## 市街地の自動運転における環境情報のフィルタリング

Filtering Environmental Information for Automatic Driving in Urban Area

北村清也 石川翔太 荒井幸代 Seiya KITAMURA Shota ISHIKAWA Sachiyo ARAI

千葉大学大学院融合理工学府都市環境システムコース

Department of Urban Environment Systems, Graduate School of Science and Engineering, Chiba University

This paper focuses on lack of explainability about the actions made by Deep reinforcement learning (DRL). DRL shows its superiority on tasks with multi-dimensional visual-input such as playing Atari games and navigation of robot. Also, DRL has been expected as an useful approach to realize a self-driving car. Though a lot of inputs would be available within urban environment, we could not specify the essential inputs to decide the appropriate action. The results derived from DRL are the tacit knowledge that is difficult to transfer to another system by means of writing it down as a symbolic way. Thus, human designer of self-driving operation may be reluctant to accept and utilize these results.

For the above reason, as a first stage to reach a symbolic representation, we propose a filtering method to specify the necessary inputs for the right actions, and show the effectiveness of it via some experiments.

## 1. はじめに

近年, 交通事故減少や渋滞解消などへの期待から, 市街地の 自動運転が積極的に研究されている.市街地の自動運転では, 交差点右折時の対向車や歩行者など,数多くの対象物を含む環 境入力に対して適切な行動出力が求められる. この自動運転 のタスクに対して,深層強化学習が期待される.深層強化学習 は,環境の状態識別に有効な深層学習を行動学習の一手法と して知られる強化学習に導入した手法であり, 文献 [Mnih 13] にその有効性が示されている.深層強化学習によって,自動運 転時に入力される周囲の環境情報から適切な行動出力が可能 となる.しかし,深層強化学習は深層学習と同様,ニューラル ネットワークによる関数近似器であり,入出力の関係が不明瞭 なブラックボックスといえる.したがって、環境の状態に対す る運転行動の妥当性が説明できない. 行動妥当性の説明は、自 動運転車設計に向けた信頼性を担保する上で必須であり, 深層 強化学習導入における課題といえる.この課題に対して,本研 究では,環境の状態と行動の関係を明らかにすることを目的と して,環境情報である入力の中から適切な行動出力に不可欠な 属性を抽出する方法を提案する. ここで、すべての環境情報か らの運転行動に必要な情報の抽出を「フィルタリング」と定義 する.フィルタリングによって深層強化学習における入出力関 係の説明が可能になるため、深層強化学習の自動運転導入にも 貢献する.

本研究では、はじめに深層強化学習を用いて市街地の典型的 なシーンに対する運転行動を獲得する.つぎに、学習後のネッ トワークを走行時に分析し、運転行動に必要な環境情報のフィ ルタリングを行う.

## 2. 問題設定

## 2.1 シミュレーション環境

運転行動の学習に用いるシミュレーション環境として TORCS(The Open Racing Car Simulator)を用いる. TORCS は自動運転タスクの研究に多く用いられているシミュ

連絡先:北村清也,千葉大学大学院融合理工学府都市環境シス テム,千葉市稲毛区弥生町1-33



図 1: 本稿における対象シーン

レータである.本稿では,制御対象の自動車をエージェントと 呼ぶ.エージェントは,周囲の環境特徴や自身の状態を表すセ ンサから得た状態入力をもとに行動の出力をする.対象とする 運転シーンは図1に示すスタートラインから400mまでの直 進走行区間である.

# 2.2 深層強化学習のモデリング 深層強化学習モデル

運転行動の学習には、Deep Deterministic Policy Gradient(以下,DDPG)[Lillicrap 15]を用いる。DDPG は Actor-Critic に深層学習を導入した手法である。Algorithm1 に,DDPG による運転行動の学習手順を示す。2.1 で定義したエージェントは、Algorithm1 における  $\mu(s|\theta^{\mu}), Q(s,a|\theta^{Q}), \mu', Q'$ から構成される学習器を備える。

エージェントは状態  $s_t$  を観測し, Actor ネットワークを介 して行動  $a_t$  を選択する.  $a_t$  を実行後,報酬  $r_t$ ,次状態  $s_{t+1}$  を 観測し,  $(s_t, a_t, r_t, s_{t+1})$  をリプレイメモリ R に格納する. Rからミニバッチをランダムサンプリングし,期待報酬  $y_i$  を用 いて Critic を更新した後, Actior を更新する.  $\mu(s|\theta^{\mu})$ の更 新には決定論的方策勾配法が,  $Q(s, a|\theta^Q)$ の更新にはベルマ ン方程式が用いられる. また,  $\mu'$ , Q' は期待報酬  $y_i$  の計算に 用いられる. 以上の手順を繰り返して学習を行う.

本研究において,  $\mu(s|\theta^{\mu})$ ,  $Q(s,a|\theta^{Q})$ はニューロンが 500 個, 100 個, 100 個の 3 層からなる中間層と入力層, 出力層で 構成される.入出力層のニューロン数は, 状態入力, 行動出力 の次元によって決定される. **Algorithm 1** DDPG の学習アルゴリズム

1:	Actor, Critic ネットワーク $\mu(s \theta^{\mu})$ , $Q(s,a \theta^{Q})$ の重み
	$ heta^{\mu},   heta^{Q}$ をランダム値で初期化
2:	ターゲットネットワーク $\mu', \ Q'$ の重み
	$\theta^{\mu'} \leftarrow \theta^{\mu}, \ \theta^{Q'} \leftarrow \theta^{Q}$ をランダム値で初期化
3:	リプレイメモリ R を初期化
4:	for $learn = 1$ to $F$ do
5:	探索ノイズ <b>N</b> を 0 で初期化
6:	初期状態 <b>s</b> 1 を観測
7:	for $t = 1$ to T do
8:	行動 $a_t = \mu(s_t  heta^{oldsymbol{\mu}}) + N_t$ を選択
9:	探索ノイズを更新: $N_t \leftarrow N_t - w_1 \cdot N_t + w_2 \cdot random$
10:	$a_t$ を実行し,報酬 $r_t$ ,次状態 $s_{t+1}$ を観測
11:	$(s_t, a_t, r_t, s_{t+1})$ を $R$ に格納
12:	Rから 128 個のミニバッチをランダムにサンプリング
13:	期待報酬 $y_i = r_i + \gamma Q'(oldsymbol{s_{i+1}}, \mu'(oldsymbol{s_{i+1}}  heta^{\mu'})  heta^{Q'})$ を
	計算
14:	誤差 $E = rac{1}{128} \sum_{i=1}^{128} (Q(\boldsymbol{s_i}, \boldsymbol{a_i}   \theta^Q) - y_i)^2$ を最小化す
	るよう Critic を更新
15:	決定論的方策勾配定理を用いてActorを更新:
16:	$rac{\partial J}{\partial  heta ^{\mu}}pprox rac{1}{128}\sum_{i=1}^{N} rac{\partial Q(m{s}_{m{i}},m{a}_{m{i}}   heta ^{arphi})}{\partial m{a}_{m{i}}} rac{\partial \mu(m{S}_{m{i}}   heta ^{\mu})}{\partial  heta ^{\mu}}$
17:	ターゲットネットワーク μ′,Q′ を更新:
18:	$ heta^{\mu'} \leftarrow  au  heta^{\mu} + (1- au)  heta^{\mu'}$
19:	$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$
20:	end for
21:	end for

#### 状態入力

状態入力は,表1に示す5個の属性を用いた5次元入力とする. これより,入力層のニューロン数は5であり,Algorithm1 では時刻*t*において5次元の状態ベクトル*st*が入力として扱われる. 状態入力は値域をすべて [-1,1] に正規化する.

表 1: ス	本実験に用いる入力次元	Loiacono	13
--------	-------------	----------	----

入力次元	值域	単位	詳細
trackPos	$(-\infty, +\infty)$	-	車と道路中心との距離
angle	$[-\pi, \pi]$	rad	道路軸に対する車の角度
speedX	$(-\infty, +\infty)$	km/h	車の縦軸における速度
speedY	$(-\infty, +\infty)$	km/h	車の横軸における速度
speedZ	$(-\infty, +\infty)$	km/h	車の鉛直方向における速度

#### 行動出力

	表 2:	本実験に	用いる	出力次元	Loiacono	13
--	------	------	-----	------	----------	----

出力次元	值域	単位	詳細
steering	[-1,1]	-	-1 が右, 1 が左への最大舵角を表す 舵角値

行動出力は、表2に示す1個の属性を用いた1次元出力とする. これより、出力層のニューロン数は1であり、Algorithm1では時刻tにおいて1次元の行動ベクトル $a_t$ が出力として扱われる.

#### 報酬

図2に、エージェントへ与える報酬 rt を示す.本実験では、 報酬をペナルティに限定する.エージェントは、ペナルティを なるべく受け取らないための行動を学習する.はじめに、最も 望ましくない状態としてコースアウトが考えられるため、コー スアウト時に最大のペナルティ-20を与える.つぎに、エージェ ントは直進走行時において、道路の中心付近を道路と平行に走

コースアウト時  

$$r_t = -20$$
  
毎ステップ  
if |trackPos| < 0.2:  
pos\_penalty = 0  
else: pos\_penalty = -10|trackPos|<sup>2</sup>  
if |angle| < 0.01:  
ang\_penalty = 0  
else: ang\_penalty = -|angle|  
 $r_t$  = pos\_penalty + ang\_penalty  
if speedX < 60.0:  
 $r_t = r_t - 5$ 

行することが望ましいと考えられるため、入力次元の trackPos と angle に応じたペナルティpos\_penalty, ang\_penalty を毎ス テップ与える.しかし、人間の運転熟練者でも trackPos と angle が常に 0 であることはない.したがって、多少の余裕を考慮 し、 |trackPos| < 0.2 oDとき pos\_penalty = 0、 |angle| < 0.01のとき ang\_penalty = 0 とする.また、他車両が存在しない 直線道路では、エージェントは徐行運転をする必要がなく、低 速走行はしないため、入力次元の speedX が 60.0km/h 以下の ときに-5 のペナルティを加える.

## 3. 勾配を用いたフィルタリング

本章では、提案法であるフィルタリングについて説明する. 2.2 より、フィルタリングは、状態ベクトルsにおける全要素 からの行動出力に必要な属性の抽出である.本研究では、出 力に対する入力の勾配(以下,勾配)を用いてフィルタリング を行う.勾配は入力値の微小変化に対する出力値の変化量であ る.勾配の大きい入力次元は微小変化による出力変化が大き く、影響度も大きいと考えられるため、エージェントに重要視 されているといえる.勾配は入力次元の数だけ要素が存在する ベクトルであり、本稿では式(1)で表される.フィルタリング は勾配ベクトルにおける各要素の大きさを比較して行う.

$rac{\partial \mu(oldsymbol{s}  heta^{\mu})}{\partial oldsymbol{s}} =$	[	$\tfrac{\partial \mu(\boldsymbol{s} \theta^{\mu})}{\partial s_1}$	$\frac{\partial \mu(\boldsymbol{s} \theta^{\mu})}{\partial s_2}$	<u>,</u>	$\frac{\partial \mu(\boldsymbol{s} \theta^{\mu})}{\partial s_{3}},$	$\frac{\partial \mu(\boldsymbol{s} \theta^{\mu}}{\partial s_{4}}$	<u>)</u> ,	$\frac{\partial \mu(\boldsymbol{s} \theta^{\mu})}{\partial s_{5}}$	$]^T$
								(1	L)

#### **3.1** フィルタリング手順

Algorithm2 に示すフィルタリング手順によって、入力次元 をn次元から行動出力に必要なm次元へフィルタリングでき る(n > m).以下、Algorithm2 について説明する.

#### I.L 体の安定走行を示すエージェントの獲得

深層強化学習を用いて n 次元入力により安定した直進走行を 学習させる.エージェントの安定走行を判断する基準として, 横加速度を用いる.ここで横加速度を,車の加速度ベクトルの うち進行方向と垂直な成分と定義する.通常,自動車の走行中 は横加速度が  $0.2G \sim 0.3G$  であり,これ以上大きくなると搭乗 者は不快感や恐怖感を抱くといわれている [Nasukawa 08] た め, Algorithm2 中で 0.3G という基準値を設けている.エー ジェントが 400m を走行可能かつ stability = good のとき,安 定走行とする. Algorithm 2 勾配を用いたフィルタリング 1: (n,m) = (存在する次元,必要な次元) (n > m) 2. i = 0I.L 体の安定走行を示すエージェントの獲得 3: while i < L do Algorithm1 を用いて n 次元入力で学習 4: 学習終了後,再度走行 5:6:  $(accel, ave\_accel, SD\_accel) =$ (観測される横加速度, accel の平均, accel の標準偏差) G:重力加速度 7:  $over\_accel = 0$ 8: for t = 1 to race\_finish do 9: accel を観測 10: if accel > 0.3G then 11:  $over\_accel = over\_accel + accel/0.3G$ 12:end if 13:if  $t = race_finish$  then 14: 走行距離 dist を観測 15:end if 16: end for 17:if  $(ave\_accel + SD\_accel + over\_accel) < 0.3G$  then 18: stability = good19: 20: end if 21: if  $dist = 400 \land stability = good$  then 安定走行 22: $i \leftarrow i + 1$  $23 \cdot$  $agent_i$ としてエージェントを保存 24:end if 25:26: end while Ⅱ. 勾配評価と評価値の更新 27: X = {x<sub>i</sub>|j は n 以下の自然数 }: 各入力次元の評価値 28: *R* = {*r<sub>j</sub>*|*j* は *n* 以下の自然数 }: 各入力次元のランク変数 29:  $X \leftarrow 0$ 30: for k = 1 to *L* do *agent<sub>k</sub>* における n 次元の勾配を分析 31: 分析した勾配の大きさを降順にソート,順位を R に代入 32: 以下の更新式により X を更新: 33:  $x_j \leftarrow x_j(n-r_j+1)$ 34: 35: end for Ⅲ. 走行テストに用いるエージェントの選択 36: *M* = *X* における上位 *m* 個の入力次元 37: for l = 1 to L do Mall:agentl における勾配の走行時平均上位 m 個の 38: 入力次元 *M*<sub>5m</sub>:agent<sub>l</sub> における勾配の走行開始 5m 平均上位 39. m 個の入力次元 if  $M = M_{all} \wedge M = M_{5m}$  then 40: agent<sub>l</sub> を採用 41: end if  $42 \cdot$ 43: end for **IV.** *M* を入力とした走行テスト 44: if  $dist = 400 \land stability = qood$  then フィルタリング成功 45: 46: **else** フィルタリング失敗 47:  $48 \cdot$ 2 行目に戻る

## 49: end if

## Ⅱ. 勾配評価と評価値の更新

獲得された L 体のエージェントを用いて勾配評価をする. こ こで, n 個の要素をもつ X はフィルタリングに用いる評価値 であり,各要素が各入力次元の評価値をそれぞれ表す. R はラ ンク変数であり, X と同じ要素数をもつ. R の各要素 r<sub>i</sub> は x<sub>i</sub> の入力次元に対応しており,勾配評価で明らかとなった勾配の 大きさによる各入力次元の順位が整数で格納される. 勾配は エージェントを n 次元入力で再度走行させることにより求め, 走行区間における平均値の大きさを比較して評価する. フィル タリング開始時に X は初期化され,更新式によって L 回更新 される.

Ⅲ. 走行テストに用いるエージェントの選択

L 回の更新終了後,評価値 X の値が大きい上位 m 次元を エージェントが重要視する次元として選択し、これを M とす る. 続いて, 走行テストを行うエージェントを選択する. 走 行テストを行うエージェントは、L個の中から2つの条件に より選択する.  $M_{all}$  を  $agent_l$  における勾配の走行時平均上 位 m 個の入力次元, M<sub>5m</sub> を agent<sub>l</sub> における勾配の走行開始 5m 平均上位 m 個の入力次元とすると、M は、勾配の走行時 平均による順位を用いた評価値更新結果に基づいているため,  $M = M_{all}$  である必要がある.また、L回の更新結果で得ら れた入力次元の順位は、エージェントが重要視している入力次 元の順位として信頼できると考える.しかし, M<sub>5m</sub> を満たさ ずに5次元入力で安定走行をするエージェントの場合,走行 開始時において本来重要視すべき入力次元が重要視されてい ないことになる. そのようなエージェントが安定走行である場 合,過学習が考えられる. TORCS は必ず同じ位置から停止状 態で走行を開始するため、走行開始時の過学習が起こりやすい と考えられる.本研究では、エージェントの学習と学習後の走 行は同じコースを用いている. すなわち, 訓練データとテスト データが同一であり、過学習による安定走行が考えられる.し たがって, $M = M_{5m}$ を満たす必要がある.

## **IV.** M を入力とした走行テスト

Ⅲ. で採用されたエージェントを用いて, *M* を入力とする走 行テストを行う. 採用されたすべてのエージェントが走行テス トにおいて安定走行を示した場合,フィルタリングは成功と なる.

## 4. 計算機実験

本実験では、学習に用いる入力次元 n = 5、フィルタリン グによって獲得する必要な次元 m = 2、獲得エージェント数 L = 10とする.

## 4.1 フィルタリング

#### I.L 体の安定走行を示すエージェントの獲得

直進走行を学習させ,安定走行を示すエージェントを10体 獲得した.図3(a)に示す自動車の軌跡は,獲得された10個の エージェントから1個を用いたものであり,青線が自動車の走 行軌跡,両端の黒線が道路端,赤線の内側がtrackPosによるペ ナルティが0になる範囲を表している.図3(a)は,Algorithm2 で定義した安定走行ができたことを示す.

## Ⅱ. 勾配評価と評価値の更新

獲得したエージェントの勾配を求め,各入力次元の評価値 を更新する.図3(b)は400mの走行時における勾配の推移を 示しており,trackPos, angleの勾配が常に大きいことがわか る.図3(b)中の表はグラフの平均値をとったものである.こ れにより各入力次元に順位がつけられるため,ランク変数 R の各要素に各入力次元の順位が格納され,0で初期化した評価



表 3: 評価値 X の更新結果

X	入力次元	1回目	10 回目
$x_1$	trackPos	5	50
$x_2$	angle	4	39
$x_3$	speedX	2	19
$x_4$	speedY	1	22
$x_5$	speedZ	3	20

値 X が更新される.表3に,評価値 X の更新結果を示す.本 実験では,表3に示すとおりに各入力次元を X の要素に対応 させた.表3における1回目の列には,図3(b)による X の 更新結果が示されている.10回目の列は,獲得された10体の エージェントによる X の最終的な更新結果である.

#### Ⅲ. 走行テストに用いるエージェントの選択

はじめに、評価値更新の結果とm = 2より、trackPos, angle を影響度の大きい、重要視される入力次元として選択する.

続いて,2次元入力による走行テストを行うエージェントを 選択する.本実験では,10体のエージェントの中から,8体 が採用された.

#### **IV.** *M* を入力とした走行テスト

採用された8体のエージェントを用いて,trackPos, angle の2次元入力による走行テストを行う.図4に,図3のエー ジェントを用いたtrackPos, angleの2次元入力による自動車 の走行軌跡を示す.図4の青線が2次元入力時の自動車の走行 軌跡を表し,両端の黒線,赤線の内側は図3(a)と同様である. 濃桃の破線は,図3(a)に示した5次元入力時の走行軌跡であ る.図4の青線は,Algorithm2で定義した安定走行ができた ことを示している.さらに5次元入力時と比較して,より直線 的な走行をしていることがわかる.走行テストに用いた8個の エージェントは,すべて図4と同様の結果を示していた.これ らより,2次元入力で安定走行が可能であり,trackPos, angle は直進走行に必要な次元であることがわかった.したがって, 入力次元を5次元から2次元へフィルタリングできた.



#### 4.2 検証実験

4.1 より、5 次元入力から直進走行に不可欠な次元が抽出で きた.ここでは、次元を削減した際の学習性能について検証す る.必要な次元だけを入力した学習結果として安定走行が確認 できた場合、提案法により特定された入力次元だけで安定走行 が獲得可能であり、提案法の有効性が示せる.

図5に示すのは、trackPos, angleの2次元入力による直 進走行の学習で獲得されたエージェントを用いた自動車の走 行軌跡である.図5の青線で表される自動車の走行軌跡は、 Algorithm2に定義した安定走行ができたことを示している. また、図3(a)と比較すると、より直線的な走行であることが わかる.以上により、提案法により抽出された次元だけで安定 走行が獲得可能であることが示された.

#### 5. 考察

ここでは、フィルタリングされた trackPos, angle の 2 次 元入力による走行が 5 次元入力時よりも直線的であった理由 について考察する.

本実験では、提案法のフィルタリングによって直進走行の学 習に用いた5次元入力から不可欠な次元を抽出できた. 直進走 行を5次元入力、1次元出力によって学習したときの最適行動 は、tracPos, angle がともに0の状態で舵角 steering が0で 出力されることである. 深層強化学習ではノードの値と、ノー ド間の重みとの積によって値が伝播、出力されるため、重みが 大きいノードほど出力に大きく影響を与えると考えられる. し たがって、多次元入力における不要な次元にかかる重みがノイ ズのような役割をしており、本実験では、5次元入力によるノ イズを除くことでより直線的な走行、すなわち最適行動に近い 挙動を示したと考えられる.

## 6. まとめ

本研究では、自動運転に対して深層強化学習を適用する際の課題とされる行動妥当性を説明可能にするための第一歩として、行動出力に必要な環境情報の属性を抽出するためのフィルタリングを提案した.提案法では、深層強化学習を用いて直進走行を学習したエージェントの入力の勾配を用いた.提案法によって抽出された、勾配が大きい入力次元だけを入力としても安定走行が確認できたことから、提案法によるフィルタリングの有効性を示した.

#### 参考文献

- [Mnih 13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, Martin Riedmiller: Playing Atari with Deep Reinforcement Learning, arXiv:1312.5602 [cs.LG], (2013)
- [Lillicrap 15] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, Daan Wierstra: Continuous control with deep reinforcement learning, arXiv:1509.02971 [cs.LG], (2015)
- [Nasukawa 08] 茄子川捷久, 宮下義孝, 汐川満: 自動車の走行 性能と試験法, 東京電機大学出版局, (2008)
- [Loiacono 13] Daniele Loiacono, Luigi Cardamone, Pier L. Lanzi: Simulated Car Racing Championship Competition Software Manual, arXiv:1304.1672v2 [cs.AI], (2013)