
Regret Minimization Experience Replay in Off-Policy Reinforcement Learning

Xu-Hui Liu*, **Zhenghai Xue***, **Jing-Cheng Pang**, **Shengyi Jiang**, **Feng Xu**, **Yang Yu†**

National Key Laboratory for Novel Software Technology

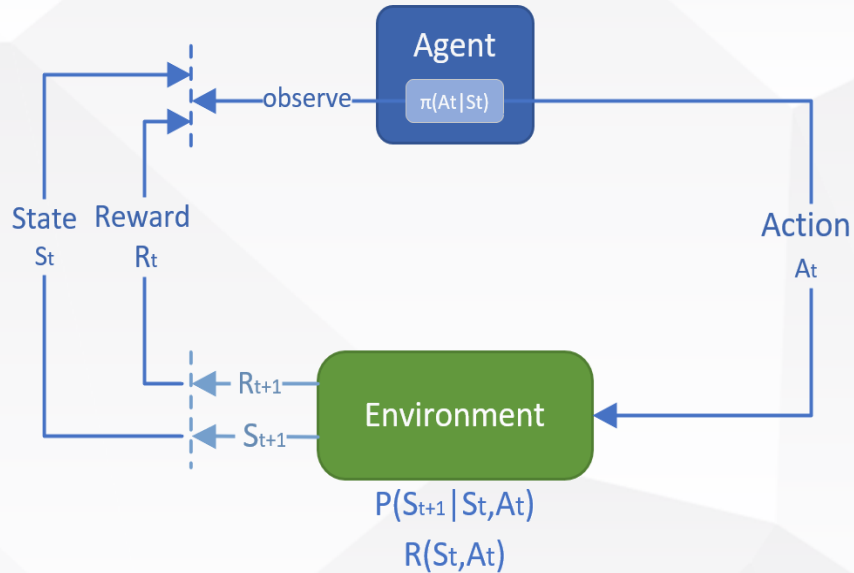
Nanjing University, Nanjing 210023, China

liuxh@lamda.nju.edu.cn, xuezh@smail.nju.edu.cn

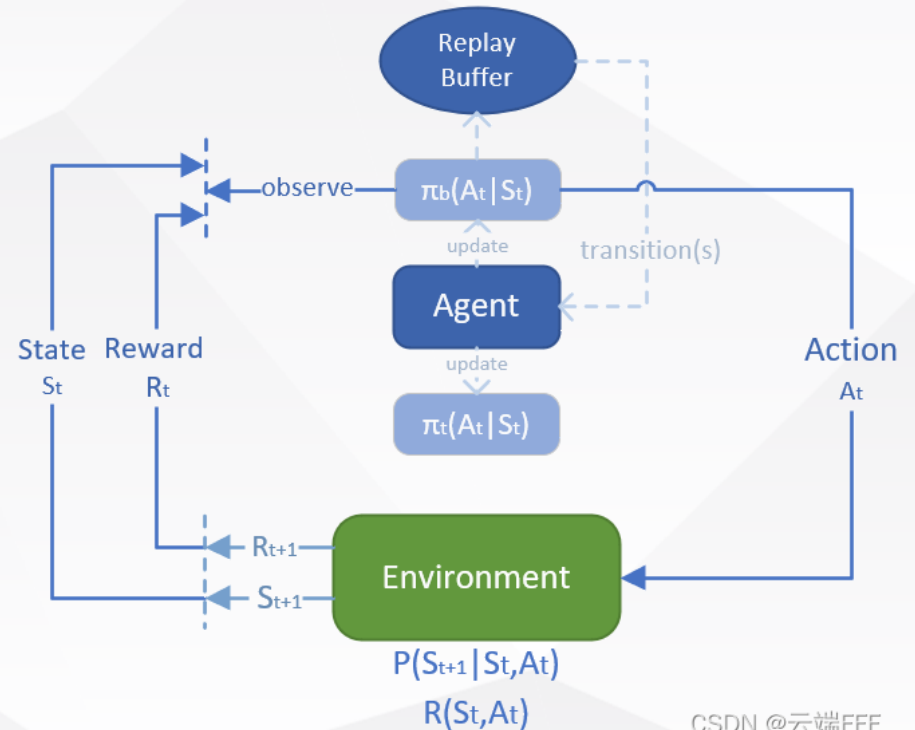
{pangjc, jiangsy, xufeng}@lamda.nju.edu.cn, yuy@nju.edu.cn

NIPS 2021

Experience Replay



CSDN @云端FFF



CSDN @云端FFF

- eliminate circular dependencies
- higher data efficiency
- better data distribution (i.i.d)

Non-uniform Experience Replay

- Key Idea: RL agent can learn **more effectively** from **some transitions** than from others
- Any **loss function evaluated with non-uniformly sampled** data can be transformed into **another uniformly sampled loss function with the same expected gradient**

importance sampling ratio

$$\underbrace{\mathbb{E}_{i \sim \mathcal{D}_1} [\nabla_Q \mathcal{L}_1(\delta(i))]}_{\text{expected gradient of } \mathcal{L}_1 \text{ under } \mathcal{D}_1} = \mathbb{E}_{i \sim \mathcal{D}_2} \left[\frac{p_{\mathcal{D}_1}(i)}{p_{\mathcal{D}_2}(i)} \nabla_Q \mathcal{L}_1(\delta(i)) \right].$$

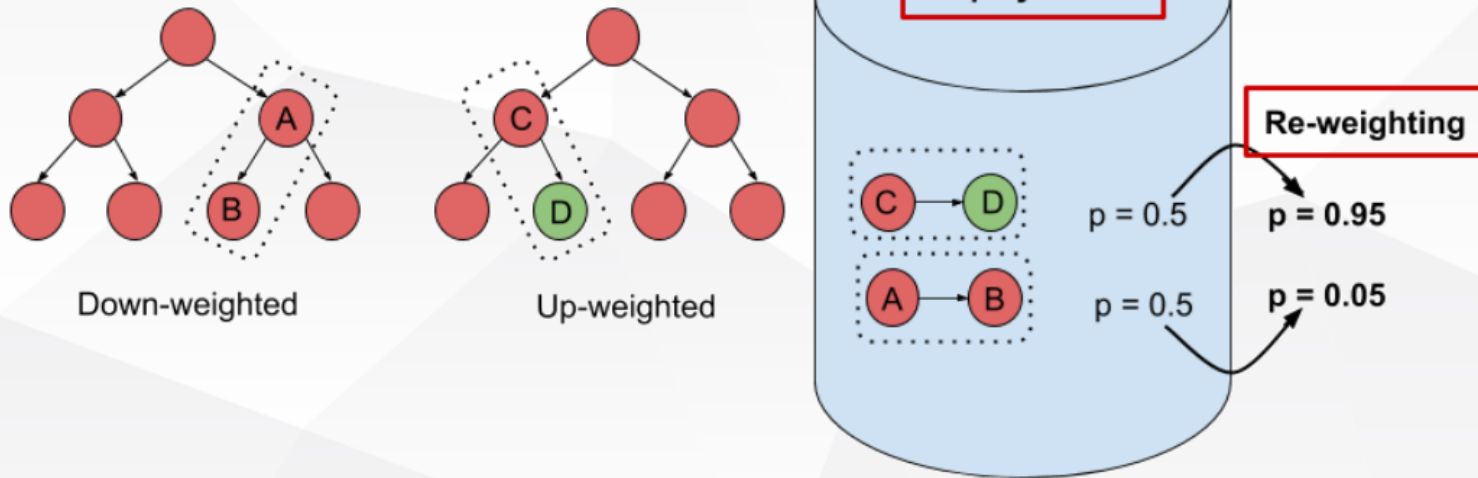
$$\nabla_Q \mathcal{L}_2(\delta(i)) = \frac{p_{\mathcal{D}_1}(i)}{p_{\mathcal{D}_2}(i)} \nabla_Q \mathcal{L}_1(\delta(i))$$

$$\mathbb{E}_{\mathcal{D}_1} [\nabla_Q \mathcal{L}_1(\delta(i))] = \mathbb{E}_{\mathcal{D}_2} [\nabla_Q \mathcal{L}_2(\delta(i))]$$

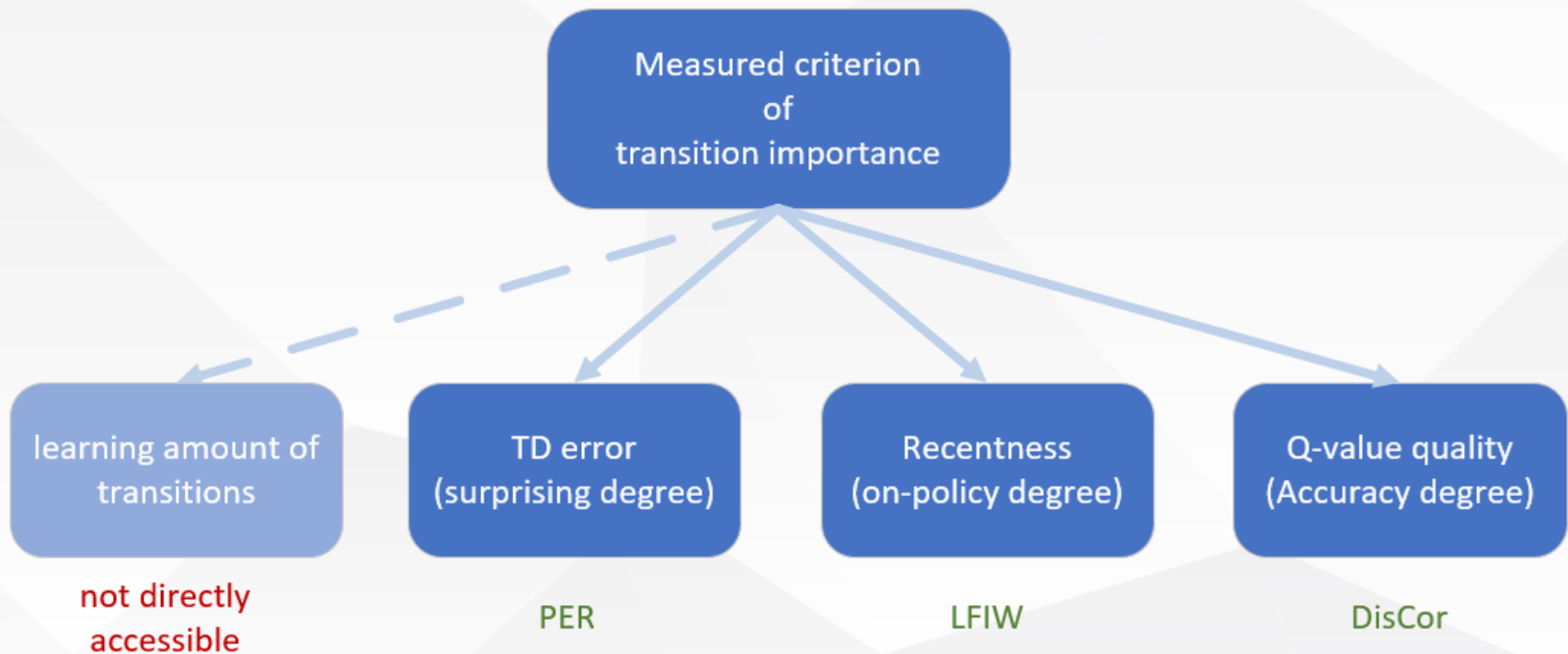
Non-uniform Experience Repaly

$$\theta \leftarrow \arg \min_{\theta} \mathbb{E}_{s,a \sim \mathcal{D}} \left[(Q_{\theta}(s,a) - (r(s,a) + \gamma \mathbb{E}_{s'|s,a} [\max_{a'} \bar{Q}(s',a')]))^2 \right]$$

$$Q_k \leftarrow \arg \min_Q \frac{1}{N} \sum_{i=1}^N w_i(s,a) \cdot (Q(s,a) - [r(s,a) + \gamma Q_{k-1}(s',a')])^2$$



Non-uniform Experience Repaly



PER

the priority of transition i

$$\left\{ \begin{aligned} p(i) &= \frac{|\delta(i)|^\alpha + \epsilon}{\sum_j (|\delta(j)|^\alpha + \epsilon)} \\ p(i) &= \frac{\left(\frac{1}{\text{rank}(i)}\right)^\alpha}{\sum_j \left(\frac{1}{\text{rank}(j)}\right)^\alpha} \end{aligned} \right.$$

between uniform sampling and greedy sampling

power-law distribution with exponent α (more robust)

$$\mathcal{L}_{\text{PER}} = w(i) \mathcal{L}(\delta(i))$$

$$\mathcal{L}_i = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a, \theta_i))^2]$$

$$w_i = \left(\frac{1}{N} \cdot \frac{1}{p(i)}\right)^\beta \leftarrow \text{annealing to 1}$$

$$w(i) = \frac{w_i}{\max_j w_j} \quad \text{for stability reasons}$$

LFIW

Bellman equation

$$Q^\pi(s, a) = \mathcal{B}^\pi Q^\pi(s, a)$$

↑
Bellman operator

$$\mathcal{B}^\pi Q(s, a) := r(s, a) + \gamma \mathbb{E}_{s', a'} [Q(s', a')]$$

loss for Q-network

$$L_Q(\theta; \mathcal{D}) = \mathbb{E}_{(s, a) \sim \mathcal{D}} \left[(Q_\theta(s, a) - \hat{\mathcal{B}}^\pi Q_\theta(s, a))^2 \right]$$

↑
replay buffer

↑
Tend to \mathcal{B}^π

introduce prioritization

$$L_Q(\theta; d, w) = \mathbb{E}_d \left[w(s, a) (Q_\theta(s, a) - \mathcal{B}^\pi Q_\theta(s, a))^2 \right]$$

↑
sampling distribution (buffer distribution)

objective

$$\arg \min_{\theta} L_Q(\theta; d, w) = \arg \min_{\theta} L_Q(\theta; d^w)$$

select a favorable priority distribution $d^w \propto d \cdot w$.

$$\boxed{d^w = d^\pi}$$

- Estimate the density ratio **only rely on samples** (e.g. from the replay buffer)

Lemma 1 ([27]). Assume that f has first order derivatives f' at $[0, +\infty)$. $\forall P, Q \in \mathcal{P}(\mathcal{X})$ such that $P \ll Q$ and $w : \mathcal{X} \rightarrow \mathbb{R}^+$,

$$D_f(P||Q) \geq \mathbb{E}_P[f'(w(\mathbf{x}))] - \mathbb{E}_Q[f^*(f'(w(\mathbf{x})))]$$
 (9)

where f^* denotes the convex conjugate and the equality is achieved when $w = dP/dQ$.

$$D_f(p||q) = \int q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

f-divergences

$$w(s, a) := d^\pi(s, a) / d^D(s, a)$$

- Two types of replay buffer
 - smaller (faster) replay buffer \longrightarrow smaller size, more on-policiness
 - regular (slow) replay buffer \longrightarrow bigger size, more off-policiness

LFIW

- Estimate the density ratio via minimizing the follow objective over **network** $w_\psi(x)$

$$L_w(\psi) := \mathbb{E}_{\mathcal{D}_s}[f^*(f'(w_\psi(s, a)))] - \mathbb{E}_{\mathcal{D}_f}[f'(w_\psi(s, a))]$$

the outputs $w_\psi(s, a)$ are **forced to be non-negative via activation functions**

- self normalization** with temperature hyperparameter T

$$\tilde{w}_\psi(s, a) := \frac{w_\psi(s, a)^{1/T}}{\mathbb{E}_{\mathcal{D}_s}[w_\psi(s, a)^{1/T}]}$$

- The final objective for TD learning over Q is then

$$L_Q(\theta; d^\pi) \approx L_Q(\theta; \mathcal{D}_s, \tilde{w}_\psi) := \mathbb{E}_{(s,a) \sim \mathcal{D}_s} [\tilde{w}_\psi(\mathbf{x}) (Q_\theta(s, a) - \hat{B}^\pi Q_\theta(s, a))^2]$$

estimate via MC

DisCor

$$\mathcal{L}(Q) = \mathbb{E}_{s \sim \beta(s), a \sim \pi_k(a|s)} [|Q_k(s, a) - Q^*(s, a)|].$$

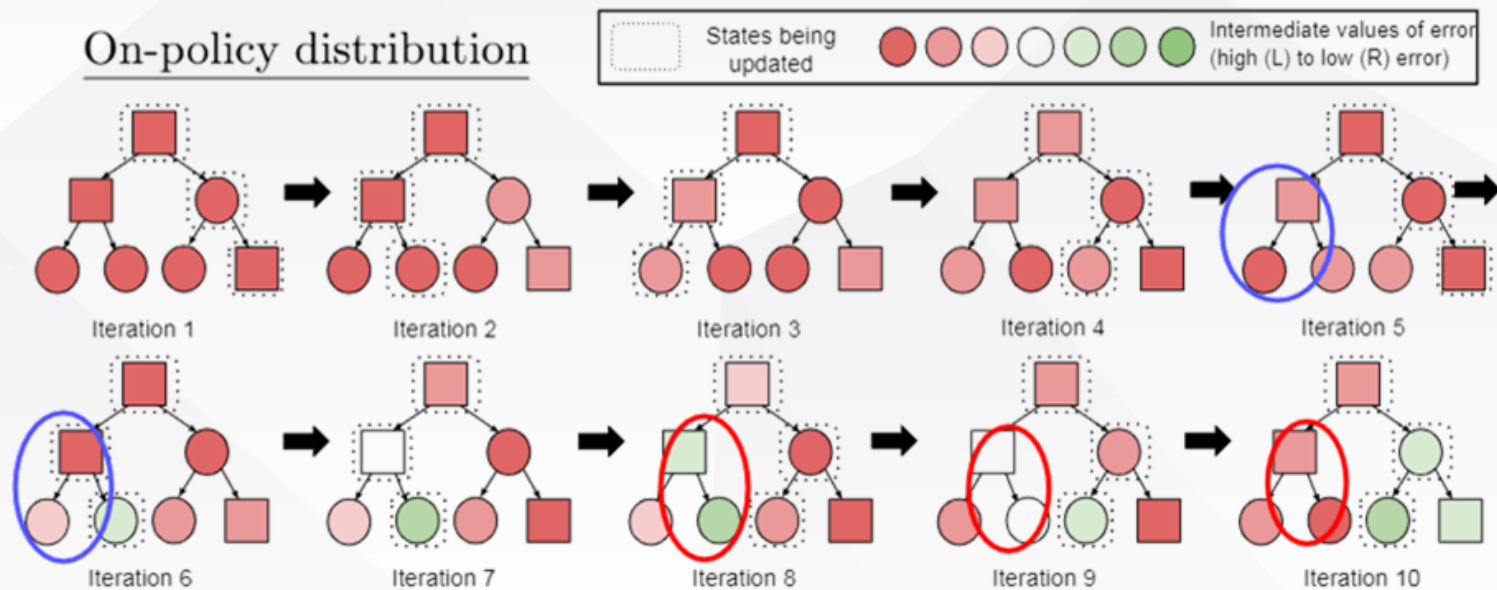
$$\mathcal{L}(Q) = \mathbb{E}_{s \sim \beta(s), a \sim \pi_k(a|s)} [|Q_k(s, a) - \mathcal{B}^* Q_k(s, a)|]$$

can be a **wrong target**

precise target

function approximator make things **worse**

DisCor



- leaf state: rarely visit, provide incorrect TD target
- root state: frequently visit, fit to incorrect target
- state with similar features affect each other

DisCor

$$\min_{p_k} \mathbb{E}_{d^{\pi_k}} [|Q_k - Q^*|]$$

$$\text{s.t. } Q_k = \arg \min_Q \mathbb{E}_{p_k} [(Q - \mathcal{B}^* Q_{k-1})^2], \quad \sum_{s,a} p_k(s,a) = 1, \quad \forall s,a \quad p_k(s,a) \geq 0$$



$$p_k(s,a) \propto \exp(-|Q_k - Q^*|(s,a)) \frac{|Q_k - \mathcal{B}^* Q_{k-1}|(s,a)}{\lambda^*}$$

$$w_k(s,a) = \frac{p_k(s,a)}{\mu(s,a)} \quad \text{densities } \mu(s,a) \text{ are unknown}$$



$$q_k^* = \arg \min_{q_k} -\mathbb{E}_{q_k} [\log p_k] + (\tau) D_{\text{KL}}(q_k \| \mu)$$



$$q_k^*(s,a) \propto (\mu_k) \cdot \exp\left(\frac{\log p_k(s,a)}{\tau}\right)$$

$$\therefore \frac{q_k^*}{\mu_k} \propto \exp\left(\frac{-|Q_k - Q^*|(s,a)}{\tau}\right) \frac{|Q_k - \mathcal{B}^* Q_{k-1}|(s,a)}{\lambda^*}$$

DisCor

$$\frac{q_k^*}{\mu_k} \propto \exp\left(\frac{-|Q_k - Q^*|(s, a)}{\tau}\right) \frac{|Q_k - \mathcal{B}^* Q_{k-1}|(s, a)}{\lambda^*}$$

$$\Delta_k(s, a) + \sum_{i=1}^k \gamma^{k-i} \alpha_i \geq |Q_k - Q^*|(s, a),$$

$$\Delta_k(s, a) = |Q_k(s, a) - (\mathcal{B}^* Q_{k-1})(s, a)| + \gamma (P^{\pi_{k-1}} \Delta_{k-1})(s, a).$$

$$\forall s, a \quad c_1 \leq |Q_k - \mathcal{B}^* Q_{k-1}|(s, a) \leq c_2$$

where

$$c_1 = \min_{s, a} |Q_{k-1} - \mathcal{B}^* Q_{k-2}|,$$

$$c_2 = \max_{s, a} |Q_{k-1} - \mathcal{B}^* Q_{k-2}|$$

↓ lower bound

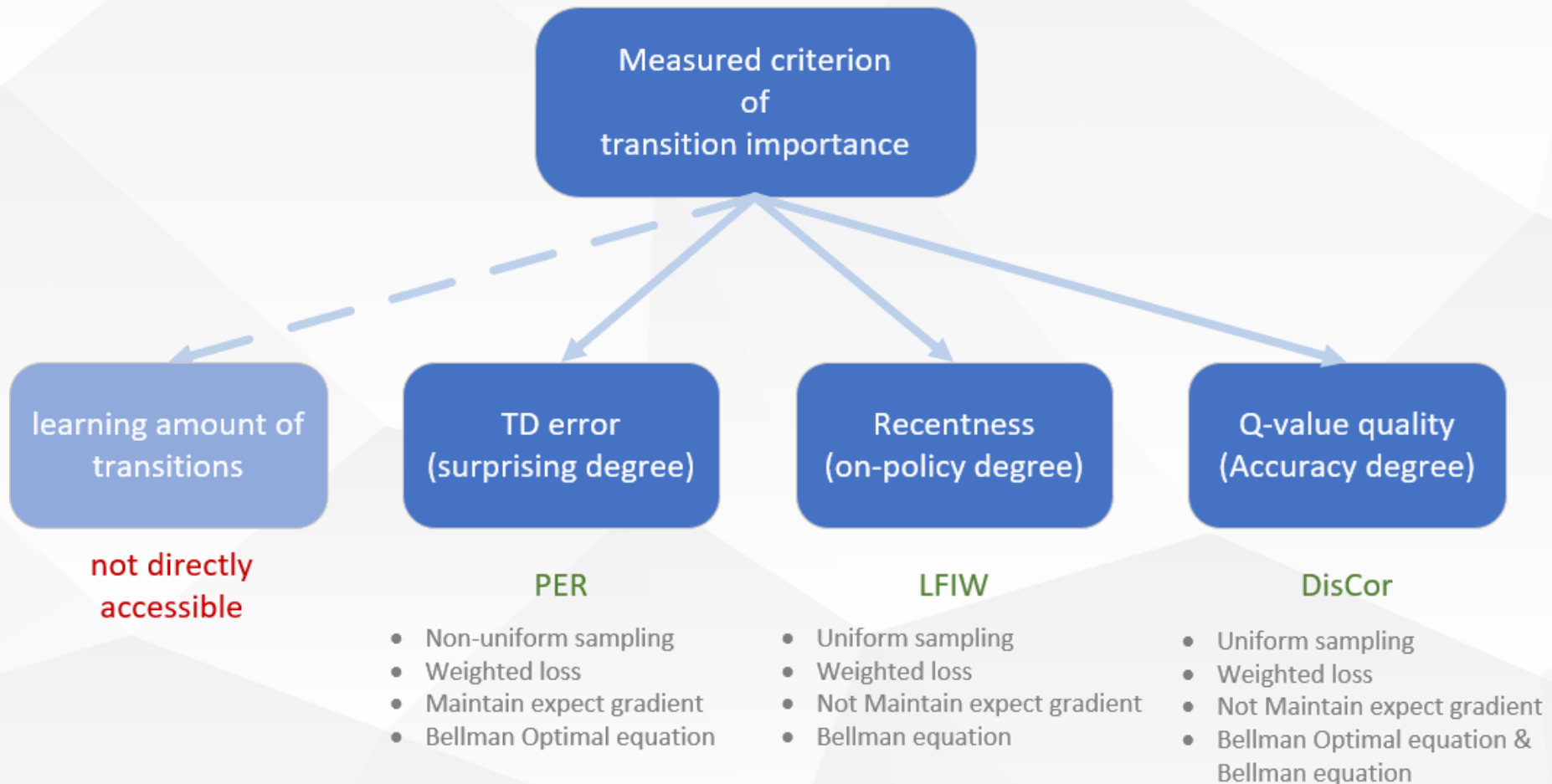
$$w_k \propto \exp\left(\frac{-c_2 - \gamma [P^{\pi_{k-1}} \Delta_{k-1}](s, a)}{\tau}\right) \frac{c_1}{\lambda^*}$$

↓

$$w_k(s, a) \propto \exp\left(-\frac{\gamma [P^{\pi_{k-1}} \Delta_{k-1}](s, a)}{\tau}\right).$$

$$P^\pi Q(s, a) := \mathbb{E}_{s' \sim P(\cdot|s, a), a' \sim \pi(a|s)} [Q(s', a')]$$

Non-uniform Experience Replay



Regret Minimization Experience Replay in Off-Policy Reinforcement Learning

Xu-Hui Liu*, **Zhenghai Xue***, **Jing-Cheng Pang**, **Shengyi Jiang**, **Feng Xu**, **Yang Yu[†]**

National Key Laboratory for Novel Software Technology

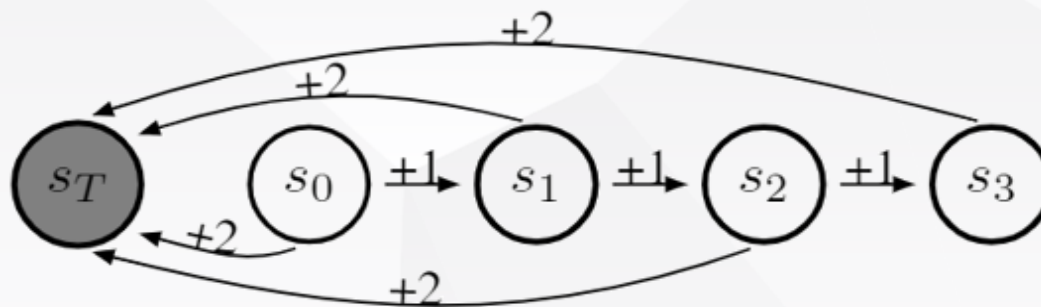
Nanjing University, Nanjing 210023, China

liuxh@lamda.nju.edu.cn, xuezh@smail.nju.edu.cn

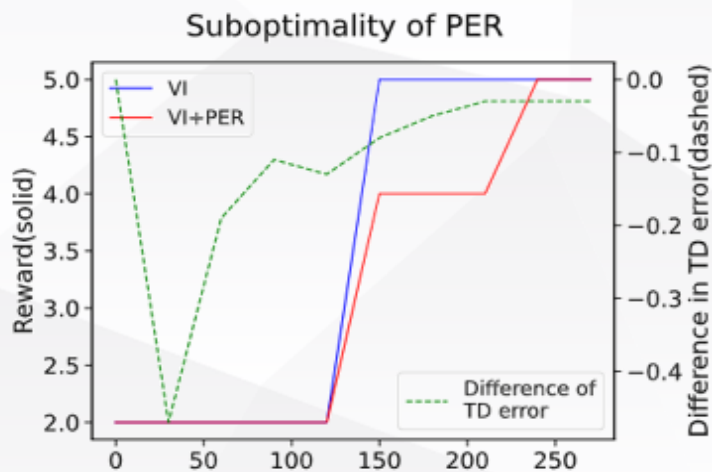
{pangjc, jiangsy, xufeng}@lamda.nju.edu.cn, yuy@nju.edu.cn

NIPS 2021

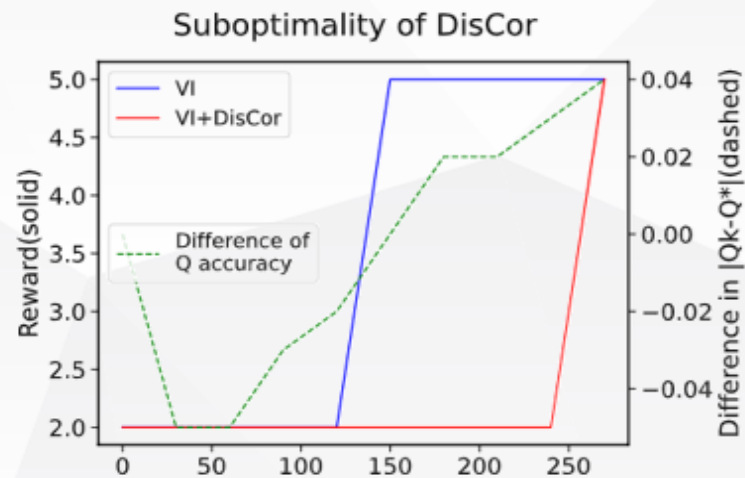
Bad case



(a)



(b)



(c)

Preliminaries

- Expect return of policy

$$\eta(\pi) = \mathbb{E}_\pi[\sum_{t \geq 0} \gamma^t r(s_t, a_t)], \text{ 其中 } s_0 \sim \rho_0, a_t \sim \pi(\cdot | s_t), s_{t+1} \sim T(\cdot | s_t, a_t)$$

$$\downarrow \lim_{t \rightarrow \infty} (\gamma^0 + \gamma^1 + \gamma^2 + \dots + \gamma^t) r(s, a) = \lim_{t \rightarrow \infty} \frac{\gamma^0 (\gamma^{t+1} - 1)}{\gamma - 1} r(s, a) = \frac{1}{1 - \gamma} r(s, a)$$
$$\eta(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{d^\pi(s, a)}[r(s, a)]$$

- Optimal policy

$$\pi^* = \arg \max_\pi \eta(\pi)$$

- Regret: the expected loss in return by following policy π instead of the optimal policy

$$\text{Regret}(\pi) = \eta(\pi^*) - \eta(\pi)$$

- Boltzmann exploration policy

$$\pi_k(s) = \frac{\exp(Q_k(s, a))}{\sum_{a'} \exp(Q(s, a'))}$$

Formalize the problem

$$\begin{aligned} \min_{w_k} \quad & \eta(\pi^*) - \eta(\pi_k) \\ \text{s.t.} \quad & Q_k = \arg \min_{Q \in \mathcal{Q}} \mathbb{E}_{\mu} [w_k(s, a) \cdot (Q - \mathcal{B}^* Q_{k-1})^2(s, a)], \\ & \mathbb{E}_{\mu} [w_k(s, a)] = 1, \quad w_k(s, a) \geq 0, \end{aligned}$$

data distribution of the replay buffer

$$\begin{aligned} \min_{p_k} \quad & \mathbb{E}_{d^{\pi_k}} [|Q_k - Q^*|] \\ \text{s.t.} \quad & Q_k = \arg \min_Q \mathbb{E}_{p_k} [(Q - \mathcal{B}^* Q_{k-1})^2], \quad \sum_{s,a} p_k(s, a) = 1, \quad \forall s, a \quad p_k(s, a) \geq 0 \end{aligned}$$

Analyze computationally

Theorem 1 (Informal). *Under mild conditions, the solution w_k to a relaxation of the optimization problem 1 in MDPs with discrete action spaces is*

$$w_k(s, a) = \frac{1}{Z_1^*} (E_k(s, a) + \epsilon_{k,1}(s, a)). \quad (2)$$

In MDPs with continuous action spaces, the solution is

$$w_k(s, a) = \frac{1}{Z_2^*} (F_k(s, a) + \epsilon_{k,2}(s, a)). \quad (3)$$

where

$$E_k(s, a) = \underbrace{\frac{d^{\pi_k}(s, a)}{\mu(s, a)}}_{(a)} \underbrace{(2 - \pi_k(a|s))}_{(b)} \underbrace{\exp(-|Q_k - Q^*|(s, a))}_{(c)} \underbrace{|Q_k - \mathcal{B}^*Q_{k-1}|(s, a)}_{(d)}$$

$$F_k(s, a) = 2 \underbrace{\frac{d^{\pi_k}(s, a)}{\mu(s, a)}}_{(a)} \underbrace{\exp(-|Q_k - Q^*|(s, a))}_{(c)} \underbrace{|Q_k - \mathcal{B}^*Q_{k-1}|(s, a)}_{(d)}$$

Z_1^*, Z_2^* are normalization factors, $\epsilon_{k,1}(s, a)$ and $\epsilon_{k,2}(s, a)$ satisfy $\max \left\{ \frac{\epsilon_{k,1}(s, a)}{E_k(s, a)}, \frac{\epsilon_{k,2}(s, a)}{F_k(s, a)} \right\} \leq \epsilon_{\pi_k}$. CSDN @云端FFF

DisCor

$$\frac{q_k^*}{\mu_k} \propto \exp \left(\frac{-|Q_k - Q^*|(s, a)}{\tau} \right) \frac{|Q_k - \mathcal{B}^*Q_{k-1}|(s, a)}{\lambda^*}$$

ReMERN

$$w_k(s, a) \propto \frac{d^{\pi_k}(s, a)}{\mu(s, a)} \exp(-\gamma [P^{\pi_{k-1}} \Delta_{k-1}](s, a))$$

Algorithm 1 ReMERN

1: Initialize Q-values $Q_\theta(s, a)$, a replay buffer μ , an **error model** $\Delta_\phi(s, a)$, and a **weight model** κ_ψ .

2: **for** step k in $\{1, \dots, N\}$ **do**

3: Collect M samples using π_k , add them to replay buffer μ , sample $\{(s_i, a_i)\}_{i=1}^N \sim \mu$.

4: Evaluate $Q_\theta(s, a)$, $\Delta_\phi(s, a)$ and $\kappa_\psi(s, a)$ on samples (s_i, a_i) .

5: Compute target values for Q and Δ on samples:

$$y_i = r_i + \gamma \max_{a'} Q_{k-1}(s'_i, a').$$

$$\hat{a}_i = \arg \max_a Q_{k-1}(s'_i, a).$$

$$\hat{\Delta} = |Q_\theta(s, a) - y_i| + \gamma \Delta_{k-1}(s'_i, \hat{a}_i).$$

6: **Optimize** κ_ψ using

$$L_\kappa(\psi) := \mathbb{E}_{\mathcal{D}_s} [f^*(f'(\kappa_\psi(s, a)))] - \mathbb{E}_{\mathcal{D}_f} [f'(\kappa_\psi(s, a))].$$

7: **Compute** w_k using

$$w_k(s, a) \propto \frac{d^{\pi_k}(s, a)}{\mu(s, a)} \exp(-\gamma [P^{\pi_{k-1}} \Delta_{k-1}](s, a)).$$

8: Minimize Bellman error for Q_θ weighted by w_k .

$$\theta_{k+1} \leftarrow \operatorname{argmin}_\theta \frac{1}{N} \sum_i w_k(s_i, a_i) (Q_\theta(s_i, a_i) - y_i)^2.$$

9: **Minimize ADP error** for training ϕ .

$$\phi_{k+1} \leftarrow \operatorname{argmin}_\phi \frac{1}{N} \sum_{i=1}^N (\Delta_\phi(s_i, a_i) - \hat{\Delta}_i)^2.$$

10: **end for**

ReMERT

- Estimate the upper bound of $|Q_k - Q^*|$ by **temporal Structure**

$$|Q_k - Q^*| \leq |Q_k - \mathcal{B}^* Q_{k-1}| + |\mathcal{B}^* Q_{k-1} - Q^*|$$

projection error error accumulate through the MDP

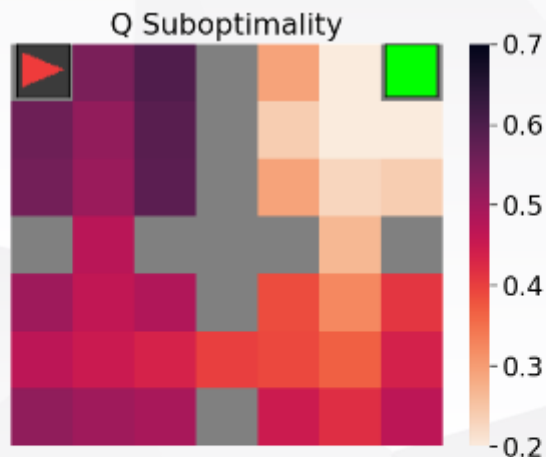


Figure 2: The visualized error of target Q value in a GridWorld Environment. The Q error is visualized by the color of the grid.

ReMERT

Definition 2 (Distance to End) : 给定 MDP \mathcal{M} , 策略 π 在 \mathcal{M} 中交互得到轨迹 $\tau = \{s_t, a_t\}_{t=0}^T$, 定义该轨迹中 (s_t, a_t) 的“到终点的距离”为

$$h_\tau(s_t, a_t) = T - t$$

作者的直觉表明 $|Q_k - Q^*|$ 的值和“到终点的距离”有关

Theorem 2 (Informal). *Under mild conditions, with probability at least $1 - \delta$, we have*

$$\begin{aligned} & |Q_k(s, a) - Q^*(s, a)| \\ & \leq \mathbb{E}_\tau \left(f(h_\tau^{\pi k}(s, a))(L_{Q_{k-1}} + c) + \gamma^{h_\tau^{\pi k}(s, a)+1} c \right) + g(k, \delta) \end{aligned} \quad (8)$$

where $c = \max_{s, a} (Q^*(s, a^*) - Q^*(s, a))$, $f(t) = \frac{\gamma - \gamma^t}{1 - \gamma}$, $L_{Q_{k-1}} = \mathbb{E}[|Q_{k-1} - \mathcal{B}^* Q_{k-2}|]$ and $g(k, \delta)$ decreases exponentially as k increases.

CSDN @云端FFF

$|Q_k - Q^*|$ increases with $h_\tau^{\pi k}$, and the remaining terms are negligible as k increases.

$$\begin{aligned} |Q_k(s, a) - Q^*(s, a)| & \approx \mathbb{E}_\tau \text{TCE}_c(s, a) \\ & = \mathbb{E}_\tau \left(f(h_\tau^{\pi k-1}(s, a))(L_{Q_{k-1}} + c) + \gamma^{h_\tau^{\pi k-1}(s, a)+1} c \right) \end{aligned}$$

ReMERT

$$w_k(s, a) \propto \frac{d^{\pi_k}(s, a)}{\mu(s, a)} \exp(-\mathbb{E}_\tau \text{TCE}_c(s, a))$$

Algorithm 2 ReMERT

- 1: Initialize Q-values $Q_\theta(s, a)$, a replay buffer μ , and a **weight model** κ_ψ .
- 2: **for** step k in $\{1, \dots, N\}$ **do**
- 3: Collect M samples using π_k , add them to replay buffer μ , sample $\{(s_i, a_i)\}_{i=1}^N \sim \mu$.
- 4: Evaluate $Q_\theta(s, a)$ and $\kappa_\psi(s, a)$ on samples (s_i, a_i) .
- 5: Compute target values for Q on samples:
 $y_i = r_i + \gamma \max_{a'} Q_{k-1}(s'_i, a')$.
 $\hat{a}_i = \arg \max_a Q_{k-1}(s'_i, a)$.
- 6: **Optimize** κ_ψ **using**

$$L_\kappa(\psi) := \mathbb{E}_{\mathcal{D}_s} [f^*(f'(\kappa_\psi(s, a)))] - \mathbb{E}_{\mathcal{D}_f} [f'(\kappa_\psi(s, a))].$$

- 7: **Compute** w_k **using**

$$w_k(s, a) \propto \frac{d^{\pi_k}(s, a)}{\mu(s, a)} \exp\left(-\mathbb{E}_{q_{k-1}(\tau)} \text{TCE}_c(s, a)\right).$$

- 8: Minimize Bellman error for Q_θ weighted by w_k .
 $\theta_{k+1} \leftarrow \underset{\theta}{\operatorname{argmin}} \frac{1}{N} \sum_i w_k(s_i, a_i) (Q_\theta(s_i, a_i) - y_i)^2$.

- 9: **end for**

Experiments

Deterministic env (MoJoCo & DMC)

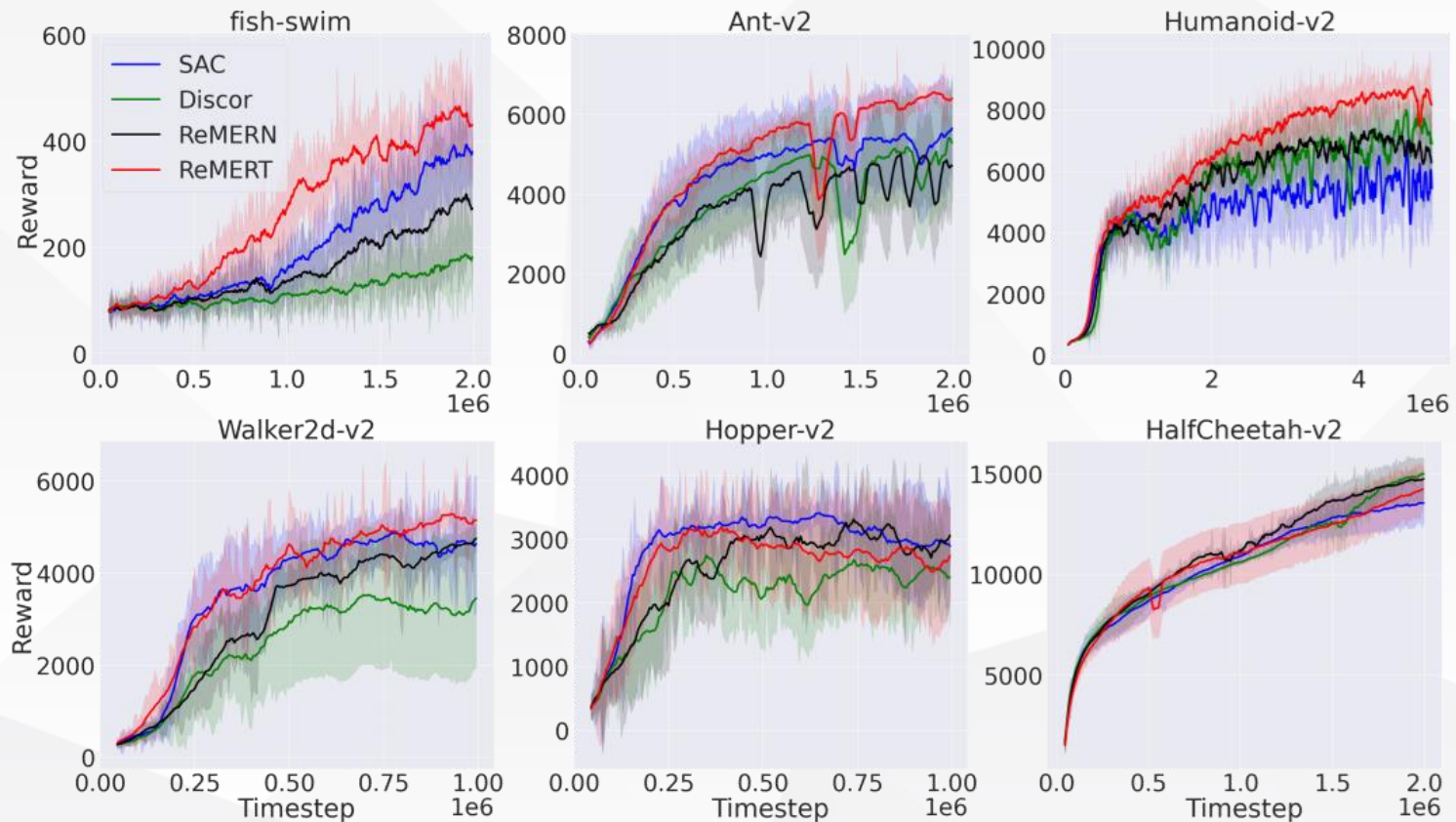
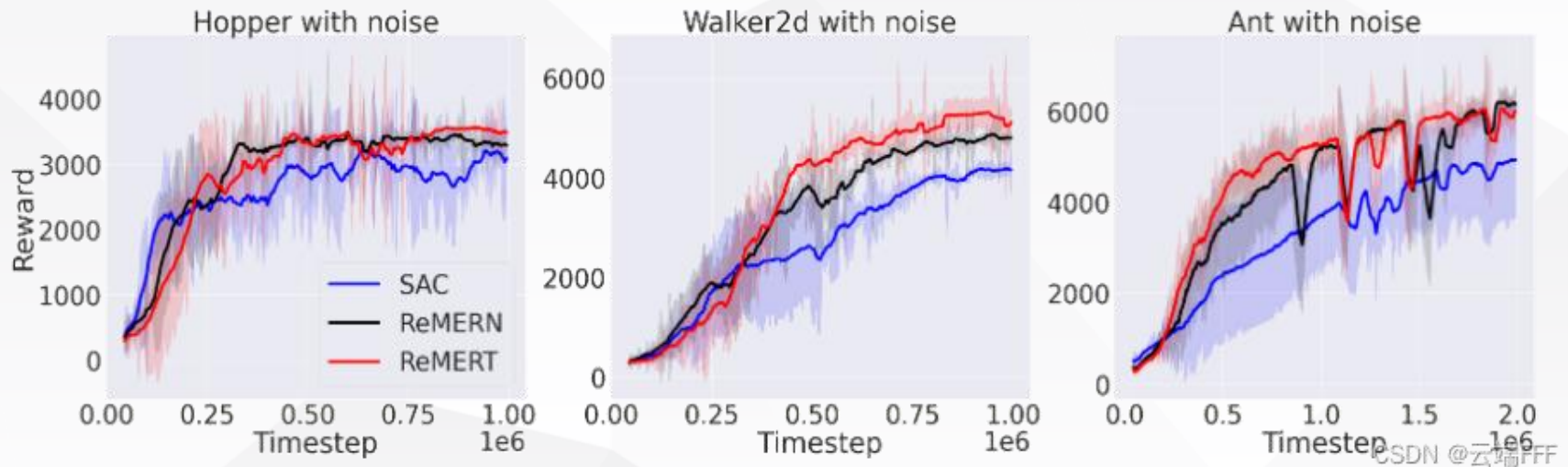


Figure 3: Performance of ReMERT, ReMERN with SAC and DisCor as baselines on continuous control tasks.

Experiments

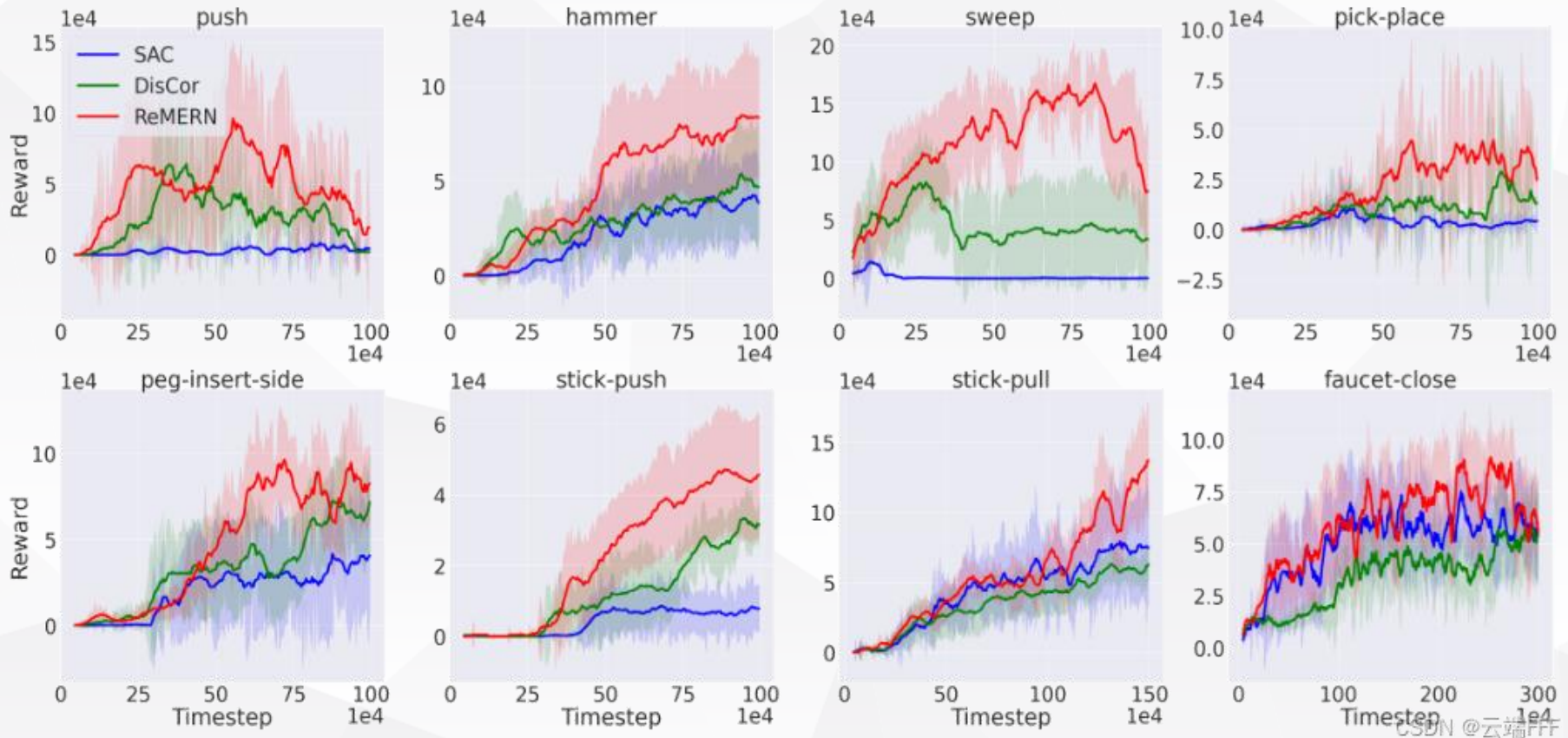
MoJoCo with noisy reward



CSDN @云瑞FF

Experiments

Meta-World with high randomness



CSDN @云端FF



— • **The End** • —

thanks