

モンテカルロ法とニューラルネットワークによる $\alpha \beta$ 法の効率化

日大生産工 (院) ○平林 佑一
日大生産工 松田 聖

1 まえがき

モンテカルロ法は乱数により近似解を求める手法である。評価関数を使用することなく、ゲームが終了するまでシミュレーションすることによって評価値が得られるため、コンピュータ囲碁のAIにしばしば用いられる。しかし、単体で適用した場合には探索空間が大きいほどシミュレーションを多く施行しなければ尤もらしい評価値を得ることができず、時間を多大に費やしてしまう。そこで筆者は嘗てニューラルネットワークと組み合わせることにより、ニューラルネットワークの出力(評価値)が良い手に対してモンテカルロ法を実行する手法を提案した。これにより探索効率上昇が見込めたが再考の余地が残る結果となった。

本研究では $\alpha \beta$ 法にニューラルネットワークを組み合わせることにより尤もらしい手から探索させることで枝刈り効率上昇を目的としている。また、3手先読み $\alpha \beta$ 法の葉ノードの評価値をモンテカルロ法で求めるため、枝刈りによるモンテカルロ法の効率上昇も兼ねている。なお、ニューラルネットワークの学習や実際のシミュレーションをする際の例として9路盤の囲碁を用い、ゲーム中盤から始めることとする。

2 モンテカルロ法

モンテカルロ法とは乱数を用いてシミュレーションを繰り返し、近似解を求める計算手法である。本研究では囲碁の盤面評価をモンテカルロ法で行う。

囲碁に適用する場合、囲碁のルールに則った全ての着手可能な場所(ただし眼は除く)に乱数で石を黑白交互に打つ。そしてそれを終局まで何度もシミュレーションする。終局時に最終的な勝率を見て最適手を選ぶ。

合法手 i のプレイアウトを行った回数を s_i 、プレイアウトによって得られた報酬の和を X_i とする。報酬はプレイアウトの結果、勝利であれば1、敗北であれば0を与える。

このとき、合法手の勝率は

$$\overline{X_i} = \frac{X_i}{s_i} \quad (1)$$

によって表わされる。この値は合法手を打った直後の盤面の評価値とする。なお、この評価値は後述するニューラルネットワークの教師信号として利用する。

Fig1にモンテカルロ法の構成を示す。

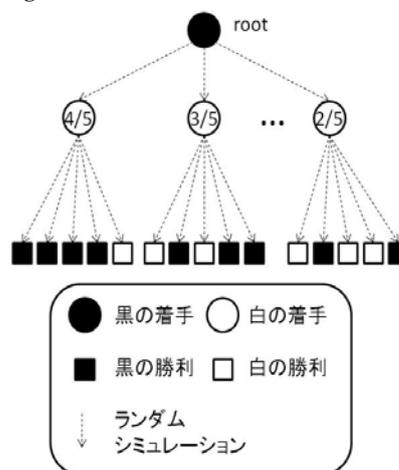


Fig1.モンテカルロの構成

このモンテカルロ法ではプレイアウトによって得られた結果のみを用いるため、プレイアウトの回数に依存して評価の精度が変わる。囲碁では一回の探索での時間が数秒~数分しかないため、評価値の精度を高くするにはプレイアウトによる合法手の選択を完全ランダムではなく、良さそうな手を重点的に探索する必要がある。そのため、本研究では現在の状態から次の手を選択する際に $\alpha \beta$ 法とニューラルネットワークを用いて学習することで不要なプレイアウトを省くことを目的としている。

3 $\alpha \beta$ 法

ミニマックス法は得られた評価値から相手の手番には最小値、自分の手番には最大値を取るように木を逆に辿ることで手を選ぶ手法である。このミニマックス法の探索時間を省略した方法を $\alpha \beta$ 法という。

Efficient $\alpha \beta$ Method by Monte Carlo Method and Neural Network

Yuichi HIRABAYASHI and Satoshi MATSUDA

