

## モンテカルロ+UCTにおける探索木のだまし構造

池田 心<sup>†1</sup>      橋本隼一<sup>†1</sup>      土井佑紀<sup>†1</sup>

マルチアームバンディット問題 (MAB) を対象に発展した  $\epsilon$ -B 戦略は、探索と収穫のジレンマに対する一つの有力な回答である。近年、 $\epsilon$ -B 戦略をゲームの木探索に応用した  $\epsilon$ -B が盛んに研究され、モンテカルロ法と組み合わせた囲碁プログラムが登場している。しかし、報酬を得られる確率が静的な MAB と異なり、 $\epsilon$ -B 原理が支配する二人ゲームでは  $\epsilon$ -B 値に基づく探索が非効率である可能性もある。本稿では、囲碁において  $\epsilon$ -B が好ましくない挙動を指し、その本質を抜き出してベンチマーク化する。

### Deceptive Structure of Search Tree for UCT with Monte-Carlo

IKEDA KOKOLO,<sup>†1</sup> HASHIMOTO JUNICHI<sup>†1</sup>  
and DOI YUUKI<sup>†1</sup>

B strategy which has been developed for solving multi-armed bandit problem (MAB) is a major solution for the exploitation-exploration dilemma. It is the application of  $\epsilon$ -B for a tree search, and is intensively studied especially for Go with Monte Carlo method. However, in contrast to the MAB, the ratio for rewards is not static but dynamic in two-player game such as Go. In this paper, we present a failure case of  $\epsilon$ -B for Go, generalize it and build a benchmark tree search problem with deceptive structure.

#### 1. はじめに

マルコフ決定過程 (MDP) などの繰り返し環境にエージェントが置かれ、不完全情報下でできるだけ多くの報酬を得ようとする場合、現在分かっている範囲で高い報酬を得る“収穫

(exploitation)”とより高い報酬が得られる可能性を模索する“探索 (exploration)”はしばしば両立せず、探索と収穫のジレンマと呼ばれる。このジレンマを簡潔に表したモデルにマルチアームバンディット問題 (MAB) がある。これは 1 状態 (局面) 複数行動の MDP であり、各行動  $a_i$  後に一定の確率  $p_i$  で報酬 1 を得られるというものである。

UCB(Upper Confidence Bound)<sup>1)</sup>とは、行動 (指手) 選択回数  $n$ 、行動  $a_i$  の選択回数  $n_i$  と報酬 (勝率) の期待値  $e_i$  に対して、 $e_i + C\sqrt{\frac{\ln(n)}{n_i}}$  で定義される値であり ( $C$  は設計パラメータ)、 $n$  に対して  $n_i$  が相対的に小さいときに報酬の期待値が“上がりうる幅”を組み込んだものである。UCB 値が最大となる行動を選択していくことで、期待値の高い行動の exploitation と発展の見込みのある行動の exploration を制御することを可能にしており、前述の MAB など良い性能を示している。

UCT(Upper Confidence bound for Tree search)<sup>4)</sup>は、ゲームなどの木探索 (Tree Search) に UCB の考え方を導入したものである。例えば囲碁や将棋などの多くのゲームでは、状態分岐数 (取り得る着手の数) や木の深さ (終局までの手数) が非常に大きく、厳密解法を用いることは事実上不可能である。そのため、良さそうな着手を辿りながら深くまで探索し (exploitation)、かつ可能性のあるそれ以外の分岐にも探りを入れる (exploration) が必要になる。UCT では、探索したノードそれぞれに訪問回数と報酬の平均値を保持し、子ノードの UCB 値が最も高くなるものを選択することでこのバランスを取る。

多くの場合、UCT は優れた探索性能を示す、すなわち、深く読むべきところに多くの探索資源を割きつつ、二番手三番手の着手にも現在の報酬平均に応じた資源を投じる。盤面の評価をランダムゲームによって求めるモンテカルロ法との組み合わせは、多くの囲碁プログラムで用いられているところである<sup>3)5)</sup>。しかしながら、UCB 戦略が開発された MAB とは異なり、minimax 原理に支配された二人ゲームでは各行動を取った場合の報酬の期待値は静的ではない。UCB 値における補正項  $C\sqrt{\frac{\ln(n)}{n_i}}$  は報酬の期待値が静的であることを前提としており、動的に変化する場合には不十分である可能性がある。

本稿では、囲碁で UCT が好ましくない挙動を示す場合を発見したことを報告したうえで、その本質を抜き出してベンチマーク化する。木探索を学問として研究していくためには、具体的なゲームを題材に勝率を云々すると同時に、容易に共有・調整・実験ができるベンチマーク問題の整備が必須であると考えている。

<sup>†1</sup> 北陸先端科学技術大学院大学  
JAIST Hokuriku

## 2. 囲碁における問題点の発見

### 2.1 nomitan プロジェクト

我々のチームでは、2年前に野口・松井両氏によって開発された囲碁プログラム **nomitan** をベースに囲碁プログラムの研究を進めている。**nomitan** はUCTに基づく探索と、ランダムゲームによるモンテカルロ評価を組み合わせたプログラムであり、ランダムゲームでは勾配法で棋譜学習した係数による着手の選択確率計算、UCTではUCB1-Tuned値による探索ノード選択<sup>1)</sup>を行っている。探索終了後は、現局面で最も訪問回数が高い着手(子ノード)を打つが、これは最も有望な着手を最も丹念に読んでいるはずだという信念に基づく<sup>2)</sup>。

### 2.2 次の一手問題

囲碁プログラムの強さは、通常、プログラム同士の対戦によって計測する。しかし、より詳細にその挙動を知りたい場合には、局面を与えて次の一手を探索させることが有効であるし、また低コストである。

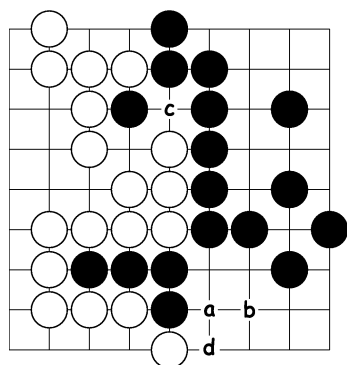


図1 次の一手問題の例：コミ3目半。b(7,8)が正解

図1は次の一手問題の例である。この問題ではa(6,8)が普通の手に見えるが、後手を引いてc(5,3)に回られてしまうので失敗であり、b(7,8)がc,dを見合いにした正解である。なお初手dはaに切られて失敗である。この問題は人間にとっては5級程度の問題であり、決して難しい部類のものではない。

### 2.3 好ましくない探索結果

nomitanに80万回のプレイアウト(ランダムゲーム)を与えて探索をさせたところ、有

力な手a(6,8),b(7,8)について、その勝率の期待値と訪問回数は図2のように変化した。なお、nomitanの場合80万プレイアウトの探索には通常のPCでおおよそ1分を要する。

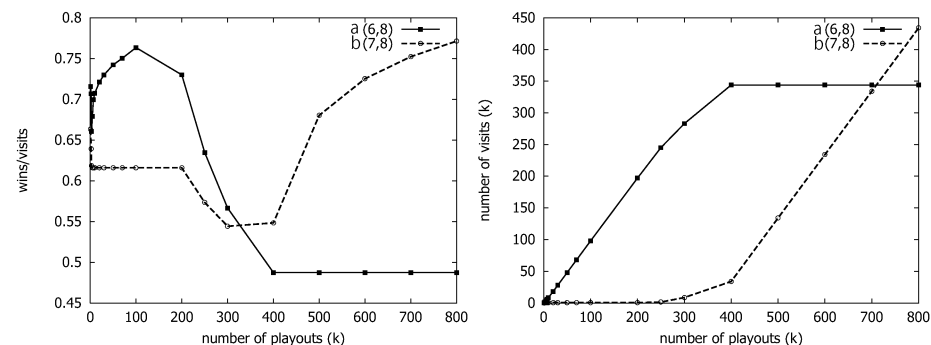


図2 nomitanによる探索、着手a(6,8),b(7,8)の勝率の期待値(左)と訪問回数(右)

探索の序盤20万プレイアウトくらいまでは、aの勝率がbよりかなり高く(左図)、殆どの探索資源がaに割かれている(右図)。中盤20万~40万プレイアウトでaの勝率は大きく下がり、bにも探索資源が割かれるようになる。終盤40万~70万プレイアウトでは逆にbの勝率が上昇し、aには探索資源が割かれない。

このような挙動は別段奇妙なものではない。有望そうな手を読んでいるうちに相手側の好手に気づくこと、逆に駄目そうな手を読み返して自分側の好手に気づくこと、これらは囲碁に限らず二人ゲームには頻繁に生じる現象である。この探索について問題を挙げるとすれば、(1)序盤20万プレイアウトまでbは殆ど探索されないこと、(2)勝率が逆転したあとも「訪問回数最大の手を打つ」というルールに従えば50万~70万プレイアウトのような明らかにbを選ぶべき場合でもaを選んでしまうこと、である。これらの問題点を整理すべく、次章ではより簡潔なベンチマーク問題を提案する。

## 3. だまし構造のあるベンチマーク探索木

確率的最適化の分野では、“一見悪そうな領域に最適解があり、一見良さそうな領域に探索を重点化すると失敗する”ことをだまし構造があると言う<sup>6)</sup>。前章で紹介した例も、着手aは正解ではないが一見良く見え、着手bは正解なのに一見悪く見えることが問題であった。そこで、このような探索木もだまし構造があるということにする。

また、クラス分類問題等の呼び方になぞらえて、一見良く見えることを **false positive** (偽陽性)、一見悪く見えることを **false negative** (偽陰性) と呼ぶことにする。例えば、医療ではしばしば簡易検査を行い陽性のものを精密検査するが、このような場合、偽陽性は精密検査のコストがかかる、偽陰性は病気を見逃すリスクがあるということになる。探索においては、偽陽性は本来探索すべきでない箇所を探索するコストがかかる、偽陰性は正解を見逃すリスクがある、と読み替えることができるだろう。

### 3.1 ベンチマーク

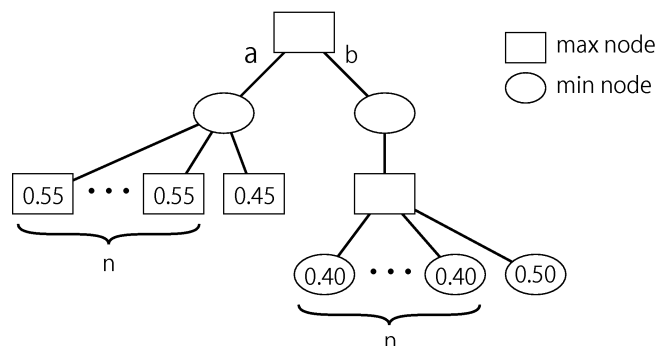


図 3 だまし構造のあるベンチマーク問題の例

図 3 は、だまし構造のある二人ゲーム探索木の例である。末端ノードでは、書かれた数字の確率に従って報酬 1 を返す (自分が勝つ確率であると考え)。木の構造があらかじめ分かっているならば、着手 a の勝率は相手が min を選択するので 0.45、着手 b の勝率は自分が max を選択するので 0.50 であり、b が正解である。しかし、着手 a 側には多くの勝率 0.55 の手があり、相手が正しく選択しない序盤では全体として有望に見える (false positive)。一方着手 b 側には多くの勝率 0.40 の手があり、自分が正しく選択しない序盤では全体として悪く見える (false negative)。

図 4 は、UCT (tuned ではない) で本ベンチマーク ( $n = 5$ ) を解かせた場合の、着手 a, b の勝率の期待値と訪問回数である。前章の実験と同様、序盤は着手 a 側の勝率期待値が高く (左) 探索も重点化されている (右)。中盤、勝率が逆転し、終盤は着手 b 側にだけ探索が集中している。そして、45000~65000 回訪問のあたりでは、訪問回数最大の手 (a) が正解でない、という点も囲碁のケースと同じである。

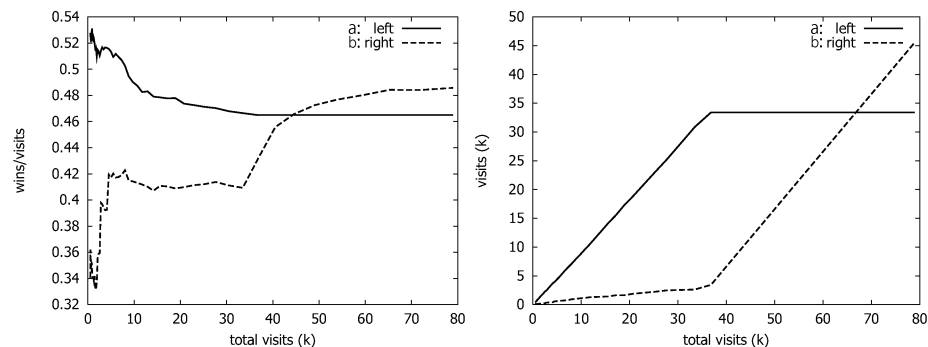


図 4 UCT による探索、勝率の期待値 (左) と訪問回数 (右)

図 5 は、100 回同様の探索を行い、40000 回訪問時点での状況をプロットしたものである。横軸は、着手 a の勝率期待値 - 着手 b の勝率期待値を表し、これが正の場合は、勝率の推定・比較に失敗しているということになる。縦軸は、着手 a の訪問回数 - 着手 b の訪問回数を表し、これが正の場合は正解を打てないということになる。図を見ると、第一象限 (勝率推定も打つ手も不正解)、第三象限 (勝率推定も打つ手も正解) に加え、第二象限 (勝率推定には成功しているが、打つ手は不正解) にも相当数の点があることがわかる。これもまた、訪問回数最大を選択することの弱点を証明している。

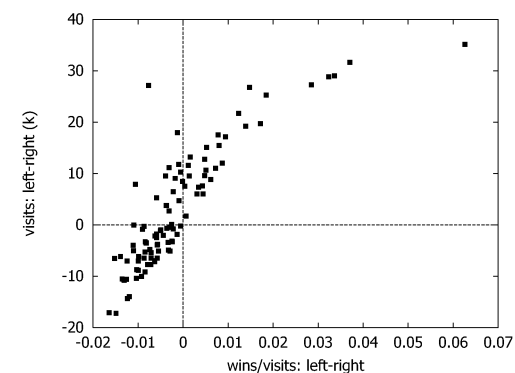


図 5 40000 回訪問時点での、勝率差 (横軸) と訪問回数差 (縦軸)、100 試行

### 3.2 バリエーション

本稿で提示した木 (図 3) は, だまし構造のあるベンチマークの一例に過ぎない. 他にも重要なものとしてだまし構造が positive negative の片側しかなかったり, 0.3 や 0.7 などの不純物的なノードが付いていたり, だまし構造が 2 手目と 4 手目に二重になっていたりなど, 多くのものが考えられる. 一方で, 現実に存在しないようなケースをあこれ想像して極端なベンチマークを作ることもあまり生産的とは言えない. 今後は, 現実の問題を分析しつつ, 系統的なベンチマークの整備を行う予定である.

### 4. おわりに

本稿では, だまし構造が現実の問題に存在することを示し, その本質を抜き出してベンチマーク化した, それに対処する方法については述べてこなかった. ここでは, 幾つかの方向性を提示することで結びとし, 具体的な提案と実験は別稿に譲ることとする.

まず, UCB 値の補正項パラメータ  $C$  を大きくすることは, 有望そうでないノードにもある程度の探索資源を割くことで, だまされる程度を小さくすることができると考えられる. しかし一方, 本当に有望でないノードにも資源を割くため, だましがいない場合の性能の低下が深刻であると予測する.

次に, 勝率の推定には成功しているのに着手を誤る点については, 訪問回数最大の手を選択するのではなく, LCB 値 (UCB 値の補正項の符号をマイナスにしたもの) 最大の手を選択するなどの方法で回避できる. ただし, これはあくまで探索終了後の話であり, 探索中にだまされることの改善にはならない.

一見良い自分の手に対して, 重点探索ののち相手の好手が発見された場合, 相手はその手を選ぶ可能性が高い訳であり, 今まで探索した他の手の勝率を尊重し平均化して親ノードに渡すことは適切でない. このような場合, 過去の記憶を忘れる, あるいは  $\min(\max)$  値を親ノードに渡すなどの手法も有望であろう. ただしこの手法も, だましがいない場合の性能低下とのトレードオフは避けられない.

いずれにせよ, 各着手の勝率が静的なものであるとして開発された UCB をそのまま二人ゲームに適用することには限界があると考えられる. さまざまなベンチマーク問題を作り, それらに対応できる手法を新規に考える必要があるのではないかな.

### 参考文献

- 1) Auer, P., Cesa-Bianchi, N. and Fischer, P.: Finite-time Analysis of the Multiarmed Bandit Problem, *Machine Learning*, Vol.47, pp.235–256 (2002).
- 2) Coulom, R.: Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search, *CG06* (2006).
- 3) Gelly, S., Wang, Y., Munos, R. and Teytaud, O.: Modification of UCT with Patterns in Monte-Carlo Go, Technical Report RR-6062, Inria (2006).
- 4) Kocsis, L. and Szepesvári, C.: Bandit based Monte-Carlo Planning, *European Conference on Machine Learning*, pp.282–293 (2006).
- 5) Tom, D. and Müller, M.: A Study of UCT and Its Enhancements in an Artificial Game, *ACG 2009*, pp.55–64 (2009).
- 6) 池田 心, 小林重信: GA の探索における UV 現象と UV 構造仮説, 人工知能学会論文誌, Vol.17, No.3G, pp.239–246 (2002).