



南京航空航天大学

Nanjing University of Aeronautics and Astronautics



模式分析与机器智能
工业和信息化部重点实验室

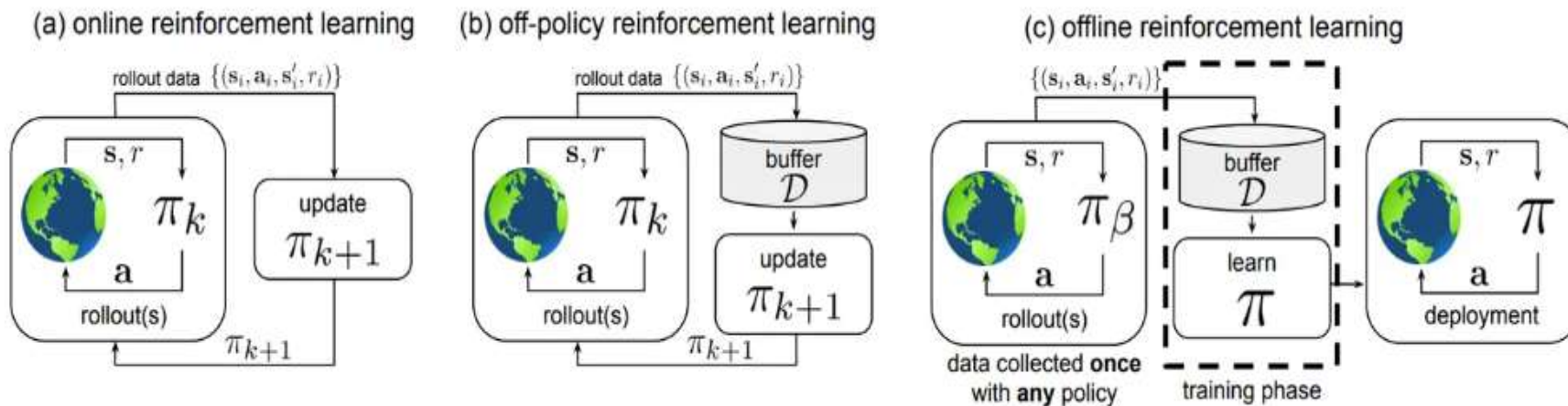
MIT Key Laboratory of
Pattern Analysis & Machine Intelligence

Beyond Uniform Sampling: Offline Reinforcement Learning with Imbalanced Datasets

NIPS | 2023

Background

- Reinforcement learning



- Offline Reinforcement learning

$$J(\pi) = \mathbb{E}_{\tau \sim p_{\pi}(\tau)} \left[\sum_{t=0}^H \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] \xrightarrow{\text{constraint}} \max_{\pi_{\theta}} J(\pi_{\theta}) - \alpha \mathbb{E}_{(s,a) \sim \mathcal{D}} [\mathcal{C}(s, a)],$$

Motivation



- **Problem Statement**

多数offlineRL模型趋向于选择数据集中出现的动作，当数据集本身不平衡或者次优数据较多，结果会很差。

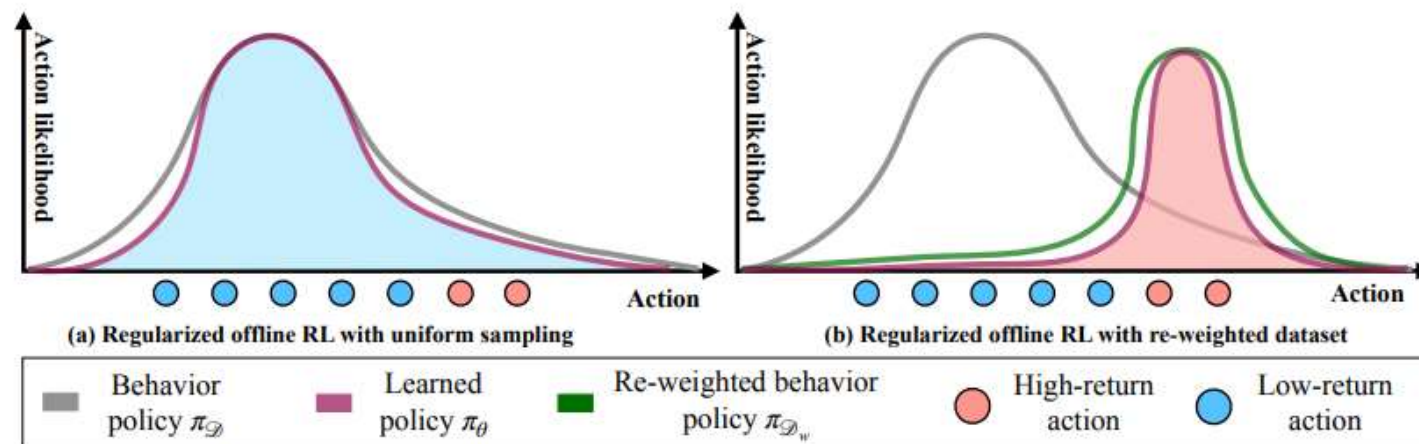


Figure 1: The dots represent actions in the dataset, where imbalanced datasets have more low-return actions. (a) Regularized offline RL algorithms [22, 7, 18] equally regularize the policy π_{θ} on each action, leading to imitation of low-return actions and a low-performing π_{θ} . The color under the curves shows the policy's performance $J(\pi_{\theta})$, with red indicating higher performance and blue indicating lower performance. (b) Re-weighting the dataset based on actions' returns allows the algorithm to only regularize on actions with high returns, enabling the policy π_{θ} to imitate high-return actions while ignoring low-return actions.

Background

Definition 3.1 (Dataset imbalance). RPSV of a dataset, $\mathbb{V}_+[G(\tau_i)]$, corresponds to the second-order moment of the positive component of the difference between trajectory return: $G(\tau_i) := \sum_{t=0}^{T_i-1} \gamma^t r(s_t^i, a_t^i)$ and its expectation, where τ_i denote trajectory in the dataset:

$$\mathbb{V}_+[G(\tau_i)] \doteq \mathbb{E}_{\tau_i \sim \mathcal{D}} \left[(G(\tau_i) - \mathbb{E}_{\tau_i \sim \mathcal{D}}[G(\tau_i)])_+^2 \right] \quad \text{with} \quad x_+ = \max\{x, 0\}, \quad (3)$$

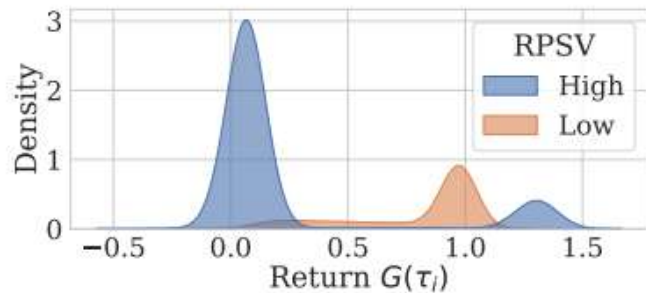


Figure 2: Return distribution of datasets with high and low RPSV. Low RPSV datasets have returns centered at the mean, while high RPSV datasets have a wider distribution extending towards higher returns. See Appendix A.4 for details.

- unnecessarily conservative

Imbalance dataset + constraint

Methodology



- Mitigating Unnecessary Conservativeness By Weighting Samples

Adding more experiences from high-performing policies

Importance sampling

$$w(s, a) = \frac{\mathcal{D}_w(s, a)}{\mathcal{D}(s, a)}, \quad (\mathcal{D}_w \text{ unknown})$$

$$\begin{aligned} \max_{\pi} \hat{J}_{\mathcal{D}}^{\gamma}(\pi) - \alpha \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} [w(s, a) \mathcal{C}(s_t, a_t)] \\ \Updownarrow \\ \max_{\pi} \hat{J}_{\mathcal{D}}^{\gamma}(\pi) - \alpha \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}_w} [\mathcal{C}(s_t, a_t)]. \end{aligned}$$

Difficulty:

The key challenge is determining the weights w since the state-action distribution of the better policy $\pi \mathcal{D}_w$ is initially unknown

off-policy evaluation techniques (optidice)

to determine if a weighting function corresponds to a high-return policy

Methodology



- OptiDICE

provides a framework for improving policy optimization using offline datasets in reinforcement learning

$$\max_{d \geq 0} \min_{\nu} \mathbb{E}_{(s,a) \sim d} [R(s,a)] - \alpha D_f(d || d^D) \quad (5)$$

$$\begin{aligned} & + \sum_s \nu(s) \left((1 - \gamma) p_0(s) + \gamma (\mathcal{T}_* d)(s) - (\mathcal{B}_* d)(s) \right), \\ & = (1 - \gamma) \mathbb{E}_{s \sim p_0} [\nu(s)] + \mathbb{E}_{(s,a) \sim d^D} [-\alpha f(w(s,a))] \\ & + \mathbb{E}_{(s,a) \sim d^D} [w(s,a) (e_\nu(s,a))] =: L(w, \nu). \end{aligned} \quad (7)$$

$$D_f(d^\pi || d^D) := \mathbb{E}_{(s,a) \sim d^D} \left[f\left(\frac{d^\pi(s,a)}{d^D(s,a)}\right) \right] \quad F = x \log x \text{ 时为 KL 散度}$$

$$(\mathcal{B}_* d)(s) := \sum_{\bar{a}} d(s, \bar{a})$$

$$(\mathcal{T}_* d)(s) := \sum_{\bar{s}, \bar{a}} T(s | \bar{s}, \bar{a}) d(\bar{s}, \bar{a})$$

DICE简单来说可以理解
为初始状态转移分布和
状态转移分布，在结合
重要度采样所构造的优
化目标

Methodology: Density ratio Weighting (DW)



Our goal is to discover a weighting function w that can emulate drawing state-action samples from a better dataset \mathcal{D}_w that is collected by an alternative behavior policy $\pi_{\mathcal{D}_w}$ with higher return than the behavior policy $\pi_{\mathcal{D}}$ that collected the original dataset \mathcal{D}

- **Optimizing the Weightings**

$$J^\gamma(\pi_{\mathcal{D}_w}) \approx \mathbb{E}_{(s,a) \sim \mathcal{D}_w} [r(s, a)] = \mathbb{E}_{(s,a) \sim \mathcal{D}} [w(s, a) r(s, a)].$$

when the the dataset \mathcal{D}_w represents a stationary state-action distribution:

$$\mathcal{D}_w(s') = (1 - \gamma)\rho_0(s') + \gamma \sum_{s,a} \mathcal{T}(s'|s, a) \mathcal{D}_w(s, a) \quad \forall s' \in \mathcal{S}, \quad \mathcal{D}_w(s') := \sum_{a' \in \mathcal{A}} \mathcal{D}_w(s', a') \quad (7)$$



$$\mathcal{D}_w(s') = \sum_{s,a} \mathcal{T}(s'|s, a) \mathcal{D}_w(s, a) \quad \forall s' \in \mathcal{S}.$$

Methodology: Density ratio Weighting (DW)



- 将数据集 \mathcal{D}_w 写成 $w(s,a)$ 的形式，则有下列的约束成立：

$$\mathcal{D}(s')w(s') = \sum_{s,a} \mathcal{T}(s'|s,a)w(s,a) \quad \forall s' \in \mathcal{S}, \quad w(s) := \sum_{a \in \mathcal{A}} \frac{\mathcal{D}_w(s,a)}{\mathcal{D}(s,a)} \quad (9)$$

$$J^\gamma(\pi_{\mathcal{D}_w}) \approx \mathbb{E}_{(s,a) \sim \mathcal{D}_w} [r(s,a)] = \mathbb{E}_{(s,a) \sim \mathcal{D}} [w(s,a)r(s,a)]. \quad (6)$$

将（6）式的无折扣版本和（9）式的约束合并，则有以下带约束的优化目标：

$$\begin{aligned} \max_w J(\pi_{\mathcal{D}_w}) &= \mathbb{E}_{(s,a) \sim \mathcal{D}} [w(s,a)r(s,a)] \\ \text{subject to } &\mathbb{E}_{(s,a,s') \sim \mathcal{D}} [w(s') - w(s,a) \mid s'] = 0 \quad \forall s' \in \mathcal{S}. \end{aligned} \quad (10)$$

Implementation



进行参数化，并给出状态边际分布的计算方式

$$w(s, a) \doteq \frac{\mathcal{D}_w(s, a)}{\mathcal{D}(s, a)} = \frac{\mathcal{D}_w(s) \pi_{\mathcal{D}_w}(a|s)}{\mathcal{D}(s) \pi_{\mathcal{D}}(a|s)} = \frac{\mathcal{D}_w(s)}{\mathcal{D}(s)} \times \frac{\pi_{\mathcal{D}_w}(a|s)}{\pi_{\mathcal{D}}(a|s)}.$$

$$w_{\phi, \psi}(s, a) = \exp \phi(s) \exp \psi(s, a), \quad w_{\phi}(s) = \exp \phi(s) \quad (12)$$

ϕ and ψ are neural networks

(two-layer Multilayer Perceptron (MLP) with 256 neurons and ReLU activation in each layer.)

$$\max_{\phi, \psi} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[\underbrace{w_{\phi, \psi}(s, a) r(s, a)}_{\text{Return}} - \lambda_F \underbrace{(w_{\phi}(s') - w_{\phi, \psi}(s, a))^2}_{\text{Bellman flow conservation penalty}} \right] - \lambda_K \underbrace{D_{KL}(\mathcal{D}_w || \mathcal{D})}_{\text{KL regularization}}, \quad (13)$$

↓
最终优化目标

↓
加快收敛 提高稳定性

↓
数据集覆盖范围有限

Methodology: DW & IQL



Algorithm 1 Implicit Q-learning

Initialize parameters $\psi, \theta, \hat{\theta}, \phi$.

TD learning (IQL):

for each gradient step **do**

$$\psi \leftarrow \psi - \lambda_V \nabla_{\psi} L_V(\psi) \rightarrow L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^{\tau}(Q_{\hat{\theta}}(s,a) - V_{\psi}(s))].$$

$$\theta \leftarrow \theta - \lambda_Q \nabla_{\theta} L_Q(\theta) \rightarrow L_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} [(r(s,a) + \gamma V_{\psi}(s') - Q_{\theta}(s,a))^2]$$

$$\hat{\theta} \leftarrow (1 - \alpha)\hat{\theta} + \alpha\theta$$

end for

Policy extraction (AWR):

for each gradient step **do**

$$\phi \leftarrow \phi - \lambda_{\pi} \nabla_{\phi} \bar{L}_{\pi}(\phi) \rightarrow L_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta(Q_{\hat{\theta}}(s,a) - V_{\psi}(s))) \log \pi_{\phi}(a|s)].$$

end for

Methodology: DW & IQL



IQL. We reweight state-value function (V), state-action value function (i.e., Q-function, Q), and policy (π) in IQL as follows:

$$\min_V \mathbb{E}_{(s,a) \sim \mathcal{D}} [w_{\phi}(s) L_2^{\tau}(Q(s, a) - V(s))] \quad (20)$$

$$\min_Q \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[w_{\phi,\psi}(s, a) (r(s, a) + \gamma V(s') - V(s))^2 \right] \quad (21)$$

$$\max_{\pi} \mathbb{E}_{(s,a) \sim \mathcal{D}} [w_{\phi,\psi}(s, a) \exp(\beta (Q(s, a) - V(s))) \log \pi(a|s)], \quad (22)$$

where L_2^{τ} denotes the upper expectile loss [18] and β denotes the temperature parameter for IQL. Our implementation is adapted from the official implementation⁴ for implicit Q-learning (IQL) [18].

Algorithm



Algorithm 1 Density-ratio weighting with generic offline RL algorithms (details in Appendix A.3)

- 1: **Input:** Dataset \mathcal{D}
 - 2: Initialize policy π and weighting function $w_{\phi, \psi}$
 - 3: **while** not converged **do**
 - 4: Sample a batch \mathcal{B} of tuples of states, actions, rewards, and next states (s, a, r, s') from \mathcal{D}
 - 5: Update $w_{\phi, \psi}$ with batch \mathcal{B} using Equation 13
 - 6: Update policy π and value function with an offline RL training objective with weights $w_{\phi, \psi}(s, a)$ and \mathcal{B} (e.g., [22, 7, 18])
 - 7: **end while**
-

Summary



目标： 本文希望优化目标不是聚焦于整个数据集，而是仅聚焦于数据集中好的部分
实施：

1. 通过重要度采样将原数据集 D 和未知数据集 D_w 联系起来
2. 通过off-policy evaluation方法，将数据集分布 D_w 和策略 π_w 的累积回报联系起来
3. 通过对 w 进行参数化，并加入KL散度，最终形成优化目标。
4. 与现有离线强化学习算法结合。

Experiments



- **Two stages**

- (i) Trajectories with similar initial states
- (ii) Trajectories with diverse initial states

- **interquartile mean (IQM)**

discards the bottom and top 25% of the runs and calculates the mean score of the remaining 50% runs (=N M/2c for N runs each on M tasks)

表示统计资料中各变量分散情形，但四分差更多为一种稳健统计

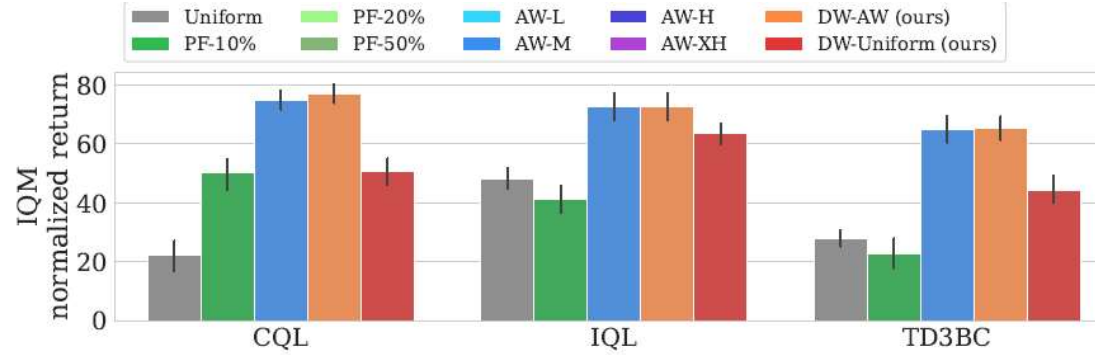
- **Baseline**

advantage-weighting (AW) (provide a better initial sampling distribution to train the weightingfunction in DW)
percentage-filtering (PF)

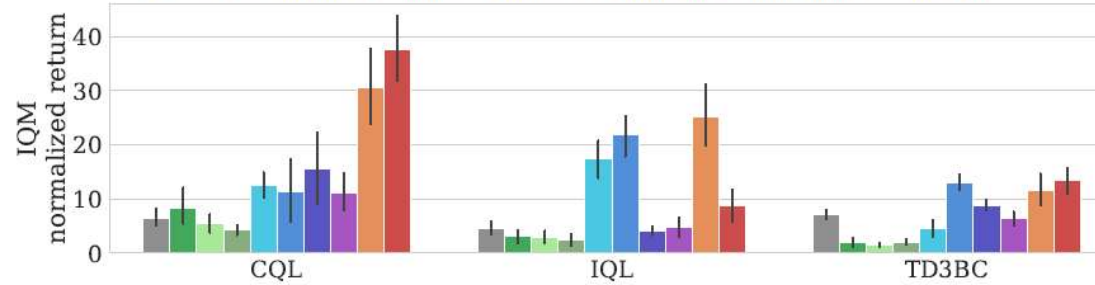
$$\mathcal{P}_{AW}(s_t^i, a_t^i) \propto \exp((G(\tau_i) - V_0(s_0^i))/\eta) \quad (\text{Advantage-weighting}) \quad (14)$$

$$\mathcal{P}_{PF}(s_t^i, a_t^i) \propto \mathbb{1}[G(\tau_i) \geq G_{K\%}] \quad (\text{Percentage-filtering}), \quad (15)$$

Experiments



(a) Results on imbalanced datasets of trajectories with similar initial states (Section 5.1).



(b) Results on smaller version of datasets used in Figure 3a.

Figure 3: (a) Our methods, DW-AW and DW-Uniform, achieve higher return than Uniform, indicating that DW can enhance the performance of offline RL algorithms on imbalanced datasets. Note that our methods in IQL, although not surpassing AW and PF-10% in performance, ours can be applied to offline RL dataset that are not curated with trajectories. (b) Our methods outperform Uniform in CQL, IQL, and TD3BC, indicating no significant overfitting in smaller datasets. DW-AW demonstrates superior returns compared to AW and PF, particularly in CQL, indicating our method effectively leverages limited data. IQL shows limited gains likely due to its difficulties in utilizing data from the rest of low-return trajectories in the dataset (see Section 5.2).

Experiments

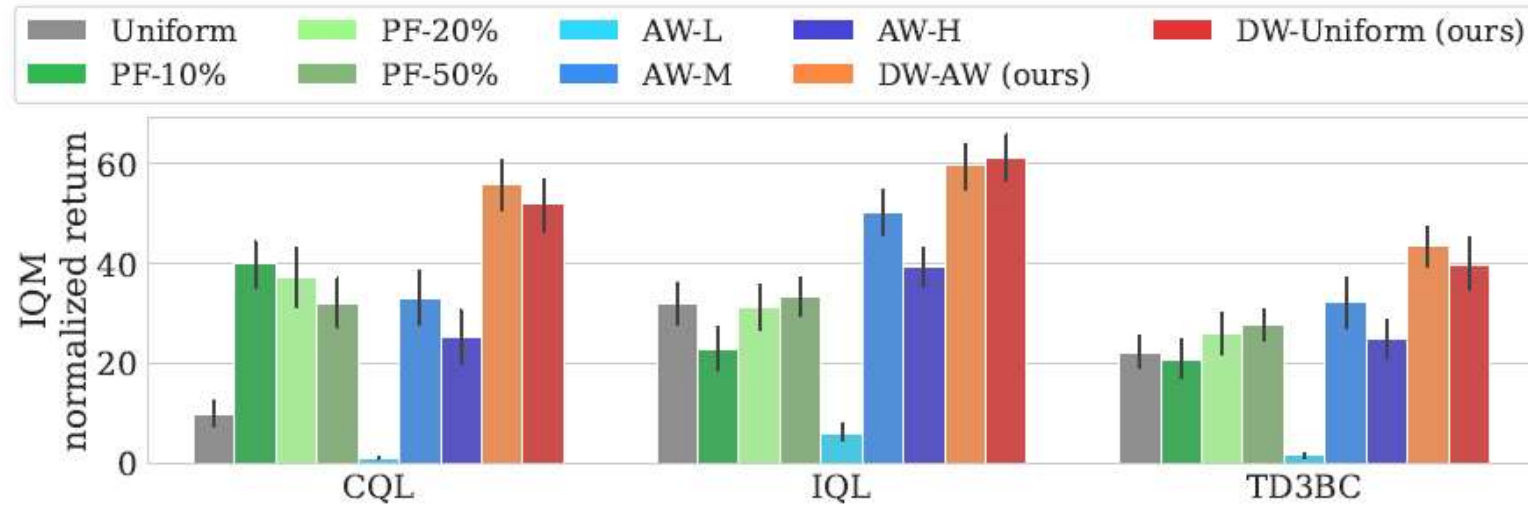


Figure 4: Results on imbalanced datasets with trajectories starting from diverse initial states (Section 5.3). Compared to Figure 3a, the performance of uniform sampling and AW decrease, showing that diverse initial states exacerbate the issue of imbalance. Our methods, DW-AW and DW-Uniform, achieve higher return than all the baselines, which suggests DW is advantageous in broader types of imbalanced datasets.