

方策勾配法による局面評価関数とシミュレーション方策の学習

五十嵐治一^{†1}, 森岡祐一, 山本一将^{†2}

本論文では強化学習の一手法である方策勾配法をコンピュータ将棋に適用する方法を考察した。方策勾配法は、報酬や方策にマルコフ性の制限なく自由に設計することができるという大きなメリットがある。本論文では、最初に全 leaf 局面の局面評価値をその局面への遷移確率値で重み付けた期待値を用いた指し手評価方式を提案する。これをベースに、探索木の各ノードにおける指し手の選択法として Boltzmann 分布に基づくソフトマックス戦略を採用した場合の局面評価関数に含まれるパラメータの学習則を導出した。しかし、探索や学習時の計算量が膨大となるため、3つの近似計算法を考案した。次に、探索時にシミュレーション方策を用いてモンテカルロ探索を行う場合や、探索の深さを制御する場合のために、局面評価関数とシミュレーション方策の両者を同時に学習する学習則を方策勾配法により導出した。さらに、この方策勾配の計算法を利用すると、局面ごとに正解手が既知の場合の教師付学習も可能であることを示し、実際に学習則を導出した。

Learning Positional Evaluation Functions and Simulation Policies by Policy Gradient Algorithm

HARUKAZU IGARASHI^{†1} YUICHI MORIOKA KAZUMASA YAMAMOTO^{†2}

This paper applies policy gradient reinforcement learning to *shogi*, a traditional Japanese board game that resembles chess. First, we propose a move evaluation function, which is defined by the expectation of the values of all leaf nodes produced by the move in a search tree that is weighted by the transition probabilities to the leaf nodes from the root node produced by the move. Since policy gradient reinforcement learning does not require Markovian properties of reward functions and policies, system designers can create the rewards functions and policies more freely than when using other reinforcement learning methods that must be applied in Markov decision processes. The learning rules of the parameters in the positional evaluation function can be calculated recursively when the Boltzmann distribution function gives the probabilities of taking branches in a search tree. We also consider three approximation methods to reduce the computation time for tree searching and parameter learning. Second, we derived the learning rules for both positional evaluation functions and simulation policies for Monte-Carlo simulation search and controlling the search depth by the policy gradient algorithm. This approach can also be applied to supervised learning problems of a teacher's moves in a given position.

1. はじめに

近年、コンピュータ将棋の実力はプロ棋士に迫るものがある[1]。例えば、2013年に行われた第2回電王戦では、プロ棋士5名とコンピュータ将棋ソフトのトップ5との対局が行われ、コンピュータ側の3勝1敗1分けという結果であった[2]。この一因となっているのが、将棋ソフトBonanzaで提案された評価関数の自動学習である[3]。現在では、評価関数をプロ棋士の棋譜データベースを利用した教師付き学習により構築することが主流となっている。

一方、教師付き学習ではなく強化学習により評価関数を学習する方法も考えられている。その代表的な強化学習法としてはTD(λ)法とTDLeaf(λ)法がある。TD(λ)法はバックギャモンでは大成功を収めており[4]、TDLeaf(λ)法はチェスにおいて棋力向上への有効性が確認されている[5]。しかし、将棋ではまだそれほど良い適用結果が報告されていない。

そこで本論文ではこれまでのTD(λ)法やTDLeaf(λ)法ではなく、“方策勾配法”と呼ばれる別の強化学習法を適用する

ことを考えた。方策勾配法は報酬を自由に設定することができるので、棋力向上だけでなく棋風の学習など様々な学習目的に対して幅広く適用できる。

しかしながら、将棋のように着手決定の際に探索(読み)を要する問題では、方策関数の中に探索処理を取り入れるための工夫が必要である。そこで、本論文ではまず、着手決定の方策として、“指し手評価の期待値”を用いた確率の方策を提案する。この方策は、探索木のleaf局面の局面評価関数を用いて再帰的に表現されるので、最適方策の学習は局面評価関数中の特徴量パラメータの学習に帰着する。この学習則を導出した後、近似計算法としてPGLeaf法を初めとするいくつかの学習法を提案する。

次に、指し手評価の期待値を計算する際に、leaf局面への遷移確率の値を着手選択の方策とは別の“シミュレーション方策”(simulation policy)を用いて計算する場合を考える。これは、モンテカルロ探索のように局面や指し手の評価値を何らかのシミュレーションで決定する場合である。また、この遷移確率の値は“激指”[6]の“実現確率”[a]

^{†1} 芝浦工業大学工学部情報工学科
Shibaura Institute of Technology
^{†2} 株式会社コスモ・ウェブ
Cosmoweb Co., Ltd.

a) 激指の実現確率の計算においてベースとなっている「指し手の遷移確率」は、その手を指すか指さないかの2値選択の選択確率であり、合法手の間での選択確率とは異なる。本論文では後者の場合を考えている。

に相当し、探索木の前向き枝刈り処理に用いることも可能である。本論文では、このシミュレーション方策中のパラメータと局面評価関数中の特徴量パラメータの両方を同時に学習できる強化学習則を導出する。さらに、強化学習ではなく、その局面での正解手を与える教師付き学習の場合の学習則も同様な方法で導出できることを示す。

2. 方策勾配法による強化学習

2.1 方策勾配法とは

強化学習では、Q学習のように行動価値関数を通して、あるいはTD法のように状態価値関数を通して、間接的に方策を学習する価値ベースの強化学習法 (value-based algorithm) がよく知られている[4]。一方、方策中にパラメータを入れておき、パラメータ空間内での期待報酬関数の最急勾配を計算することにより、方策を直接学習する強化学習法がある。WilliamsのREINFORCEアルゴリズム[7]や、部分観測マルコフ決定過程 POMDP (Partially Observable Markov Decision Processes) 環境における木村らの確率的傾斜法[8]などである。また、Q値を用いて上記の勾配関数を表現する方式[9][10][11]や、自然勾配を利用する方式[12]も考案されている。これら一連の方策ベースの強化学習法は、“方策勾配法”(policy gradient method)と呼ばれ、例えば、Petersらの文献[13]中に簡潔にまとめられている。

本研究では、五十嵐らが提案している方策勾配法[14]を用いる。この方式は、Williamsのエピソード単位の学習方式 (episodic REINFORCE algorithm) [7]に基づいており、環境モデル (状態遷移確率と報酬) と方策に関する単純マルコフ性を必要としない。これまでにマルチエージェント学習の標準問題として知られている追跡問題や、粒子群を用いた最適化手法であるPSO (Particle Swarm Optimization)等へ適用され、有効性が確認されている[15][16]。

2.2 方策勾配法による学習

t 回目 ($t=1,2,\dots,L_a$) の手番局面 u_t において学習エージェント A が指し手 a_t を選択する確率 (方策) を

$$\pi_a(a_t|u_t;\omega) = \exp(E_a(a_t, u_t; \omega)/T_a) / Z_a \quad (1)$$

$$Z_a \equiv \sum_{x \in A(u_t)} \exp(E_a(x, u_t; \omega)/T_a) \quad (2)$$

とする。ただし、 ω は評価関数中の学習パラメータ、 T_a は温度パラメータで、 $A(u_t)$ は A の手番局面 u_t における全ての合法手の集合である。(1)の右辺は Boltzmann 分布と呼ばれる確率分布関数であり、 $E_a(a_t, u_t; \omega)$ は手番局面 u_t における指し手 a_t の評価を表す指標であり”目的関数”と呼ぶ。一方、対戦エージェント B の方策 $\pi_b(b_t|v_t)$ は既知であるとする。ただし、 v_t は対戦エージェントの t 回目 ($t=1,2,\dots,L_b$) の手番局面であり、 b_t はそのときの指し手を表している。

一局の指し手と出現局面との時系列データ (棋譜) を”エピソード”と定義する。エピソード終了後に勝敗等を考

慮して学習エージェントに報酬 r を与える。一般に、両対局者の指し手の決定は確率的方策によるものとする。したがって、学習エージェント A の指し手数 (≡エピソード長 L_a) や報酬 r の観測値もエピソードごとに変動する。

文献[14]の方策勾配法を適用して、一局当たりの期待報酬値 $E[r]$ を極大化することを考える。そこでは、

$$\partial E[r] / \partial \omega = E \left[r \sum_{t=1}^{L_a} e_\omega(t) \right] \quad (3)$$

$$e_\omega(t) \equiv \partial \ln \pi_a(a_t|u_t; \omega) / \partial \omega \quad (4)$$

と表されることから、学習則として

$$\Delta \omega = \varepsilon \cdot r \sum_{t=1}^{L_a} e_\omega(t) \quad (5)$$

が用いられている。 ε は学習係数で小さな正数にとる。(5) は報酬と実現手の選択確率の勾配との積 (相関) に比例させて各パラメータを更新 (強化) している。今、方策が(1)の場合、(4)の特徴的適正度 $e_\omega(t)$ は、

$$e_\omega(t) = (1/T_a) \left[\partial E_a(a_t, u_t; \omega) / \partial \omega - \sum_{x \in A(u_t)} \pi_a(x|u_t; \omega) \partial E_a(x, u_t; \omega) / \partial \omega \right] \quad (6)$$

と表される。本論文で用いる方策勾配法のアルゴリズムをまとめると次の様になる。

【方策勾配法の学習アルゴリズム】

- step1: 各時刻 t において方策 $\pi_a(a_t|u_t; \omega)$ により指し手 a_t を選択し、特徴的適正度 $e_\omega(t)$ を計算する。
- step2: エピソード終了後に報酬 r を与える。
- step3: 学習則(5)により更新量 $\Delta \omega$ を計算する。
- step4: ω を $\omega + \Delta \omega$ と変更し、step1 から繰り返す。ただし、終了条件を満たせば終了する。

3. 探索と局面評価関数による指し手の評価

指し手の評価は、読み (探索木の展開) を伴う方がより正確と考えられる。そこで、(1)の目的関数 $E_a(a_t, u_t; \omega)$ を、着手後の局面 $v=v(a_t, u_t)$ ではなく、探索木 $G_D(a_t, u_t)$ の末端の局面 (以下、leaf 局面) の評価値を用いた関数とする。ここで、 $G_D(a_t, u_t)$ は局面 $v(a_t, u_t)$ をルートとする深さ D の探索木である。ただし、学習エージェント A の手番から次の A の手番までを深さの1単位とし、出現局面 u_t を深さ0の、 $G_D(a_t, u_t)$ の leaf 局面を深さ D の A の手番局面とする。

本研究では、以下の“指し手評価の期待値”

$$E_a^*(a_t, u_t; \omega) \equiv \sum_{u \in U_D(a_t, u_t)} P(u|a_t, u_t; \omega) E_a^*(u; \omega) \quad (7)$$

を(1)の目的関数 $E_a(a_t, u_t; \omega)$ として用いることを提案する。

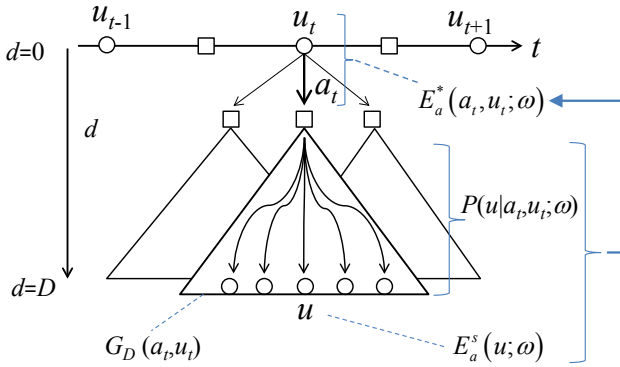


図 1 指し手評価の期待値 $E_a^*(a_t, u_t; \omega)$ と leaf 局面での局面評価値 $E_a^s(u; \omega)$, 遷移確率 $P(u|a_t, u_t; \omega)$ の関係を表す。
 Figure 1 Expected evaluation function $E_a^*(a_t, u_t; \omega)$ of move a_t , positional evaluation function $E_a^s(u; \omega)$ of leaf node u , and transition probability $P(u|a_t, u_t; \omega)$.

ただし, $U_D(a_t, u_t)$ は探索木 $G_D(a_t, u_t)$ の全 leaf 局面の集合を, $E_a^s(u; \omega)$ は leaf 局面 u での静的局面評価関数の値を表している. $P(u|a_t, u_t; \omega)$ は確率的な方策により leaf 局面 u へ遷移する確率である. これらの概念図を図 1 に示す.

また, 着手決定のためのアルゴリズムは以下のように表される:

【着手決定のアルゴリズム】

- step1: 手番局面 u_t において全ての合法手 a_t を生成する.
- step2: 各 a_t に対して $E_a^*(a_t, u_t; \omega)$ を計算する.
- step3: 方策 $\pi_a(a_t|u_t; \omega)$ により着手を選択する.

ただし, (7) の $E_a^*(a_t, u_t; \omega)$ は, 4.2 で後述するように再帰により厳密に計算できる ((12) 参照). ただし, 再帰の最下層では leaf 局面の評価値が呼び出される. また, 方策 $\pi_a(a_t|u_t; \omega)$ は, (1) の Boltzmann 分布を念頭に置いている.

通常, ゲーム木探索における指し手評価では, 探索木 $G_D(a_t, u_t)$ に対して min-max 探索法や $\alpha\beta$ 探索法により得られた leaf 局面での局面評価値を指し手 a_t の評価値とする. これは, (7) の右辺の計算において期待値計算を厳密に行わないで, 探索木の最善応手手順(principal variation)の leaf 局面(principal leaf) $u^*_D(a_t, u_t)$ の局面評価値 $E_a^s(u^*_D(a_t, u_t); \omega)$ で代表するという一種の近似計算に相当する.

また, ヒューリスティクスを用いた探索木の枝刈りも, (7) の右辺の探索過程において leaf 局面への遷移確率をゼロとおくことに相当する. 例えば, 激指チームの“実現確率”(= “親の実現確率” × “指し手の遷移確率”) による枝刈り [5] も同様である. これについては 5.2 で考察する.

さらに, 近年, 囲碁などで盛んなモンテカルロ探索 [17] は, 局面評価のために多数回のプレイアウトを行う. これは, (7) の右辺の期待値操作を, あるシミュレーション方

策により生成した leaf 局面の評価値の単純平均操作で置き換えたと見なすことができる. シミュレーション方策を用いた近似計算法については, 改めて 6. で詳細に検討する.

本論文で指し手の評価として (7) のような期待値を提案した理由は次の 2 つである. まず, 上記のように様々な指し手探索法を特殊ケースとして導くことが可能で理論的な見通しが良いことである. 次に, 最善応手手順だけを利用すると, 読みの深さや評価関数の精度に限界がある場合には, 最善応手手順以外の変化手順をも十分考慮して指し手の評価を行う方が, 読みの深さが有限であることと評価値の誤差に起因する探索の揺らぎに対して頑健な評価法を与えるのではないかと考えたからである.

4. 探索と方策勾配法による評価関数の学習

4.1 学習則

3. では (1) の目的関数として出現局面 u_t における指し手 a_t の直接的な評価値 $E_a(a_t, u_t; \omega)$ ではなく, (7) に示した a_t 以下の全 leaf 局面の評価値 $\{E_a^s(u; \omega)\} (u \in U_D(a_t, u_t))$ と各 leaf 局面への遷移確率 $P(u|a_t, u_t; \omega)$ とを用いて計算することを提案した. したがって, 学習エージェント A の方策 (1), (2) は,

$$\pi_a(a_t|u_t; \omega) = \exp(E_a^*(a_t, u_t; \omega)/T_a) / Z_a \quad (8)$$

$$Z_a \equiv \sum_{x \in A(u_t)} \exp(E_a^*(x, u_t; \omega)/T_a) \quad (9)$$

と表される. このときの学習則は, (5), (6) より,

$$\Delta\omega = \varepsilon \cdot r \sum_{t=1}^{L_a} e_\omega(t) \quad (10)$$

$$e_\omega(t) = (1/T_a) \left[\partial E_a^*(a_t, u_t; \omega) / \partial \omega - \sum_{x \in A(u_t)} \pi_a(x|u_t; \omega) \partial E_a^*(x, u_t; \omega) / \partial \omega \right] \quad (11)$$

と表される.

(8)~(11) は, $E_a^*(a_t, u_t; \omega)$ と $\partial E_a^*(a_t, u_t; \omega) / \partial \omega$ の値が局面 u_t における合法的指し手 a についてすべて分かれば計算できる. ただし, これらの値は局面 u_t において指し手 a を指した局面以下の部分木 $G_D(a, u_t)$ の全 leaf 局面 $u \in U_D(a, u_t)$ に依存する. したがって, 2.2 で述べた通常の方策勾配法の適用方式では, 出現局面 u_t 以下の深さ 1 の局面に含まれる特徴量パラメータのみが更新対象となるが, 探索を伴う本方式では全 leaf 局面に含まれる特徴量パラメータすべてが更新対象となり, 一対局あたりの学習の効率化が期待できる.

4.2 指し手評価の期待値とその勾配の再帰計算

(8)~(11) の $E_a^*(a_t, u_t; \omega)$ と $\partial E_a^*(a_t, u_t; \omega) / \partial \omega$ は再帰的に計算できることを示す. まず, 深さ d ($0 \leq d \leq D-1$) における学習エージェント A の手番局面 u^d において, 指し手 a^d により生成される対戦相手 B の手番局面を $v^d = v(a^d, u^d)$, その局面から B が指し手 b^d を指して得られた学習エージェント A

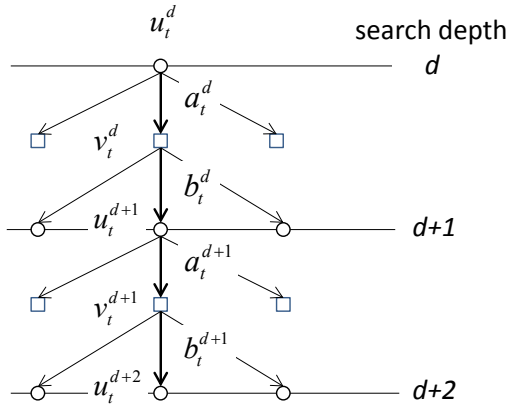


図 2 探索の深さ, 手番局面, 指し手の関係.

Figure 2 Search depth, positions and moves.

の手番局面を $u^{d+1} = u(b_t^d, v_t^d)$ とする (図 2). ただし, 対戦相手のエージェント B の方策 π_b は既知とする.

この時, 探索の深さ d における $E_a^*(a_t^d, u_t^d; \omega)$ と $\partial E_a^*(a_t^d, u_t^d; \omega) / \partial \omega$ は次のように再帰的に書ける.

$$E_a^*(a_t^d, u_t^d; \omega) = \sum_{b_t^d} \pi_b(b_t^d | v_t^d) \cdot \sum_{a_t^{d+1}} \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) \cdot E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) \quad (12)$$

$$\partial E_a^*(a_t^d, u_t^d; \omega) / \partial \omega = (\partial / \partial \omega) \sum_{b_t^d} \pi_b(b_t^d | v_t^d) \cdot \sum_{a_t^{d+1}} \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) \cdot E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) \quad (13)$$

$$= \sum_{b_t^d} \pi_b(b_t^d | v_t^d) \sum_{a_t^{d+1}} \partial \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) / \partial \omega \cdot E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) + \sum_{b_t^d} \pi_b(b_t^d | v_t^d) \sum_{a_t^{d+1}} \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) \cdot \partial E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) / \partial \omega \quad (14)$$

$$= \sum_{b_t^d} \pi_b(b_t^d | v_t^d) \sum_{a_t^{d+1}} \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) \cdot [e_\omega(t, d+1) E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) + \partial E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) / \partial \omega] \quad (15)$$

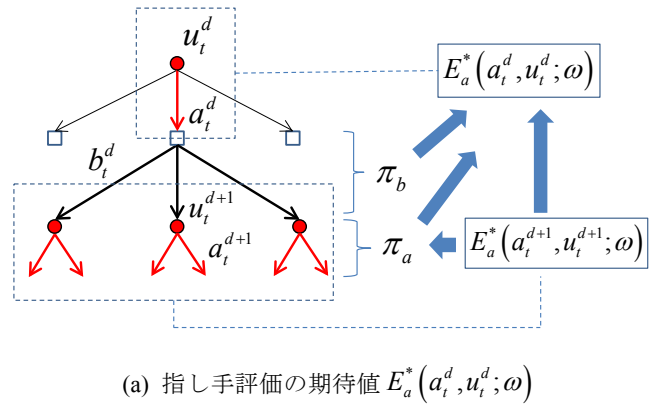
ただし,

$$e_\omega(t, d+1) \equiv \partial \ln \pi_a(a_t^{d+1} | u_t^{d+1}; \omega) / \partial \omega \quad (16)$$

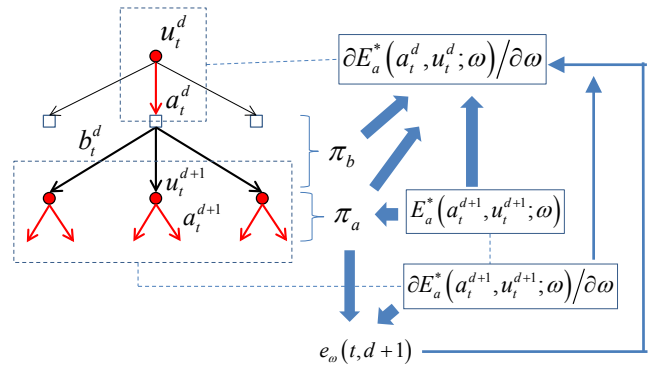
$$= (1/T_a) \left[\partial E_a^*(a_t^{d+1}, u_t^{d+1}; \omega) / \partial \omega - \sum_{x \in A(u_t^{d+1})} \pi_a(x | u_t^{d+1}; \omega) \partial E_a^*(x, u_t^{d+1}; \omega) / \partial \omega \right] \quad (17)$$

である.

また, (12), (13) における再帰の終端は, もし, u^{d+1} が leaf 局面, すなわち, $d=D-1$ ならば,



(a) 指し手評価の期待値 $E_a^*(a_t^d, u_t^d; \omega)$



(b) 1 階微係数 $\partial E_a^*(a_t^d, u_t^d; \omega) / \partial \omega$

図 3 PG 行動期待値法の再帰計算における依存関係 : (a) 指し手評価の期待値, (b) 1 階微係数の値.

Figure 3 Recursive calculation in “PG expectation algorithm”: (a) Expected evaluation function $E_a^*(a_t^d, u_t^d; \omega)$, and (b) its first derivative $\partial E_a^*(a_t^d, u_t^d; \omega) / \partial \omega$.

$$E_a^*(a_t^d, u_t^d; \omega) = \sum_{b_t^{D-1}} \pi_b(b_t^{D-1} | v_t^{D-1}) E_a^s(u_t^D; \omega) \quad (18)$$

$$\partial E_a^*(a_t^d, u_t^d; \omega) / \partial \omega = \sum_{b_t^{D-1}} \pi_b(b_t^{D-1} | v_t^{D-1}) \cdot \partial E_a^s(u_t^D; \omega) / \partial \omega \quad (19)$$

と書ける. 図 3 に上記の依存関係を表した模式図を示す.

なお, 本論文では 2.2 で述べた出現局面における指し手評価値を用いた方策勾配法を “PG 法” または単に方策勾配法, 4.1 と 4.2 で提案した全 leaf 局面に基づく指し手評価の期待値を用いた方策勾配法を “PG 行動期待値法” (PG expectation algorithm) と呼んで区別することにする.

5. 計算量削減のための近似手法のアイデア

5.1 min-max 探索または $\alpha\beta$ 探索の適用 : PGLeaf 法

3. の指し手評価には, (12) の再帰計算で探索木の最下段での全 leaf 局面の局面評価値を知る必要がある. さらに, 4. の学習においては, 全 leaf 局面での勾配値も必要である. したがって, 指し手決定と学習にかかる計算時間は膨大となる可能性が予想される. そこで, 計算量を削減するための近似手法に関するアイデアを本章では述べる.

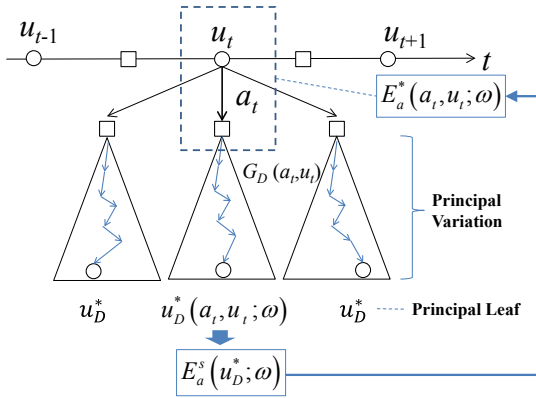


図 4 PGLeaf 法

Figure 4 PGLeaf algorithm.

まず、対戦相手 B の方策 π_b として min 探索を用いる。この近似に加えて、学習エージェント A の方策として max 探索を行う ((8) で $T_a \rightarrow 0$ と置くことに相当する)。すなわち、min-max 探索、あるいは $\alpha\beta$ 探索を行い、最善応手手順だけを考える。これは(7)の遷移確率において、

$$P(u|a_t, u_t; \omega) = \begin{cases} 1 & \text{if } u = u_D^*(a_t, u_t; \omega) \\ 0 & \text{otherwise} \end{cases} \quad (20)$$

と置いたと解釈できる。このように指し手評価の期待値 $E_a^*(a_t, u_t; \omega)$ を principal leaf $u_D^*(a_t, u_t; \omega)$ の局面評価値 $E_a^s(u_D^*(a_t, u_t; \omega); \omega)$ で置き換えた指し手決定法と学習法を“PGLeaf 法”と呼ぶことにする。PGLeaf 法では学習時に(12)~(17)のような再帰計算は不要で、通常の $\alpha\beta$ 探索アルゴリズムをそのまま利用できる。PGLeaf 法 の概念を 図 4 に示す。

5.2 反復深化法の適用

探索時に反復深化法を適用する方法が考えられる。ある深さ D を設定し、leaf 局面 u_D の集合とそれらの局面評価値 $E_a^s(u_D; \omega)$ を用いて leaf 局面までの遷移確率 $P(u_D|a_t, u_t; \omega)$ と指し手評価の期待値 $E_a^*(a_t, u_t; \omega)$ を計算する。ただし、遷移確率 $P(u_D|a_t, u_t; \omega)$ が閾値以下であればそれ以下の部分木はカットする。次に D を 1 だけ増やしてこの操作を繰り返す。

(12)~(19)での指し手評価の期待値の再帰的計算や学習時には、カットされないで残った枝の leaf 局面だけを用いる。この場合、残った枝の leaf 局面に含まれる特徴量パラメータすべてが更新される。図 5 に模式図を示す。

5.3 異なる評価関数の principal leaf による期待値の計算法

この近似法は 5.1 で述べた PGLeaf 法の合議制バージョンに相当する。まず、 N 個の異なる評価関数を持った探索アルゴリズム k によりそれぞれ min-max 探索を行い、 N 個の principal leaf $\{u_{D,k}^*\} (k=1, 2, \dots, N)$ を求める。次に、各 principal leaf における局面評価値 $E_a^s(u_{D,k}^*)$ を計算し、信頼度 α_k を重み係数とする線形和により、指し手評価の期待値を

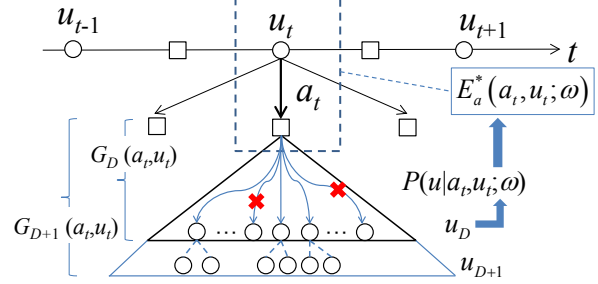


図 5 PG 行動期待値法への反復深化法の適用

Figure 5 An iterative-deepening search applied to PG expectation algorithm.

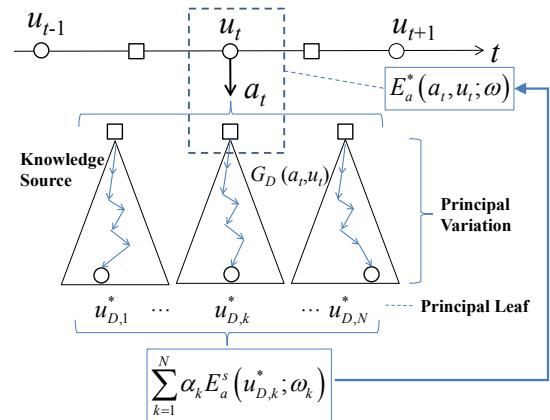


図 6 異なる評価関数による期待値操作

Figure 6 Expectation with different positional evaluation functions.

$$E_a^*(a_t, u_t; \omega) \approx \sum_{k=1}^N \alpha_k E_a^s(u_{D,k}^*(a_t, u_t; \omega_k); \omega_k) \quad (21)$$

と近似する。学習時には(21)を(11)へ代入して得られる特徴的適正度を用いる。探索アルゴリズム k は自らが探索した principal leaf $u_{D,k}^*(a_t, u_t; \omega_k)$ に含まれている特徴量に関するパラメータを更新する。これは、複数の探索アルゴリズム (知識エージェント) によるある種の“合議”による指し手決定[18]と、各探索アルゴリズムの評価関数の学習方法を与えており並列処理向きである。この際、異なる探索アルゴリズムの生成法として、評価関数にランダムノイズを付加する方法も考えられる。図 6 にこれらの考えをまとめた説明図を示す。なお、(21)の信頼度 α_k は学習パラメータと考えて本学習方式の枠組みで学習することも可能である。

6. 探索時にシミュレーション方策を用いる場合の学習

6.1 着手選択時の方策とシミュレーション方策

(8)で定義した方策 $\pi_a(a_t|u_t; \omega) = \exp(E_a^*(a_t, u_t; \omega)/T_a)/Z_a$ は、エ

エージェント A が手番で着手を選択する際の方策（以下，“着手決定方策”）である。その際には，(7)で定義された”指し手評価の期待値” $E_a^*(a_i, u_i; \omega)$ の値が必要であった。この期待値は，深さ D での leaf 局面 $u \in U_D(a_i, u_i)$ の局面評価値とその leaf 局面 u への遷移確率 $P(u|a_i, u_i)$ とから(7)で計算される。

4. で述べた”PG 行動期待値法”では，この $P(u|a_i, u_i)$ の計算にも着手決定方策(8)を用いていた。したがって，深さ D での全ての leaf 局面 $u \in U_D(a_i, u_i)$ の評価値を知る必要があり，厳密に求めようとする計算量が膨大となる。

そこで，本章では着手決定方策 $\pi_a(a_i|u_i; \omega)$ とは別に，leaf 局面の評価値を使用しない方策を探索用として用意する。通常，このような探索木生成のための方策は，“シミュレーション方策” (simulation policy) と呼ばれており，モンテカルロ探索や，実現確率を用いて前向きの枝刈りを行う場合の遷移確率の計算に用いられている[6]。

本章では，シミュレーション方策が局面評価関数 $E_a^s(u; \omega)$ 中の ω に依存しない場合を考える。つまり，激指のように，探索中の指し手選択の方法と leaf 局面の評価法とは独立であるとする。さらに，探索木中の現在のノード局面の情報だけを用いて指し手を選択する場合を考える。つまり，シミュレーション方策がマルコフ性を持っているとする。

今，シミュレーション方策を $\pi'_a(a|u; \theta)$ で表す。ただし， u は探索木中のノード局面（エージェント A の手番）， θ はシミュレーション方策に含まれるパラメータの総称であり ω とは異なる。これらを用いると，手番局面 u_i から指し手 a を経由した leaf 局面 u への遷移確率 $P(u|a, u_i)$ は，

$$P(u|a, u_i; \theta) = \pi'_a(a|u_i; \theta) \pi_b(b_i^0|v_i^0) \pi'_a(a^1|u_i^1; \theta) \pi_b(b_i^1|v_i^1) \cdots \pi'_a(a^{D-1}|u_i^{D-1}; \theta) \pi_b(b_i^{D-1}|v_i^{D-1}) \quad (22)$$

$$= \prod_{d=0}^{D-1} \pi'_a(a^d|u_i^d; \theta) \pi_b(b_i^d|v_i^d) \quad (23)$$

と表すことができる。ただし，上付きの添え字は探索木の深さを表し， $a^0=a$ ， $u_i^0=u_i$ ， $u_i^D=u$ とする。したがって，(7)の指し手評価の期待値は，

$$E_a^*(a_i, u_i; \omega, \theta) \equiv \sum_{u \in U_D(a_i, u_i)} P(u|a_i, u_i; \theta) E_a^s(u; \omega) \quad (24)$$

と2種類のパラメータ ω ， θ を含む。着手決定方策 π_a も同様であり，以下本章では $\pi_a = \pi_a(a|u; \omega, \theta)$ と記す。

次に，上記2種類のパラメータに関する学習則は以下のようにまとめることができる。

$$\Delta \omega = \varepsilon \cdot r \sum_{t=1}^{L_a} e_\omega(t) \quad (25)$$

$$e_\omega(t) \equiv \partial \ln \pi_a(a_i|u_i; \omega, \theta) / \partial \omega \quad (26)$$

$$= (1/T_a) \left\{ E_{\pi'_a} \left[\partial E_a^s(u; \omega) / \partial \omega | a_i, u_i \right] \right.$$

$$\left. - E_{\pi_a} \cdot E_{\pi'_a} \left[\partial E_a^s(u; \omega) / \partial \omega | x, u_i \right] \right\} \quad (27)$$

$$\Delta \theta = \varepsilon \cdot r \sum_{t=1}^{L_a} e_\theta(t) \quad (28)$$

$$e_\theta(t) \equiv \partial \ln \pi_a(a_i|u_i; \omega, \theta) / \partial \theta \quad (29)$$

$$= (1/T_a) \left\{ E_{\pi'_a} \left[E_a^s(u; \omega) \sum_{d=0}^{D-1} e'_\theta(d) | a_i, u_i \right] \right.$$

$$\left. - E_{\pi_a} \cdot E_{\pi'_a} \left[E_a^s(u; \omega) \sum_{d=0}^{D-1} e'_\theta(d) | x, u_i \right] \right\} \quad (30)$$

ただし，

$$E_{\pi_a} [z|u_i] \equiv \sum_{x \in A(u_i)} z \cdot \pi_a(x|u_i; \omega, \theta) \quad (31)$$

$$E_{\pi'_a} [z|a, u_i] \equiv \sum_{u \in U_D(a, u_i)} z \cdot P(u|a, u_i; \theta) \quad (32)$$

$$E_{\pi_a} \cdot E_{\pi'_a} [z|x, u_i] \equiv E_{\pi_a} [E_{\pi'_a} [z|x, u_i] | u_i] \quad (33)$$

$$= \sum_{x \in A(u_i)} \pi_a(x|u_i; \omega, \theta) \left[\sum_{u \in U_D(x, u_i)} z \cdot P(u|x, u_i; \theta) \right] \quad (34)$$

$$e'_\theta(d) \equiv \partial \ln \pi'_a(a^d|u_i^d; \theta) / \partial \theta \quad (35)$$

と定義した。なお， $E_{\pi_a}[\cdot|u]$ は局面 u において着手決定方策 $\pi_a(a|u; \omega, \theta)$ による期待値操作を， $E_{\pi_a}[\cdot|a, u]$ は局面 u で a を選択し，それ以降はシミュレーション方策 $\pi'_a(a|u; \theta)$ を用いて指し手選択を行ってシミュレーションした場合の期待値計算を表している。

(27)は，高報酬を得た対局での着手の選択確率を高めるためにその手の評価値を高めたい。そこで，シミュレーション方策は固定しておいて，その手から始まるシミュレーションにより得られる leaf 局面 u の局面評価値 $E_a^s(u; \omega)$ を高めるように ω を $E_a^s(u; \omega)$ の増加する勾配方向へ動かすと解釈できる。

一方，(30)では，同様の目的ではあるが，今度は局面評価関数を固定しておいてシミュレーション方策の方を調整する。すなわち，シミュレーション方策中のパラメータ θ の更新量 $\Delta \theta$ を，高評価の leaf 局面への遷移確率 $P(u|a_i, u_i; \theta)$ の値を増加させるように調節している。このときの調整法は， u_i を初期状態， a_i を初期行動とする方策勾配法による強化学習を行っていると思なすことができる。実際，(30)の π'_a による期待値 $E_{\pi'_a} [E_a^s(u) \cdots | a_i, u_i]$ は，一般的な方策勾配法の学習則(3)~(5)において， $\pi'_a(a|u; \theta)$ を方策 π ， $E_a^s(u)$ を報酬 r とする期待報酬値の勾配 $\partial E_{\pi} [r] / \partial \theta = \partial E_{\pi} [E_a^s(u) | a_i, u_i] / \partial \theta$ になっている。つまり，シミュレーション方策の強化学習を方策勾配法を用いて指し手の探索シミュレーション内で行

うことができることを表している。

なお、シミュレーション方策がマルコフ性を持っておらず、探索のルート局面 u_t から現ノード局面 u_t^d までの状態行動履歴 $h_t^d \equiv \{u_t, a_t, u_t^1, a_t^1, \dots, u_t^{d-1}, a_t^{d-1}\}$ に依存する場合も、 $\pi'_a(a|u_t^d; \theta)$ を $\pi'_a(a|u_t^d, h_t^d; \theta)$ と置き換えると、(25)~(35)はそのまま成り立つ。例えば、直前の手に応じて指し手を変化させる場合などがこれにあたる。

シミュレーション方策の例としては、次の Boltzmann 分布を考える。

$$\pi'_a(a_t^d | u_t^d, h_t^d; \theta) = \exp(E'_a(a_t^d, u_t^d, h_t^d; \theta) / T'_a) / Z'_a \quad (36)$$

$$Z'_a \equiv \sum_{x \in A(u_t^d)} \exp(E'_a(x, u_t^d, h_t^d; \theta) / T'_a) \quad (37)$$

ただし、目的関数 $E'_a(a_t^d, u_t^d, h_t^d; \theta)$ 中の特徴量は、囲碁では石の局所的な配置パターンなどであり、将棋では激指などでの指し手選択のための特徴量、例えば、「王手かどうか」「ひもをつける手かどうか」[6]などである。これらは、人間の将棋では、「手筋」「型」のような経験的で断片的なミニ知識による指し手、あるいは、「直観」、「第一感」などの「深い読み」を伴わない処理で指される手と考えられる。

実際の対局における探索時には、このようなシミュレーション方策を用いて、探索のルートノードから探索木の各ノードへの遷移確率の値を次のように近似的に計算できる。

$$P(u_t^d | a, u_t; \theta) = \pi'_a(a^{d-1} | u_t^{d-1}; \theta) \pi_b(b_t^{d-1} | v_t^{d-1}) P(u_t^{d-1} | a, u_t; \theta) \quad (38)$$

上記の遷移確率値は、モンテカルロ探索による指し手評価や、 $\alpha\beta$ 探索の際の前向き枝刈り処理に用いることができる。例えば、モンテカルロ探索においては、leaf局面に対して遷移確率値を計算すれば、(24)の指し手評価の期待値 $E^*_a(a, u; \omega, \theta)$ を近似的に計算することができる。したがって、膨大な回数の試行を行って局面評価値の平均操作を行う必要がなく、試行回数の削減に役立つ。また、 $\alpha\beta$ 探索においては、探索木の途中のノードへの遷移確率値が閾値以下になればそのノード以下の探索を打ち切るなどの処理が考えられる。この方法には、激指の実現確率による枝刈り処理とは異なり、兄弟手の良し悪しの度合いにより探索の深さを制御できる利点がある。例えば、飛車を取る手があれば端歩を突く手は殆ど読む必要はないが、他に有力な手がない局面では深く読む必要があると言うような場合に有効であると考えられる[19]。

6.2 シミュレーション方策と局面評価関数の教師付き学習

6.1 ではシミュレーション方策を用いた方策勾配法による強化学習を考察した。本節では強化学習ではなく、ある局面において正解手が与えられた場合の学習、すなわち、教師付き学習を考える。6.1 で用いた方策勾配法では、エピソードごとの期待報酬値の勾配を計算したが、学習システムの正解手に対する選択確率値の勾配を同じように計算

することができる。なお、簡単のために本節ではシミュレーション方策がルート局面からそのノード局面までの指し手履歴によらないマルコフ性のある場合を考える。そうでない場合も全く同様に導出できる。

通常、局面評価関数に教師付き学習を適用する際は、プロ棋士の棋譜データベース等から、局面とそこで指された指し手を唯一の正解手として局面・指し手ペアの訓練データを作成する。しかし、ここではより一般的な場合を扱う。すなわち、正解手を1つに限定せずに、正解と思われる複数の指し手に対してそれらを選択する確率分布を学習させることにする。そこで、今、正解の着手決定方策を π^* 、学習システムの着手決定方策を π_a とし、シミュレーション方策 π'_a は(36)を仮定する。次の誤差関数 U_{err} を考える。

$$U_{err}(\pi^*, \pi_a) \equiv \sum_{a \in A(s)} \pi^*(a|s) \ln[\pi^*(a|s) / \pi_a(a|s; \omega, \theta)] \quad (39)$$

$U_{err}(\geq 0)$ は、正解の方策 π^* と学習システムの方策 π_a との距離を表すカルバック・ライブラー情報量 (Kullback-Leibler divergence) である。ただし、6.1 と同じく、

$$\pi_a(a|s; \omega, \theta) = \exp(E_a^*(a, s; \omega, \theta) / T_a) / Z_a \quad (40)$$

$$Z_a \equiv \sum_{x \in A(s)} \exp(E_a^*(x, s; \omega, \theta) / T_a) \quad (41)$$

$$E_a^*(a, s; \omega, \theta) = \sum_{u \in U_D(a, s)} P(u|a, s; \theta) E_a^s(u; \omega) \quad (42)$$

と仮定する。このとき、 ω に関する勾配ベクトルは、

$$\partial U_{err} / \partial \omega = - \sum_{a \in A(s)} \pi^*(a|s) \cdot \partial \ln \pi_a(a|s; \omega, \theta) / \partial \omega \quad (43)$$

と表されるが、右辺中の対数微分の項は、(26),(27)において $a_t = a$, $u_t = s$ と置き換えた式で表される。すなわち、(43)を用いて局面評価関数 $E_a^s(u; \omega)$ 中の ω の更新量は、

$$\Delta \omega = -\varepsilon \cdot \partial U_{err} / \partial \omega \quad (44)$$

と計算すればよい。

また、 θ に関する勾配ベクトルは、

$$\partial U_{err} / \partial \theta = - \sum_{a \in A(s)} \pi^*(a|s) \cdot \partial \ln \pi_a(a|s; \omega, \theta) / \partial \theta \quad (45)$$

と表されるが、右辺の対数微分の項は、(29),(30)において $a_t = a$, $u_t = s$ と置き換えた式で表される。したがって、(45)の値を用いて、シミュレーション方策 $\pi'_a(a|s; \theta)$ 中のパラメータ θ の更新量は次のように計算すればよい。

$$\Delta \theta = -\varepsilon \cdot \partial U_{err} / \partial \theta \quad (46)$$

なお、探索(読み)にシミュレーション方策を用いないで着手決定方策を用いる場合でも局面評価関数中の ω の教師付き学習は可能である。その場合は、(43)の右辺の対数微分は(11)と同一であり、4.2 や 5.での手法が使える。

6.3 シミュレーション方策の教師付き学習に関する先行研究との関係

激指は探索時にその局面における指し手の選びやすさである実現確率を用いた枝刈りを行い、探索の深さの制御を積極的に行っている。この実現確率の計算はシミュレーション方策を用いた状態遷移確率の近似計算の一種とみなすことができる。激指ではシミュレーション方策中のパラメータ θ を、人間の棋譜データベースから統計的処理により求めている。さらに、Bonanza の学習法をベースにした局面評価関数中の ω の教師付学習も行っているが、これら2つの学習は全く独立しており直接的な関係はない。

また、方策勾配を利用したシミュレーション方策の教師付き学習の先行研究として、Policy-gradient simulation balancing法が囲碁の場合に提案されている[20]。そこでは、訓練局面 s に対して正解となる局面評価値 $V^*(s)$ が必要となり、かつ、着手決定方策として局面評価関数を利用していないので局面評価関数の学習は行っていない。

さらに、着手方策中の局面評価関数をTD(λ)法で求めている、その学習結果として得られたヒューリスティクスをUCT探索におけるシミュレーション方策へ組み込む将棋の研究がある[21]。しかし、そこでの着手決定方策の学習は探索のない学習であり、かつ、シミュレーション方策の学習結果が着手決定方策の学習に反映されることはない。

本方式では、正解の方策 π^* に対して、着手決定方策で用いられる評価関数 $E_a^s(u; \omega)$ 中のパラメータ ω と、シミュレーション方策 $\pi_a'(a|u; \theta)$ 中のパラメータ θ とが同時に連携して(39)の誤差を減らすように学習を行うことができる。

7. おわりに

本論文では、これまでに強化学習の手法としてコンピュータ将棋に用いられてきた TD(λ)法や TDLeaf(λ)法ではなく、方策勾配法を適用する手法についての理論的な検討を行った。その結果、最初に、将棋のように着手決定に探索を要する場合については、“指し手評価の期待値”による確率の方策を用いた“PG 行動期待値法”と呼ぶ着手決定方式を提案した。次に、その近似計算法として“PGLeaf法”など3つの方法を提案した。

さらに、この着手決定のための方策とは別に、探索時にシミュレーション方策を用いる場合への方策勾配法の適用についても考察し、両方の方策に含まれるパラメータの学習則を導出することができた。学習後のシミュレーション方策は、モンテカルロ探索における期待値の計算や、探索木中のノード局面への遷移確率値の計算に用いることができる。したがって、モンテカルロ探索の試行回数の削減や、探索時の深さ制御のための前向き枝刈り処理の精度向上に役立つと考えられる。

最後に、ここまでに用いた方策勾配の計算は、強化学習だけでなく、局面ごとに正解手の方策を確率分布の形で与

える教師付き学習問題へも適用することができることを示した。この場合のシミュレーション方策と局面評価関数に関する学習則も導出することが出来た。

今後は、本論文で展開した学習則や探索法を実装し、実験により有効性を検証、評価して行く予定である。

参考文献

- 1) 松原仁 編著：コンピュータ将棋の進歩⑥プロ棋士に並ぶ、共立出版(2012)。
- 2) 第2回電王戦の公式ページ：<http://ex.nicovideo.jp/denousen2013/>
- 3) 保木邦仁：局面評価の学習を目指した探索結果の最適制御，第11回ゲームプログラミングワークショップ，pp.78-83(2006)。
- 4) Sutton, R. S. and Barto A. G. : Reinforcement Learning, The MIT Press, Massachusetts (1998)。
- 5) Baxter, J., Tridgell, A., and Weaver, L., : KnightCap: A chess program that learns by combining TD(λ) with game-tree search, Proceedings of the Fifteenth International Conference (ICML '98), pp.28-36 (1998)。
- 6) 鶴岡慶雅：「激指」の最近の改良について，松原仁 編著：コンピュータ将棋の進歩⑥プロ棋士に並ぶ，第4章，共立出版(2012)。
- 7) Williams, R. J. : Simple Statistical Gradient- Following Algorithms for Connectionist Reinforcement Learning, Machine Learning, Vol.8, pp.229-256 (1992)。
- 8) 木村元，山村雅幸，小林重信：部分観測マルコフ決定過程下での強化学習-確率的傾斜法による接近，人工知能学会誌，Vol.11, No.5, pp761-768 (1996)。
- 9) Sutton, R.S., McAllester, D., Singh, S. and Mansour, Y. : Policy Gradient Methods for Reinforcement Learning with Function Approximation, Proc. of Advances in Neural Information Processing Systems 12 (NIPS' 99), pp.1057-1063 (2000)。
- 10) Konda, V. R. and Tsitsiklis, J. N.: Actor-Critic Algorithms, Proc. of Advances in Neural Information Processing Systems 12 (NIPS' 99), pp.1008-1014 (2000)。
- 11) 阿部健一：強化学習一価値関数推定と政策探索”，計測と制御，第41巻，第9号，pp.680-685 (2002)。
- 12) Kakade, S.: A natural policy gradient, Proc. of Advances in Neural Information Processing Systems 14 (NIPS'01), pp.1531- 1538 (2002)。
- 13) Peters, J., and Schaal, S.: Policy Gradient Methods for Robotics, Proc. of the IEEE International Conference on Intelligent Robotics Systems (IROS 2006), pp.2219-2225(2006)。
- 14) 五十嵐治一，石原聖司，木村昌臣：非マルコフ決定過程における強化学習一特徴的適正度の統計的性質一，電子情報通信学会論文誌 D, Vol. J90-D, No.9, pp.2271-2280 (2007)。
- 15) 石原聖司，五十嵐治一：マルチェージェント系における行動学習への方策こう配法の適用-追跡問題-，電子情報通信学会論文誌 D-I, Vol. J87-D1, No.3, pp.390-397 (2004)。
- 16) 五十嵐 治一，半田 雅人，石原 聖司，篠埜 功：マルチェージェントシステムにおける行動制御—PSOにおける重み係数の強化学習—，電子情報通信学会論文誌 D, Vol. J94-D, No. 10, pp. 1612-1621 (2011)。
- 17) 美添一樹：モンテカルロ木探索-コンピュータ囲碁に革命を起こした新手法-，情報処理，Vol.49, No.6, pp.686-693 (2008)。
- 18) 伊藤毅志：コンピュータ将棋における合議アルゴリズム，人工知能学会誌，Vol.26, No.5, pp.525-539 (2011)。
- 19) 一丸貴則：WCSC21 ツツカナアピール文書，http://www.computer-shogi.org/wcsc21/appeal/tsutsukana/WCSC21_tsutsukana_20110327.pdf
- 20) Silver, D., Tesauro, G.: Monte-Carlo simulation balancing. In: Bottou, L., Littman, M. (eds.) Proceedings of the 26th International Conference on Machine Learning, Montreal, Canada, pp. 945-952. Omnipress (June 2009)。
- 21) 燧 暁彦，三輪 誠，鶴岡慶雅，近山 隆：TD(λ)学習を用いた Ms. Pac-Man AI のモンテカルロ木探索の方策の学習，情報処理学会研究報告，Vol.2013-GI-29, No.2, pp.1-8 (2013)。