

# On/off-policy のハイブリッド深層強化学習とシミュレーション環境での制御問題への応用

## On/off-policy Hybrid Deep Reinforcement Learning and Simulation in Control Tasks

王 伯楠<sup>\*1</sup>   河合 新<sup>\*1</sup>   延原 肇<sup>\*1</sup>  
Bonan Wang   Shin Kawai   Hajime Nobuhara

<sup>\*1</sup>筑波大学大学院システム情報工学研究科知能機能システム専攻

Intelligent Interaction Technologies, Graduate School of Systems and Information Engineering, University of Tsukuba

Recently, deep reinforcement learning with neural network shows great performance in tasks such as game AI and robotics control tasks. However, on-policy and off-policy reinforcement learning methods proposed in related works have problems such as slow exploration speed. To solve these problems, we propose a hybrid deep reinforcement learning method which combines on-policy and off-policy reinforcement learning in this paper. The comparison experiment shows that the proposed method outperforms classic DDPG and DPPO method with an obvious advantage.

### 1. はじめに

近年、ニューラルネットワークを用いた深層強化学習は様々な分野で成果を出している [Lillicrap 15]。特にゲーム AI や Robotics コントロール問題では素晴らしい性能を示している [OpenAI 18b]。ニューラルネットワークは高い汎用性を持つ近似器として幅広い、かつ、複雑な問題に適用される。

強化学習アルゴリズムは主に on-policy と off-policy の二種類に分けられる。On-policy のアルゴリズムでは探索するエージェントと訓練するエージェントが同一である一方、off-policy のアルゴリズムではそれらは異なっている。On-policy のアルゴリズムは off-policy のアルゴリズムと比べて訓練が速く、安定性も高い。一方で、on-policy アルゴリズムは訓練が進むと、学習と探索の速度が遅くなる。Off-policy のアルゴリズムは過去の経験から学習するため、この問題を緩和できる。しかし off-policy のアルゴリズムは on-policy のアルゴリズムと比べて安定性が悪い。

本稿では、従来の on-policy と off-policy の深層強化学習の問題を解消するため、1) ハイブリッドなエージェントと 2) 訓練アルゴリズムを提案する。具体的には、1) ハイブリッドなエージェントとして、汎用性と効率的な学習のため、Actor-Critic 型エージェントを採用する。さらに、2) その訓練アルゴリズムとして、長期経験と短期経験の両方を活用するため、データを複数の短い時系列になるよう分割して、エージェントを学習させる。また、この時系列での訓練アルゴリズムに対応するために、ネットワークとして LSTM(Long Short-Term Memory) を用いる。提案手法では、目的関数に DDPG の価値関数と DPPO(Distribute Proximal Policy Optimization)[Schulman 17] の価値関数の両方を用いる。

本稿は次のように構成される。第 2 章では関連する従来研究を紹介する。第 3 章では提案手法を説明する。第 4 章では評価実験をその結果を示す。最後に第 5 章で結論を述べる。

### 2. 関連研究

この章では、主に関連研究の DDPG と DPPO アルゴリズムを説明する。

連絡先: 王 伯楠, 筑波大学大学院システム情報工学研究科知能機能システム専攻, wang@cmu.iit.tsukuba.ac.jp

### 2.1 DDPG(Deep Deterministic Policy Gradient)

DDPG アルゴリズム [Lillicrap 15] は典型的な off-policy の Actor-Critic アルゴリズムである。DDPG エージェントは行動を決めるネットワーク  $\pi$  と評価ネットワーク  $Q$  で構成されている。 $\pi$  は状態の観測値  $s \in \mathbb{R}^m$  を入力とし、行動  $a \in \mathbb{R}^n$  あるいは行動の分布を出力する。 $m, n$  はそれぞれ状態空間と行動空間の次元数である。 $Q$  ネットワークは状態  $s$  と対応する行動  $a$  を評価し、評価値  $q \in \mathbb{R}$  を出力する。

訓練する時、まず  $Q$  を更新する。 $Q$  を訓練した後、それを教師とし、 $\pi$  を訓練する。 $\pi$  は  $q$  を最大化するように  $\theta^\pi$  を更新する。DDPG は Replay Buffer から経験データをランダムでサンプリングして訓練する。よって、DDPG アルゴリズムは自由に探索できるが、動作の変化が激しかったり、収束が遅かったりするなどの問題がある。

### 2.2 DPPO(Distribute Proximal Policy Optimization)

DPPO アルゴリズム [Schulman 17] は典型的な on-policy のアルゴリズムである。DPPO エージェントは行動ネットワーク  $\pi$  と価値ネットワーク  $V$  で構成されている。 $\pi$  ネットワークは状態の観測値  $s$  を入力とし、行動の確率分布  $N^n$  を出力する。 $V$  ネットワークは状態の観測値ベクトル  $s$  を入力とし、価値の予測値  $v \in \mathbb{R}$  を出力する。

DPPO アルゴリズムは最近の経験データしか使えない。よって、DPPO アルゴリズムは常に新しいデータが必要である。DPPO アルゴリズムは学習がはやく、かつ安定性も高い。しかし、DPPO は訓練が進むと探索が遅くなる。

### 3. 提案手法

#### 3.1 ネットワーク構造

従来法の問題を解消するため、本論文は短期経験と長期経験の両方を活用した on/off-policy のハイブリッドエージェントと強化学習アルゴリズムを提案する。提案手法では Actor-Critic 型エージェントを採用する。短期と長期経験の特徴に対応して、行動策略部分の Actor は LSTM ネットワークを採用する。 $\pi$  は一連の状態の観測値  $s_t, s_{t-1}, \dots, s_{t-K+1}$ 、LSTM 初期細胞状態  $c_{init}$  と初期出力  $h_{init}$  を入力とし、行動の確率分布を出力する。ただし  $t \in \{1, 2, \dots, T\}$  は時刻あるいはステップ

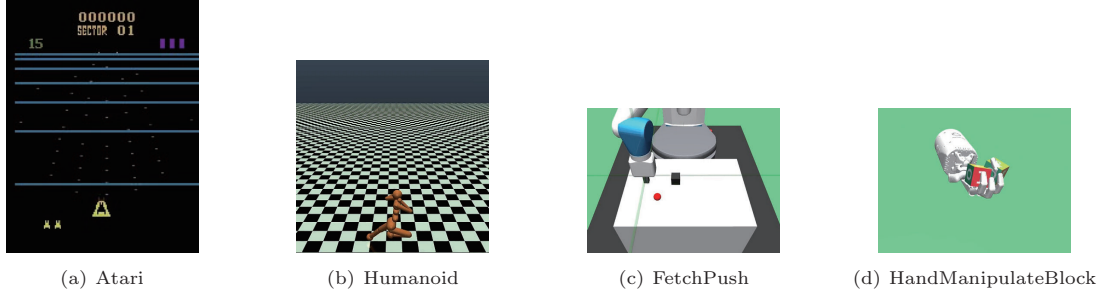
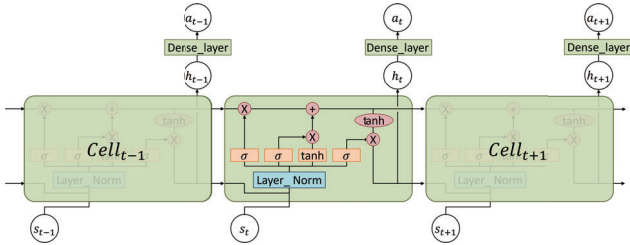


図 1: Environments Example

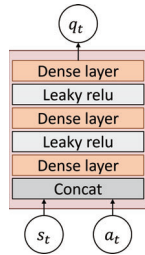
数を表し、下添え字  $t$  はステップ  $t$  における値を表すこととする。また、 $T$  は最後のステップを表す。 $k \in \{1, 2, \dots, K\}$  は LSTM セルの順番数を表し、 $K$  は LSTM のセル数を表す。

$$\pi_t(s_t, s_{t-1}, \dots, s_{t-K+1}, c_{init}, h_{init}) \sim N^n. \quad (1)$$

LSTM の各ゲートの活性化関数の前に Layer Normalization を用いる。各セルの出力に全結合層をつけて最後の出力にする。

図 2:  $\pi$  ネットワーク構造

価値関数部分の Critic ネットワークは MLP を採用し、 $Q$  とする。DDPG の Critic ネットワークと同じく  $s$  と  $a$  を入力とし、価値の予測値を出力する。負数の報酬値に対応するため、活性化関数は Leaky Relu を用いる。

図 3:  $Q$  ネットワーク構造

安定性のため、ネットワークは main ネットワークと target ネットワークで構成されている。main ネットワークは学習によって更新され、target ネットワークは main ネットワークのパラメーターを用いて更新される。

### 3.2 目的関数

長期経験と短期経験の両方を活用するため、それぞれ  $Q$  と  $\pi$  の目的関数を設計した。短期経験を使う on-policy 的な目的関数は  $L^{vf}(\theta^Q)$  と  $L^{pg}(\theta^\pi)$  があり、長期経験を使う off-policy 的な目的関数は  $L^Q(\theta^Q)$  と  $L^\pi(\theta^\pi)$  がある。ただし、 $\theta^\pi, \theta^Q$  はそれぞれのネットワークの重みなどのパラメーターを表している。アルゴリズムはこの 4 つの目的関数を最適化する。

on-policy の価値目的関数  $L^{vf}(\theta^Q)$  :

$$L^{vf}(\theta^Q) = \mathbb{E}(Q(s_t, a_t) - Q'_t)^2 \quad (2)$$

$$\text{where} \quad (3)$$

$$Q'_t = \sum_{i=t}^T (\gamma^{i-t} r_i) + Q(s_{T+1}, \pi(s_{T+1})). \quad (4)$$

on-policy の行動目的関数  $L^{pg}(\theta^\pi)$  :

$$L^{pg}(\theta^\pi) = \mathbb{E} \left( \frac{\pi(a_t | s_t)}{\pi'(a_t | s_t)} * \hat{A}_t \right), \quad (5)$$

$$\hat{A}_t = \delta_t + \gamma \lambda \delta_{t+1} + \dots + (\gamma \lambda)^{T-t} \delta_T, \quad (6)$$

$$\delta_t = r_t + \gamma * Q_{t+1} - Q_t. \quad (7)$$

ただし、 $\pi'$  は target ネットワークの  $\pi$  を表す。off-policy の価値目的関数  $L^Q(\theta^Q)$  :

$$L^Q(\theta^Q) = \mathbb{E}(Q(s_t, a_t) - c_1 * Q'_1 - c_2 * Q'_2)^2, \quad (8)$$

$$Q'_1 = r_t + \gamma * Q(s_{t+1}, \pi(s_{t+1})), \quad (9)$$

$$Q'_2 = r_k + \gamma * r_{k+1} + \dots + \gamma^{K-k} * r_K + \quad (10)$$

$$\gamma^{K-k+1} * Q(s_{K+1}, \pi(s_{K+1})). \quad (11)$$

ただし、 $c_1$  と  $c_2$  は重み引数である。

off-policy の行動目的関数  $L^\pi(\theta^\pi)$  :

$$L^\pi(\theta^\pi) = \mathbb{E}(Q(s_t, \pi(s_t))). \quad (12)$$

$L^Q(\theta^Q)$  と  $L^{vf}(\theta^Q)$  を最小化し、 $L^\pi(\theta^\pi)$  と  $L^{pg}(\theta^\pi)$  を最大化する。

### 3.3 訓練アルゴリズム

訓練する時、先ず複数のエージェントを使って探索し、データを収集する。各エージェントは環境の中で一定のステップ数  $T$  を実行し、これを一つの iteration とする。収集したデータ  $(s_t, a_t, r_t, s_{t+1}, c_t, h_t)$  を Replay Buffer  $\mathcal{M}$  に保存する。探索が効率よく進むため、各エージェントは探索の途中、小さい確率  $\epsilon$  でランダムな行動を取る。 $\epsilon$  はエージェントによって違う。安定のため、target ネットワークを用いて探索を進め、main ネットワークを更新する。

## 4. 評価実験

### 4.1 実験条件

評価実験は深層強化学習においてベンチマーク的なタスク Humanoid タスクを用いる。Humanoid タスクの目的はシミュ

**Algorithm 1** on/off-policy algorithm

---

```

1: initialize Replay Buffer  $\mathcal{M}$ ;
2: initialize Network  $\pi_{main}, \pi_{target}, Q_{main}, Q_{target}$ ;
3: for  $i$  in  $Max\_Iterations$  do
4:   for  $w$  in  $N\_Workers$  do
5:     reset Iteration Buffer  $\mathcal{I}$ ;
6:     for  $t$  in  $T$  do
7:        $a_t, c_t, h_t \leftarrow \pi_{target}(s_t)$ ;
8:        $a_t \leftarrow \text{Random Action with } \epsilon$ ;
9:        $s_{t+1}, r_t \leftarrow \text{Environment}(a_t)$ ;
10:      store  $(s_t, a_t, r_t, s_{t+1}, c_t, h_t)$  in  $\mathcal{I}$ ;
11:    end for
12:    store  $\mathcal{I}$  in  $\mathcal{M}$ ;
13:    update  $\theta_{main}^Q$  with  $L^{vf}(\theta^Q)$ ;
14:    update  $\theta_{main}^\pi$  with  $L^{pg}(\theta^\pi)$ ;
15:     $\theta_{target}^Q \leftarrow (1 - \tau) * \theta_{target}^Q + \tau * \theta_{main}^Q$ ;
16:     $\theta_{target}^\pi \leftarrow (1 - \tau) * \theta_{target}^\pi + \tau * \theta_{main}^\pi$ ;
17:  end for
18:  sample random transition  $Batch$ ;
19:  update  $\theta_{main}^Q$  with  $L^Q(\theta^Q)$ ;
20:  update  $\theta_{main}^\pi$  with  $L^\pi(\theta^\pi)$ ;
21:   $\theta_{target}^Q \leftarrow (1 - \tau) * \theta_{target}^Q + \tau * \theta_{main}^Q$ ;
22:   $\theta_{target}^\pi \leftarrow (1 - \tau) * \theta_{target}^\pi + \tau * \theta_{main}^\pi$ ;
23: end for

```

---

レーション環境の中の人間を前進することを学習させることである。エージェントは人間の各関節の動作を決める。

表 1: 実験条件

|                       |                  |
|-----------------------|------------------|
| Deep Learning Library | Tensorflow 1.9.0 |
| Simulation Library    | openAI gym       |
| Env                   | Humanoid-v2      |

訓練する時、Worker ごとに on-policy の目的関数を用いてエージェントを更新する、五つの Worker ごとに off-policy の目的関数を用いてエージェントを更新する。

表 2: Hyper Parameters

|                     |                  |
|---------------------|------------------|
| $T$                 | 1024             |
| N_Worker            | 5                |
| $\pi$ learning rate | 0.00004-0.000001 |
| $Q$ learning rate   | 0.00008-0.000002 |
| batch_size          | 32               |
| epoch               | 5                |
| sample_size         | 640              |
| $\tau$              | 0.2              |

**4.2 実験結果**

本研究はベンチマークタスク Humanoid 環境を用いて、従来の DDPG、DPPO と提案手法に対して比較実験を行いました。実験の結果は図 4 に示す。DDPG アルゴリズムはランダムな短い経験だけ用いるため、長い時系列の中の行動の関連性を無視している。よって、エージェントは最初から行動が固定し、前に進まない。DPPO アルゴリズムは学習ができるが、学習の速度が遅い。

提案手法では、安定かつ高速な学習を実現している。提案手法の on-policy 部分は安定な学習を実現し、off-policy 部分はその上探索を進ませ、学習を加速させたと思われる。

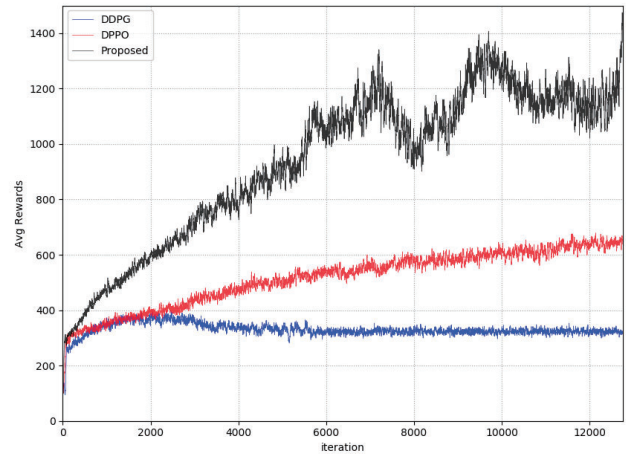


図 4: Performance Comparison of DDPG, DPPO and Proposed Method on benchmark task Humanoid

**5. おわりに**

本研究は従来の深層強化学習アルゴリズムの問題点を解決するため、短期経験と長期経験の両方を活用する on/off-policy のハイブリッド深層強化学習アルゴリズムを提案した。Humanoid タスクにおいて、提案手法は従来手法より優れた性能を示した。

Humanoid タスクのような、状態空間と行動空間は連続で、報酬値も丁寧に設計されているタスクに対して、提案手法は有効であると思われる。しかし現実の問題において、報酬値は二進値のタスクは多数存在している。エージェントの汎用性を向上させるため、これからは HER[Andrychowicz 17] 手法の投入と Robotics タスク [OpenAI 18a] への実装を検討している。

**参考文献**

- [Andrychowicz 17] Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, P., and Zaremba, W.: Hindsight Experience Replay, *CoRR*, Vol. abs/1707.01495, (2017)
- [Lillicrap 15] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D.: Continuous control with deep reinforcement learning, *CoRR*, Vol. abs/1509.02971, (2015)
- [OpenAI 18a] OpenAI, : Ingredients for Robotics Research (2018)
- [OpenAI 18b] OpenAI, : Learning Dexterous In-Hand Manipulation, *CoRR*, Vol. abs/1808.00177, (2018)
- [Schulman 17] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O.: Proximal Policy Optimization Algorithms, *CoRR*, Vol. abs/1707.06347, (2017)