グラフ畳み込み構造を持つ予測状態表現を用いた 深層強化学習による移動ロボットナビゲーション

○松本耕平(九州大学) 河村晃宏(九州大学) 安琪(東京大学) 倉爪亮(九州大学)

1. はじめに

生活環境で活動するサービスロボットの実現には、歩 行者などが存在する動的な環境での自律移動が不可欠 である.歩行者の行動は様々な影響によって変化する 可能性があり、ロボットの行動にも影響を受ける可能性 がある.本研究で提案する自律移動ロボットナビゲー ションのための深層強化学習手法は、エージェントの 行動後の環境の変化を予測できる予測状態表現(PSR) に基づいて、ロボットの行動の影響による歩行者の行 動の変化に対応する. さらに、従来の PSR モデルでは 考慮できない、歩行者同士の影響による歩行者の行動 の変化に対応するために、グラフ畳み込み構造を用い た新しい PSR モデルを提案し. 従来手法と比較して有 効であることをシミュレーションによって示す. さら に、歩行者数の変化に対応するために、PSR の状態を 統合する2つの手法を提案し、この2つの手法の性能 を比較する.

2. Predictive State Representation

予測状態表現 (Predictive State Representation : PSR)[1] は、考え得る全てのパターンの観測・行動から 予想される結果が全て分かっていれば、動的システム を完全に把握できているという考え方に基づいている.

基本的な PSR モデルは離散的な観測・行動からなる システムのみに適用できる手法であるが,これまでに, PSR を連続的なシステムに対応できるようにした手法 が提案されている [2, 3].本稿では,これらの手法を総称して Recurrent PSR (RPSR) と呼ぶ.

RPSR の状態の更新は2つの手順で行われる.

• Extension:状態 q_t に線形写像 W_{ext} を適用し,拡張状態 p_t を得る.拡張状態 p_t は,拡張された行動 $a_{t:t+k}$ によって条件付けられた拡張された観測 $o_{t:t+k}$ の条件付き分布である.また, W_{ext} は学習 によって最適化されるパラメータである.

$$e_t = W_{\text{ext}} q_t \tag{1}$$

Conditioning: 時刻 *t* における行動 *a_t* と観測 *o_t* に, 既知の条件付け関数 *f_{cond}* により,以下のように状態が更新される.

$$q_{t+1} = f_{\text{cond}}\left(e_t, a_t, o_t\right) \tag{2}$$

q_t と *e_t* は条件付き確率であり, *f*_{cond} はベイズ則を適 用する.連続的なシステムに応用するために,これら に対して,分布のヒルベルト空間埋め込み [2] とカーネ ルベイズ則 [4] を用いる.

本研究では、RPSR のモデルとして RFF-PSR [3] を 用いる. RFF-PSR では、観測と行動のデータに RBF カーネルによる写像を適用後、ランダムフーリエ特徴 [5] を抽出し、ランダム主成分分析 [6] を用いて次元削



図1 歩行者による動的環境における移動ロボットナビ ゲーションに対する PSR のモデリングの概念図

減したものを、それぞれ観測と行動データの特徴量として用いる.この特徴量を抽出する関数を ϕ で表す.これを用いて、観測予測関数 f_{pred} によって以下のように時刻 tにおける観測が予測される.

$$\hat{o} = f_{\text{pred}} \left(q_t, \phi \left(a_t \right) \right)$$
$$= W_{\text{pred}} \left(q_t \otimes \phi \left(a_t \right) \right) \tag{3}$$

ここで, *W*_{pred} は学習で最適化される線形写像であり, ⊗ はクロネッカー積を表す.

3. 提案手法

3.1 歩行者による動的環境における移動ロボットナ ビゲーションへの PSR の適用

本研究では、PSR に基づく深層強化学習手法を, 複数の歩行者が存在する動的環境における移動ロボット ナビゲーションに適用した.移動ロボットのナビゲー ションにおける PSR の要素は,以下のように構成する ことができる.

- 観測:歩行者とロボットの位置と速度データを観 測として扱う.各歩行者,ロボットの観測はベク トル $(p_x^i, p_y^i, v_x^i, v_y^i)$ であり,i番目の歩行者または ロボットにおいて, (p_x^i, p_y^i) は位置, (v_x^i, v_y^i) は速 度を表す.
- 行動:本研究では、ホロノミックな全方位移動ロボットを想定し、2次元空間におけるロボットのx軸方向の入力速度 v_x と y 軸方向の入力速度 v_y からなる 2 次元ベクトル (v_x, v_y)を用いる.

これによって構成される PSR は、ロボットの指令速度 を入力することで、歩行者の位置と次のステップの速 度を予測することができるモデルとなる. PSR のモデ ルの概念図を図1に示す.

3.2 提案手法のアーキテクチャ

提案手法のアーキテクチャを図2に示す.提案手法 はFeature extractor, State updater, Observation predictor, State integrator, Value estimator により構成 される.

• Feature extractor:入力情報から特徴を抽出する.



図2 提案手法のアーキテクチャ

- State updater:観測と行動から抽出された特徴量 を用いて、状態を更新する.
- Observation predictor:行動から抽出された状態 特徴量を用いて、次の観測を予測する.
- State integrator:各歩行者に対応する PSR の状態を統合する.
- Value estimator:統合された状態とロボットの特徴から、価値を推定する.

3.3 Graph convolutional RPSR



図3 GC-RPSR の状態更新の概念図

RPSR モデルは、エージェントの行動が環境に与え る影響を考慮することができ、これにより、ロボット の行動が各歩行者に与える影響を考慮することができ る.しかし、エージェントの行動に起因しない影響は 考慮されていないため、歩行者同士の影響を考慮する ことはできない.この問題に対処するために、Graph Convolutional RPSR (GC-RPSR)を提案する.

本アーキテクチャは、RPSR の状態更新時において、 歩行者間の相互作用を表現するためのグラフ畳み込み 構造を有している. この GC-RPSR の状態更新処理の 概念図を図 3 に示す. グラフ畳み込みに用いる隣接行 列 A_t のカーネルは、式 (4) に示すように、各位置の差 分の L_2 ノルムの逆数に基づく [7].

$$a_t^{ij} = \begin{cases} 1/\|p_t^i - p_t^j\|_2 & \text{if } \|p_t^i - p_t^j\|_2 \neq 0\\ 0 & \text{Otherwise} \end{cases}$$
(4)

さらに,式 (5) に示す正規化処理を行い \tilde{A}_t を実際の 計算に用いる.

$$\tilde{A}_t = \Lambda_t^{-\frac{1}{2}} \hat{A}_t \Lambda_t^{-\frac{1}{2}} \tag{5}$$

ここで, $\hat{A}_t = A_t + I$ であり, Λ_t は次数行列を表す. グ ラフ畳み込みと式 (1), (2) を用いることで, GC-RPSR の状態更新は以下のように表される.

$$\tilde{q}_t = f_{gc} \left(q_t, \tilde{A}_t \right) \tag{6}$$

$$e_t = W_{\text{ext}}\tilde{q}_t \tag{7}$$

$$q_{t+1} = f_{\text{cnd}}\left(e_t, a_t, o_t\right) \tag{8}$$

ここで, f_{qc} はグラフ畳み込み関数を表す.

3.4 状態統合手法

各歩行者に対応する状態から,歩行者数の変化に対 して,スケーラブルに価値評価を行うための統合処理 が必要である.これに対して.グラフの畳み込みを用 いる手法と占有地図に基づく手法の2種類を提案する.

3.4.1 グラフ畳み込みによる統合

各歩行者に対応する PSR の状態を統合するために, グラフ畳み込みを用いる. Value estimator はグラフ畳 み込みが適用されたロボット特徴量のみを受け取る. 隣接行列に用いるカーネルは,式(4)で表される GC-RPSR の状態更新と同じである.この手法の概念図を 図4に示す.



図4 グラフ畳み込みによる状態統合の概念図

3.4.2 占有地図による統合

各歩行者に対応する PSR の状態を統合するために, 占有地図を使用する.歩行者の状態は,歩行者の位置 に応じて対応するセルに格納される.同じセルに複数 の歩行者の状態を格納する必要がある場合は平均値を 格納する.この手法の概念図を図5に示す.



図5 占有地図による状態統合の概念図

3.5 行動生成

行動は、Value estimator f_v , PSR モデル f_p の State updater f_p^q , Observation predictor f_p^o を用いて、得ら れる報酬が最大になるように、行動空間 \mathcal{A} から式 (9) に従って選択される.ここで、 γ は割引率、 s_t は時刻 tにおける統合された状態、 p_t は時刻 t におけるロボッ トの位置を表し、 $\hat{o}_t = f_p^o(q_t, a_t)$ である.

$$a_t \leftarrow \operatorname{argmax}_{a_t \in \mathcal{A}} R\left(\hat{o}_t\right) + \gamma^{\Delta t} f_v\left(s_t, p_t\right) \tag{9}$$

さらに, $R(o_t)$ は時刻 t における報酬関数であり,式 (10) で表される. d_t はロボットと歩行者間の最小距離 を表し, p_q はロボットの目標位置を示す.

$$R(o_t) = \begin{cases} -0.25 & \text{if } d_t < 0\\ -0.1 + d_t/2 & \text{else if } d_t < 0.2\\ 1 & \text{else if } p_t = p_g\\ 0 & \text{otherwise} \end{cases}$$
(10)

さらに, d-step planning[8] を利用し, 式 (11) に従っ

て、 d ステップ先を考慮して行動を生成する.

$$\begin{aligned}
f_{v}^{d}(s_{t}, p_{t}) &= \\
\begin{cases}
f_{v}(s_{t}, p_{t}) & \text{if } d = 1 \\
\frac{1}{d}f_{v}^{1}(s_{t}, p_{t}) + \frac{d-1}{m}\max_{a_{t}}(& (11) \\
R(\hat{o}_{t+1}) + \gamma f_{v}^{d-1}(\hat{s}_{t+1}, \hat{p}_{t+1})) & \text{otherwise}
\end{aligned}$$

3.6 学習手法

学習は2つのステップで行われる. 最初のステップ では, PSR パラメータを事前学習する. この事前学習 は, ORCA [9] に従った探索方針でデータを収集した 後, two-stage regression [10] によって行われる. その 後, 収集したデータを用いてアルゴリズム1に従い, 強 化学習によって Value estimator の学習させ, 教師あり 学習で PSR モデルを学習させる.

Algorithm 1: 提案手法の学習

```
PSR f_p と Value estimator f_v を two-stage regression と模倣
 学習によって事前学習
Value estimator のターゲットネットワーク \hat{f}_v を初期化
for i = 1 to E do
     探索方策に従い at を選択し,報酬 rt, 観測 ot, ロボットの位
     置 p_t を取得
エピソード終了後,軌跡 (o_t, a_t, r_t, p_t) をバッファ B に保存
     バッファ B から M セットの軌跡をサンプルする
     for j = 1 to M do
           以下の軌跡を取得する
                 観測 \mathbf{o}^j = \left\{ o_1^j, o_2^j, \dots, o_T^j \right\},
                 ロボットの位置 \mathbf{p}^{j} = \left\{ p_{1}^{j}, p_{2}^{j}, \dots, p_{T}^{j} \right\},
                行動 \mathbf{a}^j = \left\{ a_1^j, a_2^j, \dots, a_T^j \right\}
           以下の軌跡を計算する
                 状態 \mathbf{q}^j = \left\{ q_1^j, q_2^j, \dots, q_T^j \right\},
                 価値の目標値 \mathbf{y}^j = \left\{ y_1^j, y_2^j, \dots, y_T^j \right\},
                 統合された状態 \mathbf{s}^j = \left\{ s_1^j, s_2^j, \dots, s_T^j \right\}
     end
     L_{\text{pred}} = \text{MSE}(f_p^o(\mathbf{q}, \mathbf{a}), \mathbf{o})を最小化するように f_pを更新
     L_{\text{value}} = \text{MSE}(f_v(\mathbf{s}, \mathbf{p}), \mathbf{y})を最小化するように f_vを更新
     if i \mod d then
      \hat{f}_v \leftarrow f_vによって\hat{f}_vを更新
     \mathbf{end}
```

end

4. シミュレーション実験

4.1 シミュレーション環境

本研究では、CrowdNav 環境での circle crossing シ ナリオを利用する [11, 12]. このシナリオでは、ロボッ トは初期位置 (x, y) = (0, -4) からゴール地点 (x, y) =(0, 4) を目指して進む. 歩行者は ORCA 従って行動し、 初期化時に半径 4m の円上にランダムに配置される. ロ ボットと歩行者は互いに影響し合い、位置や速度によっ て挙動が変化する. 評価は 500 パターンのテストケー スを用いて行う.

4.2 GC-RPSR と RPSR の比較

GC-RPSR と RPSR を比較することで,その提案手法の有効性を確認した.成功率,衝突率,実行時間,平均収益を比較する.この実験では State integrator には単純な連結を使用した.結果を表1に示す.

表1 RPSR と GC-RPSR の比較

	成功率 [%]	衝突率 [%]	実行時間 [s]	、 平均収益
J 124	/x·/J [/0]	四八十 [70]		1
RPSR	57.8	33.8	8.84 ± 1.32	0.258
GC-RPSR	64.0	11.0	8.73 ± 1.31	0.356

この表から, GC-RPSR はすべての評価項目で RPSR を上回っていることがわかる.したがって,提案する PSR モデルは基本的なモデルである RPSR と比較し て,移動ロボットナビゲーションタスクにおいて優れ ていると言える.

4.3 ベースラインと状態統合手法を組み合わせた提 案手法の比較

図6にグラフ畳み込みを用いた状態統合を行った提 案モデル(GGC-RPSR)の結果の軌跡と,占有地図を 用いた状態統合を行った提案モデル(OGC-RPSR)の 結果の軌跡を示す.各図において,黒い軌跡はロボッ ト,その他の色の軌跡は歩行者,数字は時間ステップ を表している.どちらの手法も,歩行者を避けながら ロボットをゴールまで誘導できることが分かる.しか し,GGC-RPSRの方がより短い経路を生成できるた め,効率が良いと言える.



図 6 GGC-RPSR により生成される軌跡 (左)と OGC-RPSR により生成される軌跡 (右)

また, 関連研究で提案された手法 [12] をベースライン として提案モデルとの比較を行う.提案する 2 つの手 法とベースライン手法の数値比較を表 2 に示す.GGC-RPSR はすべての評価項目で他の手法を上回っている. この結果は,GC-RPSR と,グラフ畳み込みを用いた 状態統合を用いた提案手法の有効性を示している.こ れは,PSR の構造がより最適な経路計画を実現するた めに有効であったためであると考えられる.また,グ ラフ畳み込みを用いた状態統合手法の方が最適化性能 が高いため,学習時とテスト時の歩行者数が同じ本実 験では,GGC-RPSR の方が OGC-RPSR より性能が 高かったと推察される.

表 2 提案手法とベースラインの比較						
手法	成功率 [%]	衝突率 [%]	実行時間 [s]	平均収益		
RGL [12]	92.0	7.0	9.09 ± 1.31	0.560		
GGC-RPSR	94.6	5.4	7.79 ± 0.10	0.587		
OGC-RPSR	91.4	6.8	11.54 ± 1.03	0.512		

4.4 歩行者数が学習時とテスト時で異なる場合の比較

歩行者の数が学習時とテスト時で異なる状況で,2つ の提案手法を比較した.両モデルとも歩行者数が5人 の環境で学習し,歩行者数が1~10人の環境でテストを 行った.各モデルの歩行者5人の場合での結果との成 功率の差を図7に示す.GGC-RPSRの成功率の最大 低下率は43.4%であり,これに対しOGC-RPSRの最 大低下率は7.4%である.この結果から,OGC-RPSR は学習時とテスト時で歩行者数を変化させた場合によ り安定していることがわかる.

両者を比較すると、GGC-RPSR は歩行者数によら ず、各歩行者に対応する状態を1つの特徴量に統合し ている.そのため、未学習の歩行者数ではうまく機能 しない可能性があると推察される.一方、OGC-RPSR は、各セルの歩行者に対応する状態を統合しているた め、歩行者の数に関わらず、各セルに統合される状態 の数が一定になることが多いのではないかと考えらる. そのため、OGC-RPSR の方が安定した性能を発揮し たと考察する.



図7 歩行者が5人の場合での結果との成功率の差の比較

4.5 歩行者数を変えながらの学習

4.4 節で不安定な結果であった GGC-RPSR におい て、歩行者数を 1~10 人で変化させながら学習を行い、 歩行者数 5 人で学習した場合と成功率を比較した.結 果を図 8 に示す.この結果より歩行者数を 1~10 人で 変化させながら学習を行った場合、成功率の最大値が 下がるものの、各歩行者数においての成功率の差は小 さくなっていると言える.各歩行者数での成功率の平 均及び標準偏差を見ると、歩行者数 5 人で学習した場 合 78.2 ± 15.3%であるが、歩行者数を変化させながら 学習を行った場合 85.7±7.7%であり、より安定した結 果になっていることが分かる.



図 8 歩行者数を変化させながら学習を行った場合と歩 行者数 5 人で学習した場合の成功率の比較

5. 結論

本研究では、歩行者による動的環境における移動ロ ボットナビゲーション課題において、PSR に基づく深 層強化学習法を用いた.グラフ畳み込みを利用するこ とで歩行者間の相互作用を考慮した、新しい PSR の構 造を提案した.また、学習時とテスト時の歩行者数の 変化に対応するため、各歩行者に対応した PSR の状態 を統合する方法を提案し、シミュレーション実験によ り提案モデルの有効性を確認した.また、訓練時とテ スト時で歩行者数が異なるシナリオでは、占有地図を 用いた状態統合を行う OGC-RPSR がより安定してい ることを確認した.さらに、GGC-RPSR については、 歩行者数を変化させながらの学習によって、安定性が 増すことを確認した. 今後は、本研究で提案した2種類の状態統合手法の より詳細な解析を行い、本研究の考察が妥当であるか を確認する.加えて、占有地図ベースの手法について も、4.5節と同様の条件で学習を行い、今回の結果との 比較を行う.さらに、本研究の実験は全てシミュレー ション環境で行ったが、今後、実際のロボットへの適 用も進める予定である.

参考文献

- S. Singh, M. James, and M. Rudary, "Predictive State Representations: A New Theory for Modeling Dynamical Systems," in *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 512–519, 2004.
- [2] B. Boots, A. Gretton, and G. J. Gordon, "Hilbert space embeddings of predictive state representations," in *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 92–101, 2013.
- [3] A. Hefny, C. Downey, and G. Gordon, "An efficient, expressive and local minima-free method for learning controlled dynamical systems," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pp. 3191–3198, 2018.
- [4] K. Fukumizu, L. Song, and A. Gretton, "Kernel Bayes' rule: Bayesian inference with positive definite kernels," *Journal of Machine Learning Research*, vol. 14, pp. 3753–3783, 2013.
- [5] "Random features for large-scale kernel machines," in Advances in Neural Information Processing Systems, 2008.
- [6] N. Halko, P. G. Martinsson, and J. A. Tropp, "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions," *SIAM Review*, vol. 53, no. 2, pp. 217–288, 2011.
- [7] A. Mohamed, K. Qian, M. Elhoseiny, and C. Claudel, "Social-STGCNN: A Social Spatio-Temporal Graph Convolutional Neural Network for Human Trajectory Prediction," in *Proceedings of the IEEE/CVF Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14412–14420, 2020.
- [8] J. Oh, S. Singh, and H. Lee, "Value prediction network," in Advances in Neural Information Processing Systems (NeurIPS), pp. 6119–6129, 2017.
- [9] J. Van Den Berg, S. J. Guy, M. Lin, and D. Manocha, "Reciprocal n-body collision avoidance," in *Proceed*ings of the International Symposium of Robotic Research, pp. 3–19, 2011.
- [10] A. Hefny, C. Downey, and G. J. Gordon, "Supervised learning for dynamical system learning," in Advances in Neural Information Processing Systems (NeurIPS), pp. 1963–1971, 2015.
- [11] C. Chen, Y. Liu, S. Kreiss, and A. Alahi, "Crowdrobot interaction: Crowd-aware robot navigation with attention-based deep reinforcement learning," in *Proceedings of the IEEE International Conference* on Robotics and Automation (ICRA), pp. 6015–6022, 2019.
- [12] C. Chen, S. Hu, P. Nikdel, G. Mori, and M. Savva, "Relational graph learning for crowd navigation," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10007–10013, 2020.