

# 連続空間における信頼度予測を用いた強化学習

蔡 詩祐<sup>†</sup>

平原 誠<sup>†</sup>

<sup>†</sup> 法政大学大学院 理工学研究科

## 1. はじめに

強化学習を扱う上での課題の一つに、価値関数に基づく行動選択と、価値関数の推定値を改良する探索的行動とのバランスが挙げられる[1].

学習の進捗が進むにつれて探索確率は下げていくのが理想であるが、問題に応じて学習進捗の推移が異なるため、探索確率の適切な変更スケジュールを事前に設定するのは困難である。この問題に対して、“信頼度”という指標に基づいて探索確率を制御する提案がされている [2]. 本研究ではこの信頼度による探索確率の制御を連続空間における時系列課題に適用する。

## 2. 先行研究

### 2.1 離散空間上での信頼度による探索確率の制御[2]

離散空間における  $Q$  学習の一般式は

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha \delta_t \quad (1)$$

$$\delta_t = r_t + \gamma \max_a \{Q_t(s_{t+1}, a)\} - Q_t(s_t, a) \quad (2)$$

で表される。式(2)は  $TD$  誤差  $\delta_t$  と呼ばれる。

状態  $s$  に対する信頼度の逆数  $R^2(s)$  は

$$R_{t+1}^2(s_t) = R_t^2(s_t) + \alpha_R (\delta_t^2 + \gamma_R R_t^2(s_{t+1}) - R_t^2(s_t)) \quad (3)$$

で更新される。ここで  $\alpha_R, \gamma_R$  は 0 以上 1 未満の係数である。 $R^2$  値は、二乗  $TD$  誤差  $\delta_t^2$  が大きい場合に増加し、小さい場合は減少することになる。

この  $R^2$  を用いて、状態  $s$  で行動  $a$  を選択する確率  $P(s, a)$  を

$$P(s, a) = \frac{\exp\left(\eta \frac{Q(s, a)}{R(s)}\right)}{\sum_a \exp\left(\eta \frac{Q(s, a)}{R(s)}\right)} \quad (4)$$

により定める。softmax 法の温度パラメータ部分を  $R$  とすることで、学習の進捗に応じて  $P(s, a)$  のランダム性を調整することを表している。

### 2.2 連続空間における時系列課題での強化学習[3]

連続空間における時系列課題に対して、強化学習を可能にするため、Elman 型リカレントネット(図 1)を用いて  $Q$  学習の価値関数近似を行う手法が提案されている[3]. 行動の数だけ出力素子を設け、各素子は入力(状態  $s_t$ )に対して行動価値関数  $Q^N$  を出力する。時刻  $t$  で選択した行動を  $a_t$  と表すと、その行動価値関数  $Q_{a_t}^N(s_t)$  に対する教師

信号  $T_{a_t, t}$  は

$$T_{a_t, t} = r_{t+1} + \gamma \max_a Q_a^N(s_{t+1}) \quad (5)$$

で与えられる。ここで  $\gamma$  は 0 以上 1 未満の係数、 $r$  は報酬である。この教師信号を用いて BPTT 法により学習を行う。

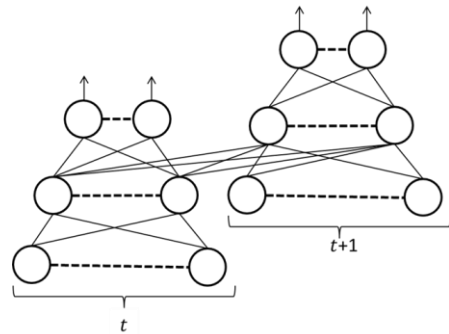


図 1 Elman 型リカレントネット

## 3. 提案手法

信頼度による探索確率の制御を連続空間上で実現するため、 $R^2$  値を、単出力の Elman 型リカレントネットで近似する。出力を  $R^2$  値とし、 $R^2$  値に対する教師信号  $T^R$  を

$$T_{s_t}^R = (\delta_t^N)^2 + \gamma_R R_t^2(s_{t+1}) \quad (6)$$

$$\delta_t^N = T_{a_t, t} - Q_{a_t}^N(s_t) \quad (7)$$

により定める。 $(\delta_t^N)^2$  は、行動価値関数  $Q^N$  とその教師信号  $T$  の二乗誤差である。 $(\delta_t^N)^2$  が大きいと信頼度のネットワークの出力である  $R^2$  値は増加、小さい場合は減少する方向へ学習することになる。 $R$  を式(4)の温度パラメータとして用いることで、連続空間上で信頼度による探索確率の制御が可能であると考えられる。

## 参考文献

- [1] R.S. Sutton, and A.G. Barto, Reinforcement Learning: An Introduction, The MIT Press, 1998.
- [2] 阪口豊, 高野光雄, “内部モデルの信頼度に基づく強化学習のアルゴリズム,” 日本神経回路学会第 11 回全国大会予稿集, 2001.
- [3] Kenta Goto, and Katunari Shibata, “Emergence of Prediction by Reinforcement Learning Using a Recurrent Neural Network,” Journal of Robotics, Volume 2010.