# Revisiting Prioritized Experience Replay: A Value Perspective

**Ang A. Li** [1]   **Zongqing Lu** [1]   **Chenglin Miao** [1]

## Abstract

Experience replay enables off-policy reinforcement learning (RL) agents to utilize past experiences to maximize the cumulative reward. Prioritized experience replay that weighs experiences by the magnitude of their temporal-difference error ($|TD|$) significantly improves the learning efficiency. But how $|TD|$ is related to the importance of experience is not well understood. We address this problem from an economic perspective, by linking $|TD|$ to *value of experience*, which is defined as the value added to the cumulative reward by accessing the experience. We theoretically show the value metrics of experience are upper-bounded by $|TD|$ for Q-learning. Furthermore, we successfully extend our theoretical framework to maximum-entropy RL by deriving the lower and upper bounds of these value metrics for soft Q-learning, which turn out to be the product of $|TD|$ and "on-policyness" of the experiences. Our framework links two important quantities in RL: $|TD|$ and value of experience. We empirically show that the bounds hold in practice, and experience replay using the upper bound as priority improves maximum-entropy RL in Atari games.

## 1. Introduction

Learning from important experiences prevails in nature. In rodent hippocampus, memories with higher importance, such as those associated with rewarding locations or large reward-prediction errors, are replayed more frequently (Michon et al., 2019; Roscow et al., 2019; Salvetti et al., 2014). Psychophysical experiments showed that participants with more frequent replay of high-reward associated memories show better performance in memory tasks (Gruber et al., 2016; Schapiro et al., 2018). As accumulating new experiences is costly, utilizing valuable past experiences is a key for efficient learning (Ólafsdóttir et al., 2018).

[1]Peking University. Correspondence to: Zongqing Lu <zongqing.lu@pku.edu.cn>, Chenglin Miao <chenglin.miao@pku.edu.cn>.

Differentiating important experiences from unimportant ones also benefits reinforcement learning (RL) algorithms (Katharopoulos & Fleuret, 2018). Prioritized experience replay (PER) (Schaul et al., 2016) is an experience replay technique built on deep Q-network (DQN) (Mnih et al., 2015), which weighs the importance of samples by the magnitude of their temporal-difference error ($|TD|$). As a result, experiences with larger $|TD|$ are sampled more frequently. PER significantly improves the learning efficiency of DQN, and has been adopted (Hessel et al., 2018; Horgan et al., 2018; Kapturowski et al., 2019) and extended (Daley & Amato, 2019; Pan et al., 2018; Schlegel et al., 2019) by various deep RL algorithms. $|TD|$ quantifies the unexpectedness of an experience to a learning agent, and biologically corresponds to the signal of reward prediction error in dopamine system (Schultz et al., 1997; Glimcher, 2011). However, how $|TD|$ is related to the importance of experience in the context of RL is not well understood.

We address this problem from an economic perspective, by linking $|TD|$ to *value of experience* in RL. Recently in neuroscience field, a normative theory for memory access, based on Dyna framework (Sutton, 1990), suggests that a rational agent should replay the experiences that lead to most rewarding future decisions (Mattar & Daw, 2018). Follow-up research shows that optimizing the replay strategy according to the normative theory has advantage over prioritized experience replay with $|TD|$ (Zha et al., 2019). Inspired by (Mattar & Daw, 2018), we define the value of experience as the increase in the expected cumulative reward resulted from updating on the experience. The value of experience quantifies the importance of experience from first principles: assuming that the agent is economically rational and has full information about the value of experience, it will choose the most valuable experience to update, which leads to most rewarding future decisions. As supplements, we derive two more value metrics, which correspond to the evaluation improvement value and policy improvement value due to update on an experience.

In this work, we mathematically show that these value metrics are upper-bounded by $|TD|$ for Q-learning. Therefore, $|TD|$ implicitly tracks the value of experience, and accounts for the importance of experience. We further extend our framework to maximum-entropy RL, which augments the reward with an entropy term to encourage exploration

(Haarnoja et al., 2017). We derive the lower and upper bounds of these value metrics for soft Q-learning, which are related to |TD| and "on-policyness" of the experience. Experiments in grid-world maze and CartPole support our theoretical results for both tabular and function approximation RL methods, showing that the derived bounds hold in practice. Moreover, we show that experience replay using the upper bound as priority improves maximum-entropy RL (*i.e.*, soft DQN) in Atari games.

## 2. Motivation

### 2.1. Q-learning and Experience Replay

We consider a Markov Decision Process (MDP) defined by a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma\}$, where $\mathcal{S}$ is a finite set of states, $\mathcal{A}$ is a finite set of actions, $\mathcal{P}$ is the transition function, $\mathcal{R}$ is the reward function, and $\gamma \in [0, 1]$ is the discount factor. A policy $\pi$ of an agent assigns probability $\pi(a|s)$ to each action $a \in \mathcal{A}$ given state $s \in \mathcal{S}$. The goal is to learn an optimal policy that maximizes the expected discounted return starting from time step $t$, $G_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$, where $r_t$ is the reward the agent receives at time step $t$. Value function $v_\pi(s)$ is defined as the expected return starting from state $s$ following policy $\pi$, and Q-function $q_\pi(s, a)$ is the expected return on performing action $a$ in state $s$ and subsequently following policy $\pi$.

According to Q-learning (Watkins & Dayan, 1992), the optimal policy can be learned through policy iteration: performing policy evaluation and policy improvement interactively and iteratively. For each policy evaluation, we update $Q(s, a)$, an estimate of $q_\pi(s, a)$, by

$$Q_{\text{new}}(s, a) = Q_{\text{old}}(s, a) + \alpha \text{TD}(s, a, r, s'),$$

where TD error $\text{TD}(s, a, r, s') = r + \gamma \max_{a'} Q_{\text{old}}(s', a') - Q_{\text{old}}(s, a)$ and $\alpha$ is the step-size parameter. $Q_{\text{new}}$ and $Q_{\text{old}}$ denote the estimated Q-function before and after the update respectively. And for each policy improvement, we update the policy from $\pi_{\text{old}}$ to $\pi_{\text{new}}$ according to the newly estimated Q-function,

$$\pi_{\text{new}} = \arg\max_a Q_{\text{new}}(s, a).$$

Standard Q-learning only uses each experience once before disregarded, which is sample inefficient and can be improved by *experience replay* technique (Lin, 1992). We denote the experience that the agent collected at time $k$ by a tuple $e_k = \{s_k, a_k, r_k, s_k'\}$. According to experience replay, the experience $e_k$ is stored into the replay buffer and can be accessed multiple times during learning.

### 2.2. Value Metrics of Experience

To quantify the importance of experience, we derive three value metrics of experience. The utility of update on expe-

rience $e_k$ is defined as the value added to the cumulative discounted rewards starting from state $s_k$, after updating on $e_k$. Intuitively, choosing the most valuable experience for update will yield the highest utility to the agent. We denote such utility as the expected value of backup $\text{EVB}(e_k)$ (Mattar & Daw, 2018),

$$\begin{aligned} \text{EVB}(e_k) &= v_{\pi_{\text{new}}}(s_k) - v_{\pi_{\text{old}}}(s_k) \\ &= \sum_a \pi_{\text{new}}(a|s_k) q_{\pi_{\text{new}}}(s_k, a) \\ &\quad - \sum_a \pi_{\text{old}}(a|s_k) q_{\pi_{\text{old}}}(s_k, a), \end{aligned} \quad (1)$$

where $\pi_{\text{old}}$, $v_{\pi_{\text{old}}}$ and $q_{\pi_{\text{old}}}$ are respectively the policy, value function and Q-function before the update, and $\pi_{\text{new}}$, $v_{\pi_{\text{new}}}$, and $q_{\pi_{\text{new}}}$ are those after. As the update on experience $e_k$ consists of policy evaluation and policy improvement, the value of experience can further be separated to evaluation improvement value $\text{EIV}(e_k)$ and policy improvement value $\text{PIV}(e_k)$ by rewriting (1):

$$\begin{aligned} \text{EVB}(e_k) = \underbrace{\sum_a [\pi_{\text{new}}(a|s_k) - \pi_{\text{old}}(a|s_k)] q_{\pi_{\text{new}}}(s_k, a)}_{\text{PIV}(e_k)} + \\ \underbrace{\sum_a \pi_{\text{old}}(a|s_k)[q_{\pi_{\text{new}}}(s_k, a) - q_{\pi_{\text{old}}}(s_k, a)]}_{\text{EIV}(e_k)}, \quad (2) \end{aligned}$$

where $\text{PIV}(e_k)$ measures the value improvements due to the change of the policy, and $\text{EIV}(e_k)$ captures those due to the change of evaluation. Thus, we have three metrics for the value of experience: EVB, PIV and EIV.

### 2.3. Value Metrics of Experience in Q-Learning

For Q-learning, we use Q-function to estimate the true action-value function. A backup over an experience $e_k$ consists of policy evaluation with Bellman operator and greedy policy improvement. As the policy improvement is greedy, we can rewrite value metrics of experience to simpler forms. EVB can be written as follows from (1),

$$\text{EVB}(e_k) = \max_a Q_{\text{new}}(s_k, a) - \max_a Q_{\text{old}}(s_k, a). \quad (3)$$

Note that EVB here is different from that in (Mattar & Daw, 2018): in our case, EVB is derived from Q-learning; while in their case, EVB is derived from Dyna, a model-based RL algorithm (Sutton, 1990). Similarly, from (2), PIV can be written as

$$\text{PIV}(e_k) = \max_a Q_{\text{new}}(s_k, a) - Q_{\text{new}}(s_k, a_{\text{old}}), \quad (4)$$

where $a_{\text{old}} = \arg\max_a Q_{\text{old}}(s_k, a)$, and EIV can be written as

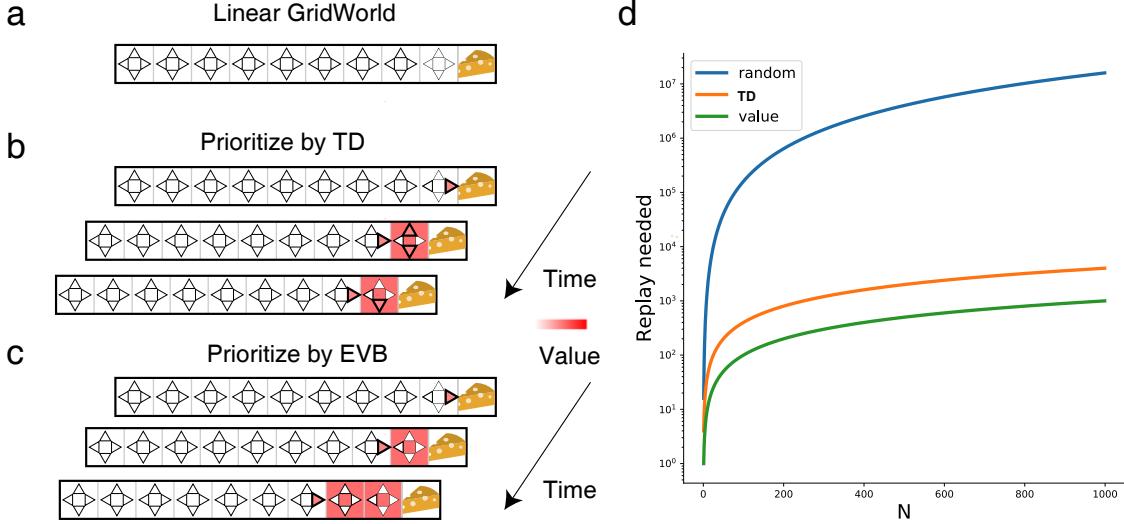$$\text{EIV}(e_k) = Q_{\text{new}}(s_k, a_{\text{old}}) - Q_{\text{old}}(s_k, a_{\text{old}}). \quad (5)$$

Figure 1: **a.** Illustration of the "Linear Grid-World" example: there are $N$ grids and 4 actions (north, south, east, west). Reward for entering the goal state (cheese) is 1; reward is 0 elsewhere. **b-c.** Examples of prioritized experience replay by $|\text{TD}|$ and value of experience (EVB). The main difference is that EVB only prioritizes the experiences that are associated with the optimal policy; while $|\text{TD}|$ is sensitive to changes in value function and will prioritize non-optimal experiences, such as those associated with north or south. Here squares represent states, triangles represent actions, and experiences associated with the highest priority are highlighted. **d.** Expected number of replays needed to learn the optimal policy, as the number of grids changes: uniform replay (blue), prioritized by $|\text{TD}|$ (orange), and EVB (green).

## 2.4. A Motivating Example

We illustrate the potential gain of value of experience in a "Linear Grid-World" environment (Figure 1a). This environment contains $N$ linearly-aligned grids and 4 actions (north, south, east, west). The rewards are rare: 1 for entering the goal state and 0 elsewhere. The solution for this environment is always choosing east.

We use this example to highlight the difference between prioritization strategies. Three agents perform Q-learning updates on the experiences drawn from the same replay buffer, which contains all the $(4N)$ experiences and associated rewards. The first agent replays the experiences uniformly at random, while the other two agents invoke the oracle to prioritize the experiences, which greedily select the experience with the highest $|\text{TD}|$ or EVB respectively. In order to learn the optimal policy, agents need to replay the experiences associated with action east in a reverse order.

For the agent with random replay, the expected number of replays required is $4N^2$ (Figure 1d). For the other two agents, prioritization significantly reduces the number of replays required: prioritization with $|\text{TD}|$ requires $4N$ replays, and prioritization with EVB only uses $N$ replays, which is optimal (Figure 1d). The main difference is that EVB only prioritizes the experiences that are associated with the optimal policy (Figure 1c), while $|\text{TD}|$ is sensitive to changes in the value function

in the value function and will prioritize non-optimal experiences: for example, the agent may choose the experiences associated with south or north in the second update, which are not optimal but have the same $|\text{TD}|$ as the experience associated with east (Figure 1b). Thus, EVB that directly quantifies the value of experience can serve as an optimal priority.

## 3. Upper Bounds of Value Metrics of Experience in Q-Learning

PER (Schaul et al., 2016) greatly improves the learning efficiency of DQN. However, the underlying rationale is not well understood. Here, we prove that $|\text{TD}|$ is the upper bound of the value metrics in Q-learning.

**Theorem 3.1.** *The three value metrics of experience $e_k$ in Q-learning ($|EVB|$, $|PIV|$ and $|EIV|$) are upper-bounded by $\alpha|TD(s_k, a_k, r_k, s'_k)|$, where $\alpha$ is a step-size parameter.*

*Proof.* From (3), $|\text{EVB}|$ can be written as

$$
\begin{aligned}
|\text{EVB}(e_k)| &= |\max_a Q_{\text{new}}(s_k, a) - \max_a Q_{\text{old}}(s_k, a)| \\
&\leq \max_a |Q_{\text{new}}(s_k, a) - Q_{\text{old}}(s_k, a)| \\
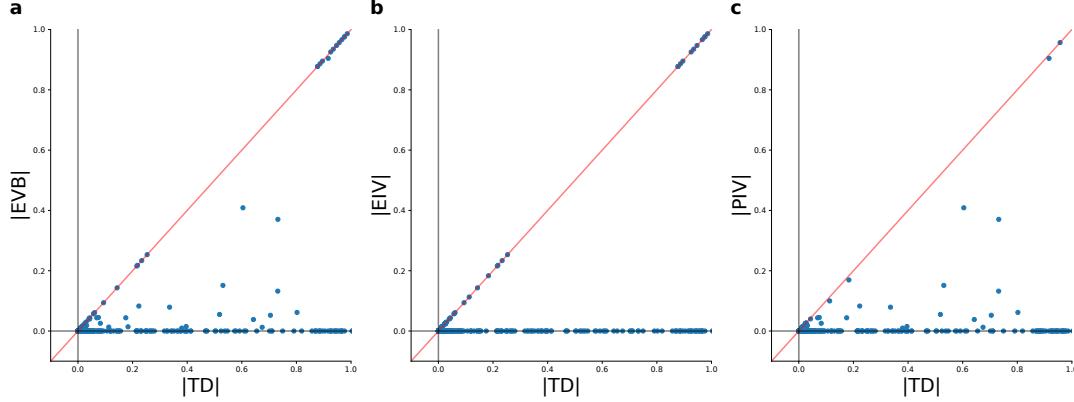&\leq \alpha|\text{TD}(s_k, a_k, r_k, s'_k)|,
\end{aligned}
\tag{6}
$$

Figure 2: The value metrics are upper-bounded by TD errors in Q-learning. **a-c.** |TD| *v.s.* |EVB| (left), |EIV| (middle) and |PIV| (right) of a tabular Q-learning agent in a grid-world maze. The red line indicates the identity line.

where the second line is from the contraction of max operator.

Proofs for the upper bounds of |PIV| and |EIV| are similar and given in Appendix A.1.  □

In Theorem 3.1, we prove that |EVB|, |PIV|, and |EIV| are upper-bounded by |TD| (scaled by the learning step-size) in Q-learning. To verify the bounds experimentally, we simulated a tabular Q-learning agent in a 5 × 5 grid-world maze [1]. The agent needs to reach the goal zone by moving one square in any of the four directions (north, south, east, west) each time (further details are described in Appendix A.4). For each transition, we record the associated TD error and value metrics. As we can see from Figure 2, all three value metrics of experience are bounded by |TD|. As our theory predicts (see Appendix A.1 for detail), |EIV| is either equal to |TD| (if the action of the experience is the optimal action before update) or 0. There is a large proportion of EVBs lies on the identity line, indicating the bound is tight. Moreover, we note that a significant proportion of value metrics lies on the x-axis. Because the value metrics are affected by the "on-policyness" of the experienced actions, and Q-learning learns a deterministic policy that makes most actions of experiences off-policy. As |TD| intrinsically tracks the evaluation and policy improvements, it can serve as an appropriate importance metric for past experiences.

## 4. Extension to Maximum-Entropy RL

In this section, we extend our framework to study the relationship between |TD| and value of experience in maximum-entropy RL, particularly, soft Q-learning.

---

[1]All the codes for the experiments are available at: https://github.com/AmazingAng/VER.

### 4.1. Soft Q-Learning

Unlike regular RL algorithms, maximum-entropy RL augments the reward with an entropy term: $r + \beta \mathcal{H}(\pi(\cdot|s))$, where $\mathcal{H}(\cdot)$ is the entropy, and $\beta$ is an optional temperature parameter that determines the relative importance of entropy and reward. The goal is to maximize the expected cumulative entropy-augmented rewards. Maximum-entropy RL algorithms have advantages at capturing multiple modes of near optimal policies, better exploration, and better transfer between tasks.

Soft Q-learning is an off-policy value-based algorithm built on maximum-entropy RL principles (Haarnoja et al., 2017; Schulman et al., 2017). Different from Q-learning, the target policy of soft Q-learning is stochastic. During policy iteration, Q-function is updated through soft Bellman operator $\Gamma^{\text{soft}}$, and the policy is updated to a maximum-entropy policy:

$$Q_{\text{new}}^{\text{soft}}(s,a) = [\Gamma^{\text{soft}} Q_{\text{old}}^{\text{soft}}](s,a) = r + \gamma V_{\text{old}}^{\text{soft}}(s')$$
$$\pi_{\text{new}}(a|s) = \text{softmax}_a(\frac{1}{\beta} Q_{\text{new}}^{\text{soft}}(s,a)),$$

where $\text{softmax}_i(x) = \exp(x_i)/\sum_i \exp(x_i)$ is the softmax function, and the soft value function $V_\pi^{\text{soft}}(s)$ is defined as,

$$V_\pi^{\text{soft}}(s) = \mathbb{E}_a\{Q_\pi^{\text{soft}}(s,a) - \log(\pi(a|s))\}$$
$$= \beta \log \sum_a \exp(\frac{1}{\beta} Q_\pi^{\text{soft}}(s,a)).$$

Similar as in Q-learning, the TD error in soft Q-learning (soft TD error) is given by:

$$\text{TD}^{\text{soft}}(s,a,r,s') = r + \gamma V_{\text{old}}^{\text{soft}}(s') - Q_{\text{old}}^{\text{soft}}(s,a).$$
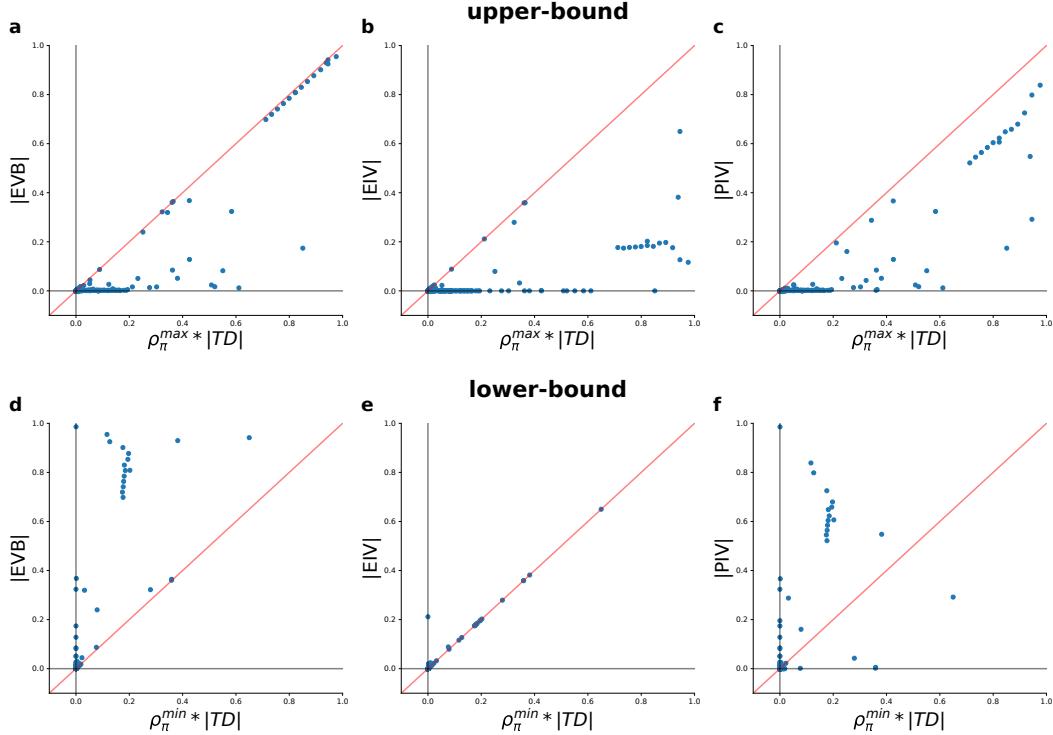
Figure 3: The value metrics and their bounds in soft Q-learning. |EVB| (left), |EIV| (middle) and |PIV| (right) as well as their theoretical lower-bound **a-c.** and upper-bounds **d-f.** of a tabular soft Q-learning agent in a grid-world maze. The red line indicates the identity line.

## 4.2. Value Metrics of Experience in Maximum-Entropy RL

Here, we extend the value metrics of experience to soft Q-learning. Similar as (1), EVB for maximum-entropy RL is defined as,

$$
\begin{aligned}
&\mathrm{EVB}^{\mathrm{soft}}(e_k) \\
&= v_{\mathrm{new}}^{\mathrm{soft}}(s_k) - v_{\mathrm{old}}^{\mathrm{soft}}(s_k) \\
&= \sum_a \pi_{\mathrm{new}}(a|s_k)\{q_{\mathrm{new}}^{\mathrm{soft}}(s_k, a) - \beta \log(\pi_{\mathrm{new}}(a|s_k))\} \quad (7) \\
&\quad - \sum_a \pi_{\mathrm{old}}(a|s_k)\{q_{\mathrm{old}}^{\mathrm{soft}}(s_k, a) - \beta \log(\pi_{\mathrm{old}}(a|s_k))\}
\end{aligned}
$$

$\mathrm{EVB}^{\mathrm{soft}}$ can be separated into $\mathrm{PIV}^{\mathrm{soft}}$ and $\mathrm{EIV}^{\mathrm{soft}}$, which respectively quantify the value of policy and evaluation improvement in soft Q-learning,

$$
\begin{aligned}
\mathrm{PIV}^{\mathrm{soft}}(e_k) = \sum_a &\{\pi_{\mathrm{new}}(a|s_k) - \pi_{\mathrm{old}}(a|s_k)\}q_{\mathrm{new}}^{\mathrm{soft}}(s_k, a) \\
&+ \beta(H(\pi_{\mathrm{new}}(\cdot|s)) - H(\pi_{\mathrm{old}}(\cdot|s_k))),
\end{aligned}
\quad (8)
$$

$$
\mathrm{EIV}^{\mathrm{soft}}(e_k) = \sum_a \pi_{\mathrm{old}}(a|s_k)[q_{\mathrm{new}}^{\mathrm{soft}}(s_k, a) - q_{\mathrm{old}}^{\mathrm{soft}}(s_k, a)]. \quad (9)
$$

Value metrics of experience in maximum-entropy RL have similar forms as in regular RL except for the entropy term,

because changes in policy lead to changes in the policy entropy and affect the entropy-augmented rewards.

## 4.3. Lower and Upper Bounds of Value Metrics of Experience in Soft Q-learning

We theoretically derive the lower and upper bounds of the value metrics of experience in soft Q-learning.

**Theorem 4.1.** *The three value metrics of experience $e_k$ in soft Q-learning ($|EVB^{soft}|$, $|PIV^{soft}|$ and $|EIV^{soft}|$) are upper-bounded by $\rho_\pi^{max} * \left|TD^{soft}\right|$, where $\rho_\pi^{max} = \max\{\pi_{old}(a_k|s_k), \pi_{new}(a_k|s_k)\}$ is a policy related term.*

*Proof.* See Appendix A.2. □

**Theorem 4.2.** *For soft Q-learning, $|EVB^{soft}|$ and $|EIV^{soft}|$ (but not $|PIV^{soft}|$) are lower-bounded by $\rho_\pi^{min} * \left|TD^{soft}\right|$, where $\rho_\pi^{min} = \min\{\pi_{old}(a_k|s_k), \pi_{new}(a_k|s_k)\}$ is a policy related term.*

*Proof.* See Appendix A.3. □

The lower and upper bounds in soft Q-learning include a policy term with |TD|. The policy related term $\rho_\pi$ quantifies the "on-policyness" of the experienced action. And the
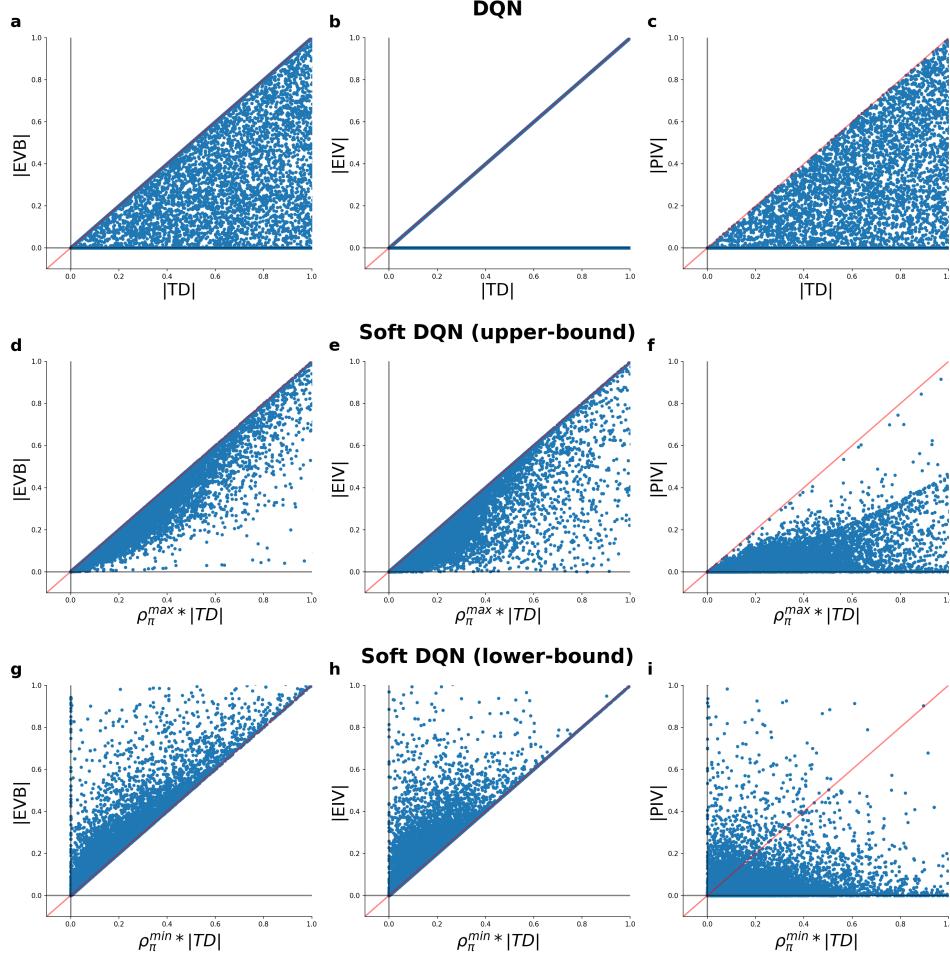
Figure 4: Results of DQN and soft DQN in CartPole. **a-c.** |TD| *v.s.* |EVB| (left), |EIV| (middle) and |PIV| (right) in DQN. **d-f.** Theoretical upper bound and (**g-i.**) lower bound *v.s.* |EVB| (left), |EIV| (middle) and |PIV| (right) in soft DQN. The red line indicates the identity line.

bounds become tighter as the difference between $\pi_{\text{old}}(a_k|s_k)$ and $\pi_{\text{new}}(a_k|s_k)$ becomes smaller. Surprisingly, the coefficient of the entropy term $\beta$ impacts the bound only through the policy term, which makes it an excellent priority even $\beta$ changes during learning (Haarnoja et al., 2018). As $0 \le \rho_{\pi}^{\max} \le 1$, the value metrics are also upper-bounded by |TD| alone, which is similar as in Q-learning. However, as $\pi(a_k|s_k)$ is usually less than 1, |TD| is a looser upper bound in soft Q-learning.

To verify the bounds in soft Q-learning experimentally, we simulated a tabular soft Q-learning agent in the grid-world maze described previously. From upper panel of Figure 3, all three value metrics of experience are upper-bounded by $\rho_{\pi}^{\max} * $|TD|. Moreover, from lower panel of Figure 3, $|\text{EVB}^{\text{soft}}|$ and $|\text{EIV}^{\text{soft}}|$ (but not $|\text{PIV}^{\text{soft}}|$) are lower-bounded by $\rho_{\pi}^{\min} * \left|\text{TD}^{\text{soft}}\right|$, supporting our theoretical analysis (Theorem 4.1 and 4.2). The proportion of non-zero values of experiences is higher in soft Q-learning than in Q-learning,

because, different from greedy policy of Q-learning, soft Q-learning learns a stochastic policy that makes experiences more "on-policy" and have non-sparse values. In summary, the experimental results support the theoretical bounds of value metrics in tabular soft Q-learning.

## 5. Extension to Function Approximation Methods

Function approximation methods, which are more powerful and expressive than tabular methods, are effective in solving more challenging tasks, such as the game of Go (Silver et al., 2016), video games (Mnih et al., 2015) and robotic control (Haarnoja et al., 2017). In these methods, we learn a parameterized Q-function $Q(s, a; \theta_t)$, where the parameters are updated on experience $e_k$ through gradient-based method,

$$\theta_{t+1} = \theta_t + \alpha \text{TD} \nabla_{\theta_t} Q(s_k, a_k; \theta_t)),$$

where $\alpha$ is the learning rate, the TD error is defined as:

$$\text{TD} = Q_{\text{target}}(s_k, a_k) - Q(s_k, a_k; \theta_t),$$

and the target Q-value $Q_{\text{target}}$ is defined as

$$Q_{\text{target}}(s_k, a_k) = r_k + \gamma \max_{a'} Q(s'_k, a'; \theta_t).$$

As $\alpha$ in function approximation Q-learning is usually very small, for each update, the parameterized function moves to its target by a small amount.

Our framework can be extended to function approximation methods by slightly modifying the definition of the value metrics of experience. Note that if we apply the original definition of EVB in (3) directly to function approximation methods, the Q-function after the update $Q_{\text{new}}(s, a) = Q(s, a; \theta_{t+1})$ involves gradient-based update, which complicates the analysis and breaks the inequalities derived in the tabular case. As a remedy, we replace $Q(s, a; \theta_{t+1})$ by the target Q-value $Q_{\text{target}}(s, a)$ in the value metrics of experience (1-5) and (7-9). The intuition is simple: the value is defined by the cause of the update (target Q-value), but not the result of the update through gradient-based method. Moreover, this modification allows our theory to apply to all function approximation methods, regardless the specific forms of the function approximator (linear function or neural networks). After the modifications, the value metrics of experience have similar form as the tabular case, and all theorems derived in the tabular case can be applied to function approximation methods.

To test whether our theoretical predictions hold in function approximation methods, we simulated one DQN (Deep-Q network) agent and one soft DQN (DQN with soft update) agent in CartPole environment, where the goal is to keep the pole balanced by moving the cart forward and backward (further details are described in Appendix A.4). From Figure 4, all value metrics of experience in DQN (Figure 4a-c) and soft DQN (Figure 4d-f) are bounded by the theoretical upper bounds. For DQN, |EVB| and |PIV| are uniformly distributed in the bounded area, while |EIV| are equal to |TD| or 0. Results are different in soft DQN, where |EVB| and |EIV| are distributed more closely towards the theoretical upper bounds, suggesting the upper bound in soft Q-learning is tighter. Moreover, (Figure 4g-i) shows the |EVB| and |EIV| are lower-bounded by $\rho_\pi^{\min} * \left|\text{TD}^{\text{soft}}\right|$, while |PIV| are not. The experimental results confirm the bounds of value metrics hold for function approximation methods.

## 6. Experiments on Atari Games

The theoretical upper bound ($\rho_\pi^{\max} * |\text{TD}|$) of value metrics balances the prediction error and "on-policyness" of the experience. To better illustrate the difference between the theoretical upper bound and |TD|, we randomly drew 50
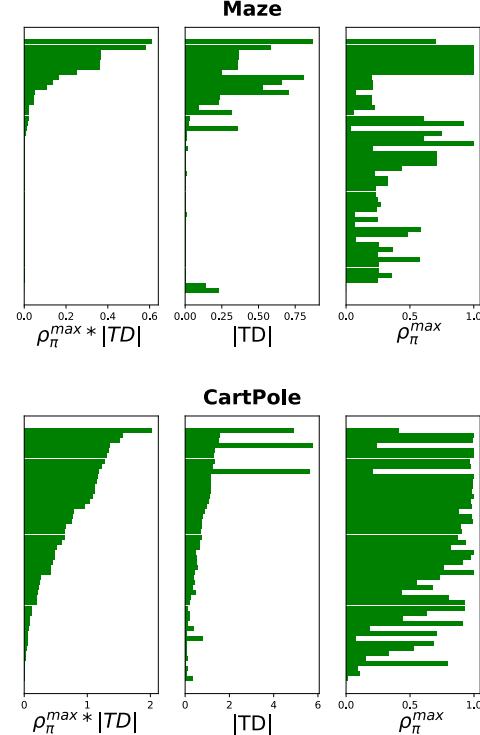


Figure 5: Illustration on the difference between |TD| and theoretical upper-bound for value metrics in soft Q-learning. Depicted are the theoretical upper bound (left), |TD| (middle), and the policy term (right) of 50 experiences from the replay buffer in the grid-world maze (upper panel) and CartPole (lower panel), ordered by the theoretical upper bound.

experiences from the grid-world maze and CartPole experiments. From Figure 5, we can see that the experiences with highest theoretical upper bounds are associated with higher |TD| and "on-policyness". To investigate whether the theoretical upper bound can serve as an appropriate priority for experience replay in soft Q-learning, we compare the performance of soft DQN with different prioritization strategies: uniform replay, prioritization with |TD| or the theoretical upper bound ($\rho_\pi^{\max} * \left|\text{TD}^{\text{soft}}\right|$), which are denoted by soft DQN, PER and VER (valuable experience replay) respectively. This set of experiments consists of 9 selected Atari 2600 games according to (Schulman et al., 2017), which balances generality of the games and limited compute power. We closely follow the experimental setting and network architecture outlined by (Mnih et al., 2015). For each game, the network is trained on a single GPU for 40M frames, or approximately 5 days. More details for the settings and hyperparameters are available in Appendix A.4.

Figure 6 shows that soft DQN prioritized by |TD| or the theoretical upper bound substantially outperforms uniform replay in most of the games. On average, soft DQN with PER
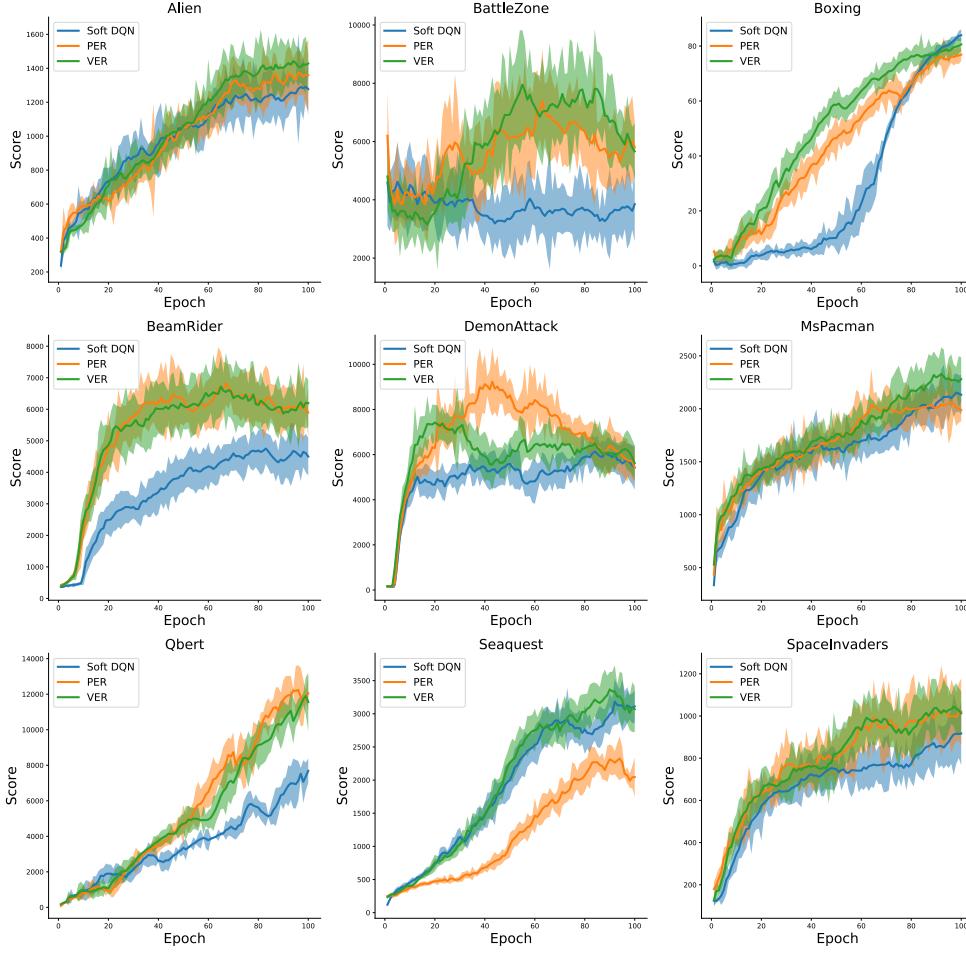
Figure 6: Learning curve of soft DQN (blue lines), and soft DQN with prioritized experience replay in term of soft TD error (PER, orange lines) and the theoretical upper bound of value metrics of experience (VER, green lines) on Atari games. Solid lines are average return over 8 evaluation runs and shaded area is the standard error of the mean.

or VER outperform vanilla soft DQN by 11.8% or 18.0% respectively. Moreover, VER shows higher convergence speed and outperforms PER in most of the games (8.47% on average), which suggest that a tighter upper bound on value metrics improves the performance of experience replay. These results suggest that the theoretical upper bound can serve as an appropriate priority for experience replay in soft Q-learning.

## 7. Discussion

In this work, we formulate a framework to study relationship between the importance of experience and $|TD|$. To quantify the importance of experience, we derive three value metrics of experience: expected value of backup, policy evaluation value, and policy improvement value. For Q-learning, we theoretically show these value metrics are upper-bounded by $|TD|$. Thus, $|TD|$ implicitly tracks the value of the experience, which leads to high sample efficiency of PER. Further-

more, we extend our framework to maximum-entropy RL, by showing that these value metrics are lower and upper-bounded by the product of a policy term and $|TD|$. Experiments in grid-world maze and CartPole support our theoretical results for both tabular and function approximation RL methods, showing that the derived bounds hold in practice. Moreover, we show that experience replay using the upper bound as a priority improves maximum-entropy RL (*i.e.*, soft DQN) in Atari games.

By linking $|TD|$ and value of experience, two important quantities in learning, our study has the following implications. First, from a machine learning perspective, our study provides a framework to derive appropriate priorities of experience for different algorithms, with possible extension to batch RL (Fu et al., 2020) and N-step learning (Hessel et al., 2017). Second, for neuroscience, our work provides insight on how brain might encode the importance of experience. Since $|TD|$ biologically corresponds to the

reward prediction-error signal in the dopaminergic system (Schultz et al., 1997; Glimcher, 2011) and implicitly tracks the value of the experience, the brain may account on it to differentiate important experiences.

# References

Daley, B. and Amato, C. Reconciling $\lambda$-returns with experience replay. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

El Ghaoui, L. Optimization models and applications. http://livebooklabs.com/keeppies/c5a5868ce26b8125, 2018. Livebook visited Spring 2018.

Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.

Glimcher, P. W. Understanding dopamine and reinforcement learning: the dopamine reward prediction error hypothesis. *Proceedings of the National Academy of Sciences*, 108(Supplement 3):15647–15654, 2011.

Gruber, M. J., Ritchey, M., Wang, S.-f., Doss, M. K., and Ranganath, C. Post-learning Hippocampal Dynamics Promote Preferential Retention of Rewarding Events. *Neuron*, 89(5):1110–1120, 2016.

Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning (ICML)*, 2017.

Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.

Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., and Silver, D. Rainbow: Combining improvements in deep reinforcement learning. *arXiv preprint arXiv:1710.02298*, 2017.

Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., and Silver, D. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.

Horgan, D., Quan, J., Budden, D., Barth-Maron, G., Hessel, M., Van Hasselt, H., and Silver, D. Distributed prioritized experience replay. In *International Conference on Learning Representations (ICLR)*, 2018.

Kapturowski, S., Ostrovski, G., Quan, J., Munos, R., and Dabney, W. Recurrent experience replay in distributed reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2019.

Katharopoulos, A. and Fleuret, F. Not all samples are created equal: Deep learning with importance sampling. *arXiv preprint arXiv:1803.00942*, 2018.

Lin, L.-J. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992.

Mattar, M. G. and Daw, N. D. Prioritized memory access explains planning and hippocampal replay. *Nature Neuroscience*, 21(11):1609–1617, 2018.

Michon, F., Sun, J.-J., Kim, C. Y., Ciliberti, D., and Kloosterman, F. Post-learning hippocampal replay selectively reinforces spatial memory for highly rewarded locations. *Current Biology*, 29(9):1436–1444, 2019.

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.

Ólafsdóttir, H. F., Bush, D., and Barry, C. The role of hippocampal replay in memory and planning. *Current Biology*, 28(1):R37–R50, 2018.

Pan, Y., Zaheer, M., White, A., Patterson, A., and White, M. Organizing experience: a deeper look at replay mechanisms for sample-based planning in continuous state domains. *arXiv preprint arXiv:1806.04624*, 2018.

Roscow, E. L., Jones, M. W., and Lepora, N. F. Behavioural and computational evidence for memory consolidation biased by reward-prediction errors. *bioRxiv*, pp. 716290, 2019.

Salvetti, B., Morris, R. G., and Wang, S.-H. The role of rewarding and novel events in facilitating memory persistence in a separate spatial memory task. *Learning & memory*, 21(2):61–72, 2014.

Schapiro, A. C., McDevitt, E. A., Rogers, T. T., Mednick, S. C., and Norman, K. A. Human hippocampal replay during rest prioritizes weakly learned information and predicts memory performance. *Nature communications*, 9(1):1–11, 2018.

Schaul, T., Quan, J., Antonoglou, I., and Silver, D. Prioritized experience replay. In *International Conference on Learning Representations (ICLR)*, 2016.

Schlegel, M., Chung, W., Graves, D., Qian, J., and White, M. Importance resampling for off-policy prediction. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

Schulman, J., Chen, X., and Abbeel, P. Equivalence between policy gradients and soft q-learning. *arXiv preprint arXiv:1704.06440*, 2017.

Schultz, W., Dayan, P., and Montague, P. R. A neural substrate of prediction and reward. *Science*, 275(5306): 1593–1599, 1997.

Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.

Sutton, R. S. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Machine learning proceedings 1990*, pp. 216–224. Elsevier, 1990.

Watkins, C. J. C. H. and Dayan, P. Q-learning. *Machine Learning*, 8(3-4):279–292, 1992.

Zha, D., Lai, K.-H., Zhou, K., and Hu, X. Experience replay optimization. *arXiv preprint arXiv:1906.08387*, 2019.

# A. Appendix

### A.1. Proof of Theorem 3.1

In this section, we derive the upper bounds of |PIV| and |EIV| in Q-learning. |PIV| can be written as

$$
\begin{aligned}
|\text{PIV}(e_k)| &= |\max_a Q_{\text{new}}(s_k, a) - Q_{\text{new}}(s_k, \arg\max_a Q_{\text{old}}(s_k, a))| \\
&= \max_a Q_{\text{new}}(s_k, a) - Q_{\text{new}}(s_k, \arg\max_a Q_{\text{old}}(s_k, a)) \\
&= \max_a Q_{\text{new}}(s_k, a) - \max_a Q_{\text{old}}(s_k, a) - \mathbf{1}_{a_{\text{old}}=a_k} \alpha \text{TD}(s_k, a_k, r_k, s_k')
\end{aligned}
\tag{10}
$$

where the second line is from that the change in Q-function following greedy policy improvement is greater or equal to 0, and the third line is from the update of Q-function. For $\text{TD}(s_k, a_k, r_k, s_k') \geq 0$, we have

$$
0 \leq \max_a Q_{\text{new}}(s_k, a) - \max_a Q_{\text{old}}(s_k, a) \leq \alpha \text{TD}(s_k, a_k, r_k, s_k').
$$

And for $\text{TD}(s_k, a_k, r_k, s_k') \leq 0$, we have

$$
\max_a Q_{\text{new}}(s_k, a) - \max_a Q_{\text{old}}(s_k, a) \leq 0.
$$

Bring above inequalities to 10, we have

$$
|\text{PIV}(e_k)| \leq \alpha |\text{TD}(s_k, a_k, r_k, s_k')|
\tag{11}
$$

Similarly, |EIV| can be written as follows,

$$
\begin{aligned}
|\text{EIV}(e_k)| &= |Q_{\text{new}}(s_k, a_{\text{old}}) - Q_{\text{old}}(s_k, a_{\text{old}})| \\
&= \mathbf{1}_{s=s_k, a_{\text{old}}=a_k} \alpha |\text{TD}(s_k, a_k, r_k, s_k')| \\
&\leq \alpha |\text{TD}(s_k, a_k, r_k, s_k')|
\end{aligned}
\tag{12}
$$

For equations (6) and (11), the equality is reached if the experienced action is the same as the best action before and after the update. For (12), the equality is met if the experienced action is the best action before update. □

### A.2. Proof of Theorem 4.1

In this section, we derive upper bounds of value metrics of experience in soft Q-learning. For soft-Q learning, $|\text{EVB}^{\text{soft}}|$ can be written as

$$
\begin{aligned}
|\text{EVB}^{\text{soft}}(e_k)| &= |\beta \log \sum_a \exp(\frac{1}{\beta} Q_{\text{new}}^{\text{soft}}(s_k, a)) - \beta \log \sum_a \exp(\frac{1}{\beta} Q_{\text{old}}^{\text{soft}}(s_k, a))| \\
&= |\beta \log \sum_a \exp\left(\frac{1}{\beta}(Q_{\text{old}}^{\text{soft}}(s_k, a) + \mathbf{1}_{a=a_k} \text{TD}^{\text{soft}})\right) - \beta \log \sum_a \exp \frac{1}{\beta} Q_{\text{old}}^{\text{soft}}(s_k, a)|.
\end{aligned}
$$

Let us define the LogSumExp function $F(\vec{x}) = \beta \log \sum_i \exp\left(\frac{x_i}{\beta}\right)$. The LogSumExp function $F(\vec{x})$ is convex, and is strictly and monotonically increasing everywhere in its domain (El Ghaoui, 2018). The partial derivative of $F(\vec{x})$ is a softmax function

$$
\frac{\partial F(\vec{x})}{\partial x_i} = \text{softmax}_i(\frac{1}{\beta}\vec{x}) \geq 0,
$$

which takes the same form as the policy of soft Q-learning. For $\epsilon < 0$, we have:

$$
\epsilon \frac{\partial F(x_1, ..., x_i, ...)}{\partial x_i} \leq F(x_1, ..., x_i + \epsilon, ...) - F(x_1, ..., x_i, ...) \leq 0.
$$

Similarly, for $\epsilon \geq 0$, we have,

$$
0 \leq F(x_1, ..., x_i + \epsilon, ...) - F(x_1, ..., x_i, ...) \leq \epsilon \frac{\partial F(x_1, ..., x_i + \epsilon, ...)}{\partial x_i}.
$$

By substituting $x_i$ by $Q_{\text{old}}^{\text{soft}}(s_k, a_k)$ and $\epsilon$ by $\text{TD}^{\text{soft}}$, and rewriting partial derivative of $F(\vec{x})$ into policy form, we have following inequalities. For $\text{TD}^{\text{soft}} \leq 0$

$$\pi_{\text{old}}(a_k|s_k)\text{TD}^{\text{soft}} \leq \beta \log \sum_a \exp(\frac{1}{\beta}Q_{\text{new}}^{\text{soft}}(s_k, a)) - \beta \log \sum_a \exp(\frac{1}{\beta}Q_{\text{old}}^{\text{soft}}(s_k, a)) \leq 0$$

Similarly, for $\text{TD}^{\text{soft}} > 0$, we have :

$$0 \leq \beta \log \sum_a \exp(\frac{1}{\beta}Q_{\text{new}}^{\text{soft}}(s_k, a)) - \beta \log \sum_a \exp(\frac{1}{\beta}Q_{\text{old}}^{\text{soft}}(s_k, a)) \leq \pi_{\text{new}}(a_k|s_k)\text{TD}^{\text{soft}}$$

Thus, we have the upper bounds of $|\text{EVB}^{\text{soft}}|$,

$$\left|\text{EVB}^{\text{soft}}(e_k)\right| \leq \max\{\pi_{\text{old}}(a_k|s_k), \pi_{\text{new}}(a_k|s_k)\} * \left|\text{TD}^{\text{soft}}\right|.$$

For $|\text{PIV}^{\text{soft}}|$, we have,

$$|\text{PIV}^{\text{soft}}(e_k)| = |\sum_a \pi_{\text{new}}(a|s)\{Q_{\text{new}}^{\text{soft}}(s_k, a) - \beta \log(\pi_{\text{new}}(a|s))\}$$
$$- \sum_a \pi_{\text{old}}(a|s)\{Q_{\text{new}}^{\text{soft}}(s_k, a) - \beta \log(\pi_{\text{old}}(a|s))\}|$$
$$= \sum_a \pi_{\text{new}}(a|s)\{Q_{\text{new}}^{\text{soft}}(s_k, a) - \beta \log(\pi_{\text{new}}(a|s))\}$$
$$- \sum_a \pi_{\text{old}}(a|s)\{Q_{\text{old}}^{\text{soft}}(s_k, a) - \beta \log(\pi_{\text{old}}(a|s))\} - \pi_{\text{old}}(a_k|s)\text{TD}^{\text{soft}}$$
$$= \beta \log \sum_a \exp \frac{Q_{\text{new}}^{\text{soft}}(s_k, a)}{\beta} - \beta \log \sum_a \exp \frac{Q_{\text{old}}^{\text{soft}}(s_k, a)}{\beta} - \pi_{\text{old}}(a_k|s)\text{TD}^{\text{soft}},$$

where the second line is because the policy improvement value is always greater than or equal to 0, and the third line is by reordering the equation.

For $\text{TD}^{\text{soft}} > 0$, we have:

$$0 \leq \beta \log \sum_a \exp \frac{Q_{\text{new}}^{\text{soft}}(s_k, a)}{\beta} - \beta \log \sum_a \exp \frac{Q_{\text{old}}^{\text{soft}}(s_k, a)}{\beta} - \pi_{\text{old}}(a_k|s)\text{TD}^{\text{soft}} \leq \pi_{\text{new}}(a_k|s_k)\text{TD}^{\text{soft}}$$

For $\text{TD}^{\text{soft}} \leq 0$, we have:

$$0 \leq \beta \log \sum_a \exp \frac{Q_{\text{new}}^{\text{soft}}(s_k, a)}{\beta} - \beta \log \sum_a \exp \frac{Q_{\text{old}}^{\text{soft}}(s_k, a)}{\beta} - \pi_{\text{old}}(a_k|s)\text{TD}^{\text{soft}} \leq \pi_{\text{old}}(a_k|s_k)\text{TD}^{\text{soft}}$$

Thus, we have the upper bounds of $\left|\text{PIV}^{\text{soft}}\right|$:

$$\left|\text{PIV}^{\text{soft}}(e_k)\right| \leq \max\{\pi_{\text{old}}(a_k|s_k), \pi_{\text{new}}(a_k|s_k)\} * \left|\text{TD}^{\text{soft}}\right|$$

Also, for $|\text{EIV}^{\text{soft}}|$, we have:

$$|\text{EIV}^{\text{soft}}(e_k)| = |\sum_a \pi_{\text{old}}(a|s)[Q_{\text{new}}^{\text{soft}}(s_k, a) - Q_{\text{old}}^{\text{soft}}(s_k, a)]|$$
$$= \pi_{\text{old}}(a|s) * \left|\text{TD}^{\text{soft}}\right|$$
$$\leq \max\{\pi_{\text{old}}(a_k|s_k), \pi_{\text{new}}(a_k|s_k)\} * \left|\text{TD}^{\text{soft}}\right|$$

There is no lower bound of the similar form for $|\text{PIV}^{\text{soft}}|$. $\qquad\square$
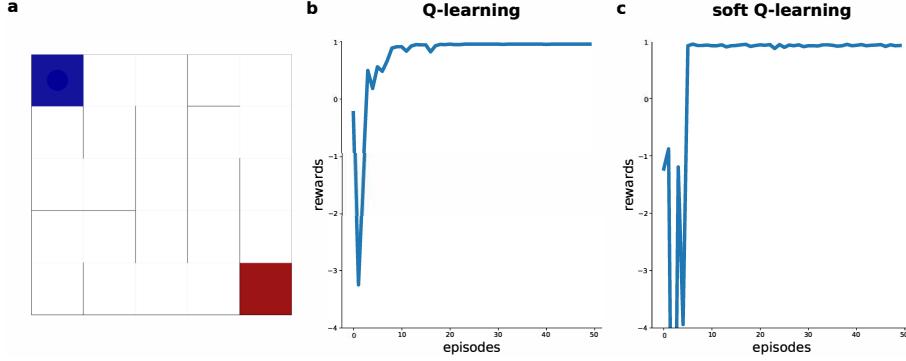
Figure 7: Grid-world maze environment and learning curves.

## A.3. Proof of Theorem 4.2

In this section, we derive lower bounds of value metrics of experience in soft Q-learning. Similar as deriving upper bounds in Appendix A.2, we derive the lower bounds for $|\text{EVB}|$ using the the LogSumExp function $F(\vec{x}) = \beta \log \sum_i \exp\left(\frac{x_i}{\beta}\right)$. For $\epsilon < 0$, we have:

$$F(x_1, ..., x_i + \epsilon, ...) - F(x_1, ..., x_i, ...) \le \epsilon \frac{\partial F(x_1, ..., x_i + \epsilon, ...)}{\partial x_i} \le 0.$$

Similarly, for $\epsilon \ge 0$, we have,

$$F(x_1, ..., x_i + \epsilon, ...) - F(x_1, ..., x_i, ...) \ge \epsilon \frac{\partial F(x_1, ..., x_i, ...)}{\partial x_i} \ge 0.$$

By substituting $x_i$ by $Q^{\text{soft}}_{\text{old}}(s_k, a_k)$ and $\epsilon$ by $\text{TD}^{\text{soft}}$, and rewriting partial derivative of $F(\vec{x})$ into policy form, we have following inequalities. For $\text{TD}^{\text{soft}} \le 0$

$$\beta \log \sum_a \exp(\frac{1}{\beta} Q^{\text{soft}}_{\text{new}}(s_k, a)) - \beta \log \sum_a \exp(\frac{1}{\beta} Q^{\text{soft}}_{\text{old}}(s_k, a)) \le \pi_{\text{new}}(a_k|s_k)\text{TD}^{\text{soft}} \le 0$$

Similarly, for $\text{TD}^{\text{soft}} > 0$, we have :

$$\beta \log \sum_a \exp(\frac{1}{\beta} Q^{\text{soft}}_{\text{new}}(s_k, a)) - \beta \log \sum_a \exp(\frac{1}{\beta} Q^{\text{soft}}_{\text{old}}(s_k, a)) \ge \pi_{\text{old}}(a_k|s_k)\text{TD}^{\text{soft}} \ge 0$$

Thus, we have the lower bounds of $|\text{EVB}^{\text{soft}}|$,

$$\left|\text{EVB}^{\text{soft}}(e_k)\right| \ge \min\{\pi_{\text{old}}(a_k|s_k), \pi_{\text{new}}(a_k|s_k)\} * \left|\text{TD}^{\text{soft}}\right|.$$

For $|\text{EIV}^{\text{soft}}|$, we have:

$$|\text{EIV}^{\text{soft}}(e_k)| = |\sum_a \pi_{\text{old}}(a|s)[Q^{\text{soft}}_{\pi_{\text{new}}}(s_k, a) - Q^{\text{soft}}_{\pi_{\text{old}}}(s_k, a)]|$$
$$= \pi_{\text{old}}(a|s) * \left|\text{TD}^{\text{soft}}\right|$$
$$\ge \min\{\pi_{\text{old}}(a_k|s_k), \pi_{\text{new}}(a_k|s_k)\} * \left|\text{TD}^{\text{soft}}\right|$$

$\square$

## A.4. Experimental Details

### A.4.1. GRID-WORLD MAZE

For the grid-world maze experiments in section 3, we use a maze environment of a $5 \times 5$ square with walls, as depicted in Figure 7a. The agent needs to reach the goal zone in the bottom-right corner. At each time step, the agent can choose

to move one square in any of the four directions (north, south, east, west). If the move is blocked by a wall or the border of the maze, the agent stays in place. Every time step, the agent gets a reward of $-0.004$ or $1$ if it enters the goal zone and the episode ends. The discount factor is $0.99$ throughout the experiments. For these experiments, we use a tabular setting for Q-learning and soft Q-learning according to section 2.1 and 4.1. For Q-learning, the behavior policy is $\epsilon$-greedy, where $\epsilon$ decays exponentially from 1 to 0.001 during training. And we set learning step size $\alpha = 1$. For soft Q-learning, the temperature parameter $\beta$ is set to 100. Total trial number is 50 for each algorithm. During training, both algorithms successfully solve the maze game, see Figure 7b-c for the learning curves.

### A.4.2. CARTPOLE

For CartPole, the goal is to keep the pole balanced by moving the cart forward and backward for 200 steps. We test our theoretical prediction on DQN and soft-DQN (DQN with soft-update). For DQN, we implement the model according to (Mnih et al., 2015), where we replace the original Q-network with a two-layer MLP, with 256 Relu neurons in each layer. The $\epsilon$ in $\epsilon$-greedy policy decays exponentially from 1 to 0.01 for the first $10,000$ steps, and remains 0.01 afterwards. For soft-DQN, all settings are the same with DQN, except for two modifications: for policy evaluation, the (soft) Q-network is updated according to the soft TD error; the policy follows maximum-entropy policy, calculated as the softmax of the soft Q values (see section 4.1). The temperature parameter $\beta$ is set to 0.5. For both algorithms, the discount factor is 0.99, the learning rate is 0.005, experience buffer size is 1000, the batch size is 16 and total environment interaction is $50,000$.

### A.4.3. ATARI GAMES

For this set of experiments, we compare the performance of vanilla soft DQN and soft DQN with PER, where we use |TD| and the theoretical upper bound as priorities (Schaul et al., 2016), respectively denoted as PER and VER (valuable experience replay). We select 9 Atari games for the experiments: Alien, BattleZone, Boxing, BeamRider, DemonAttack, MsPacman, Qbert, Seaquest and SpaceInvaders. The vanilla soft DQN is similar to that described in the above section, but the Q-network has the same architecture as in (Mnih et al., 2015). We implement PER on soft-DQN according to (Schaul et al., 2016). For all algorithms, the temperature parameter $\beta$ is 0.05, the discount factor is 0.99, the learning rate is 1e$-4$, experience buffer size is 1M, the batch size is 32, total environment interaction is $50,000$. For PER or VER, the parameters for importance sampling are $\alpha_{\text{IS}} = 0.4$ and $\beta_{\text{IS}} = 0.6$. For each game, the network is trained on a single GPU for 40M frames, or approximately 5 days.