

# Exploiting the Replay Memory Before Exploring the Environment: Enhancing Reinforcement Learning Through Empirical MDP Iteration

Hongming Zhang<sup>1</sup>, Chenjun Xiao<sup>2</sup>, Chao Gao<sup>3</sup>, Han Wang<sup>1</sup>, Bo Xu<sup>4</sup>, Martin Müller<sup>1</sup>

<sup>1</sup>Department of Computing Science and Amii, University of Alberta

<sup>2</sup>The Chinese University of Hong Kong, Shenzhen

<sup>3</sup>Edmonton Research Center, Huawei Canada

<sup>4</sup>Institute of Automation, Chinese Academy of Sciences

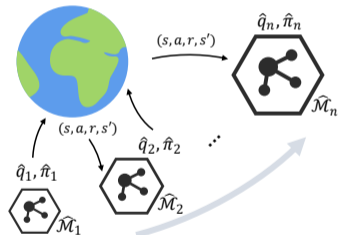


## Motivation

- Reinforcement learning is typically based on the *optimal Bellman equation* in a sampled manner.
- It suffers from *estimation errors* due to the combination of *applying the max operator to out-of-sample actions* and *bootstrapping from a function approximator*.
- How about focusing on the **currently available data**?

## Our Framework: Empirical MDP Iteration (EMIT)

- Constructs a sequence of empirical MDPs using data from the growing replay memory.
- Restrict the Bellman update to in-sample bootstrapping that uniquely solves each empirical MDP.
- Gradually expand from the empirical MDPs to the original MDP through new data collection.



An overview of our method.

## Preliminaries

**Markov Decision Process (MDP).**  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, R, P, \gamma)$ .

**Empirical MDP.**  $\hat{\mathcal{M}} := (\hat{\mathcal{S}}, \hat{\mathcal{A}}, \hat{R}, \hat{P}, \gamma)$ . Given a dataset  $\mathcal{D}$  from MDP  $\mathcal{M}$ , the empirical MDP has state space  $\hat{\mathcal{S}} = \{s | s \in \mathcal{D}\}$ , and action space  $\hat{\mathcal{A}} = \{a | (s, a) \in \mathcal{D}\}$ , with reward function  $\hat{R}(s, a) = R(s, a)$  if  $(s, a) \in \mathcal{D}$ , and  $\hat{R}(s, a) = -\infty$  otherwise.  $\hat{P}(s' | s, a) = N(s, a, s') / \sum_{s'} N(s, a, s')$  is the empirical transition dynamics based on visit counts in  $\mathcal{D}$ , and  $\gamma \in (0, 1)$  is the discount factor.

**Bellman Update.**

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a') \quad (1)$$

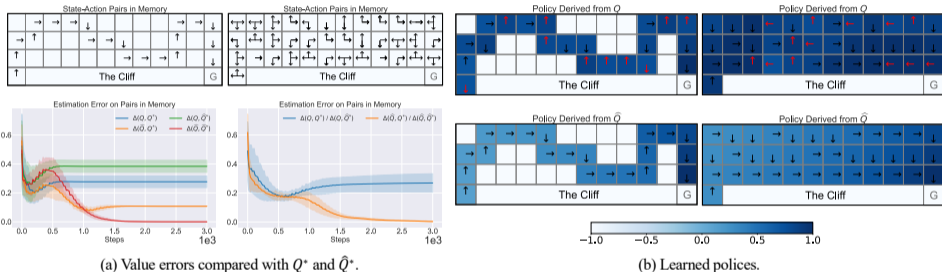
**In-Sample Bellman Update.**

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_{a': (s', a') \in \mathcal{D}} \hat{Q}(s', a') \quad (2)$$

Eq.(2) bootstraps only from in-sample actions.

# Analysis of EMIT

## Illustrative Example



(a) Value errors compared with  $Q^*$  and  $\hat{Q}^*$ .

(b) Learned policies.

**Proposition 3.2.** If the data coverage  $\mathcal{D}$  is incomplete, then the Bellman update Eq.(1) neither guarantees convergence to the optimal value  $Q^*$  for the original MDP  $\mathcal{M}$  nor to the optimal value  $\hat{Q}^*$  for the empirical MDP  $\hat{\mathcal{M}}$ , even in the limit of infinite updates.

**Proposition 3.3.** The in-sample Bellman update Eq.(2) uniquely converges to the optimal value  $\hat{Q}^*$  for the empirical MDP  $\hat{\mathcal{M}}$  in the limit of infinitely many updates. Furthermore, if the optimal trajectory for  $\mathcal{M}$  is included in  $\hat{\mathcal{M}}$ , then the greedy policy derived from  $\hat{Q}^*$  is also optimal for  $\mathcal{M}$  following the optimal trajectory.

**Proposition 3.4.** Assume  $\mathcal{M}$  has deterministic transitions. If  $\{\mathcal{D}_i\}$  are datasets collected from  $\mathcal{M}$  with  $\mathcal{D}_1 \subset \mathcal{D}_2 \subset \dots \subset \mathcal{D}_n$  and corresponding empirical MDPs  $\{\hat{\mathcal{M}}_i\}$ , then  $\hat{Q}_1^* \leq \hat{Q}_2^* \leq \dots \leq \hat{Q}_n^*$ .

## Algorithm Instantiation with EMIT

### Regularization for stable learning

$$\mathcal{L} = \text{MSE}(Q_\theta, Q_{\text{target}}) + \alpha \text{MSE}(Q_\theta, \hat{Q}) \quad (4)$$

Here,  $Q_{\text{target}} = r + \gamma \max_{a'} Q(s', a')$  and  $\alpha$  is a parameter controlling the level of regularization.

### Exploration to grow $\widehat{\mathcal{M}}$

Define the absolute difference between  $Q$  and  $\hat{Q}$  at  $(s, a)$ :

$$\delta(s, a) = |Q(s, a) - \hat{Q}(s, a)| \quad (5)$$

For algorithms in discrete action spaces such as DQN [2]:

$$\pi = \begin{cases} \text{random action,} & p = \epsilon \\ \text{argmax}_a (Q(s, a) + \delta(s, a)), & p = 1 - \epsilon \end{cases} \quad (6)$$

For algorithms in continuous action spaces such as TD3 [3]:

$$a \leftarrow \pi(s) + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \delta(s, a)), -c, c) \quad (7)$$

The added noise is clipped as in TD3.

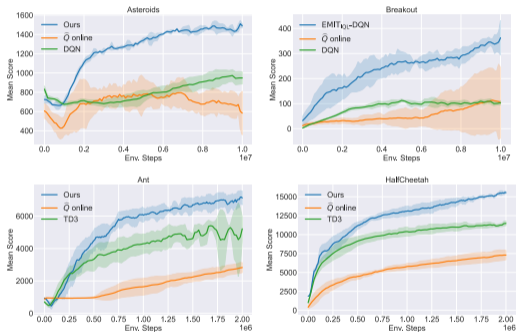
---

**Algorithm 1** Empirical MDP Iteration for Enhancing a  $Q$  Learning Algorithm Alg (EMIT-Alg)

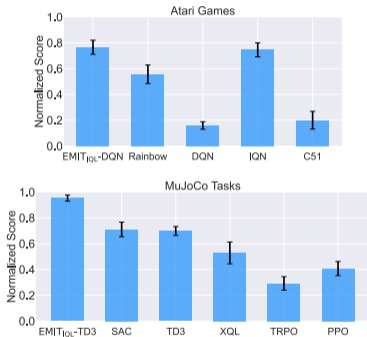
---

- 1: Alg. Initialize the replay memory  $\mathcal{D}$  and action value network  $Q_\theta$
  - 2: EMIT. Initialize  $\hat{Q}_\delta$
  - 3: Alg. Initialize the environment  $s_0 \leftarrow Env$
  - 4: **for** environment step  $t = 0$  to  $T$  **do**
  - 5:   Alg. Select an action  $a_t = \text{Alg.act}(s_t, \hat{Q}_\delta)$  as Eq. (6) or (7)  $\{\hat{Q}$  guided exploration $\}$
  - 6:   Alg. Execute  $a_t$  in  $Env$  and get  $r_t, s_{t+1}$
  - 7:   Alg. Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$
  - 8:   Alg. Sample random minibatch of transitions  $\mathcal{B}$  from  $\mathcal{D}$
  - 9:   EMIT.update( $\hat{Q}_\delta, \mathcal{B}$ ) w.r.t MSE loss derived by Eq.(2)  $\{\text{Learning of } \hat{Q}\}$
  - 10:   Alg.update( $Q_\theta, \mathcal{B}, \hat{Q}_\delta$ ) as Eq.(4)  $\{\hat{Q}$  regularized learning for  $Q\}$
  - 11: **end for**
-

## Experiments

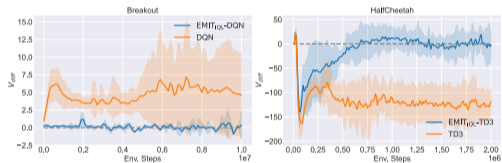


(a) Performance Enhancement of EMIT for DQN and TD3

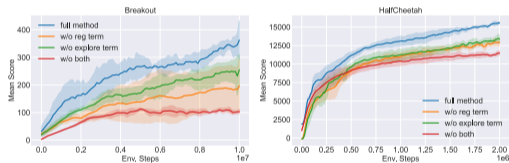


(b) Overall Performance across Various Tasks

## Experiments



(a) Estimation error of  $Q$  after regularization.



(b) The contribution of each component of EMIT.

## Conclusions

- Unlike previous methods that solely focus on solving the entire original MDP, we propose an Empirical MDP Iteration (EMIT) framework that uniquely solves each empirical MDP and incrementally approaches the original MDP through new data collection.
- Instantiate EMIT with DQN and TD3, and conduct extensive experiments on Atari and MuJoCo tasks. Demonstrate strong performance enhancements for both methods.

## Future work

- Design a learning process that is directly based on in-sample Bellman update and also take the exploration mechanism into consideration.



**THANKS !**