

ファジィ推論と方策勾配法とを融合した制御方式における逆強化学習の適用 Application of Inverse Reinforcement Learning in a Control System Integrating Fuzzy Reasoning and Policy Gradient Reinforcement Learning

小嶋 怜央[†]
Reo Kojima

石原 聖司[†]
Seiji Ishihara

1. はじめに

ファジィ推論には推論過程の理解がしやすく、計算コストが少ないという利点がある。だが、ファジィ推論を行うにはルール重みやメンバシップ関数を事前知識に基づき定義する必要がある。これに対してファジィ推論と方策勾配法とを融合した制御方式の学習則[1]が提案されており、ルール重みやメンバシップ関数が方策勾配法による学習によって獲得可能であることが自動車の速度制御問題において実験的に示されている[2][3]。ただし、当該の制御モデルの学習を行うには報酬を定義する必要がある。これに対して、逆強化学習を当該の制御モデルに適用することにより報酬関数を自動獲得しながらメンバシップ関数の学習が可能なが同様の問題において実験的に示されている[4]。ただし、文献[4]の実験では、使用するルールを先験的な知識に基づいて限定していた。また、逆強化学習を当該の制御モデルに適用した際のルール重みの学習は行われていなかった。

本研究では、使用するルールを限定せずに、文献[4]と同様の学習方法によって、ルール重みを固定した状態でのメンバシップ関数の学習結果と自動獲得された報酬関数およびメンバシップ関数を固定した状態でのルール重みの学習結果と自動獲得された報酬関数についてそれぞれ報告する。

2. 本研究での融合方式

エージェントの行動制御のファジィ制御ルールは次のように表すことができる。

Rule i :

$$\begin{aligned} &\text{if } (x_1 \text{ is } A_1^i) \text{ and } \dots \text{ and } (x_M \text{ is } A_M^i) \\ &\text{then } (y_1 \text{ is } B_1^i) \text{ and } \dots \text{ and } (y_N \text{ is } B_N^i) \text{ with } \theta_i \end{aligned} \quad (1)$$

ここで、 $i (= 1, 2, \dots, n_R)$ はルール番号、 $x_j (j = 1, 2, \dots, M)$ は状態を記述する入力変数、 $y_k (k = 1, 2, \dots, N)$ は離散的な行動を記述する出力変数である。 A_j^i / B_k^i はルール i の j / k 番目の前件部/後件部に含まれる変数 x_j / y_k を記述したファジィ表現で、それらのメンバシップ関数を $A_j^i(x_j) / B_k^i(y_k)$ と表す。 $\theta_i (\geq 0)$ はルール i の重みを表すパラメータである。

エージェントの各時刻の行動決定のために、次の目的関数

$$E(y; x, \theta, A, B) = -\sum_{i=1}^{n_R} \theta_i A^i(x) B^i(y) \quad (2)$$

の最小化を規範とした。ただし、入力値 x / 出力値 y のルール i における前件部/後件部の適合度 $A^i(x) / B^i(y)$ を

$$A^i(x) \equiv \prod_{j=1}^M A_j^i(x_j) \quad (3)$$

$$B^i(y) \equiv \prod_{k=1}^N B_k^i(y_k) \quad (4)$$

で定義している。

時刻 t における入力値 $x(t)$ に対する出力値 $y(t)$ の決定をエージェントの行動決定、すなわち、方策とし、式(2)の目的関数を持つ Boltzmann 分布による確率として次のように定義している。

$$\pi(y(t); x(t), \theta, A, B) \equiv \frac{e^{-\frac{E(y(t); x(t), \theta, A, B)}{T}}}{\sum_y e^{-\frac{E(y; x(t), \theta, A, B)}{T}}} \quad (5)$$

ただし、 T は温度と呼ばれるパラメータで、この方策で選択される出力値 $y(t)$ における目的関数 $E(y(t))$ の期待値をコントロールすることができる。すなわち、確率的な方策に基づく行動選択のランダム性の大きさを表している。

制御時には式(5)を重みとする重心 $y_G(t)$ を次のように定義して用いる。

$$\begin{aligned} y_G(t) &\equiv \sum_y y \cdot \pi(y; x(t), \theta, A, B) \\ &= \langle y \rangle_{\pi(y)} \end{aligned} \quad (6)$$

ただし、 $\langle y \rangle_{\pi(y)}$ は $\pi(y)$ による y の期待値操作である。

方策勾配法により、エピソードごとの報酬期待値を極大化する、パラメータ μ に関する学習則は次のように与えられる。

$$\Delta\mu = \varepsilon \cdot \sum_{t=0}^{L-1} r(t) \cdot e_{\mu}(t) \quad (7)$$

ただし、 L はエピソード長、 $r(t)$ は報酬関数である。また、 ε は学習係数、 $e_{\mu}(t)$ は特徴適性度

$$e_{\mu}(t) = \frac{\partial}{\partial \mu} \log \pi(y(t); x(t), \mu) \quad (8)$$

$$= -\frac{1}{T} \left[\frac{\partial E(y(t); x(t), \mu)}{\partial \mu} - \left\langle \frac{\partial E(y; x(t), \mu)}{\partial \mu} \right\rangle_{\pi(y)} \right] \quad (9)$$

である。

3. 逆強化学習

本研究では、報酬関数 $r(t)$ を次のように定義する。

$$r(t) = \boldsymbol{\tau}^T \cdot \mathbf{f}_{(x(t), y(t))} \quad (10)$$

ただし、 $\boldsymbol{\tau} = (\tau_1, \dots, \tau_u)^T$ は報酬パラメータである。入力と出力の組を $(x(t), y(t)) \in F = \{1, \dots, u\}$ と u 通りに量子化

[†] 東京電機大学 Tokyo Denki University

し、特徴ベクトル $\mathbf{f}(x(t), y(t)) = (0, \dots, f_{(x(t), y(t))=h} = 1, \dots, 0)^T$ を u 次元の one-hot ベクトルで定義する。本研究では、報酬パラメータ \mathbf{r} を最大エントロピー逆強化学習によって、

$$\Delta \mathbf{r} = \alpha \sum_{t=0}^L (\mathbf{f}_{(x(t), y(t))}^* - \mathbf{f}(x(t), y(t))) \quad (11)$$

と更新する。ここで、 α は学習係数、 \mathbf{f}^* はエキスパートの特徴ベクトルである。つまり、 $f_h < f_h^*$ では τ_h は増加し、 $f_h > f_h^*$ では τ_h は減少する。これによって報酬関数の自動獲得を行う。

4. 自動車の速度制御問題への適用

一次元的な直線道路において、後方車両速度 x_2 の制御を行うことで、一定速度 v_0 で走行する前方車両との車間距離 x_1 を適切に保持する問題を考える。 $t = 0 \sim 80$ 秒までに目標車間範囲 $[l_0, l_1]$ に入り、 $t = 110$ 秒まで目標車間範囲内を保つことができればその問題は制御成功とする。ただし、制御を行う際のアクセルとブレーキの操作量 y は、 $y = -5.0 + 0.1h$ ($h = 0, 1, \dots, 100$) の 101 通りの離散値を使用する。また、学習時のエージェントの操作量 y の決定には (5) 式の方策 $\pi(y)$ を使用し、問題が制御成功か判定する際のエージェントの操作量は (6) 式の重心 y_G を使用する。

4.1 メンバシップ関数

車間距離が“長い” / “短い”という 2 種類のファジィ表現のメンバシップ関数 $A_1(x_1 | l_0, l_1)$ と後方車両速度が“速い” / “遅い”という 2 種類のファジィ表現のメンバシップ関数 $A_2(x_2 | v_0)$ の 4 種類の前件部のメンバシップ関数をニューラルネットワークで表現する。アクセルを“強く” / “弱く”踏む、ブレーキを“強く” / “弱く”踏む、“なにもしない”の 5 種類の後件部のメンバシップ関数 $B_1(y)$ も学習可能だが、本研究では簡単のため、先験知識に基づいて定義したものを使用する。

4.2 量子化

本研究では、車間距離 x_1 を“遠い”、“やや遠い”、“やや近い”、“近い”の 4 通り、後方車両速度 x_2 を“速い”、“やや速い”、“やや遅い”、“遅い”の 4 通り、そしてアクセルとブレーキの操作量 y を $-5.5 + g \leq y < -4.5 + g$ ($g = 0, 1, \dots, 10$) のように 11 通りに量子化する。よって全体では $u = 4 \times 4 \times 11 = 176$ 通りに量子化する。

5. メンバシップ関数の学習に対する実験

学習に使用する 16 問の問題と 20 通りのファジィ制御ルールは文献 [2] と同様のものを使用する。ルール重み θ を固定した状態で前件部のメンバシップ関数 A_1, A_2 の結合荷重を (7) 式のパラメータ μ として強化学習を行う。ルール重みは先験知識に基づく固定値を使用し、既定の更新回数毎に (11) 式より報酬関数の更新を行う。全 16 問で制御成功するまで前件部のメンバシップ関数 A_1, A_2 の結合荷重の強化学習と報酬関数の更新を繰り返す。

6. ルール重みの学習に対する実験

学習に使用する 16 問の問題と 20 通りのファジィ制御ルールは文献 [2] と同様のものを使用する。前件部のメンバシップ関数 A_1, A_2 を固定した状態でルール重みを表すパラ

メータ θ を (7) 式のパラメータ μ として強化学習を行う。前件部のメンバシップ関数 A_1, A_2 は先験知識に基づいて定義したものを使用し、既定の更新回数毎に (11) 式より報酬関数の更新を行う。全 16 問で制御成功するまでルール重みを表すパラメータ θ の強化学習と報酬関数の更新を繰り返す。

7. まとめ

ファジィ推論と方策勾配法とを融合した制御方式へ逆強化学習を適用することによって報酬関数の自動獲得機能を追加し、メンバシップ関数の学習とルール重みの学習の 2 種類の実験を行った。ルール重みを固定した状態での前件部のメンバシップ関数の学習結果と自動獲得された報酬関数、前件部のメンバシップ関数を固定した状態でのルール重みの学習結果と自動獲得された報酬関数についてはそれぞれ口頭発表時に報告する。

参考文献

- [1] H. Igarashi and S. Ishihara, “An Algorithm of Policy Gradient Reinforcement Learning with a Fuzzy Controller in Policies”, *International Journal of Artificial Intelligence and Expert Systems*, Vol. 4, pp. 17-26, 2013.
- [2] 石原 聖司, 五十嵐 治一, “ファジィ制御ルールにより表現された方策を持つ方策勾配法：自動車の速度制御問題への適用”, *知能と情報*, Vol. 32, No.4, pp. 801-810, 2020.
- [3] 市毛 竣, 五十嵐 治一, 石原 聖司, “ファジィ制御と強化学習の融合 メンバシップ関数とルール重みの学習”, *人工知能学会全国大会論文集*, pp. 1D1GS203, 2022.
- [4] 小嶋 怜央, 石原 聖司, “ファジィ推論と方策勾配法とを融合した制御方式における逆強化学習による報酬関数の自動獲得”, *電子情報通信学会総合大会予稿集*, pp.148, 2024.