

## UCT を用いた訓練初期局面の多様化による TD 学習法の改善\*

三木 理斗<sup>†</sup> 三輪 誠<sup>‡</sup> 田浦 健次朗<sup>¶</sup> 近山 隆<sup>‡</sup><sup>†</sup> 東京大学工学部 <sup>‡</sup> 東京大学大学院新領域創成科学研究科 <sup>¶</sup> 東京大学大学院情報理工学系研究科

## 1 はじめに

強化学習によるコンピュータゲームプレイヤーの評価関数の学習法として、TD( $\lambda$ ) 法 [1] はさまざまなゲームで注目すべき成果を挙げている。対戦の経験から学習を行う場合には、なるべく多様な局面を展開することが望ましい。本研究では TD( $\lambda$ ) 法においてエージェントの手の選択を多様化する方法として、序盤数手を UCT 法によって決定する手法を提案する。実験としてブロックデュオの評価関数の重みを学習し、従来の  $\epsilon$ -greedy 方策や完全にランダムに序盤の手を決める手法と比較して評価し、提案手法が少ない対局数でよりよい重みを学習することを確認した。

2 TD( $\lambda$ ) 学習法2.1 ゲームにおける TD( $\lambda$ ) 法の概要

ゲームで多く用いられる評価関数は局面  $s$  のさまざまな特徴（駒の配置、持ち駒の種類など）を評価要素  $\vec{\phi}(s)$  とし、その各要素の重みつき線形和で表した次のような形をしている。

$$V(s) = \vec{\theta} \cdot \vec{\phi}(s)$$

TD( $\lambda$ ) 法では各局面  $s$  で任意の方策にしたがって手を決めて次の局面  $s'$  を決定し、誤差  $\delta$  を用いて重みを更新する。

$$\delta = r + \gamma V(s') - V(s)$$

$$\vec{e} = \gamma \lambda \vec{e} + \nabla_{\vec{\theta}} V(s)$$

$$\vec{\theta} = \vec{\theta} + \alpha \delta \vec{e}$$

$r$  は対局終了時のみ与えられる報酬で、勝ちなら 1、負けなら 0、引き分けなら 0.5 とする。 $\gamma$  は割引率と呼ばれるパラメータである。 $\vec{e}$  は適格度トレースと呼ばれる

るもので各対局の最初に零ベクトルに初期化する。 $\alpha$  は学習率と呼ばれるパラメータである。

## 2.2 学習局面の多様化

チェスのような確定ゲームで、かつ棋譜やネットワーク上の対戦サーバを用いずに学習を行う場合は、意図的に選択する手をばらつかせ、ある程度多様な局面を学習するような方策を用いる必要がある。

そのために従来用いられてきた手法に  $\epsilon$ -greedy 方策がある。これは各局面において通常はもっとも評価値の高い手を選択するが、ある確率  $\epsilon$  でランダムな手を選択する方法である。この方法では選択する手が完全にランダムであるため明らかな悪手を打つことが多いことと、ランダムな手を打った前後で TD( $\lambda$ ) におけるバックアップのつながりが途絶えてしまうという問題点がある。

## 3 提案手法

本研究では序盤 UCT 方策という手法を提案する。この手法では序盤の数手を UCT 法 [2] によって決定する。UCT 法はゲーム木を探索する手法のひとつで、ランダムシミュレーションで手を決定するモンテカルロ法に最良優先探索を組み合わせるより有望な部分木を優先的に展開するように改良した手法である。これを用いた囲碁のプログラムは世界トップクラスの成績を収めている [3]。UCT 法を用いることでランダムよりも良い手が選ばれることが期待できる。また、序盤にのみ指し手の変化をつけるため、TD( $\lambda$ ) のバックアップが切れ切れになりすぎないのでより高い精度の学習ができると思われる。

## 4 評価

評価のための実験として、ブロックデュオ (<http://www.blokus.com/>) の評価関数の学習を行った。ブロックデュオは陣取り形式でランダム要素のない二人対戦ゲームである。

$\epsilon$ -greedy 方策、序盤 UCT 方策と、比較のため序盤数手に完全ランダムに手を決定する方策 (以下、序盤

\*Improving Temporal-Difference Learning Using a Variety of Opening Moves Obtained by UCT

Ayato Miki<sup>†</sup>, Makoto Miwa<sup>‡</sup>, Kenjiro Taura<sup>¶</sup>  
and Takashi Chikayama<sup>‡</sup>

<sup>†</sup> Faculty of Engineering, The University of Tokyo

<sup>‡</sup> Graduate School of Frontier Sciences, The University of Tokyo

<sup>¶</sup> Graduate School of Information Science and Technology,  
The University of Tokyo

ランダム方策)の3方策についてそれぞれ自己対戦による学習を行った。評価要素は2007年度のGPCCのブロックデュオ大会2位のプログラムと同じものを用い、それぞれの重みの初期値は0に設定した。 $\lambda$ は0.8、割引率 $\gamma$ は1とし、学習率 $\alpha$ は次式にしたがって減衰させた[4]。

$$\alpha = \alpha_0 \frac{N_0 + 1}{N_0 + \text{Episode\#}^{1.1}}$$

今回は $\alpha_0$ は0.01、 $N_0$ は100とした。Episode#は対局回数である。 $\epsilon$ -greedy方策は1000000対局、序盤UCT方策および序盤ランダム方策は100対局の学習を行った。序盤にUCTやランダムを用いる深さは4手目までとし、UCTのシミュレーションは20000回行った。

図1に代表的な評価要素の最初の100対局での重みの変化を示す。序盤UCT方策と序盤ランダム方策は重みの学習の様子がみられるが、 $\epsilon$ -greedy方策では重みにほとんど変化がない。200対局あたりまで重みが変化していない状態が続いたが、その後変化を始め、2000対局ほどで他の方策と同じ重みに到達した。これは終局の形が大きく変わるほどの手の変化が発生しなかったためと考えられる。序盤ランダム方策も100対局で序盤UCT方策とほぼ同じ重みを学習しているが、最初の20局ほどまではやや不安定な変動が見られる。

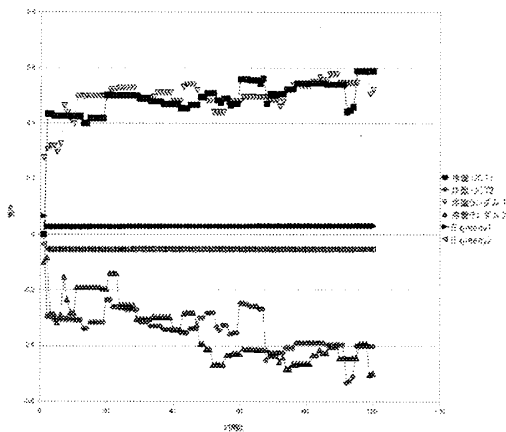


図1: 重みの変化

## 5 おわりに

序盤UCT方策を用いることで、完全にランダムな従来の手法よりも安定した学習を行うことができる。今回実装したUCT方策による学習では、学習中にUCT法を用いてシミュレーションを行うため手の決定に非常に時間がかかっており、学習時間は $\epsilon$ -greedy方策とそれほど変わらない。しかしあらかじめUCT法のシミュレーションを行って序盤の手をデータベース化しておけば、以後はその中からランダムに選択することで序盤ランダム方策と同じ速度で学習を行うことができる。本手法で得られた評価関数を大会2位のプログラムと対戦させたところ、勝率は1割程度であった。自己対戦のみで強いプレイヤーを作るにはまだ改善の必要があるだろう。

## 参考文献

- [1] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [2] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In Johannes Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou, editors, *ECML*, Vol. 4212 of *Lecture Notes in Computer Science*, pp. 282–293. Springer, 2006.
- [3] Sylvain Gelly, Yizao Wang, Rémi Munos, and Olivier Teytaud. Modification of UCT with patterns in monte-carlo go. Rapport, HAL - CCSD - CNRS, 2006.
- [4] Alborz Geramifard, Michael Bowling, Martin Zinkevich, and Richard S. Sutton. iLSTD: Eligibility traces and convergence analysis. In Bernhard Schölkopf, John Platt, and Thomas Hoffman, editors, *NIPS*, pp. 441–448. MIT Press, 2006.