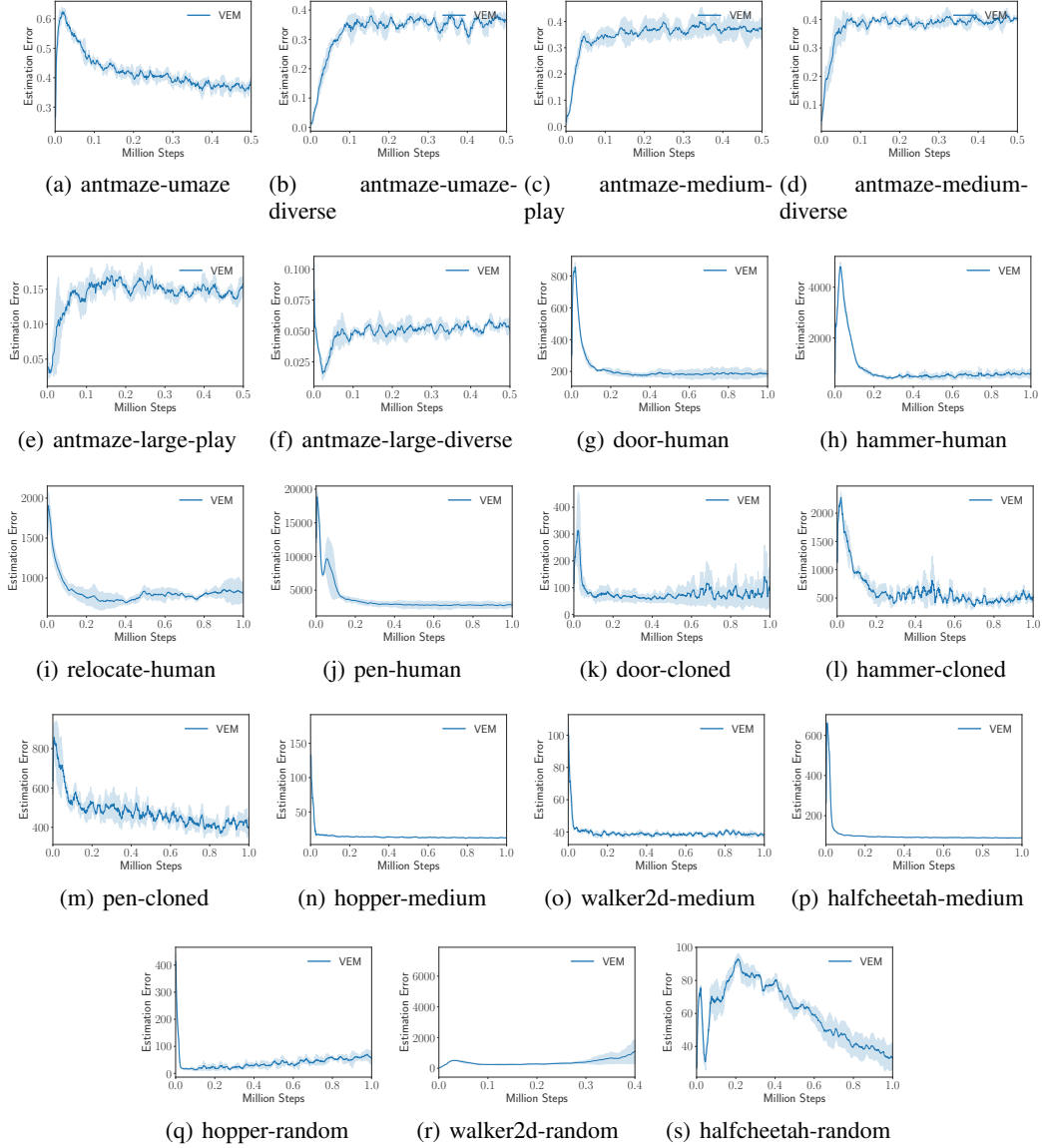


Figure 12: The value estimation error of VEM on D4RL tasks. The estimation error refers to the average estimated state values minus the average returns.