深層強化学習による施工機械の経路生成

Route Generation of Construction Machine by Deep Reinforcement Learning

田邉峻也 Shunya Tanabe 孫澤源 Zeyuan Sun 中谷優之 Masayuki Nakatani 内村裕 Yutaka Uchimura

芝浦工業大学

Shibaura institute of technology

Recently, research on artificial intelligence has been progressing in various fields. In some of the Atari 2600games, the AI player has scored higher than the skilled human players by using deep reinforcement learning techniques. In this paper, autonomous ground leveling work by a bulldozer is targeted, which are expected to optimize action of the bulldozer. In previous work, we implemented deep Q learning method by giving images of simulator as the input data for the network. However, when learning the image using the convolution layer as input using deep reinforcement learning, it requires large computational cost for learning process. This research aims to reduce the computational cost by giving smaller order of input data. This paper describes the comparison results in different order of input data. Transition of the learning sequence is also evaluated.

1. 緒言

近年,自動車の自動運転,囲碁でプロ棋士に勝利する [David 16] など様々な分野で人工知能の研究が進んでいる. Atari2600 のゲーム [Bellemare 13] をルールの教示無しに画 面の状態とスコアのみからコンピュータが繰り返し学習を行 うことで,人間よりもスコアを上回った結果が報告されてい る [Mnih 15].同論文では Deep Q Network (DQN) と呼ば れる強化学習の手法が用いられている.DQN は,環境の状況 からエージェントが採るべき行動を選択し,報酬によってよ り良い行動の選択をする Q 学習法 [Riedmiller 09] と,画像 処理の分野において大きな成果 [Krizhevsky 12] を上げてい る深層学習を組み合わせた手法である.さらに発展させた手 法として Double Q learning[Hasselt 15] や Dueling Double DQN[Ziyu 16] など様々な手法によって人工知能の性能は向上 している.

このような人工知能の応用分野として、本研究では建設作 業の自動化を対象として行っている。建設作業のなかでも不整 地の整地作業は、切り土、盛り土の状況に応じた運転経路の判 断が必要なため,最終的な整地形状(出来型)に到達するため の逐次の運転動作は、囲碁における大局観に通じるものがあ る。よって本研究では、不整地をブルドーザのような施工機械 で整地する作業の自律化を想定し、DQN の手法を用いること によりブルドーザが自律的に最適な運転経路を生成すること を目指す.本研究では、不整地をブルドーザのような施工機械 で整地する作業の自律化を想定し, DQN の手法を用いること によりブルドーザが自律的に最適な運転経路を生成すること を目指す. しかし, DQN による学習においては, ニューラル ネットワークの入力,層数,各層のニューロン数,学習率など 様々なハイパーパラメータの調整が必要であり,正しい値を設 定しないと学習時間が長くなることや、学習がうまく行かない ことがあるためハイパーパラメータ調整は重要である.ハイ パーパラメータの調整には試行錯誤を重ねていく必要がある が、畳み込み層を利用した画像を入力として学習する場合、マ ルチコアの CPU や GPU を使っても学習の過程に非常に長い

連絡先: 内村 裕, 芝浦工業大学, 東京都江東区豊洲 3-7-5, uchimura@shibaura-it.ac.jp



図 1: シミュレーション概要

時間を要することが課題である.これまでに筆者らは,DQN を整地シミュレーションに適用して経路計画する研究を行った が[孫 17,中谷 17],ネットワークへの入力が画面データのた め,入力サイズが大きく計算負荷が高くなることが課題であっ た。そこで本研究では入力サイズを小さくすることで,ニュー ラルネットワークの規模を縮減し,ニューラルネットワークに おける計算コストを減らすことで学習の高速化を図ることを目 的とする.

具体的には、ネットワークへの入力を画像データの代わり に、ブルドーザや土砂の座標データもしくは、ブルドーザと土 砂との相対座標とし、入力のサイズを小さくした場合の学習を 行い、経路の最適性と高速化について検証する.

2. 問題設定

土砂の整地作業を単純化してシミュレートするために図1に 示すように整地エリアを8×8のグリッド上に設定した.エー ジェントであるブルドーザで土砂山を移動し,図2のように地 面の凹凸を無くし整地を行うことを目的とする.本シミュレー ションでは図3のように土砂の高さを6段階で表現し,明度が 高いと土砂が高く,暗い色ほど土砂が低いことを表している. ブルドーザは上下左右の4方向に動き,土砂を押すことがで



図 2: 整地終了時



図 3: 土砂と穴の概要

きるが、引くことはできない. 基準となる地面の高さを0とし たとき、高さ –1 の地点の穴に高さ1以上の土砂を移動させる ことで整地を行う.また、図6のように、各エピソードにおい て6パターンの初期位置をランダムに与えた.報酬の設定は、 土砂の高さを –1 の地点に高さ1以上の土砂を近づけて移動 させることで +1、土砂を1エリア整地することで+10、すべ てのエリアを整地することで+50とした.すべての穴が整地 されるか、ブルドーザの行動回数が100回に達した場合に1 エピソードの終了とし、エピソードを繰り返し行うことで Q ネットワークの学習を行った.

ネットワークに与える入力の違いによる学習結果の比較のた めに以下の2つの異なる入力データを用意した.

- ケース1:ブルドーザの、土砂山・穴の座標値を入力とした場合
- ケース2:ブルドーザの周囲のセルの地面の高さを入力 とした場合

ケース2では入力をシミュレータ画面(8×8)より大きなサ イズ(15×15)のセルとしているが画面外のセルの値は0と した。ケース1とケース2のネットワークの各パラメータを 表1に示す.中間層に3層の全結合層を用いており,各層の ニューロン数は1000,250,50とした。また,ブルドーザの 前後左右方向の移動の4つのアクションを出力とした.

学習過程では、割引率 $\gamma \ge 0.99$ に設定した. ϵ -greedy 法に よるランダム率は、学習開始時では 100%とし、その後、学習 回数が増えるに従って徐々に減らし、学習回数が 100 万回以 降は 10%とした.また、学習データを保存するメモリサイズ

表 1: ケース1とケース2のパラメータ

パラメータ	ケース1	ケース2
入力サイズ	38	15×15
出力数	4	4
中間層	3 層	3 層
中間層の種類	全結合層	全結合層
ミニバッチ数	50	50
学習率	1.0×10^{-5}	1.0×10^{-5}
割引率	0.99	0.99

	Links Enter		~ `		- 0 - 1
表 2:	又献 将	§ 171	のネッ	トワーク	のバフメータ

パラメータ	文献 [孫 17]		
入力サイズ	$40 \times 40 \times 4$		
出力数	4		
中間層	7層(畳み込み3層,		
	プーリング3層,全結合1層)		
ミニバッチ数	50		
学習率	1.0×10^{-5}		
割引率	0.99		

は 10 万エピソード分のリングバッファとした.本研究では, Double Q learning を適用し,教師信号を出力するニューラル ネットワークはランダムに初期化すると同時に,1 万回の学習 ごとに学習済みのニューラルネットワークで上書きした.シ ミュレーションの1エピソードの最大行動回数は100回に設 定した.また学習回数が10万回間隔で学習したネットワーク を保存し,ランダム率を0%としたシミュレーションを行い, 行動回数と整地した穴の数を記録した.なお,シミュレータの 作成には Pygame を使用し,Double Q learning の実装にお いては,TensorFlow をフレームワークとして用いた。また, GPU として GeForce GTX 1070 (NVIDIA 製)を使用した.

さらに,ニューラルネットワークに与える入力サイズによ る計算コストの比較のために,文献 [孫 17] における画像デー タを入力とした場合の学習も行った.文献 [孫 17] では,1グ リッドを5×5ピクセルで構成した8×8のグリッドのシミュ レータの画面(40×40ピクセル,8bit 階調)のピクセルデー タを1画面として,過去3画面を含めて4画面分を入力とし ている.同ネットワークに使用したパラメータを表2に示す. 中間層には畳み込み層が3層,プーリング層が3層,全結合 層が1層の合計7層をニューラルネットワークに与えられる. しかし,ニューラルネットワークの規模が大きいため,学習時 間が長くなるのが問題点である.

3. シミュレーションによる学習結果と比較

3.1 入力データの相違による学習結果

入力層に与える入力データとしてをブルドーザの,土砂,穴 の座標値とした場合(ケース1)と,ブルドーザの周囲のセル (15×15)の地面高さを入力とした場合(ケース2)のシミュ レーションによる学習結果を図5に示す.グラフ横軸が学習回 数,縦軸が各初期位置の穴をすべて整地するまでの行動回数の 平均を示している.本結果に示すようにケース1ではブルドー ザの行動回数が減少せず,学習回数1,300万回の試行を行って も穴を全て埋めるには至らず,学習による経路生成が不調に終 わった.

ケース2では、学習を重ねるごとにブルドーザが全ての穴





(文献 [孫 17])

図 4: 生成経路の比較



図 5: ケース1とケース2の学習結果の比較

を整地するまでの行動回数が減少していることがわかる.また 穴や土砂の初期位置の異なるパターンにおいても整地を完了す る経路の最適化することができた.また,学習回数が240万 回のときに保存したニューラルネットワークを用いてランダム 率0%でブルドーザをシミュレーション上で動かしたときの6 パターンの各初期位置の平均行動回数が31.6回となった.そ の後,学習回数1,300万回まで学習を行ったが,ブルドーザは 31.6回より行動回数が少なくなることはなかった。

検証結果が示すように、ケース2は学習が進むにつれ整地 を完了するまでの行動回数を減少したが、ケース1は最大行動 回数の100回以内に全ての穴の整地を行うことができなかっ た.また、学習回数を増やしてもすべての穴を整地するまでの ブルドーザの行動回数は減少することはなかった.このとき、 エピソード開始時は、ブルドーザは上方への行動を数回選択 し、いくつかの穴を整地したが、その後全く整地に至る動作を 行わなかった.この結果から、ケース1で与えられた入力デー タでは行動価値関数の関数近似が上手く行われていないと考 えられる。ケース1はケース2と比べて入力サイズが少ない ことから情報量が少ないこと、入力データが座標データのスカ ラーの数値のみであり、ケース2のような2次元配列として の相関がないため学習がうまく進まなかったと推測できる。

表 3: ケース 2 と文献 [孫 17] の比較						
パラメータ	学習時間 [s]	最小行動回数 [回]				
ケース 2	5320	31.7				
文献 [孫 17]	12300	29.0				

3.2 ケース2と文献 [孫 17] の比較

ブルドーザの周囲のセルの地面高さを入力としたケース2 と, 文献 [孫 17] の学習結果の比較を図7に示す. 縦軸に初期 位置の異なる穴をすべて整地するまでのエージェントの行動回 数の平均を示し、横軸に学習時間を示している.本結果が示す ように、文献 [孫 17] が収束するまでの時間が約5時間かかる に対して、ケース2とした場合では約2時間で収束すること がわかる.また、ケース2と文献 [孫 17] を学習回数 100 万回 にかかる時間とすべての穴を埋めるまでのブルドーザの平均行 動回数を比較したものを表3に示す. これらの結果が示すよ うに,ケース2では文献 [孫 17] と比べて入力サイズを少なく することでニューラルネットワークの規模が小さくなり、学習 時間が約2分の1に短縮できたことがわかる.しかし,文献 [孫 17] のブルドーザの最小行動回数が 29.0 回に対して、ケー ス2では31.7回と行動回数が多くなっていた.また,同じ初 期位置のブルドーザの経路をケース2と文献 [孫 17] で比較し たものを図4に示す.矢印はすべての穴を整地するまでのブ ルドーザの軌跡であり、ブルドーザが1マス移動することで 行動回数が増えていく.この初期位置の場合ではブルドーザの 行動回数は、ケース2としたとき28回、文献 [孫 17] では22 回となった.ケース2を入力として学習された経路では、最初 に4つの土砂を運んで穴の整地を行っているが,最後の2つ の土砂を1つずつ運ぶことにより、ブルドーザの行動回数が 増えていることがわかる. 文献 [孫 17] の経路は, 最初に1つ ずつ土砂を運んでいるが,これ以上行動回数を減らすことがで きない最も効率の良い経路であった.

本研究で適用した手法 (Double DQN) では,現在の状態に おいてエージェントが最も価値の高い行動を選択するが,人間 と違い現在の状態から未来の状態を予測することができない問 題点があり,文献 [孫 17] に比べても効率が悪い経路を選択し ている.文献 [孫 17] と違い,入力数が少ないことや,ニュー ラルネットワークの規模,畳み込み層の有無などが影響して局 所解に陥ってしまった結果,最短経路に到達できなかったと考 えられる.







図 7: ケース 2 と文献 [孫 17] の学習結果の比較

4. 結言

本研究では、ブルドーザによる整地作業の自律化を目的に、 ニューラルネットワークの入力データとして、ブルドーザや土 砂の座標データ、もしくはブルドーザ周囲のセルの状態を与え て学習結果の比較を行った.ブルドーザ周囲のセルの状態を与 えることで入力サイズを小さくすることで、ネットワークの規 模も縮減し、学習の高速化を実現した.今後は経路をより最適 化することを目指すと同時に、より現実に近い問題設定におけ る学習の高速化を図る予定である.

参考文献

- [David 16] David, S., Aja, Huang. "Mastering the game of Go with deep neural networks and tree search" Nature 529 484/489, (2016)
- [Bellemare 13] Bellemare, M. G. Naddaf, Y. Veness, J. Bowling, M. :"The arcade learning environment: An evaluation platform for general agents" Journal of Artificial Intelligence Research 47 253-279 (2013)
- [Mnih 15] Mnih, V., Kavukcuoglu, K. and Silver, D. "Human level control through deep reinforcement learning" Nature 518 pp.529-533, 2015.
- [Riedmiller 09] Riedmiller, M. Gabel, T. Hafner, R. :"Reinforcement learning for robot soccer" Auton Robot (2009) 27: 55-73 DOI:10.1007/s10514-009-9120-4 (2009)
- [Krizhevsky 12] Krizhevsky, A. Sutskever, I. Hinton, G. :"ImageNet classification with deep convolutional neural networks" Adv. Neural Inf. Process. Syst.25, 1106-1114 (2012)
- [Hasselt 15] Hasselt, H., Guez, A. and Silver, D."Deep Reinforcement Learning with Double Q-learning" arXiv:1509.06461v3 cs.LG (2015)
- [Ziyu 16] Ziyu, W., Tom, S., Matteo, H., Hado, V.H., Marc, L., Nando, F. "Dueling Network Architectures for Deep Reinforcement Learning" arXiv:1511.06581v3 cs.LG (2016)
- [孫 17] 孫澤源, 中谷優之, 内村 裕, "深層強化学習による整地作業の 自律制御", ロボット学会学術講演会論文集, 3I2-04, (2017)
- [中谷 17] 中谷優之,孫澤源,内村 裕,"深層強化学習による最適経路 学習",第 60 回自動制御連合講演会,SuD1-4, (2017)