

深層強化学習によるロボットのバラ積みピッキング習得 (第1報)

町田晃平・石黒 聡・細谷 肇*

Learning of Robot Bin Picking by Deep Reinforcement Learning
(1st Report)

Kohei MACHIDA, Satoshi ISHIGURO, Hajime HOSOYA

深層強化学習 Deep Q-Learning によるロボットのバラ積みピッキングを提案する。人間の教示を組み合わせたことにより、約 10 時間の深層強化学習でピッキングに成功した。今後は、本バラ積みピッキング技術や深層強化学習技術を展開し、企業の生産性向上を支援する。

キーワード：AI、深層強化学習、ロボット

We propose robot bin picking using deep reinforcement learning called as Deep Q-Learning. We report the successes of bin picking using deep reinforcement learning in about 10 hours because of human teaching. In the future, we expand this bin picking technology and deep reinforcement learning technology to support the improvement of productivity for enterprises.

Keywords : Artificial Intelligence, Deep Reinforcement Learning, Robot

1 はじめに

工場の生産性を高めるには、ロボットの活用が必要不可欠である。しかし、現状は多くの作業を人手に頼っている。そのような作業の一つとして、バラ積みされたワークの取り出し（バラ積みピッキング）がある。

バラ積みピッキングをロボットとマシンビジョンで実現しようとした場合、バラ積みワークが複雑に重なり合っているため、従来の画像処理手法ではピッキング箇所の取得が難しいという問題があった。

一方、近年の AI（人工知能）技術の発達により、特に画像認識分野において、様々な成果が報告されている。中でも、深層強化学習は、AI が目的を達成するために自ら試行錯誤を繰り返して学習し、最適な行動方法を獲得していく手法である。深層強化学習は、事前に学習データを用意する必要がなく、AI が自動で高いスキルを獲得することができるため、産業分野への応用が期

待されている。

そこで、本研究では、この深層強化学習によるロボットのバラ積みピッキング習得を目指す。昨年度の研究では、深層強化学習による平面に置かれたワークのピッキングを実現した。本研究では、バラ積みされたワークのピッキングに取り組む。

2 研究方法

2.1 強化学習

強化学習は、エージェント（AI）が環境との相互作用により学習を行う。強化学習の枠組みを図 1 に示す。エージェントは、タイムステップ t において、観測した状態 $s_t \in S$ から行動 $a_t \in A$ を選択する。それにより、エージェントは環境から報酬 $r_{t+1} \in R$ と次の状態 s_{t+1} を受け取る。

エージェントの目標は、受け取る報酬の累積（収益）を最大化することである。収益は以下のように定義される。

ここで、 $\gamma(0 \leq \gamma \leq 1)$ は割引率と呼ばれ、

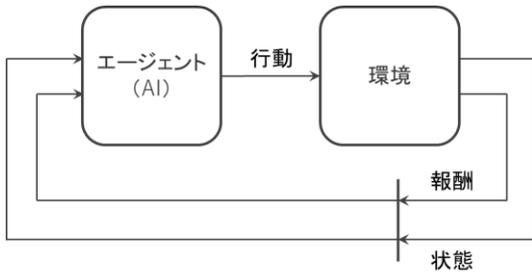


図1 強化学習の枠組み

$$R_t = \sum_{k=0}^T \gamma^k r_{t+k+1} \quad (1)$$

将来の報酬への割引を決定するパラメータである。また、 T はタイムステップの終端を表す。

方策 π のもと、状態 s において行動 a を取ることの価値を行動価値関数 $Q^\pi(s, a)$ として表す。行動価値関数は以下のように定義される。

$$\begin{aligned} Q^\pi(s, a) &= E_\pi\{R_t | s_t = s, a_t = a\} \\ &= E_\pi\left\{\sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s, a_t = a\right\} \quad (2) \end{aligned}$$

E_π は方策 π における期待値である。上式のように、行動価値関数は期待報酬として表される。エージェントに行動価値関数を学習させるアルゴリズムをQ学習と呼ぶ。Q学習では、以下のように行動価値関数を更新する。

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a) \right] \quad (3)$$

学習で獲得される行動価値関数は、使われている方策とは独立に、最適行動価値関数 Q^* を近似する。これにより、エージェントは現在の状態 s_t から行動価値関数を用いて最適な行動 $\max_a Q^*(s_t, a)$ を取ることができるようになる。

2.2 深層強化学習

本研究では、深層強化学習のアルゴリズムとして、Mnihらによって提案されたDeep Q-Learning¹⁾を用いた。Deep Q-Learningは、Q学習における行動価値関数をニューラルネットワークによって近似する手法である。Deep Q-LearningにおけるニューラルネットワークをDeep Q-



図2 実験環境

Network(DQN)と呼ぶ。DQNに状態を入力することで、各行動を行った場合の価値が出力される。

近似された行動価値関数を $Q(s, a; \theta)$ と表す。ここで、 θ はニューラルネットワークのパラメータである。Deep Q-Learningでは、このパラメータ θ を学習により更新する。

Deep Q-Learningのアルゴリズムをアルゴリズム1に示す。アルゴリズム1において、リプレイメモリ D は、遷移情報 $(\phi_t, a_t, r_t, \phi_{t+1})$ を格納しておくためのメモリである。学習の際には、このリプレイメモリ D から遷移情報のミニバッチ $(\phi_j, a_j, r_j, \phi_{j+1})$ を取り出して学習を行う。これは、Experience Replay(経験再生)と呼ばれる手法である。時系列として相関のない状態・行動・報酬の組を用いて学習した方が、学習の効率が上がるため、Deep Q-LearningではこのExperience Replayが用いられている。また、学習時は試行錯誤による探索を促すため、 ϵ -greedy方策に従って確率 ϵ でランダムな行動を選択する。それ以外の場合には、最適な行動 $\max_a Q^*(s_t, a; \theta)$ を選択する。

2.3 実験条件

実験の環境を図2に示す。ロボットは、小型のロボットアーム uArm Swift Pro (UFACTORY社)を用いた。uArm Swift Proは、pyufと呼ばれるPythonライブラリを用いて、Pythonにより制御した。カメラは、RealSense Depth Camera D435 (Intel社)を用いた。RealSense Depth Camera D435は、ステレオビジョンの深度

アルゴリズム 1 Deep Q-learning

リプレイメモリ D をキャパシティ N に初期化する

行動価値関数 $Q(s, a; \theta)$ をランダムな重み θ で初期化する

各エピソード (1~ M) に対して繰り返し:

初期画像 x_1 を初期状態 $s_1 = \{x_1\}$ として初期化し前処理 $\phi_1 = \phi(s_1)$ を行う

各タイムステップ (1~ T) に対して繰り返し:

行動 $a_t = \begin{cases} \text{random} & (\text{確率 } \varepsilon) \\ \max_a Q^*(s_t, a; \theta) & (\text{確率 } \varepsilon - 1) \end{cases}$ を選択する

行動 a_t を実行し報酬 r_t と次の画像 x_{t+1} を受け取る

次の状態 $s_{t+1} = s_t, a_t, x_{t+1}$ をセットし前処理 $\phi_{t+1} = \phi(s_{t+1})$ を行う

遷移情報 $(\phi_t, a_t, r_t, \phi_{t+1})$ をリプレイメモリ D に格納する

リプレイメモリ D から遷移情報 $(\phi_j, a_j, r_j, \phi_{j+1})$ のミニバッチを取り出す

教師データ $y_j = \begin{cases} r_j & (\phi_{t+1} \text{ が終端の場合}) \\ r_j + \gamma \max_{a'} Q(\phi_j, a'; \theta) & (\phi_{t+1} \text{ が終端でない場合}) \end{cases}$ をセットする

教師データと行動価値関数の誤差 $(y_j - Q(\phi_j, a_j; \theta))^2$ を求め誤差逆伝搬を実行する



図 3 RGB 画像

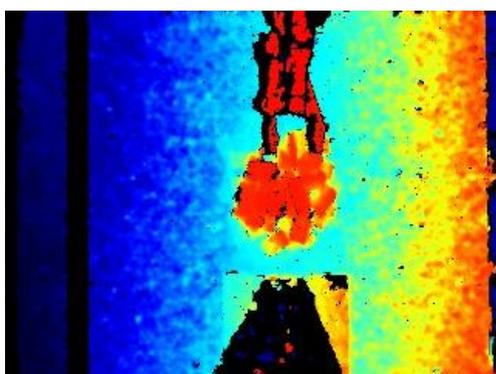


図 4 深度画像

カメラである。2つの深度センサ、RGBセンサ、IR 投射器を搭載する。本実験では、RealSense Depth Camera D435 より取得した RGB 画像 (図 3) と深度画像 (図 4) を DQN へ入力した。画像のサイズは、どちらも 320×240 の 3 チャンネルである。バラ積み部品を想定したワークは、ドミノ 60 個を使用した。ドミノは平面の上にラン

ダムに山積みした。

Deep Q-Learning のプログラムは、深層強化学習ライブラリ ChainerRL を用いて Python で実装した。

エージェントは各タイムステップにおいて、アーム水平方向の上・下・左・右移動と垂直方向の上・下移動およびピックングの 7 つの行動から選択を行う。ピックングを選択するとその場でアームを閉じ、uArm Swift Pro のピックング判定によりピックングの成功が確認されれば、エージェントに報酬+1 が与えられる。ピックングに失敗した場合やアームを移動した場合の報酬は-1 である。タイムステップの終端 T は 100 とした。タイムステップの終端 T が終わると、アームは初期位置に戻り、次のエピソードが開始される。

実験に使用した DQN の構造を図 5 に示す。3 層の畳み込み・プーリング層と 2 層の全結合層からなるニューラルネットワークである。最初の畳み込み層へは、3 チャンネルの RGB 画像と深度画像を結合した 6 チャンネルの画像を入力した。畳み込み層と全結合層の活性化関数は、Leaky ReLU を用いた。

学習の最適化手法は、SGD (確率的勾配降下法) を用いた。SGD の学習率は 0.01 とした。学習時の GPU は、NVIDIA GeForce 1080 Ti を用いた。

報酬の割引率 γ は 0.99 とした。 ε -greedy

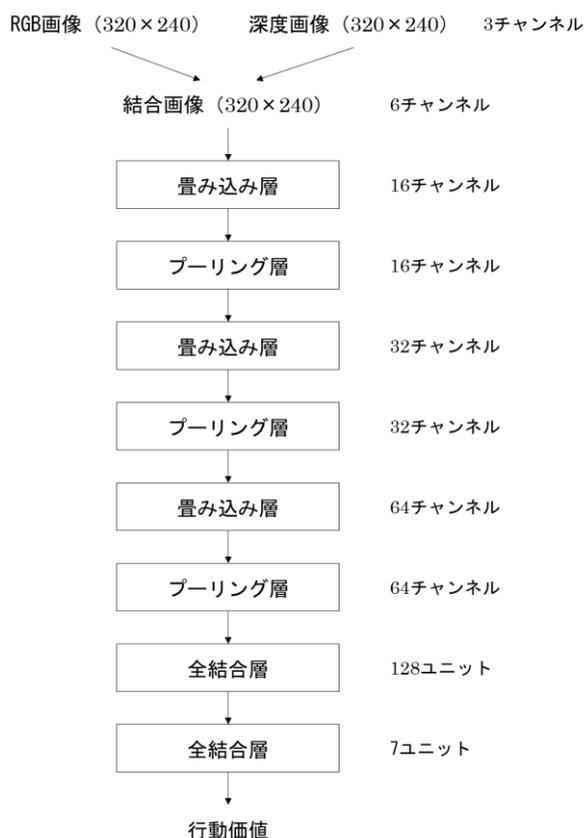


図5 Deep Q-Networkの構造

方策の確率 ϵ は、初期タイムステップで0.5、その後各タイムステップで線形に減少させ、3000タイムステップ(30エピソード)で0.3となるようにした。リプレイメモリ D のキャパシティ N は1,000,000とした。Experience Replayには、教師データと行動価値関数の誤差が大きい遷移情報を優先的に取り出して学習を行うPrioritized Experience Replay(優先順位付き経験再生)を用いた。

昨年度の研究より、事前に人間の教示を行うことで、深層強化学習の学習時間を短縮できることが確認された。そのため、本研究では、事前に約30分(500タイムステップ分)の教示を行った後に、深層強化学習を行った。教示は、各タイムステップにおいて、オペレータがロボットの取るべき行動を選択した。

3 研究結果

約30分の人間の教示後、約10時間の深層強化学習により、バラ積みされたワーク

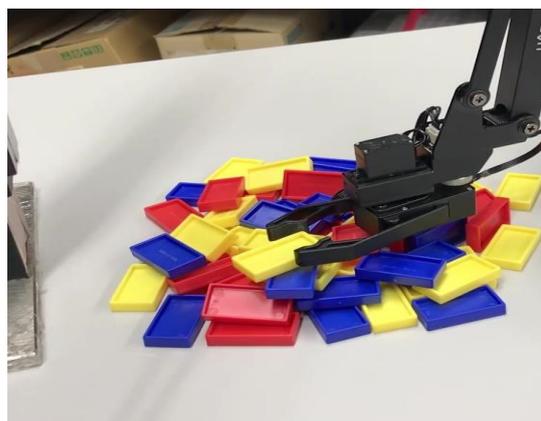


図6 ピッキング成功時(教示30分、深層強化学習10時間後)

のピッキングに成功した(図6)。

一方、人間の教示を行わなかった場合には、1日以上、深層強化学習を行ってもピッキングに成功しなかった。このことから、実用的な時間内で学習を完了するには、人間の教示が重要であることが確認された。また、本バラ積みピッキングを従来の画像処理手法等で実現しようとした場合、数日以上ソフト開発工数がかかることが想定される。このことから、従来手法と比較して、深層強化学習を用いることで開発工数の削減が期待できる。

4 まとめ

深層強化学習 Deep Q-Learning によるロボットのバラ積みピッキングを実現した。本研究の成果により、人間の教示を組み合わせることで実用的な時間内に深層強化学習を完了できることや、従来手法と比較して深層強化学習がソフトの開発工数削減につながる可能性があることが確認された。

今後は、本バラ積みピッキング技術や深層強化学習技術を県内企業に展開し、企業の生産性向上を支援する。

文献

- 1) Volodymyr Mnih, *et al.*: Human-lev elcontrol through deep reinforcement- learning, *Nature*, 518(7540):529, 02(2015)