

モンテカルロシミュレーションを用いた強化学習法の提案

大崎 泰寛[†] 柴原 一友^{††} 但馬 康宏^{††} 小谷 善行^{††}

本稿ではモンテカルロシミュレーションを用いた評価関数の強化学習手法を提案する。評価関数の自動制御の研究は古くから行われており、最近では将棋プロの棋譜を教師として兄弟局面を比較する学習手法が注目を集めている。ところが、棋譜がなく、コンピュータによる解析もまだ行われていない新しい思考ゲーム環境における強化学習の成功事例の報告は稀である。そこで、モンテカルロシミュレーションによって得られた各局面の潜在的な勝率の平均値を教師値にした新しい強化学習法を提案する。これは、学習前の評価関数の精度に依存せず、膨大な学習データとなる棋譜が必要ないことからLMS法やTD法を代表とする従来の強化学習法の諸問題を解決する可能性を示している。そして、ブロックスデュオの環境における比較実験から、提案手法が従来の学習手法よりも高い学習成果が挙げられたことを確認した。

A Proposal of Reinforcement Learning Using Monte-Carlo Simulation

Yasuhiro OSAKI[†] Kazutomo SHIBAHARA^{††} Yasuhiro TAJIMA^{††} Yoshiyuki KOTANI^{††}

In this paper, we present reinforcement learning using Monte-Carlo simulation. A study of automatic adjustment of evaluation functions has been made for a long time, recently there has been an interest in a new learning method, comparing the selected move and the others in a database of the records by professional Shogi player. However, few successful cases of reinforcement learning have been reported in new thought games which have not the games records and have not been analyzed by computers yet. Therefore, we propose two new reinforcement learning methods; the teaching values are defined as average values of the potential percentage of victory by Monte-Carlo simulation. These algorithms, not depending on the accuracy of evaluation functions before learning and not needing the database of the game records, have shown possibility of solving some problems of former reinforcement learning methods represented Least Mean Square method and Temporal Differences learning. In our methods, we confirmed that the higher learning result has been obtained than the former one in our comparison experiments by Blokusduo.

1. はじめに

計算機科学における人工知能研究の黎明期から思考ゲームの解析は行われている。思考ゲームにおいて、コンピュータの思考ルーチンが最適な着手を選定する指針の一つに評価関数がある。この評価関数とは一般的に、与えられた局面から得られる様々な評価要素とその重要性を表す重みとの線形和から局面の形勢を数値で表すものであり、精度の高い評価関数の設計がコンピュータの思考ルーチンの強化に結びつく。しかし、精度の高い評価関数の構築には適切な評価要素の重み付けが必要不可欠であり、複雑な思考ゲームの環境において評価要素の数が多ければ多いほど手作業で逐一重みを制御することが困難となる。したがって、今日に至るまで評価要素の重み付けの自動制御に関して、LMS法やTD法に代表される強化学習の

研究が取り組まれてきたが、その成功事例の報告はごく限られたものであった。

そこで本稿では、近年コンピュータ囲碁の思考ルーチンが大きな前進を遂げる引き金となったモンテカルロシミュレーション[4]を用いた強化学習法である“LMS-MC法”ならびに“TD(λ)-MC法”を提案する。

これは、各局面からモンテカルロシミュレーションをすることによってランダムに派生した末端局面から元の局面の潜在的な勝率の近似値を求め、この値に評価関数の静的な評価値が近似するように各パラメータの重みを更新するというものである。本稿では、モンテカルロシミュレーションを用いた強化学習法によって生成された評価関数と、従来の強化学習法によって生成された評価関数をブロックスデュオで直接対局させて学習成果の比較および検証を行った。これらの結果から、本手法はプロの棋譜がなく、コンピュータによる解析の歴史が浅い新しい思考ゲーム環境においても、精度の高い評価関数が学習できるという可能性を示した。

[†] 東京農工大学 工学部 情報コミュニケーション工学科
Department of Computer and Information Sciences
Tokyo University of Agriculture and Technology

^{††} 東京農工大学 共生科学技術研究院 先端情報科学部門
Division of Advanced Information Technology & Computer
Science, Institute of Symbiotic Science and Technology,
Tokyo University of Agriculture and Technology

2. 関連研究

強化学習には LMS 法や TD 法があり、それらを思考ゲームの環境に適用させて評価関数を調節する研究がされてきた[5]。ところが、この従来手法には各パラメータの重みを更新する上でそれぞれ問題点があった。

まず、思考ゲームの環境に適用された LMS 法は学習方針となる教師値が学習エピソードの末端の勝敗(勝ちなら 1, 負けなら 0)であり、この値を各局面に割り振り、局面ごとの静的な評価値を近似させるアルゴリズムである。しかし、同一の学習エピソード内では各局面の教師値がすべて均一であるため、その精度に問題があった。例えば、ある学習エピソードで結果的に勝利していれば、そこに登場したすべての局面は正しかったものとして、各局面の静的な評価値が 1 に近似するように各パラメータの重みは更新されるが、登場した局面が一様に学習されてしまうので、どの局面のどのパラメータが勝利するために重要であったのかが考慮されず、学習の精度が下がる。

また、異なる学習エピソード間で全く同じ局面が出現しても、末端の勝敗が異なれば教師値が異なるため、一般的に学習の収束に時間がかかるという問題があった[5]。

続いて、思考ゲームの環境に適用された TD 法は教師値が一つ推移した子局面で得られる静的な評価値であり、各パラメータの重みを更新することで各親子局面間の評価値の差を減らすアルゴリズムである。着手後の子局面は学習エピソードの末端に近づいた分だけ評価値の精度が増すという時間的差分に注目した学習法であるが、学習前の評価関数の精度が低い場合では教師値の精度が低いという問題があった。したがって、将棋やチェスと異なり既存の知識や棋譜がなく、コンピュータによる解析がされていない思考ゲームの環境において高い学習成果は望めない。

そこで、従来手法の問題点を解決するためモンテカルロシミュレーション(以下、モンテカルロ法)を用いた強化学習法である“LMS-MC 法”ならびに“TD(λ)-MC 法”を提案する。これらの提案手法は、大崎らによるモンテカルロ法を用いた強化学習アルゴリズム[7]の拡張にあたり、一部名称を変更したものである。提案手法では LMS 法の

ように各局面に均一な教師値が割り振られることはなく、モンテカルロ法によって得られた潜在的な勝率を割り振るため、各局面の形勢に準拠した教師値となる。そのために異なる学習エピソードに同一の局面が登場しても教師値が変動する LMS 法の問題点が解決されている。また、これらの提案手法では乱数シミュレーションによって教師値を得るために、教師値が評価関数の精度に依存しないことから TD 法の問題点が解決されている。

3. モンテカルロ法を用いた強化学習法の提案

ここ近年、モンテカルロ法やそれに改良を加えた UCT アルゴリズム[3]による着手の研究がコンピュータ囲碁の分野で関心を集め、思考ルーチンの大幅な向上に一役買う形となった。このことから、評価関数による知識表現の困難な思考ゲームに対して、解析的な見解ではなく乱数シミュレーションによる近似的な見解が有効であることの裏付けが取れた[8]。以下に LMS-MC 法の学習アルゴリズムを示す。

3.1 LMS-MC 学習アルゴリズム

モンテカルロ法を用いた強化学習法は教師値が各局面から得られる潜在的な勝率であるため、従来の LMS 法と異なり学習エピソードの最終結果を待たずして、各パラメータの重みの逐次更新が可能である。言い換えると、オンライン更新(on-line updating)が可能である。以下の図 1 に LMS-MC 法の教師値の割り振りを示す。

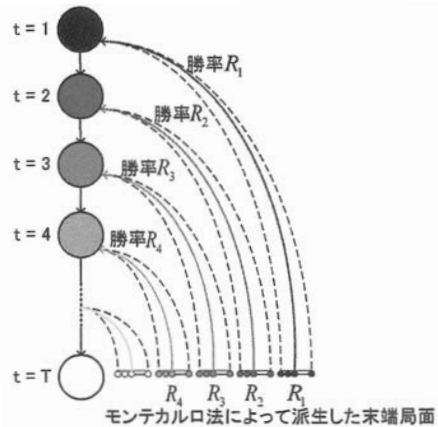


図 1 LMS-MC 法の教師値の割り振り

次に M 個の各パラメータの重み列 \bar{w} における具体的な学習手順を以下の図 2 に示す。

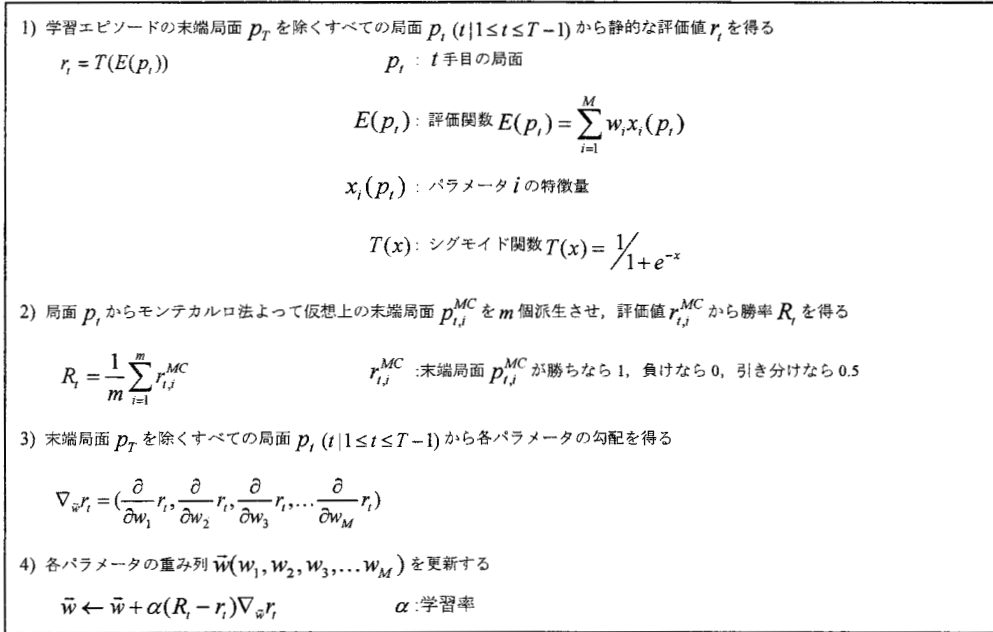


図 2 LMS-MC 法の学習アルゴリズム

t 手目の局面 p_t において、各特徴量と重みの線形和から評価値 $E(p_t)$ を求め、シグモイド関数を用いて 0 から 1 の値に正規化する。この値をモンテカルロ法による潜在的な勝率 (0 から 1 の値) に最急降下法の原理を用いて近似させた。

このように、学習エピソードの末端局面の勝敗が不要で、評価値によるブートストラップも行わず、各パラメータの重みを更新している。

LMS-MC 法の目的関数を以下の式 (1) に示す。

$$E_{LMS-MC} = \sum_{t=1}^{T-1} (R_t - r_t)^2 \quad (1)$$

目的関数の値は教師値と各局面の静的な評価値との差の 2 乗の総和であるため、値の増減幅が減少すれば学習は収束傾向にあることがいえる。特に、目的関数の値が 0 に収束すれば、対象局面からモンテカルロ法を m 回行って得られた潜在的な勝率が静的評価値として表すことができる評価関数に制御される。ここで留意すべきは、潜在的な勝率を表すために評価要素となるパラメータを十分に用意する点である。

続いて、TD 法の問題点を解決し、LMS-MC 法を“着手”の観点から改良した TD(λ)-MC 法の学習アルゴリズムを示す。

3.2 TD(λ)-MC 学習アルゴリズム

TD 法は子局面の評価値から親局面の評価値を学習するため、学習前の評価関数にある程度高い精度が求められた。提案手法である TD(λ)-MC 法では、着手後の子局面からモンテカルロ法によって得られた潜在的な勝率を教師値とすることからその問題点を解決している。

以下の図 3 に TD(λ)-MC 法の教師値の割振りを示す。

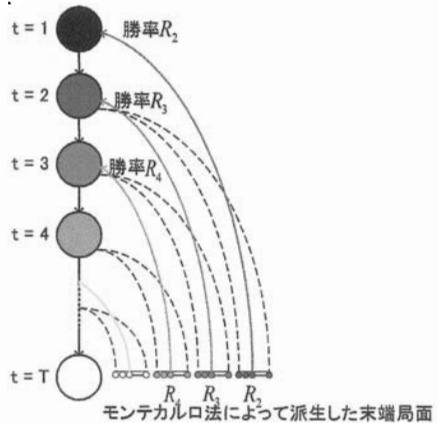


図 3 TD(λ)-MC 法の教師値の割振り

次に M 個の各パラメータの重み列 \vec{w} における具体的な学習手順を以下の図 4 に示す。

1) 学習エピソードの末端局面 p_T を除くすべての局面 p_t ($t|1 \leq t \leq T-1$) から静的な評価値 r_t を得る

$$r_t = T(E(p_t)) \quad p_t : t \text{ 手目の局面}$$

$$E(p_t) : \text{評価関数 } E(p_t) = \sum_{i=1}^M w_i x_i(p_t)$$

$$x_i(p_t) : \text{パラメータ } i \text{ の特徴量}$$

$$T(x) : \text{シグモイド関数 } T(x) = \frac{1}{1+e^{-x}}$$

2) 局面 p_t の子局面 p_{t+1} からモンテカルロ法によって仮想上の末端局面 $p_{t+1,j}^{MC}$ を m 個派生させ、評価値 $r_{t+1,j}^{MC}$ から勝率 R_{t+1} を得る

$$R_{t+1} = \frac{1}{m} \sum_{j=1}^m r_{t+1,j}^{MC} \quad r_{t+1,j}^{MC} : \text{末端局面 } p_{t+1,j}^{MC} \text{ が勝たらなら 1, 負けなら 0, 引き分けなら 0.5}$$

3) 末端局面 p_T を除くすべての局面 p_t ($t|1 \leq t \leq T-1$) から各パラメータの勾配を得る

$$\nabla_{\vec{w}} r_t = \left(\frac{\partial}{\partial w_1} r_t, \frac{\partial}{\partial w_2} r_t, \frac{\partial}{\partial w_3} r_t, \dots, \frac{\partial}{\partial w_M} r_t \right)$$

4) 各パラメータの重み列 $\vec{w}(w_1, w_2, w_3, \dots, w_M)$ を更新する $\vec{w} \leftarrow \vec{w} + \alpha(R_{t+1} - r_t) \sum_{i=1}^t \lambda^{t-i} \nabla_{\vec{w}} r_i$ λ : 過去の予定定数

図 4 TD(λ)-MC 法の学習アルゴリズム

3.1 節の LMS-MC 学習アルゴリズム同様、末端の勝敗が不要であるためオンライン更新が可能である。TD(λ)-MC 法の目的関数を以下の式 (2) に示す。

$$E_{TD(\lambda)-MC} = \sum_{t=1}^{T-1} (R_{t+1} - r_t)^2 \quad (2)$$

目的関数の値が 0 に収束すれば、対象局面の子局面からモンテカルロ法を m 回行って得られた潜在的な勝率が対象局面の静的評価値として表すことができる評価関数に制御される。LMS-MC 法との相違は、学習に着手の良し悪しが影響する点である。例えばある局面から良手を指せば、形勢の有利な子局面が派生されるのでモンテカルロ法によって得られる教師値が高くなり、悪手を指せば当然教師値が低くなる。このことから、TD(λ)-MC 法によって自動制御された評価関数では、深さ 2 の探索をせずに子局面の潜在的な勝率が静的評価から求まることがわかる。

3.3 ブロックデュオにおける評価関数の設計

ブロックデュオはモノミノからペントミノまでの計 21 種類のピースを使う二人零和有限確定完全情報ゲームである。14×14 の正方の盤上に交互

にピースを置き、その置いた累計面積で勝敗を決する思考ゲームである。自分の置いたピースの角と角をつなげて置くという制約があり、初手は決められたマス (図 5 の 15 番マスが先手、55 番マスが後手) を覆うように置く。ゲーム木の大きさの概算は 10^{80} 程度であり、これはハバースとチェスのちょうど中間に位置する。

評価関数には「盤の重み (105 種類)」「有効マスを中心とした 5×5 正方の特徴の重み (12 種類)」「盤上の累計面積 (1 種類)」の計 118 種類のパラメータを用意した。

1) 盤の重み

	1	2	3	4	5	6	7	8	9	A	B	C	D	E
1	1	2	4	7	11	16	22	29	37	46	56	67	79	92
2		3	5	8	12	17	23	30	38	47	57	68	80	93
3			6	9	13	18	24	31	39	48	58	69	81	94
4				10	14	19	25	32	40	49	59	70	82	95
5					15	20	26	33	41	50	60	71	83	96
6						21	27	34	42	51	61	72	84	97
7							28	35	43	52	62	73	85	98
8								36	44	53	63	74	86	99
9									45	54	64	75	87	100
A										55	65	76	88	101
B											66	77	89	102
C												78	90	103
D													91	104
E														105

図 5 盤の重み

コンピュータ将棋やチェスに倣って[1][2][6], 盤の重みを評価関数のパラメータとした. 図 5 は先手から見た盤の重みであり, 後手の場合は 1 番マス~105 番マスが反転する. ゲームの性質から盤の重みは線対称である.

2) 有効マスを中心とした 5×5 正方の特徴の重み

この有効マスとは, 自分の置いたピースの角のマスを目指す. ブロックデュオでは有効マスが多ければ多いほど可能手は増え, 有効マスを基準にピースを置くので重要な役割を持つものと考えた. そして, 以下の表 A, 図 6 に有効マスを中心とした 5×5 正方の特徴を示した. ただし, 死角マスとは自分がもう置けないマスを指し, 空マスとは自分にとって有効マスでも死角マスでも置いたマスでもないマスを指すものとする.

表 A 各マスの特徴の分類

マスの特徴	先手	後手
①	空マス	有効マス
②	空マス	有効マス
③	空マス	死角マス
④	有効マス	空マス
⑤	有効マス	有効マス
⑥	有効マス	死角マス
⑦	死角マス	空マス
⑧	死角マス	有効マス
⑨	死角マス	死角マス
⑩	置いたマス	置かれたマス
⑪	置かれたマス	置いたマス
⑫	盤外	盤外

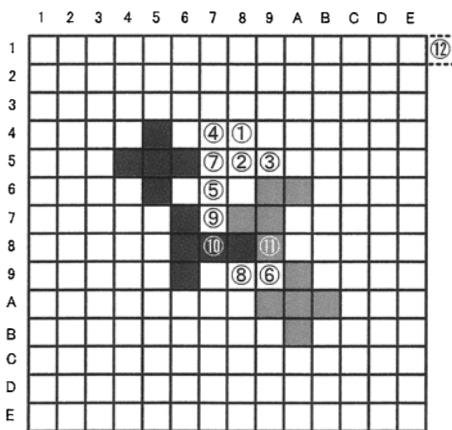


図 6 各マスの特徴の出現例

3) 盤上の累計面積

ブロックデュオは置いたピースの面積の総和

で勝敗を決するので, 盤上の累計面積は重要な評価要素であると考えた.

以上 118 種のパラメータから評価関数を作った.

4. 学習成果の比較実験

学習前の各パラメータの重みはすべて 1 とし, 従来手法である LMS 法と TD(λ) と, 提案手法である LMS-MC 法と TD(λ)-MC 法の 4 種類の学習法から得られた学習成果の比較実験を行った. 以下の図 7 のように, 自己対局で得られた学習データから各パラメータの重みを増減し, 更新された評価関数から再び学習データを得る. この繰り返しを 1000 回それぞれの学習法で行った.

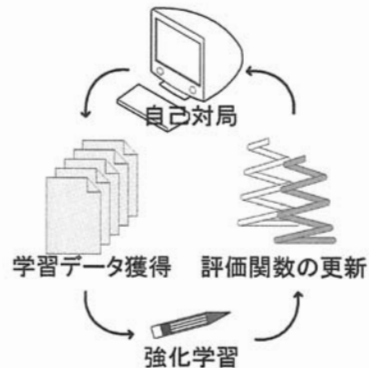


図 7 学習フロー

また, 全 118 種類のパラメータの重みをゲームの前半と後半にそれぞれ独立させて学習した. 例えば将棋の場合, 序盤は定跡, 終盤は詰めというようにゲームの進行状況によって着手の役割が異なることから, これに倣って各パラメータの重みを一辺倒に学習することを避けた. ブロックデュオの平均手数は 30 手前後であったため, 前半を初手~14 手目, 後半を 15 手目~終局とした.

その他にも, 自己対局によって得られる学習データが同一局面ばかりで偏向しないように微小の乱数を加えて最善手以外の着手も含めた.

4.1 各パラメータの重みと目的関数の推移

以下の図 8~図 11 に盤の重みを示し, 図 12, 図 13 に (空マス・空マス), (空マス・有効マス), (空マス・死角マス), (累計面積) の重みを示し, 図 14, 図 15 には目的関数の推移を示した.

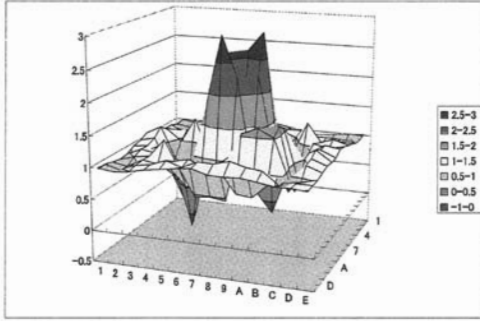


図8 LMS-MC 法による盤の重み (ゲーム前半)

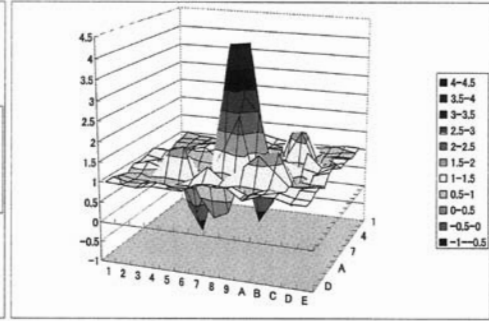


図9 TD(λ)-MC 法による盤の重み (ゲーム前半)

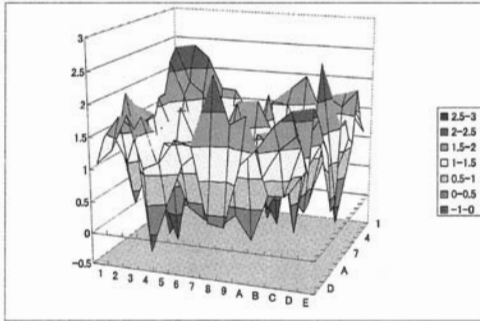


図10 LMS-MC 法による盤の重み (ゲーム後半)

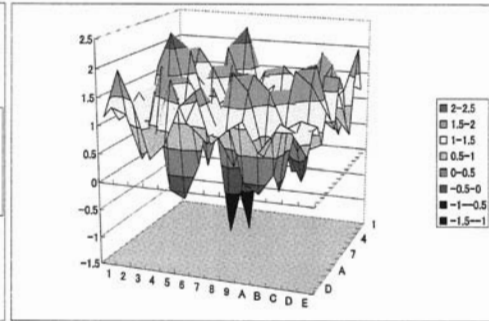


図11 TD(λ)-MC 法による盤の重み (ゲーム後半)

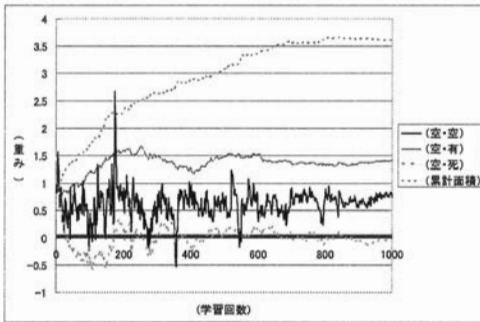


図12 LMS-MC 法によるマスの特徴の重み (ゲーム前半)

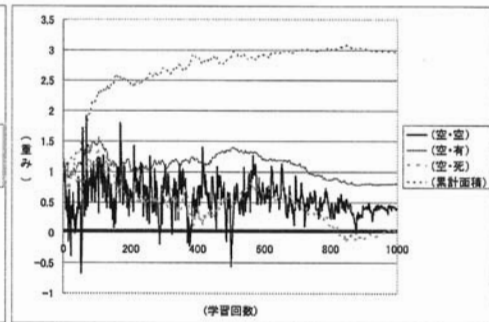


図13 TD(λ)-MC 法によるマスの特徴の重み (ゲーム前半)

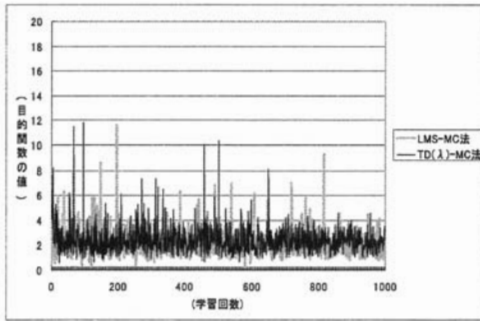


図14 LMS-MC 法とTD(λ)-MC 法の目的関数の推移

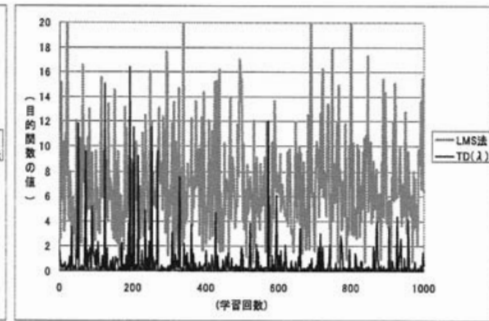


図15 LMS 法とTD(λ)の目的関数の推移

盤の重みに関しては、LMS-MC 法と TD(λ)-MC 法から得られた学習結果が良く似たものとなった。ゲーム前半では盤の中央の重みが 1 よりも高いため、ピースを盤の中央に置き可能手を増やすことが勝率を高くするということがわかった。一方、盤の中央よりもその周囲の方がゲーム後半では重みが高いことがわかった。

有効マスを中心とした 5×5 マスの特徴の重みに関しても、LMS-MC 法と TD(λ)-MC 法から得られた学習結果が良く似ていた。盤上の累計面積が学習前の重みの 3~3.5 倍の値で収束していることから、ゲーム前半では盤上に面積の大きいピースを置くことが勝率を高くするということがわかった。また、(空マス・有効マス)の重みが LMS-MC 法は 1 より上がり、TD(λ)-MC 法では 1 より下がるということがわかった。

目的関数の推移に関しては、各学習法で特徴が表れた。図 14 より、TD(λ)-MC 法よりも LMS-MC 法の値が低いことがわかった。また、学習回数に応じて値が低くなっていることから、学習が収束方向に向かっていると見える。その一方で、図 15 より、TD 法は目的関数の値から学習が収束傾向にあるが、LMS 法は目的関数の値が発散しているため、収束方向に向かっていないことがわかった。

4.2 各評価関数の対局結果

続いて、それぞれの強化学習法から得られた評価関数と学習前の評価関数を合わせた 5 種類の評価関数から対局実験を行った。公正を期すため対局は先手後手それぞれ 500 対局ずつ行った。評価関数に微小の乱数を入れて次善手以下も着手するようにした。対局結果は以下の図 16~図 18 のとおりとなった。

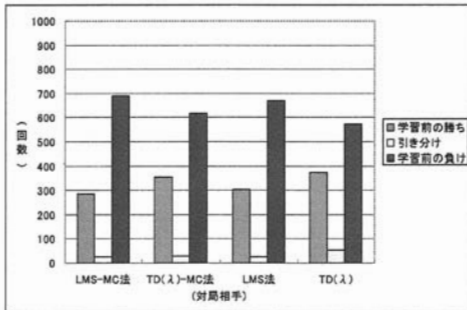


図 16 学習前の評価関数の対局結果

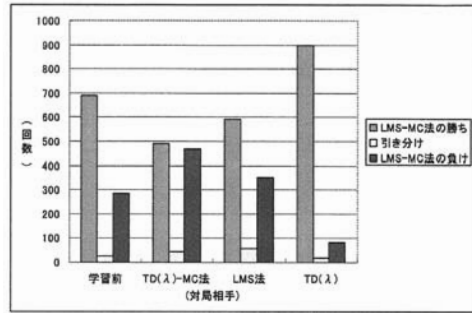


図 17 LMS-MC 法による評価関数の対局結果

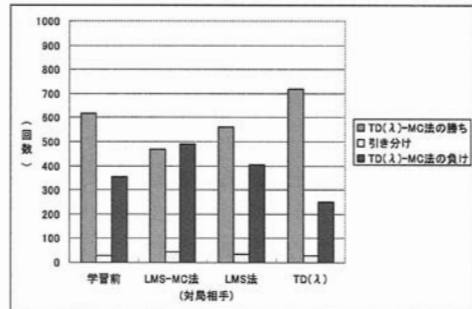


図 18 TD(λ)-MC 法による評価関数の対局結果

図 16 より、各強化学習をした評価関数が学習前の評価関数を上回る結果となった。特に、LMS-MC 法による評価関数は学習前の評価関数に対して、690 勝 284 敗 26 分であった。

図 17 より、提案手法である LMS-MC 法で学習した評価関数は他のすべての評価関数に勝ち越す結果となった。特に、従来手法の LMS 法に対して 593 勝 351 敗 56 分と 59.3% 勝つことがわかった。

図 18 より、提案手法である TD(λ)-MC 法で学習した評価関数は LMS-MC 法による評価関数を除くすべての評価関数に勝ち越す結果となった。ただし、LMS-MC 法による評価関数との差は 468 勝 489 敗 43 分とわずか 21 分であったことから、LMS-MC 法に明らかな優位性があるとは言いきれない。また、従来手法の TD 法に対して 719 勝 251 敗 30 分と 71.9% 勝つことがわかった。

4.3 乱数プログラムとの対局結果

ランダムで着手する思考ルーチンと対局することで学習前の評価関数から精度がどの程度向上したかを検証した。対局結果は以下の図 19 に示す。

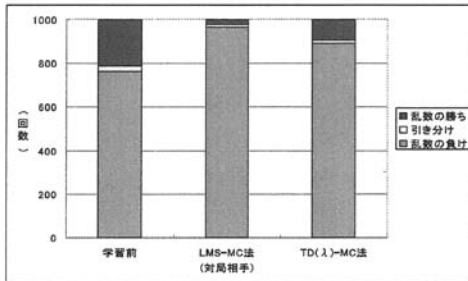


図 19 乱数プログラムとの対局結果

先後手合わせて 1000 対局行ったところ、図 19 のような対局結果が得られた。学習前の評価関数が乱数プログラムに対して 762 勝であった。一方で、LMS-MC 法による評価関数では 965 勝であり、TD(λ)-MC 法による評価関数では 893 勝であった。

これらの結果から提案手法によって学習された評価関数の精度が向上していることがわかった。

4.4 自己対局の結果

各評価関数が自己対局を 500 回ずつ行った結果を以下の図 20 に示す。

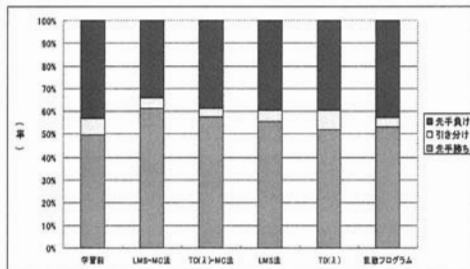


図 20 自己対局の結果

乱数プログラムは 500 対局中 52.9%先手が勝ち、学習前の評価関数の先手の勝率が最も低く 49.6%、LMS-MC 法による評価関数の先手の勝率が最も高く 61.2%であった。これらの結果から、ブロックデュオは先手が有利な思考ゲームであるといえる。

5. おわりに

本稿では、モンテカルロ法によって得られた各局面の潜在的な勝率を教師値とした学習法である LMS-MC 法および TD(λ)-MC 法を提案した。そして、プロの棋譜などの学習データが存在しないう

えにコンピュータによる解析の歴史が浅いブロックデュオにおける実験環境下では、従来手法である LMS 法や TD 法によって学習された評価関数よりも更に優れた評価関数の学習が実現した。LMS-MC 法による評価関数は LMS 法による評価関数に対し 59.3%勝ち、TD(λ)-MC 法による評価関数は TD(λ)による評価関数に対し 71.9%勝った。

今後の課題として、LMS-MC 法と TD(λ)-MC 法の優劣関係を示すことが挙げられる。教師値が異なる両者の評価関数の勝敗が本稿の実験環境下ではほぼ同じであったため、学習率や評価関数のパラメータ数を変えて比較する必要がある。また、提案手法の目的関数が収束しなかった原因に評価関数のパラメータ不足が考えられるため、パラメータの増加と目的関数の収束関係を検証する。

モンテカルロ法を用いた強化学習法は解析未だの思考ゲームにおいて従来手法よりも優れた学習成果を示したことから強化学習の更なる成果向上に大きく寄与するものと思われる。

参考文献

- [1] D.F.Beal, and M.C.Smith : Temporal Difference learning for heuristic search and game playing, ICGA Journal Vol22 No4, pp.223-235, 1999.
- [2] J.Baxter, A.Triggell, L.Weaver : Experiments in Parameter Learning using Temporal Differences, ICGA Journal Vol.21 No.2, pp.84-99, 1998.
- [3] Sylvain Gelly, Yizao Wang, Remi Munos, Olivier Teytaud: Modification of UCT with Patterns in Monte-Carlo Go, RR-6062-INRIA, pp.1-19, 2006.
- [4] Remi Coulom : Monte-Carlo Tree Search in Crazy Stone, 第 12 回ゲーム・プログラミング ワークショップ, pp.74-75, 2007.
- [5] Richard S. Sutton and Andrew G.Barto, 三上 貞寿・皆川 雅章 共訳 : 強化学習, 森北出版株式会社, 2000.
- [6] 小谷善行 : コンピュータ将棋の頭脳, サイエンス社, 2007.
- [7] 大崎泰寛 柴原一友 但馬康宏 小谷善行: TD(λ)-MC 法を用いた評価関数の強化学習 第 12 回ゲーム・プログラミング ワークショップ, pp.36-43, 2007.
- [8] 中村秋吾, 三輪誠, 近山隆 : 静的評価関数を用いた UCT の改善, 第 12 回ゲーム・プログラミング ワークショップ, pp.44-51, 2007.