

AlphaZeroを用いた 京都将棋AIの作成と評価

4 年 1 4 組 4 番石元稜

目次

- 京都将棋について
 - 京都将棋の概要
 - 京都将棋のルール
- 準備
 - 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
 - モンテカルロ木探索
 - ニューラルネットワーク
- AlphaZeroの仕組み
- 結果
- 参考文献

目次

- 京都将棋について

- 京都将棋の概要
- 京都将棋のルール

- 準備

- 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
- モンテカルロ木探索
- ニューラルネットワーク

- AlphaZeroの仕組み

- 結果

- 参考文献

京都将棋の概要

- 5 x 5 マスの盤面
- 駒は「玉」「香と」「銀角」「金桂」「飛歩」の5種。

「香と」…表「香」，裏「と」

「銀角」…表「銀」，裏「角」

「金桂」…表「金」，裏「桂」

「飛歩」…表「飛車」，裏「歩」

と表裏一体になっている。

- 駒の初期配置は右図

5	4	3	2	1	
香	と	玉	金	歩	一
					二
					三
					四
と	銀	玉	金	歩	五

図1.京都将棋の初期配置

京都将棋のルール

基本的なルール，駒の動きは将棋と同じ。

違う点は以下の通り

- ・駒は1手動くごとに裏返す
- ・とった駒は裏表どちらでも打ってよい
- ・二歩はなし、千日手は引き分け。

目次

- 京都将棋について
 - 京都将棋の概要
 - 京都将棋のルール
- 準備
 - 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
 - モンテカルロ木探索
 - ニューラルネットワーク
- AlphaZeroの仕組み
- 結果
- 参考文献

「AlphaGo」 と 「AlphaGo Zero」 と 「AlphaZero」 の違い

- AlphaGo

ハンデなしでプロに勝った初めてのコンピュータ囲碁プログラム。

- AlphaGo Zero

プロ棋士の学習データを使わず自己対戦のみで学習し、AlphaGoに100勝0敗。

- AlphaZero

「AlphaGo Zero」の改造バージョン。

囲碁だけでなくチェスや将棋も学習できるようになった。

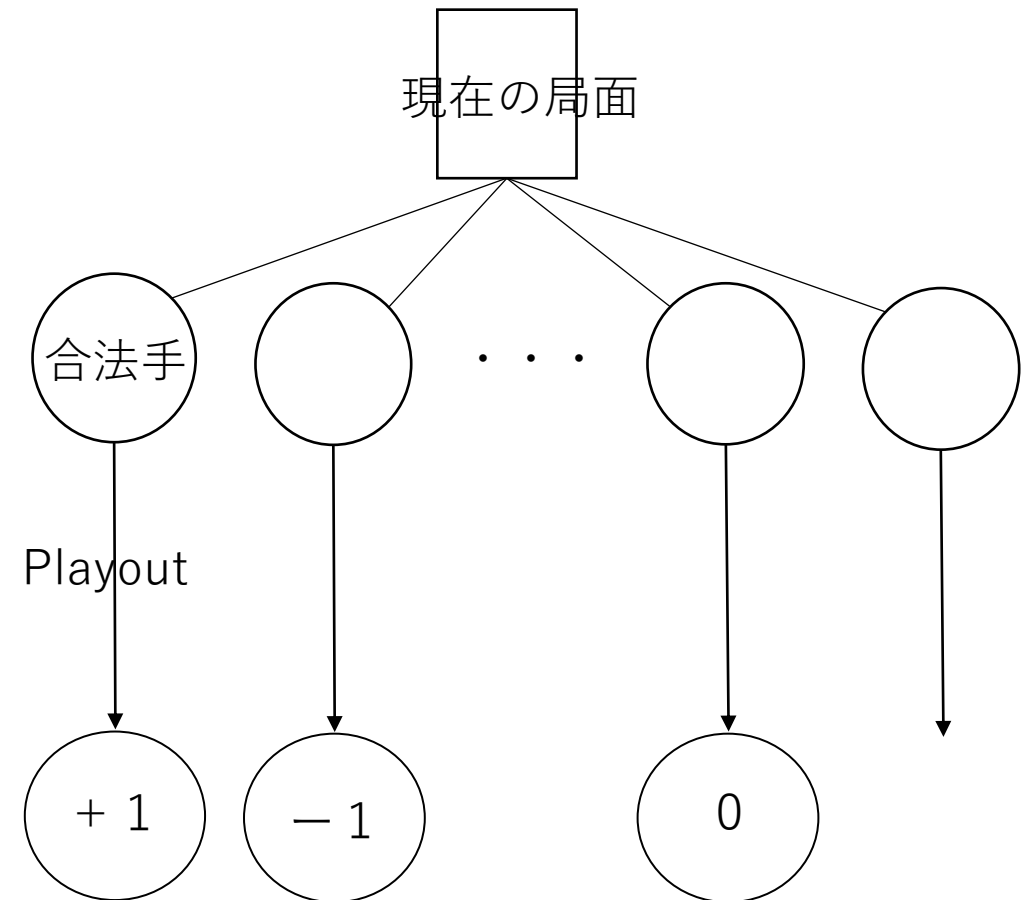
モンテカルロ木探索の前に、

• 原始モンテカルロ探索

よりよい次の一手を探索する。

合法手（現在の局面から指せる手）を
指した局面から
何度か終局までランダムにプレイし
（playout）、
勝ち(+1)負け(-1)の合計を計算する。

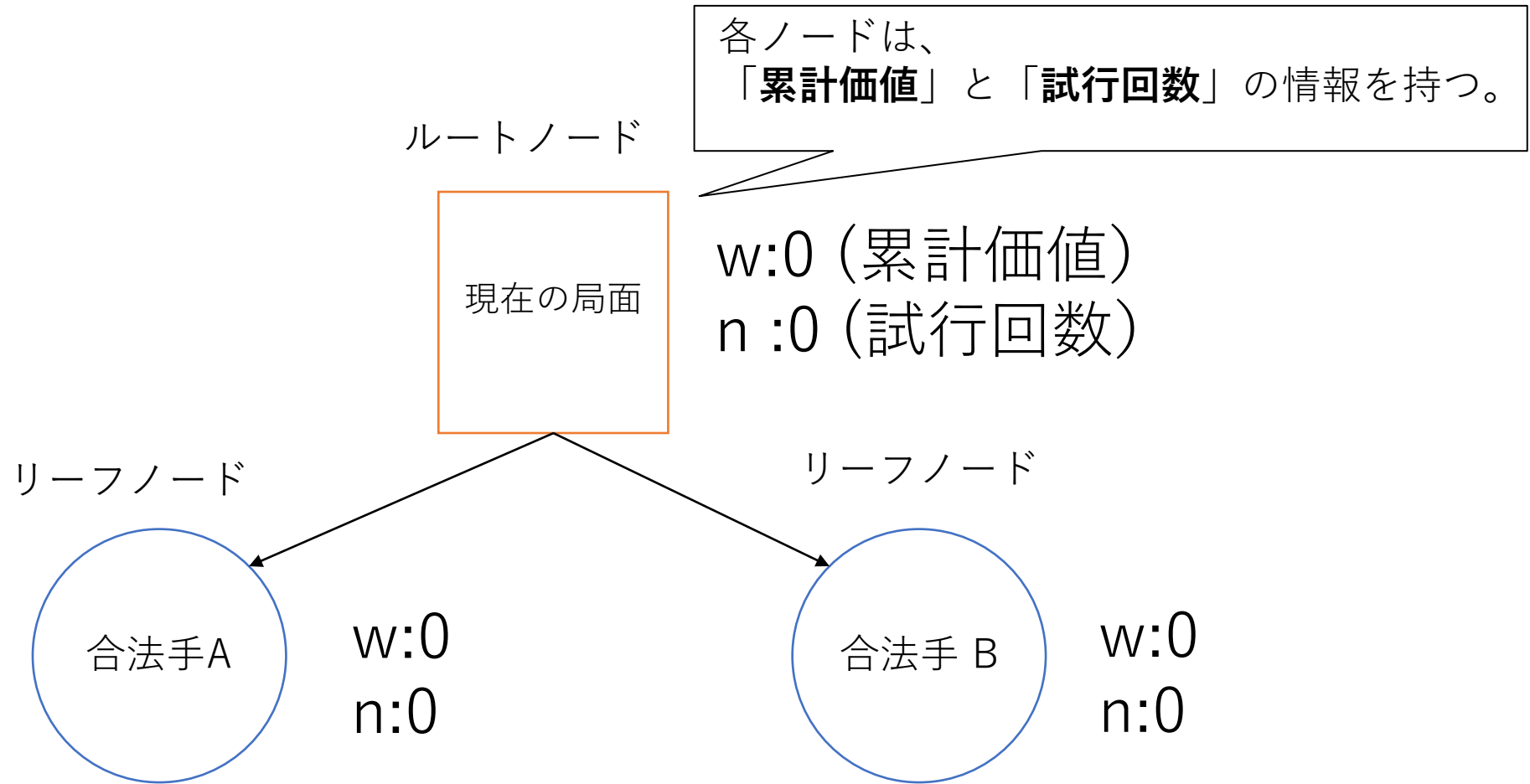
それぞれの合法手の中で
合計得点が最も大きい手を選択。



モンテカルロ木探索

モンテカルロ木探索は原始モンテカルロ探索に
「**選択**」、「**展開**」、「**評価**」、「**更新**」
の4つの操作を加えたもの。

モンテカルロ木探索ー初期状態



モンテカルロ木探索－選択

「ルートノード」から「子ノード」が存在したら
「リーフノード」に到達するまで移動する。

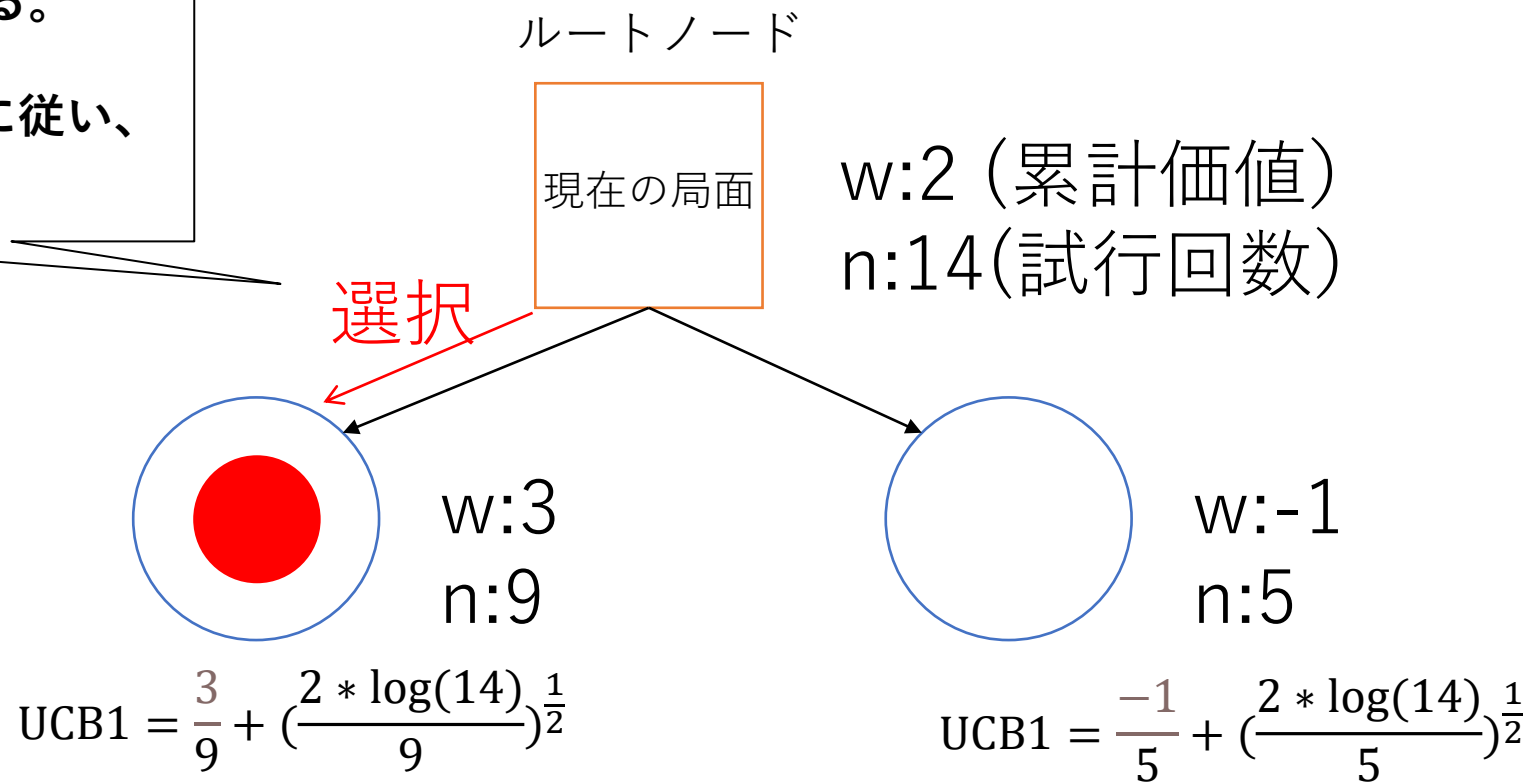
この時、「UCB1」（バイアス＋確率）に従い、
移動先を選択する。

$$\text{UCB1} = \underbrace{\frac{w}{n}}_{\text{成功率}} + \underbrace{\left(\frac{2 * \log(t)}{n}\right)^{\frac{1}{2}}}_{\text{バイアス}}$$

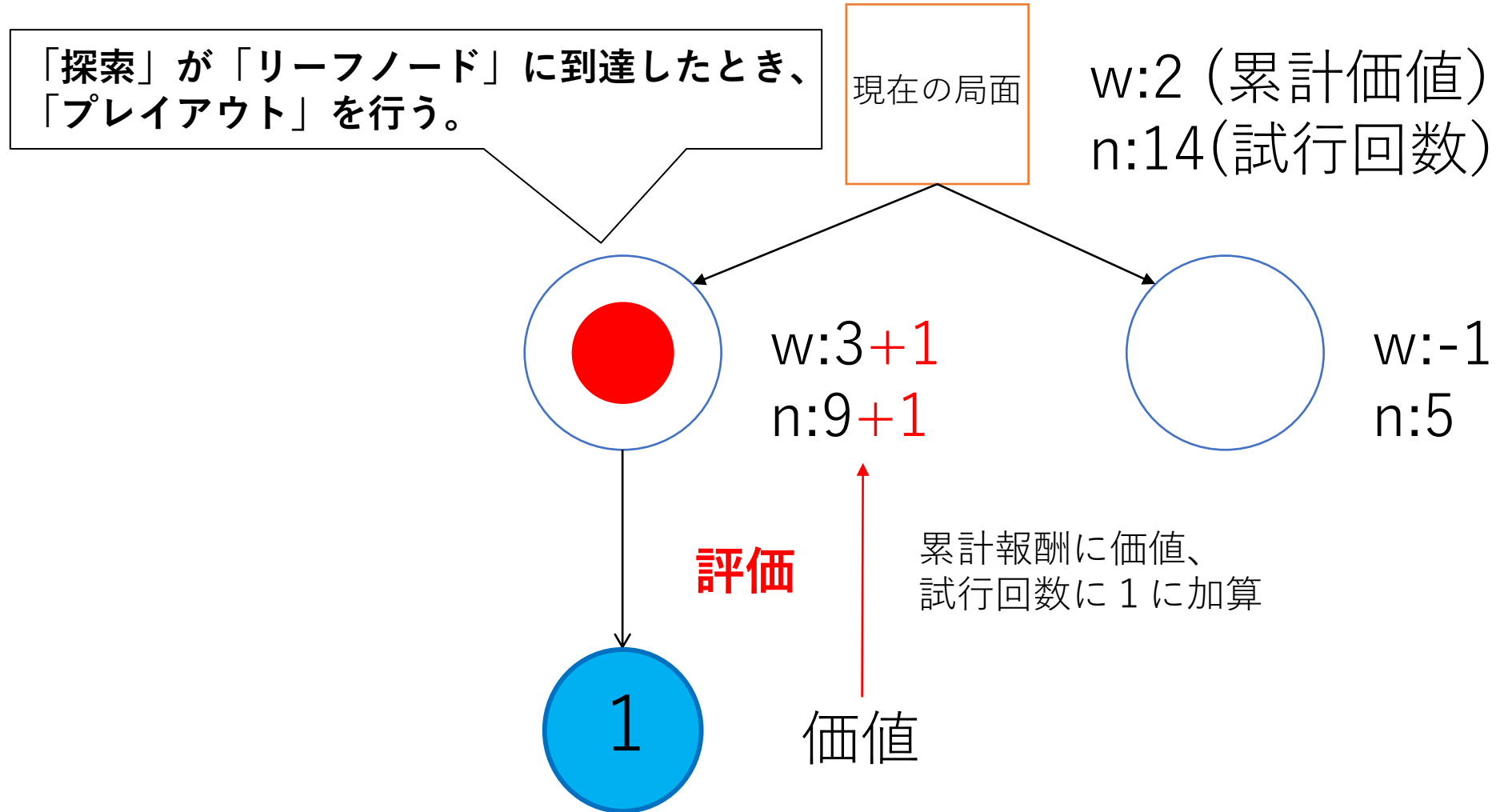
成功率

バイアス

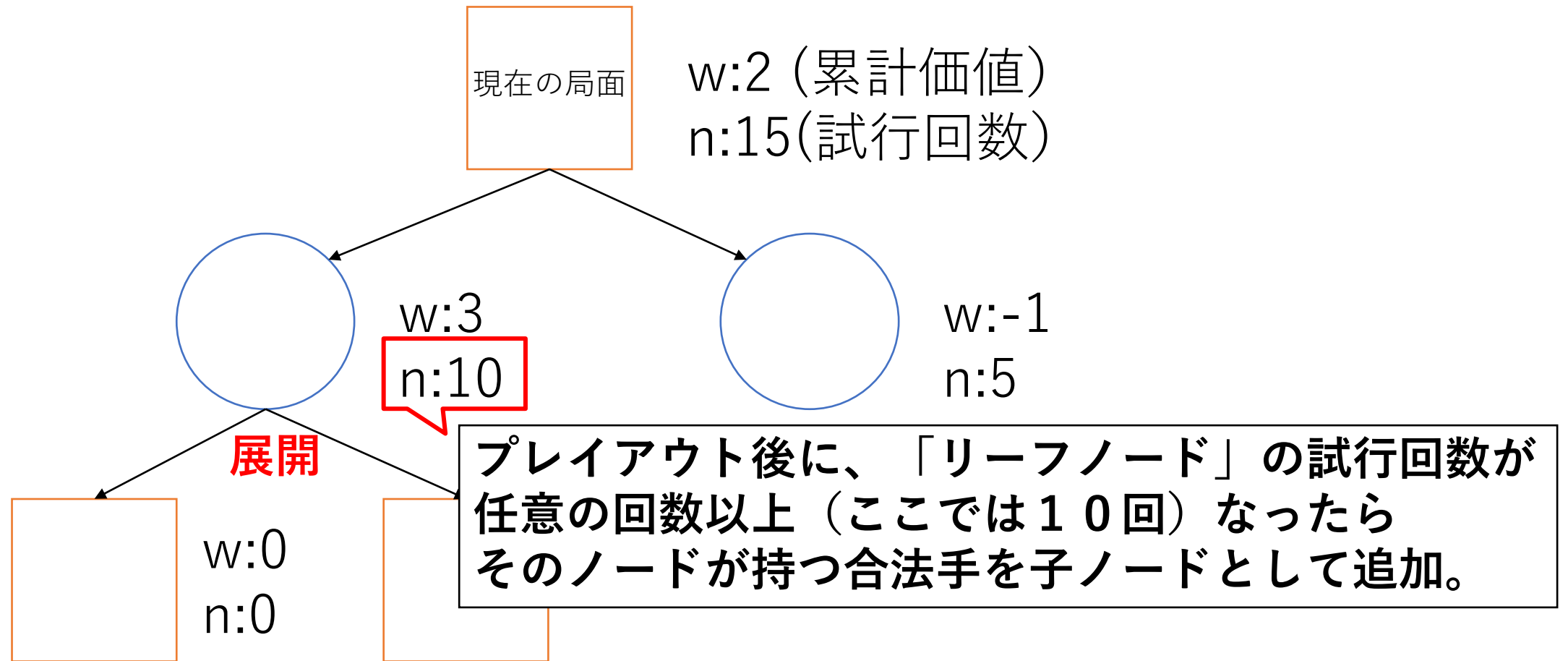
n：この行動の試行回数
w：この行動の累計価値
t：すべての行動の試行回数



モンテカルロ木探索－評価

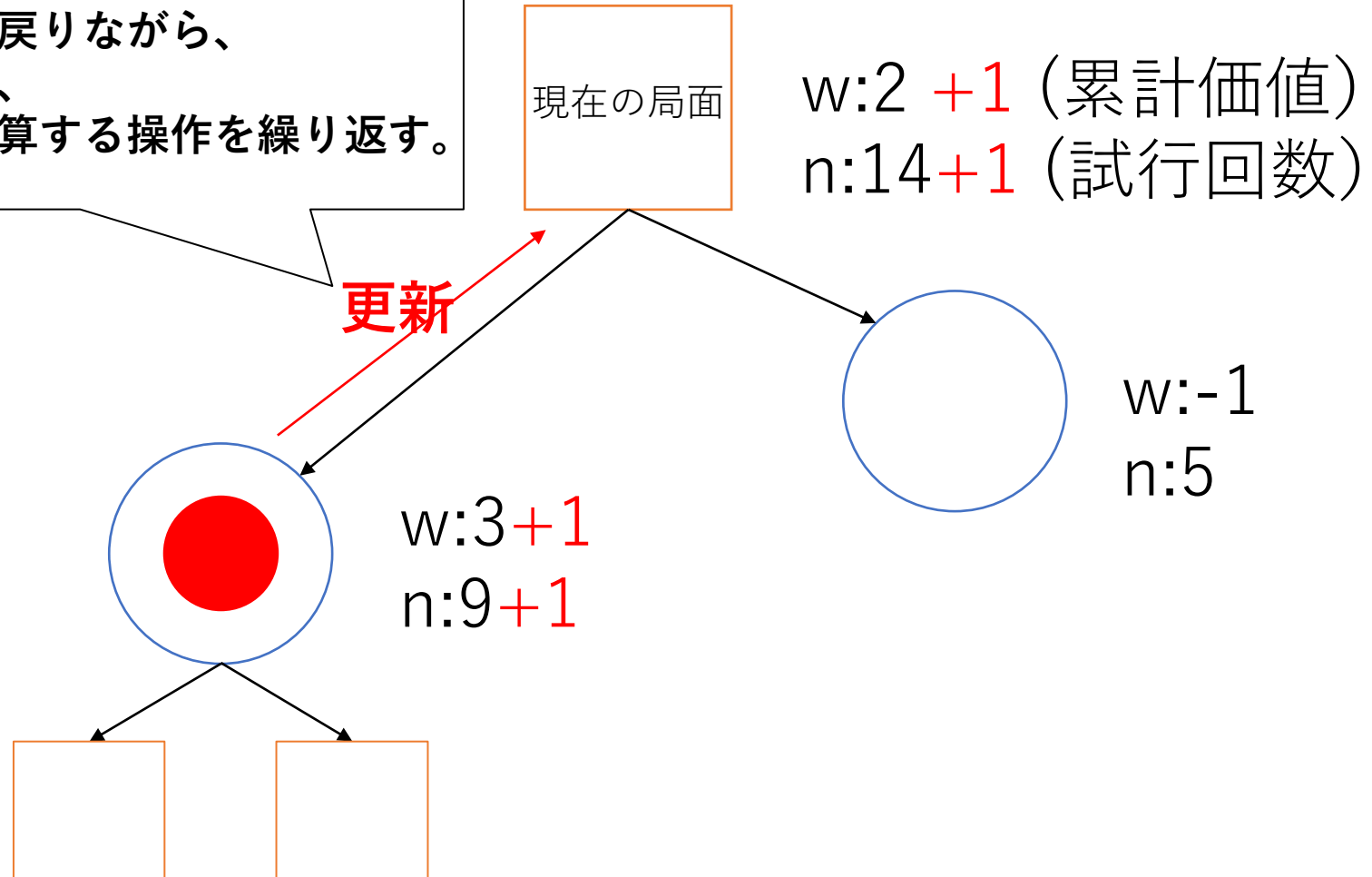


モンテカルロ木探索－展開



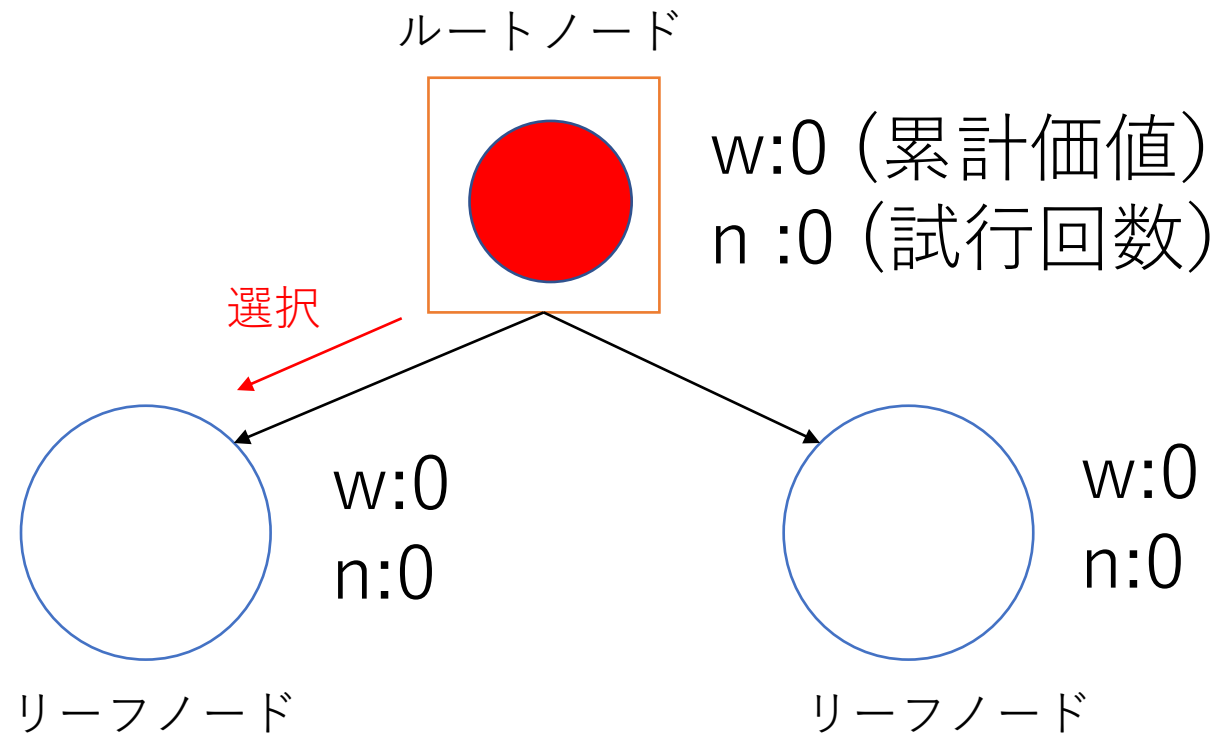
モンテカルロ木探索－更新

「プレイアウト」が終わったら、
「ルートノード」まで戻りながら、
ノードの「累計価値」、
「試行回数」に1を換算する操作を繰り返す。



モンテカルロ木探索ーシミュレーション

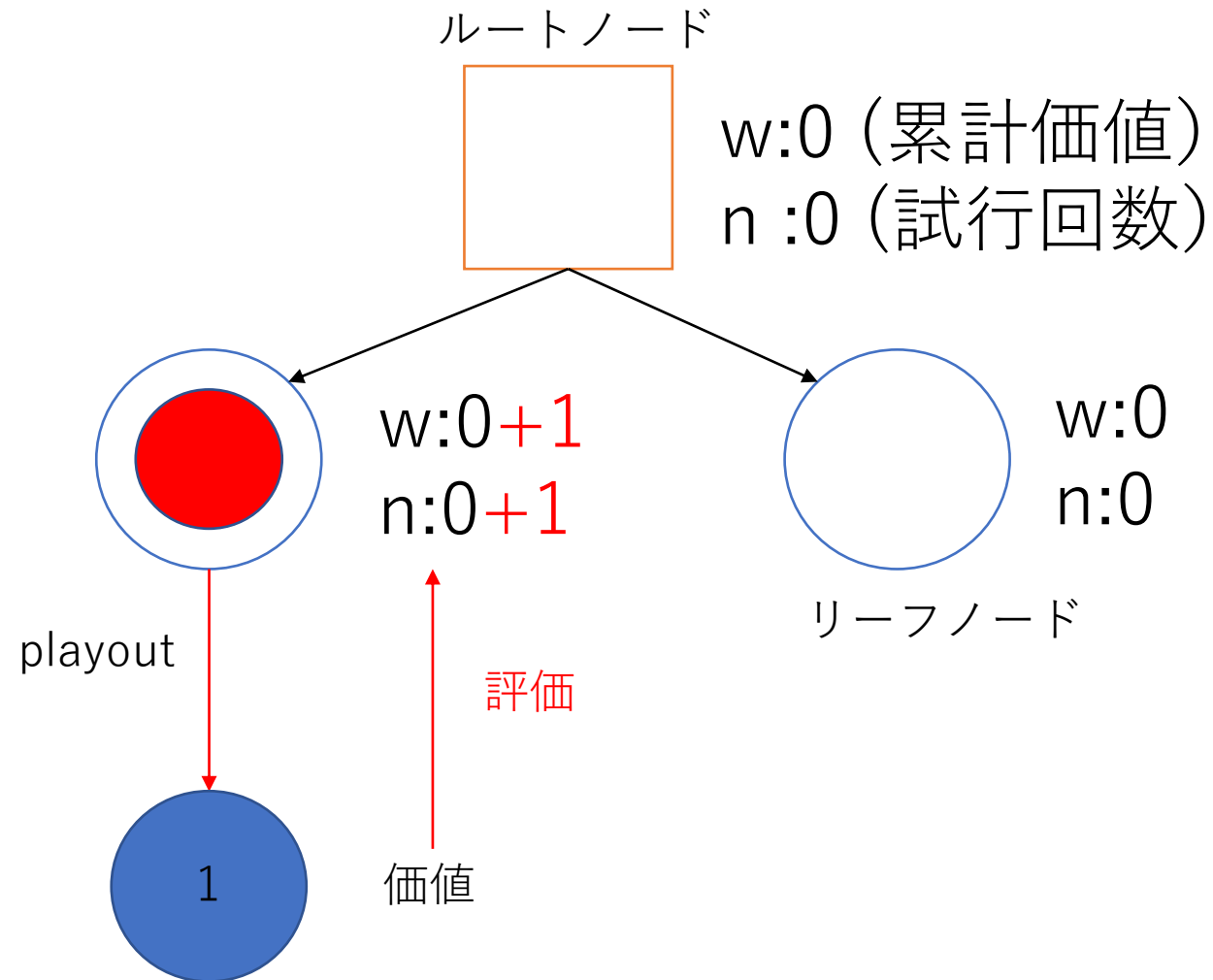
1回目



「UCB1」は試行回数が1以上じゃないと計算できないため、試行回数が0のノードを一通り「選択」する。

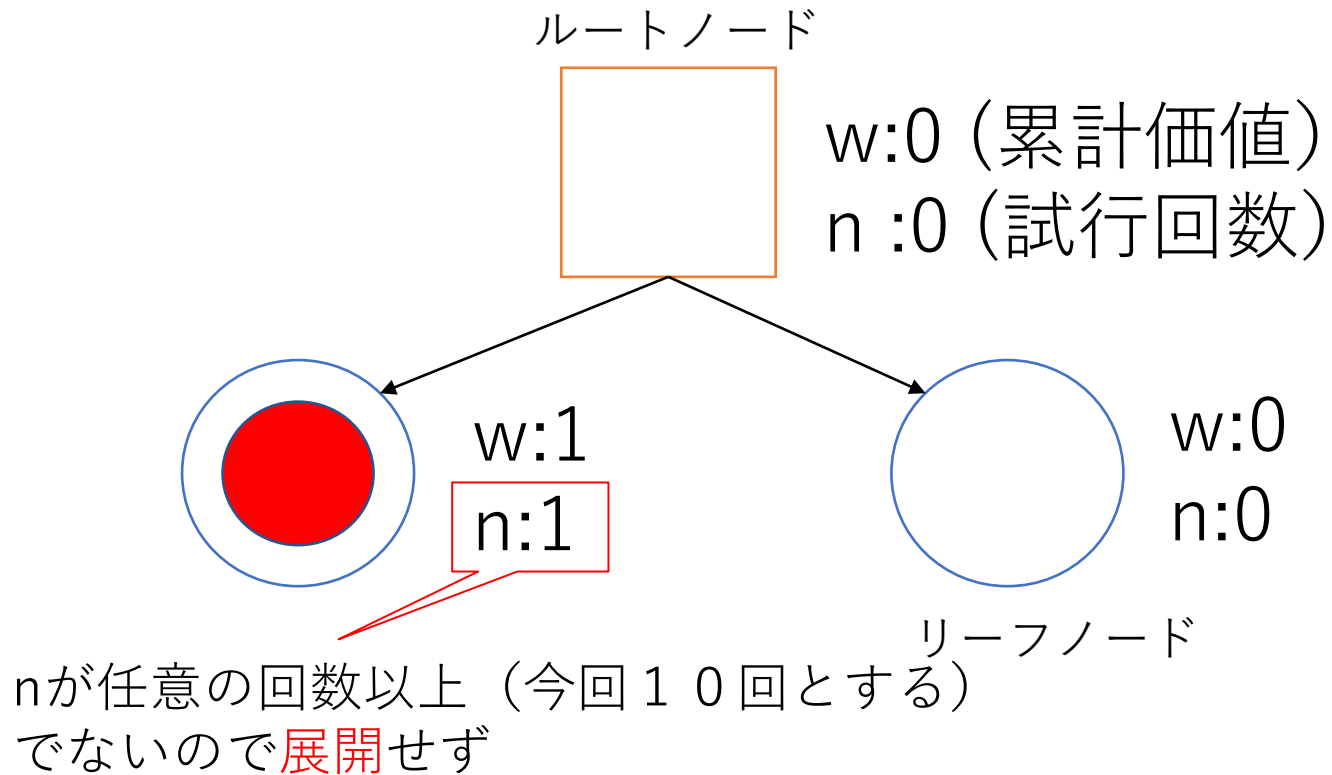
モンテカルロ木探索—シミュレーション

1回目



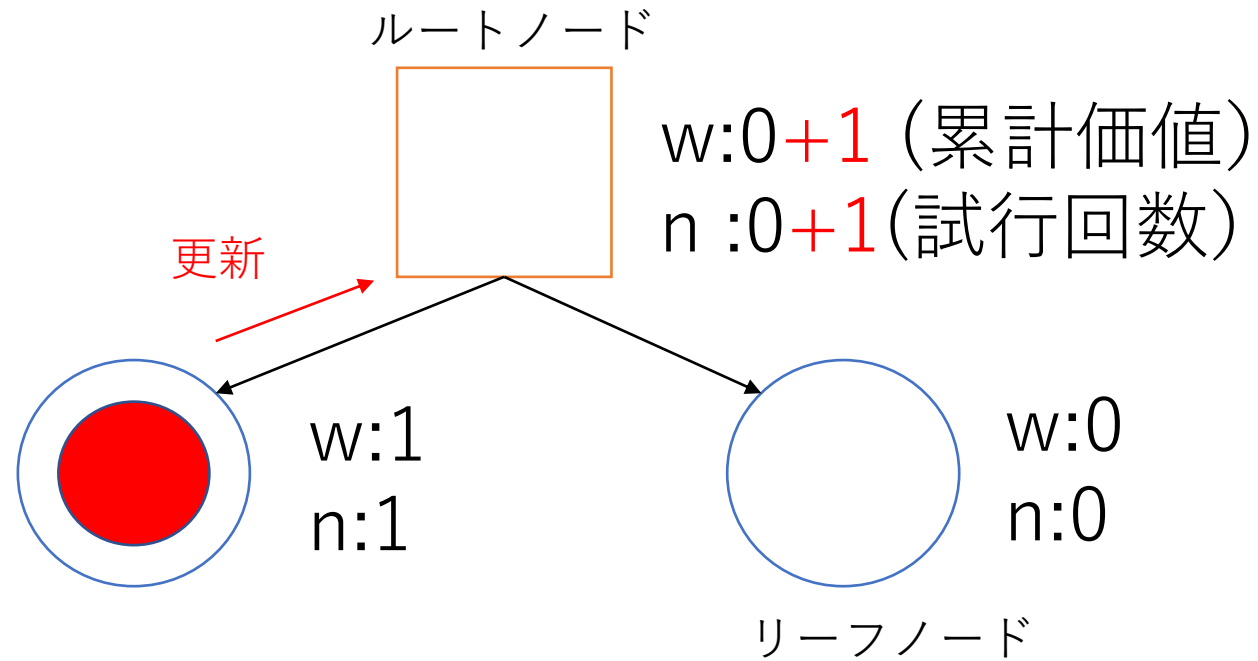
モンテカルロ木探索ーシミュレーション

1回目



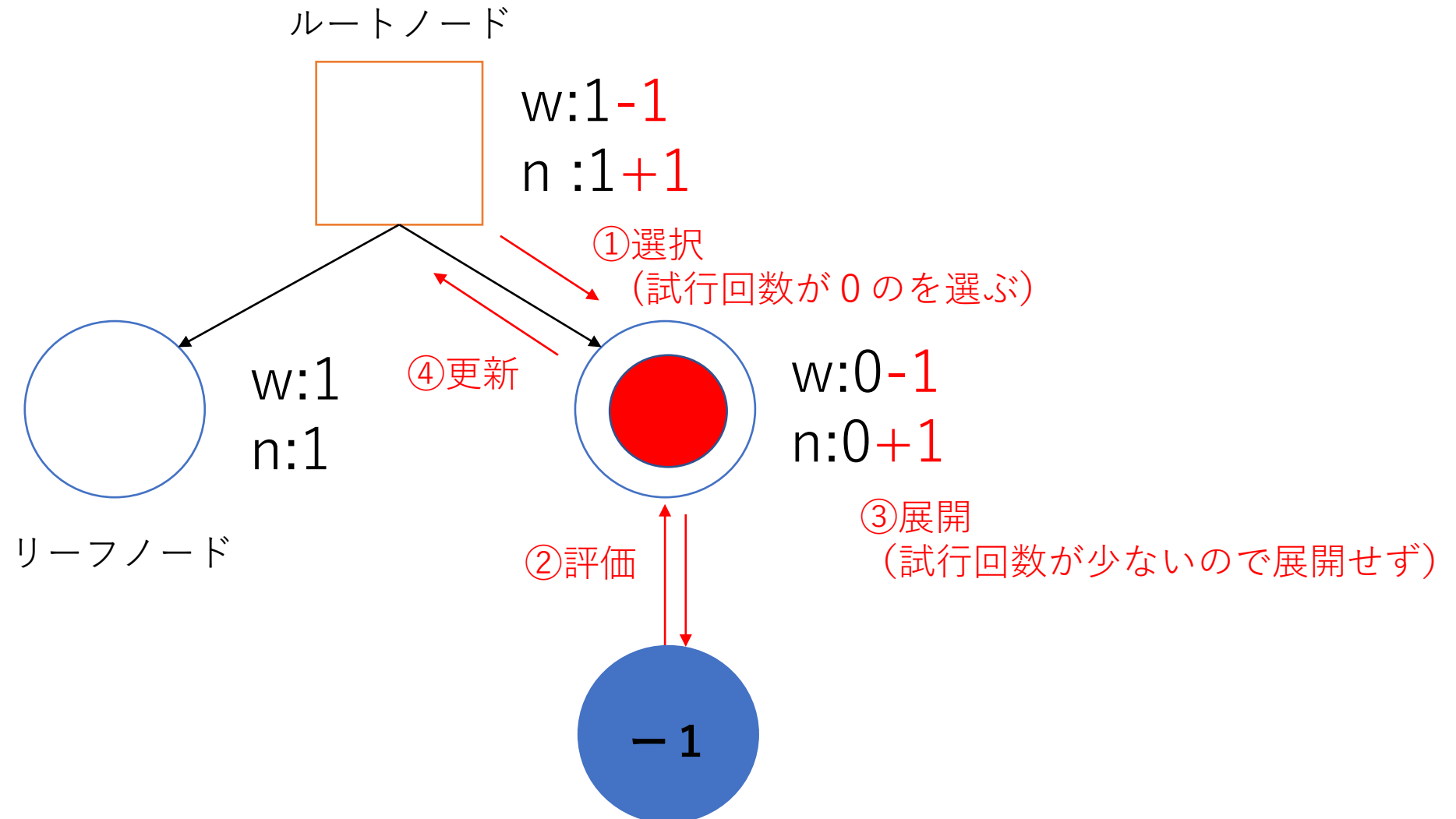
モンテカルロ木探索—シミュレーション

1回目



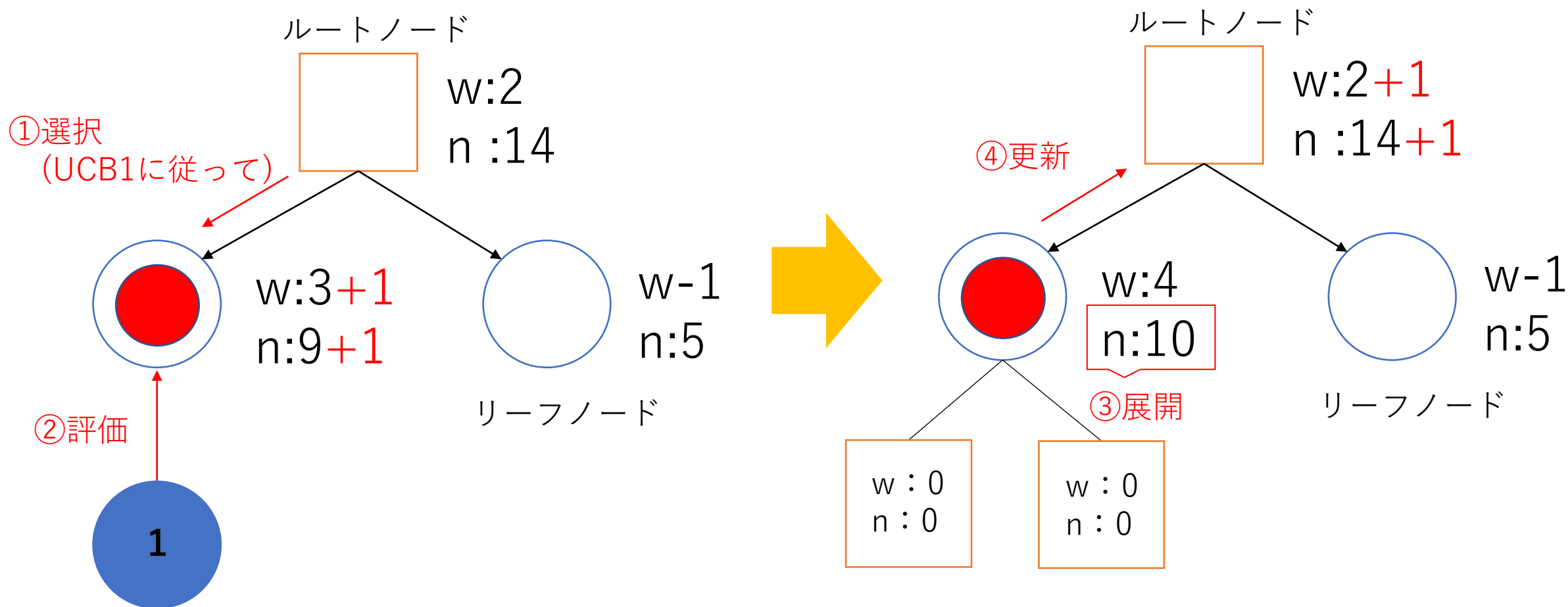
モンテカルロ木探索—シミュレーション

2回目



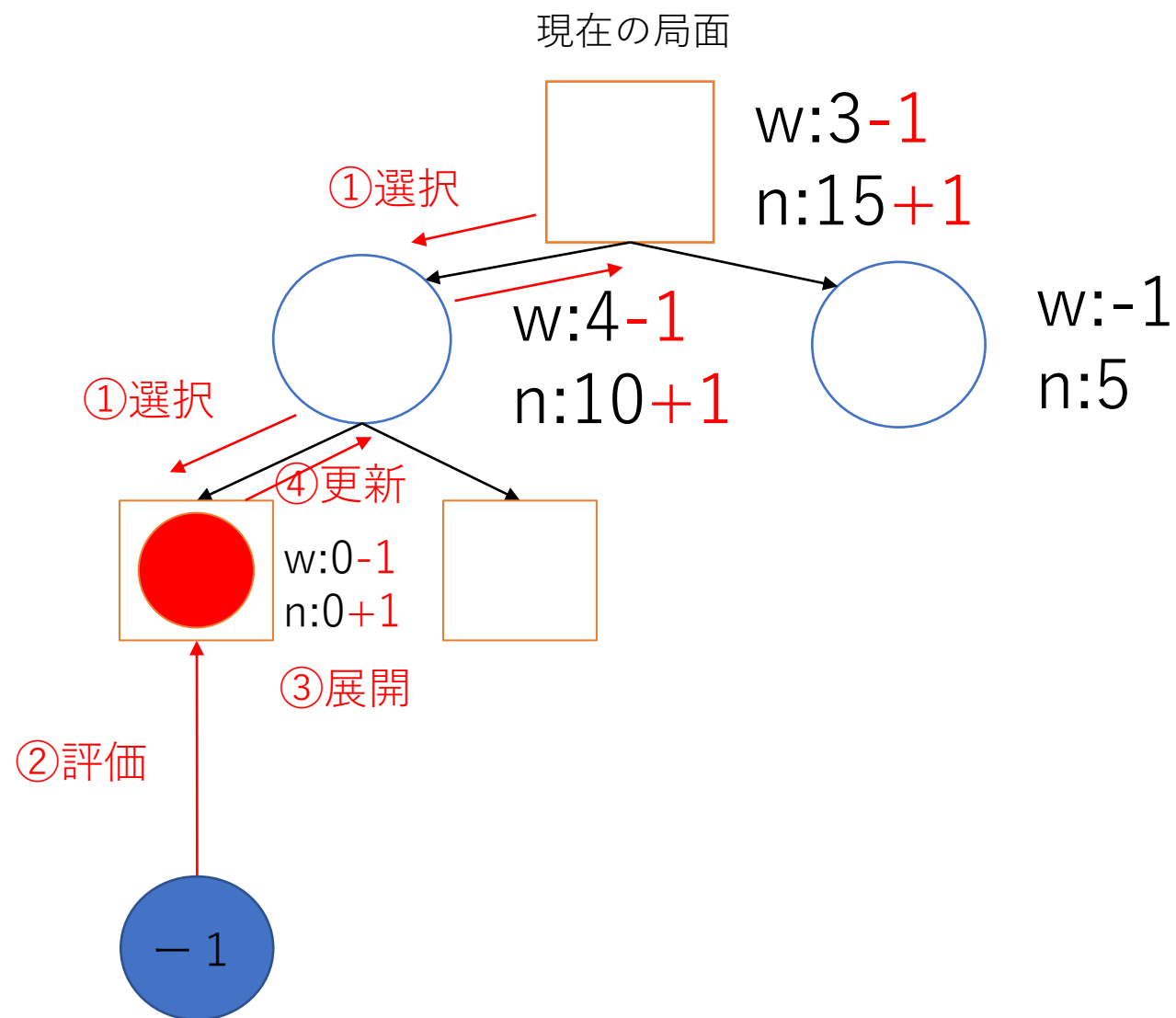
モンテカルロ木探索—シミュレーション

15回目



モンテカルロ木探索—シミュレーション

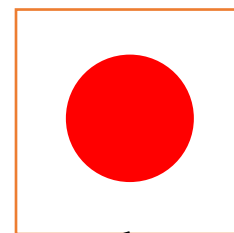
16回目



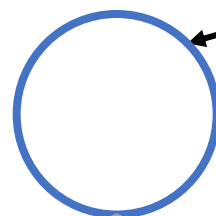
モンテカルロ木探索—行動選択

十分にシミュレーションを繰り返した後、
「**試行回数(n)**」が最大の手を「**次の一手**」
として選択する。

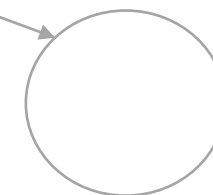
現在の局面



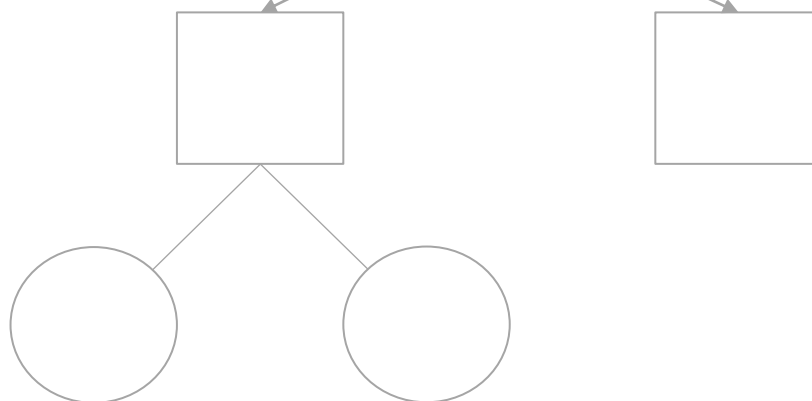
w:8 (累計価値)
n:100 (試行回数)



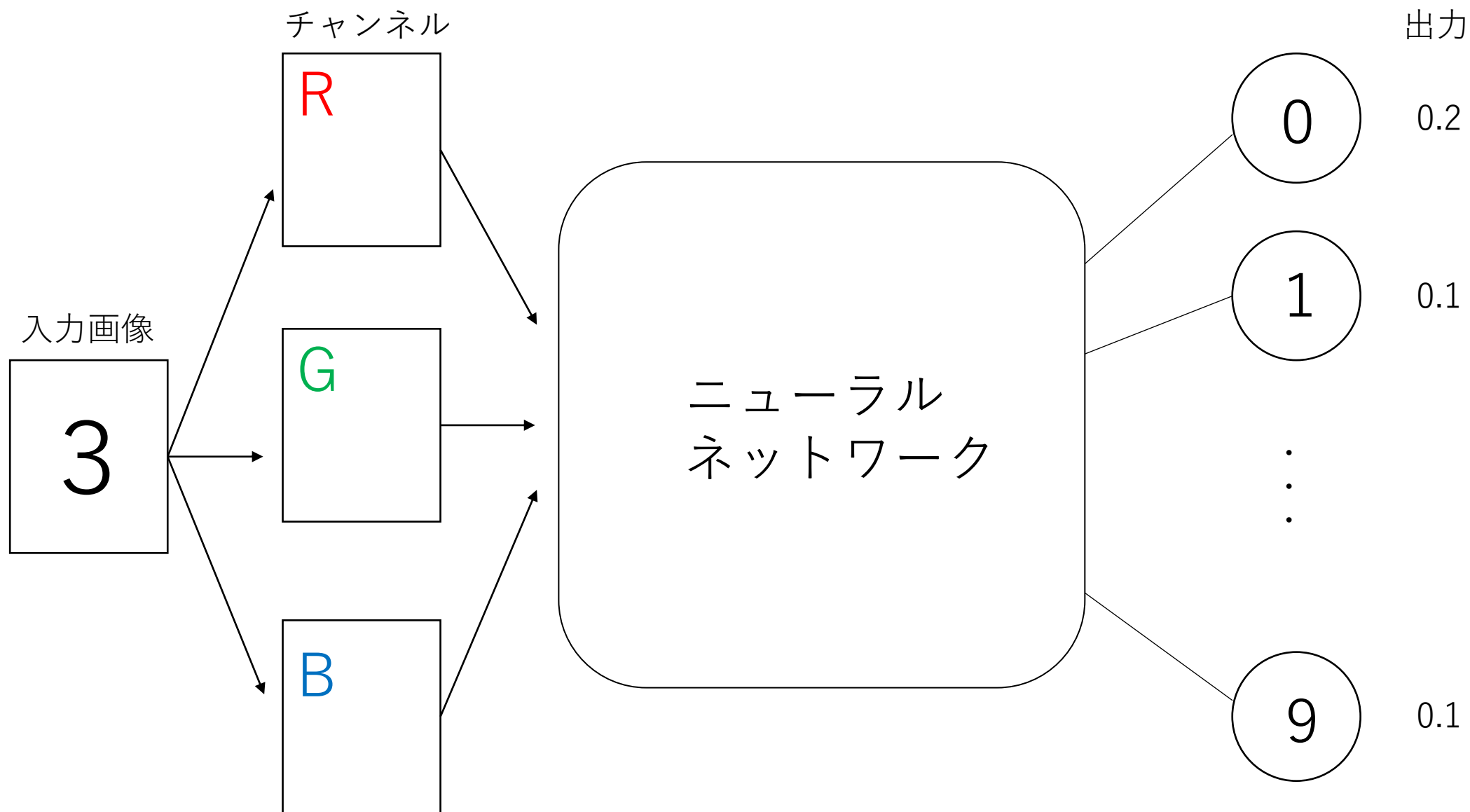
w:19
n:70



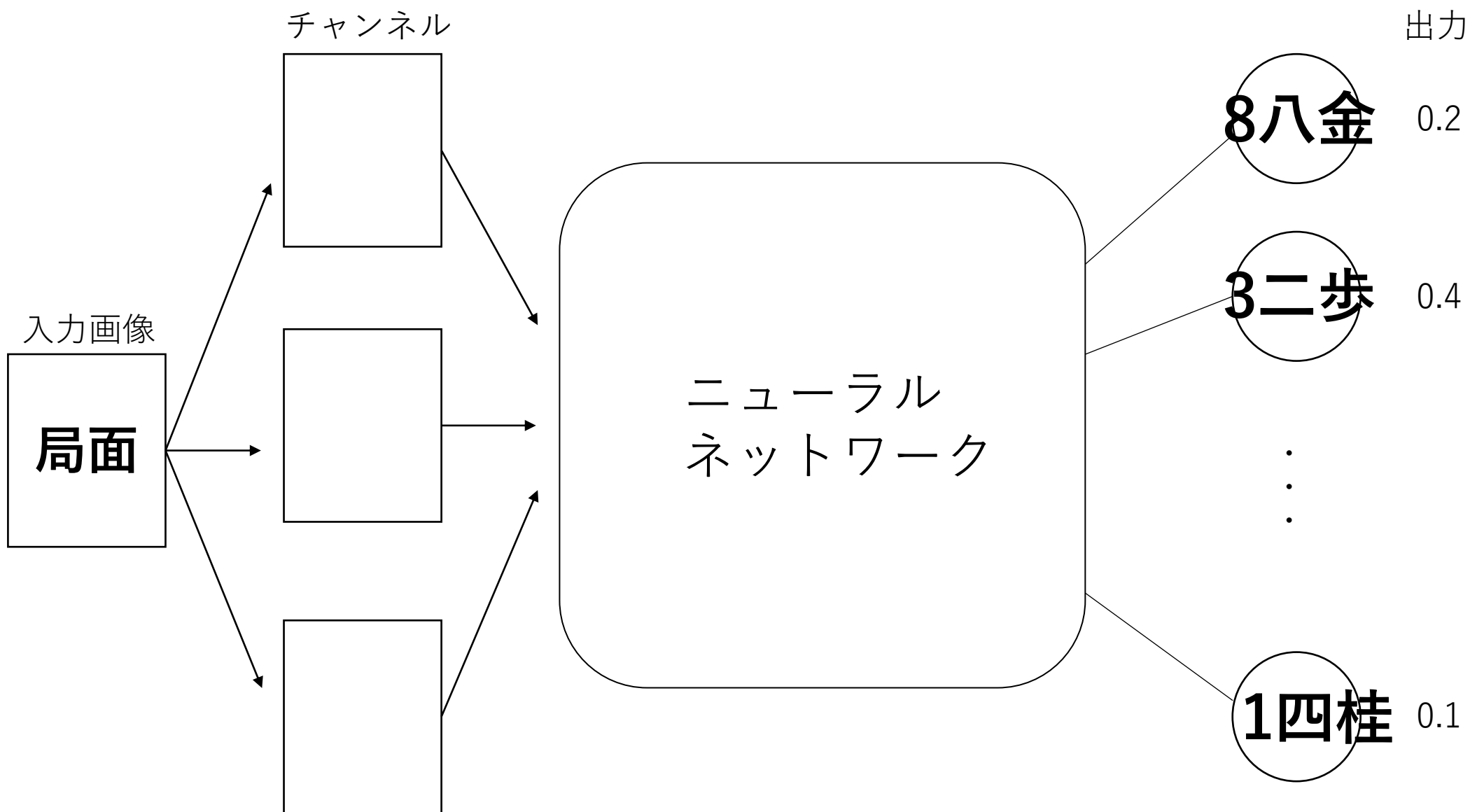
w:-11
n:30



ニューラルネットワークの画像識別の例



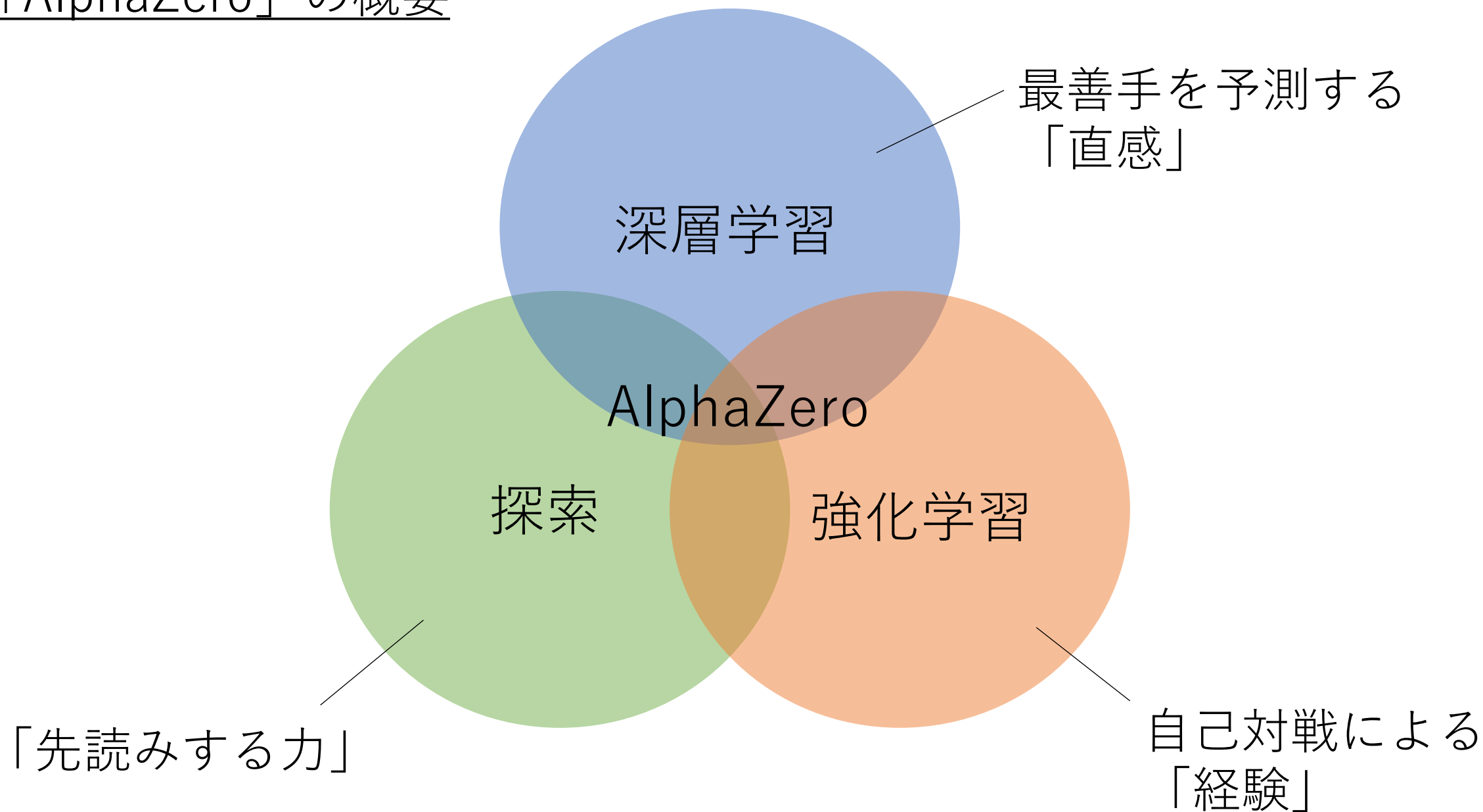
ニューラルネットワークの将棋への活用



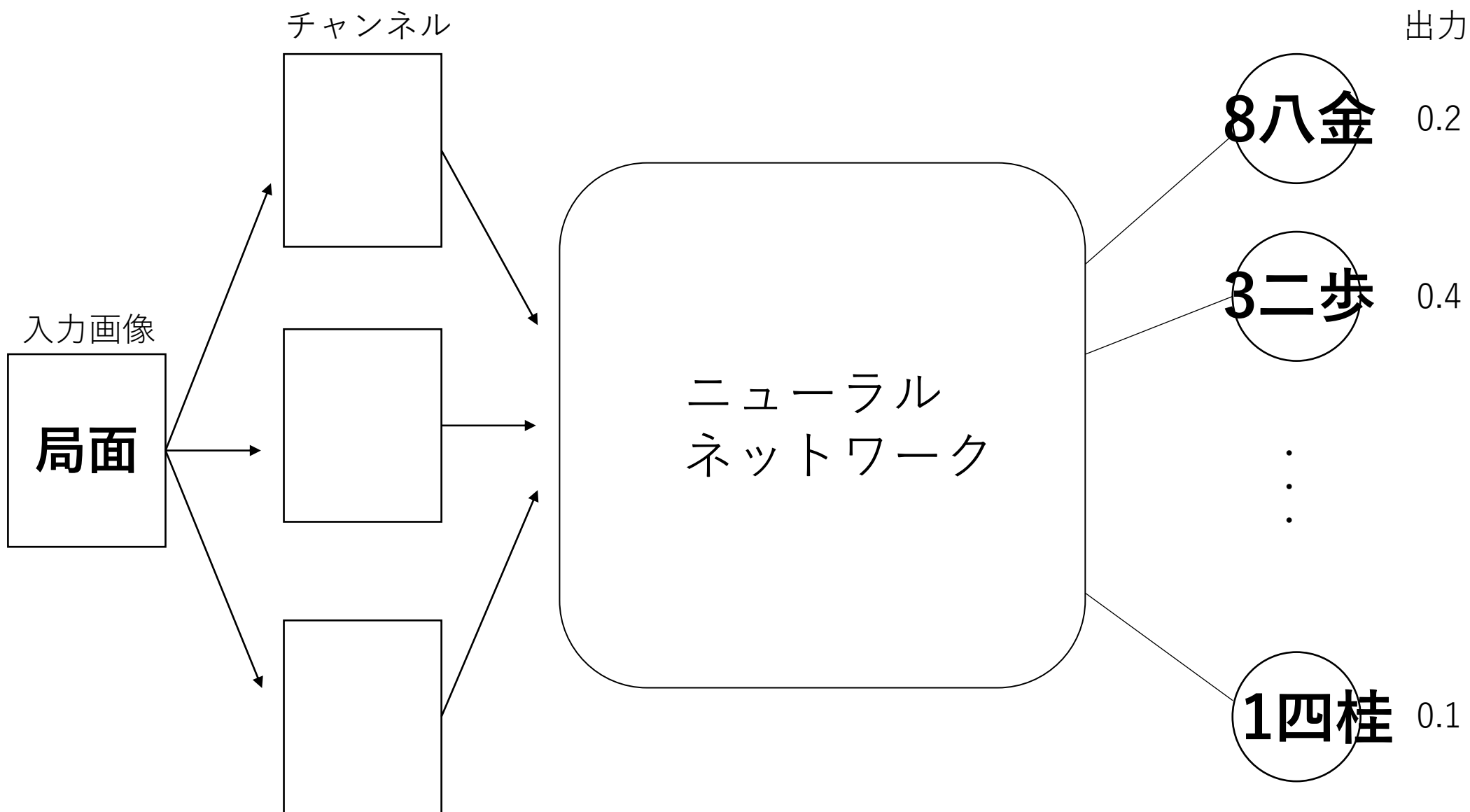
目次

- 京都将棋について
 - 京都将棋の概要
 - 京都将棋のルール
- 準備
 - 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
 - モンテカルロ木探索
 - ニューラルネットワーク
- AlphaZeroの仕組み
- 結果
- 参考文献

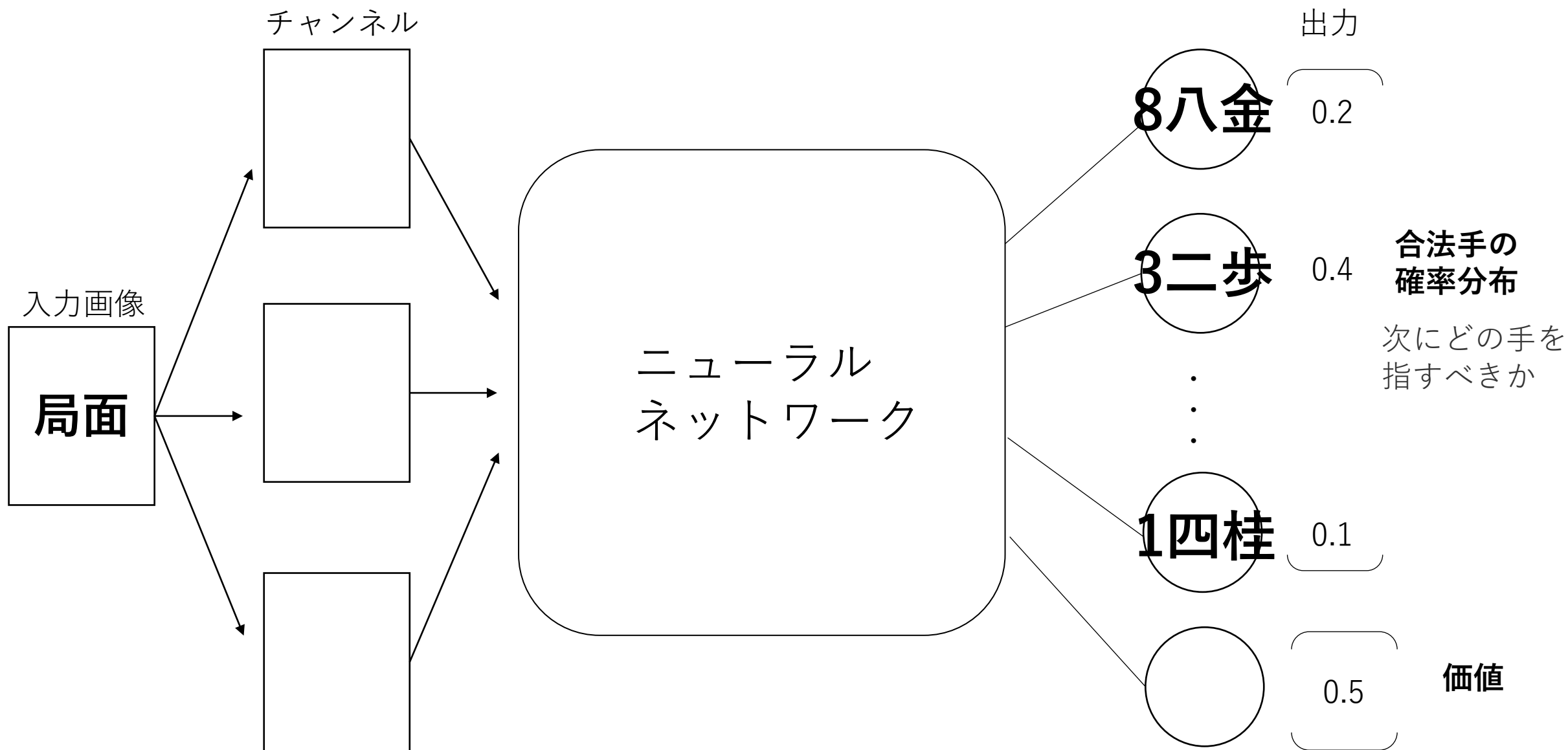
「AlphaZero」の概要



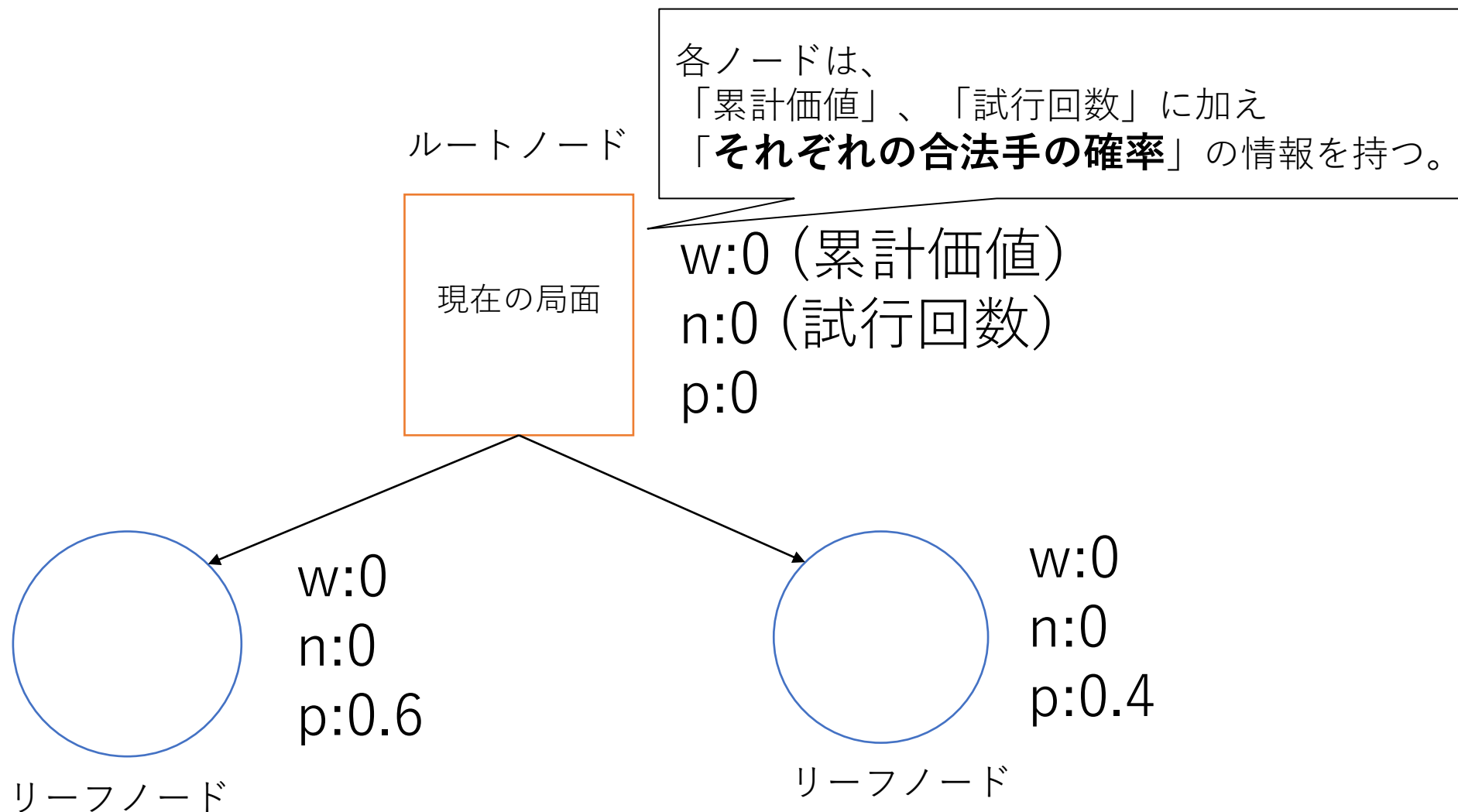
ニューラルネットワークの将棋への活用（再掲）



AlphaZeroのニューラルネットワーク



AlphaZeroのモンテカルロ木探索－初期状態



AlphaZeroのモンテカルロ木探索－選択

「ルートノード」から「子ノード」が存在したら
「リーフノード」に到達するまで移動する。

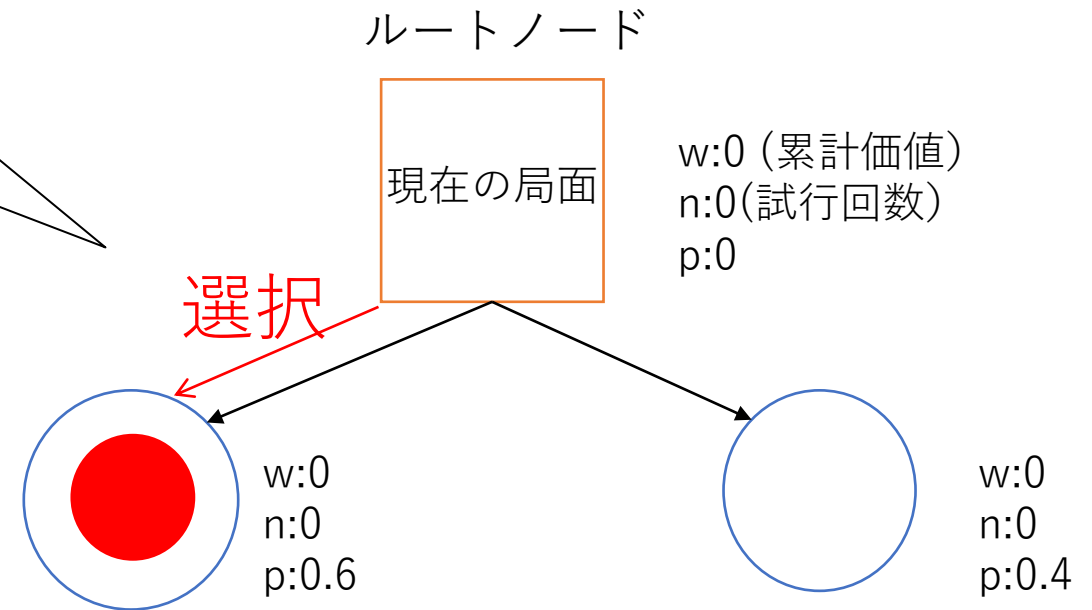
この時、「~~UCB1~~」（バイアス＋確率）に従い、
移動先を選択する。

$$\text{アーク評価値} = \underbrace{\frac{w}{n}}_{\text{成功率}} + \underbrace{c_{puct} * p}_{\text{バランス調整の定数}} * \underbrace{\frac{\sqrt{t}}{(1+n)}}_{\text{バイアス}}$$

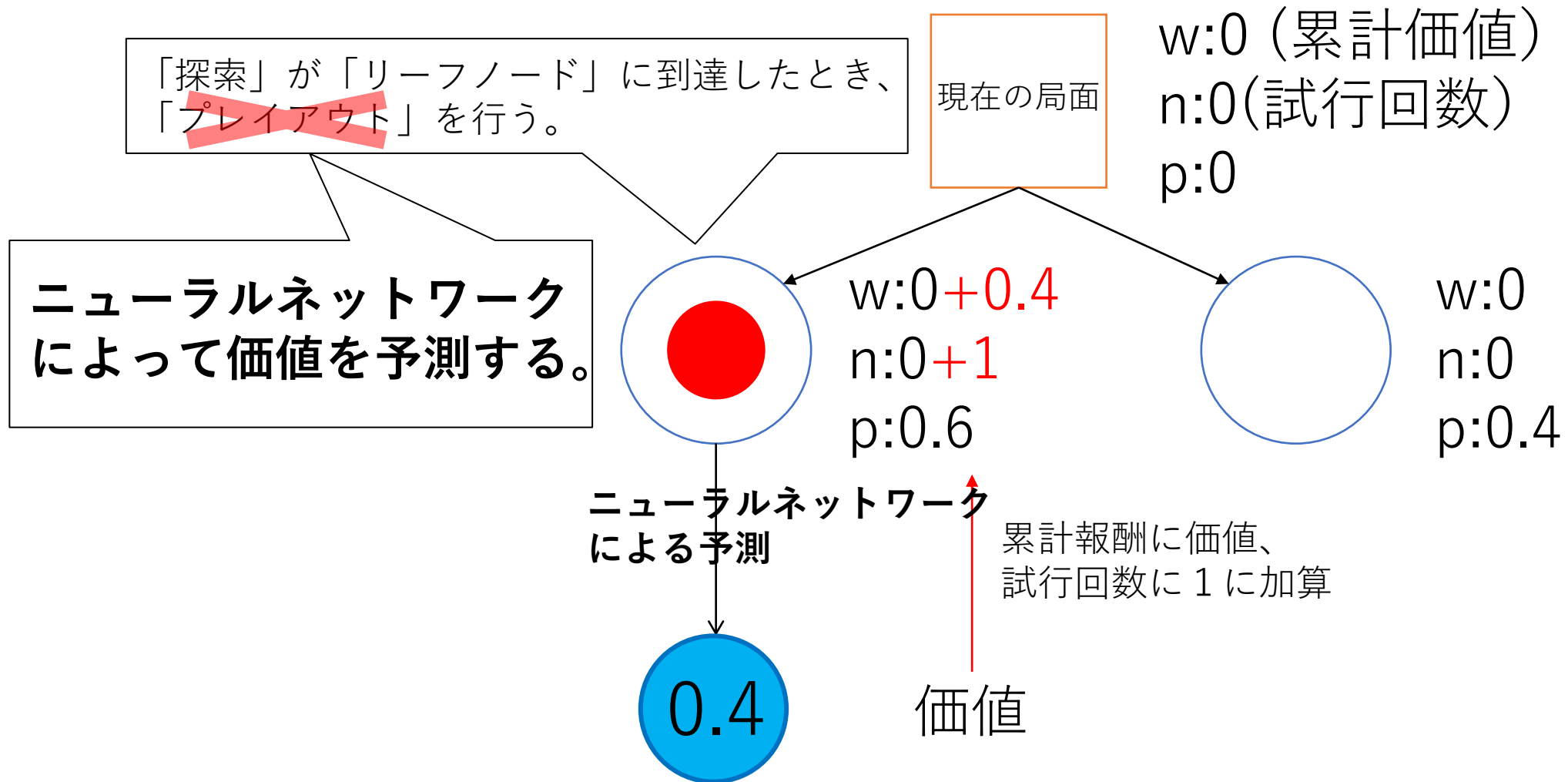
合法手の確率分布

n: このノードの試行回数, w: このノードの累計価値,
t: 累計試行回数

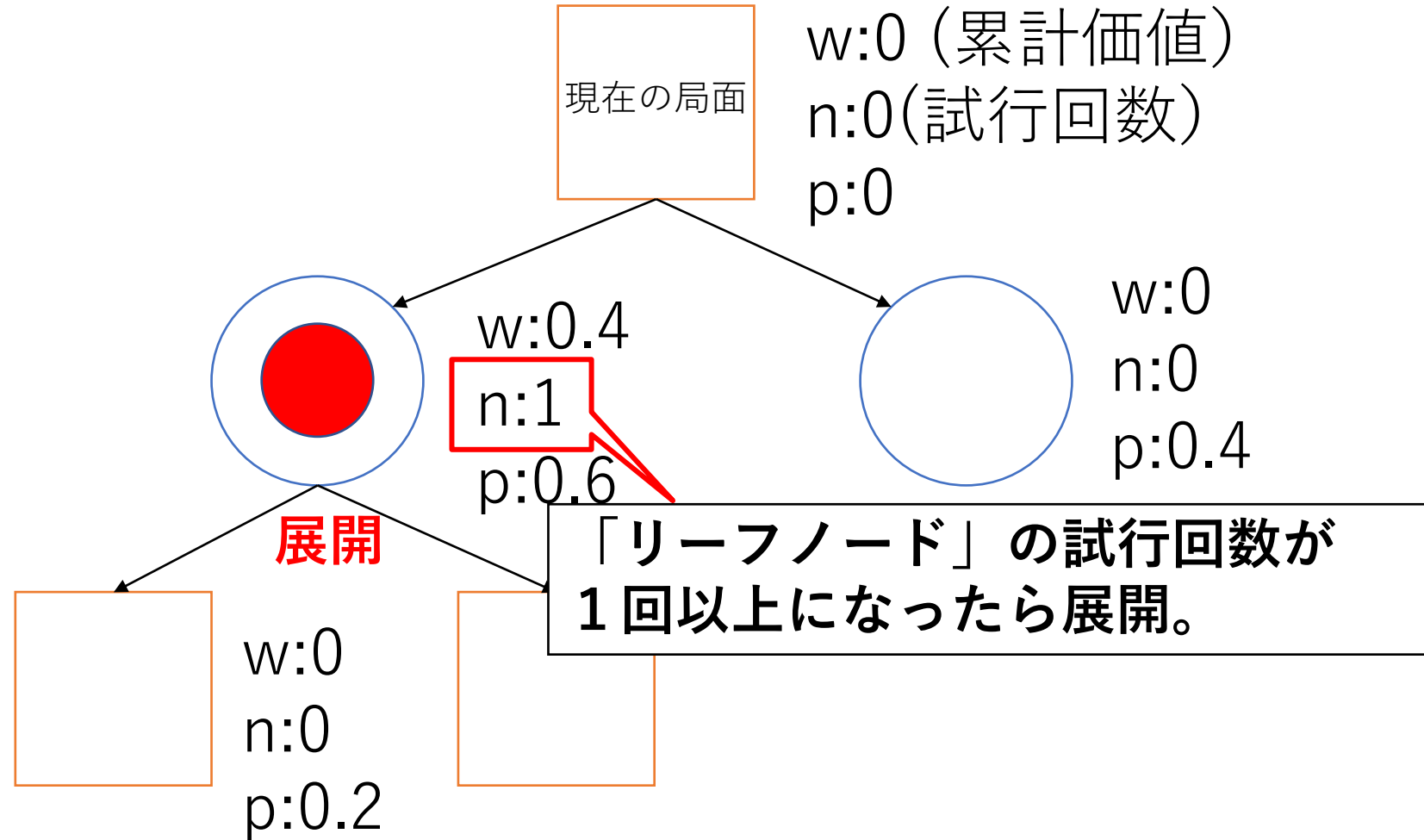
Cpuct: 「勝率」と「手の予測確率*バイアス」のバランスを調整するための定数



AlphaZeroのモンテカルロ木探索－評価

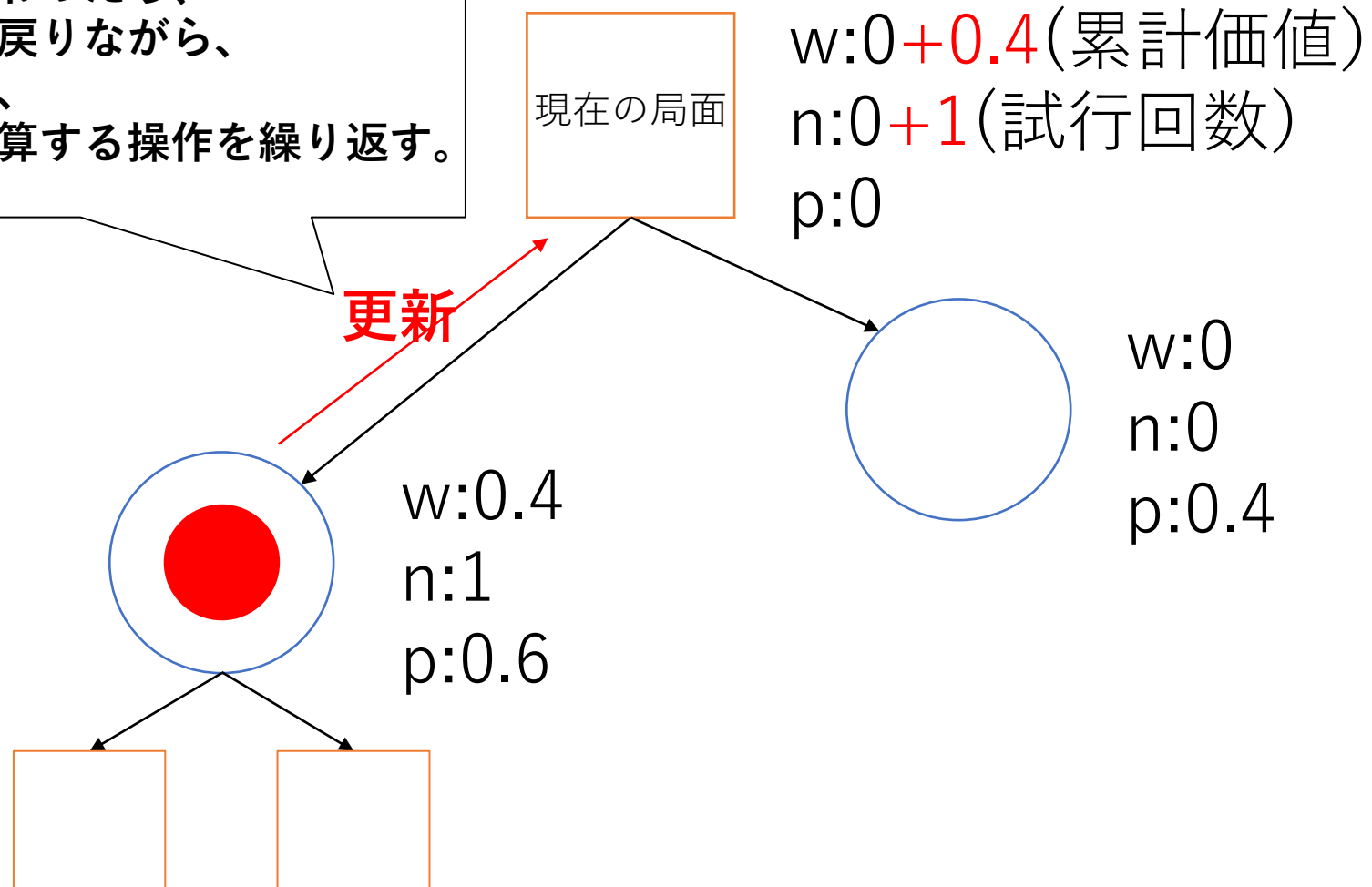


AlphaZeroのモンテカルロ口木探索ー展開



AlphaZeroのモンテカルロ木探索－更新

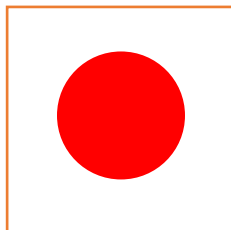
「プレイアウト」が終わったら、
「ルートノード」まで戻りながら、
ノードの「累計価値」、
「試行回数」に1を換算する操作を繰り返す。



AlphaZeroのモンテカルロ木探索—行動選択

十分にシミュレーションを繰り返した後、
「**試行回数(n)**」を「**確率分布**」に変換し、
その確率分布に従って選択する。

現在の局面

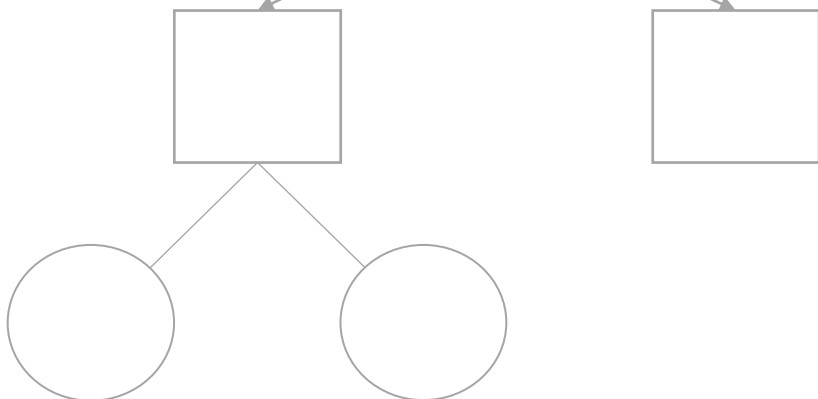


w:8 (累計価値)
n:100 (試行回数)

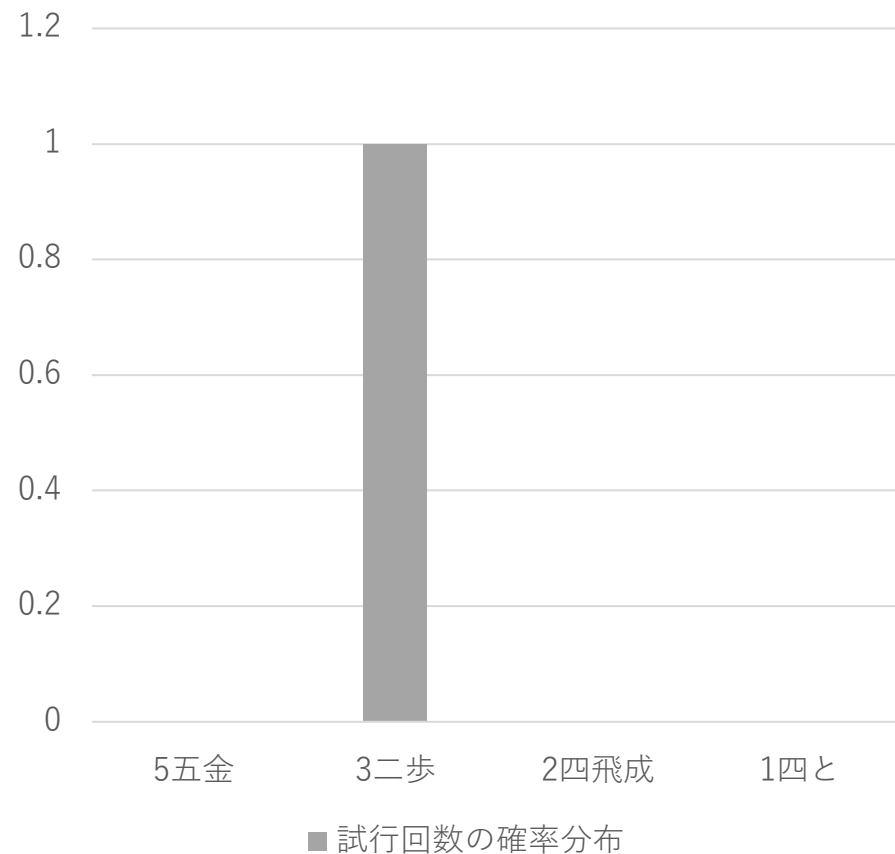
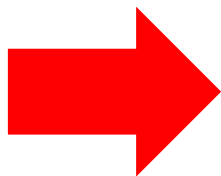
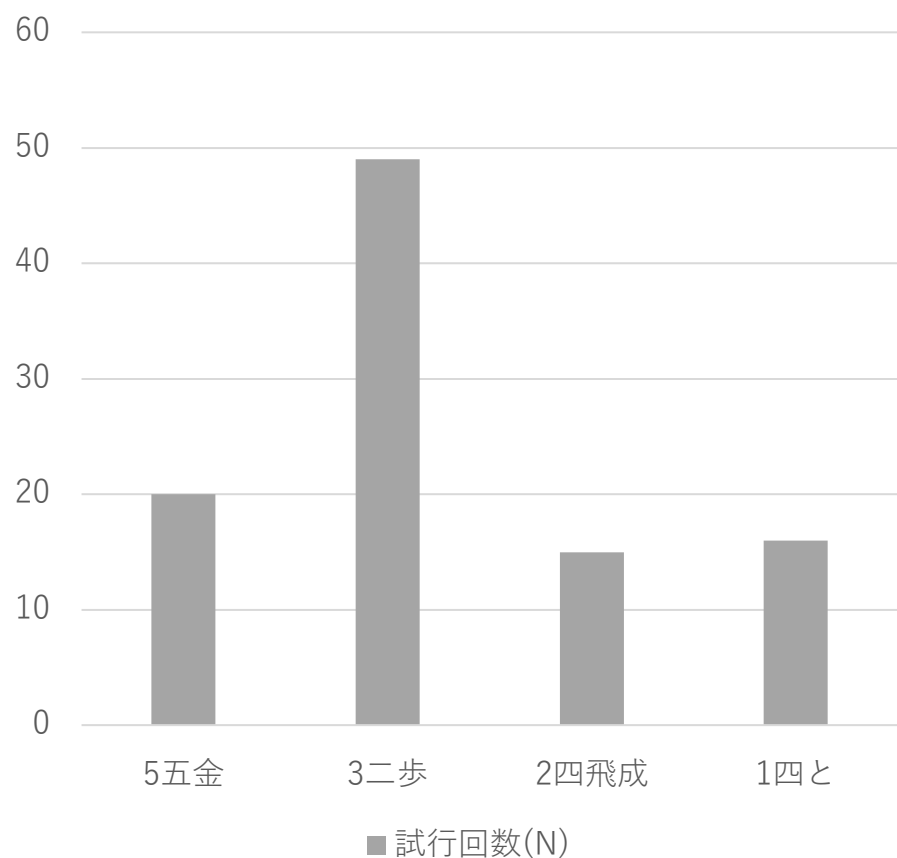
w:12.32
n:20

...

w:-11.98
n:10

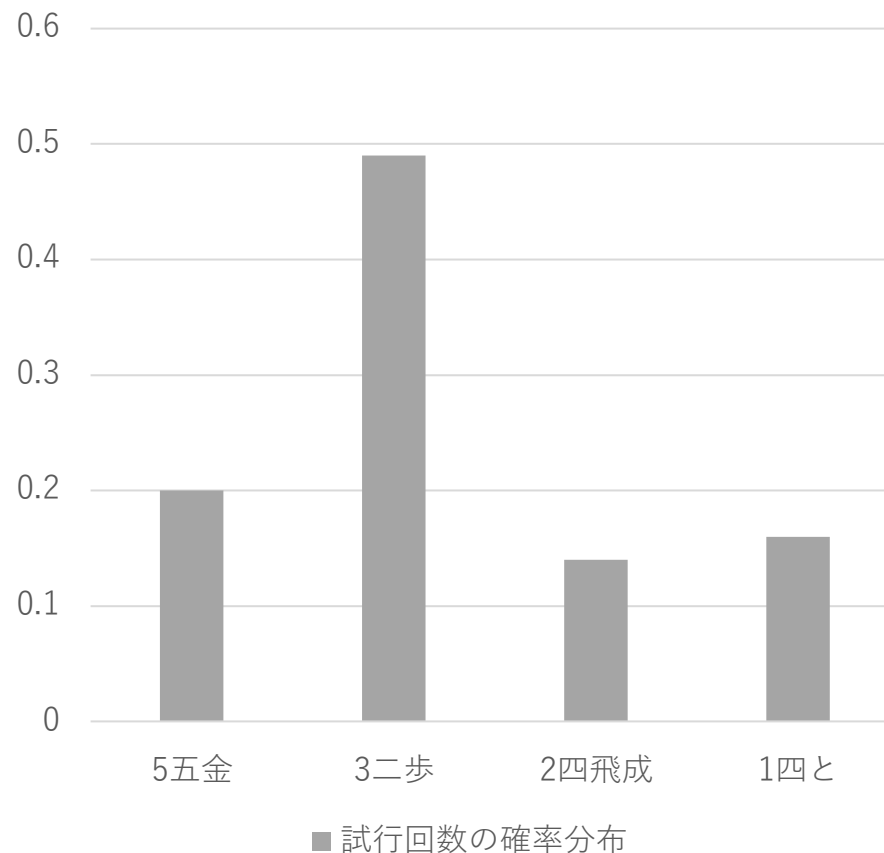
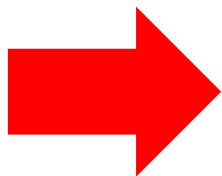
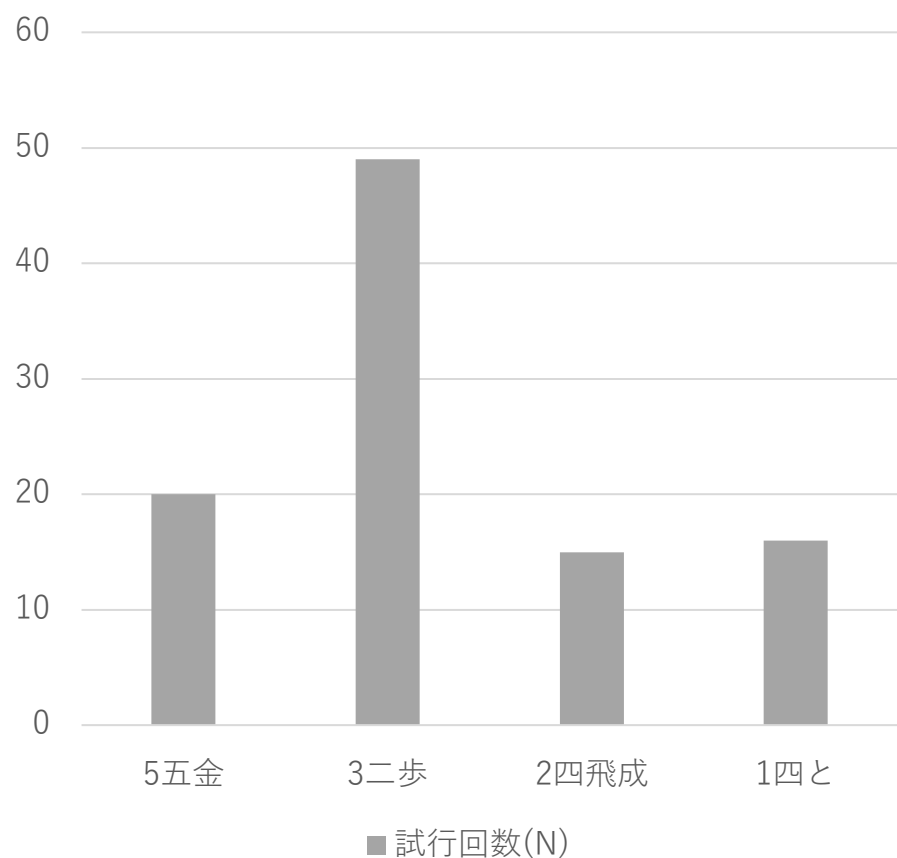


AlphaZeroのモンテカルロ木探索一行動選択 (本番時)



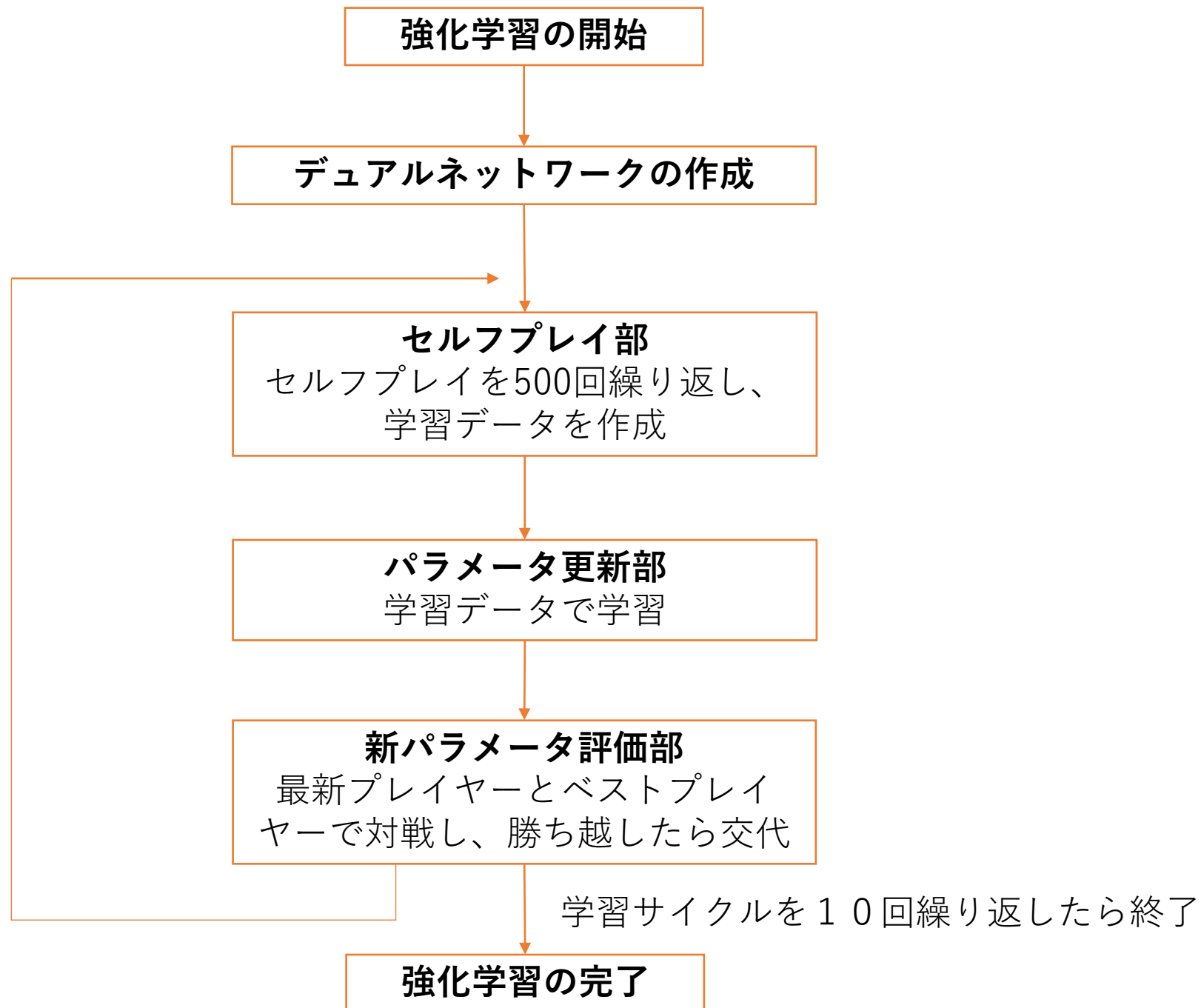
この確率分布に従って1手を選ぶ

AlphaZeroのモンテカルロ木探索一行動選択 (学習時)

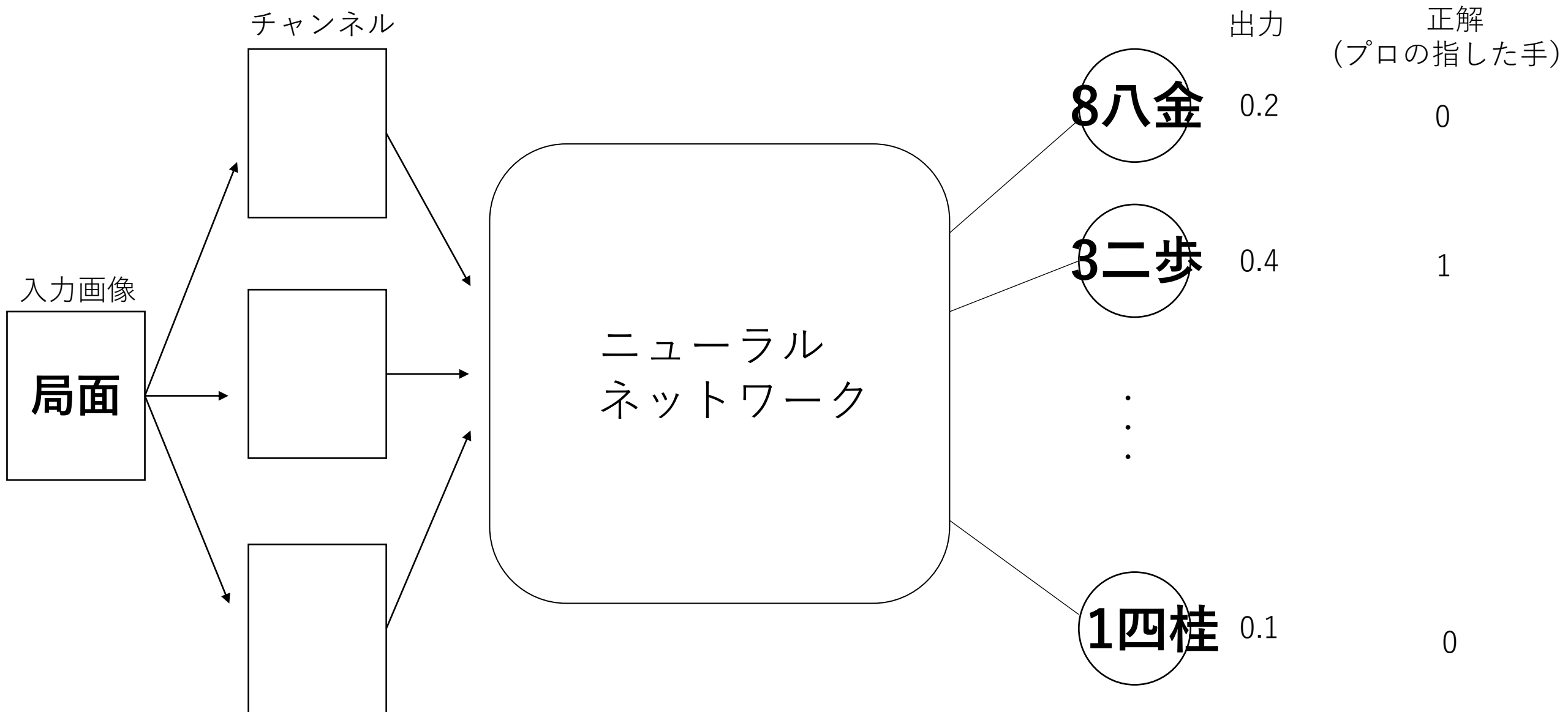


この確率分布に従って1手を選ぶ

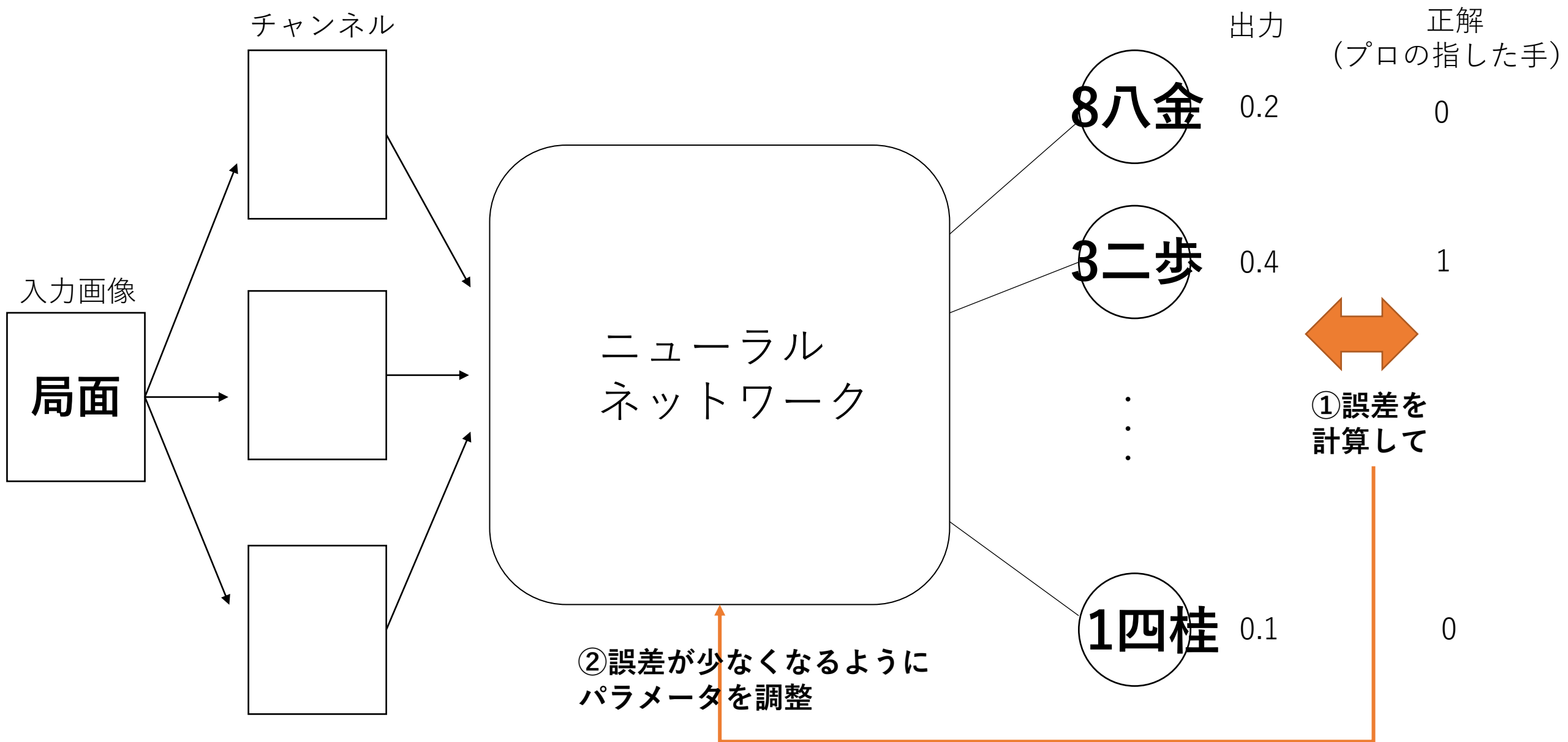
学習



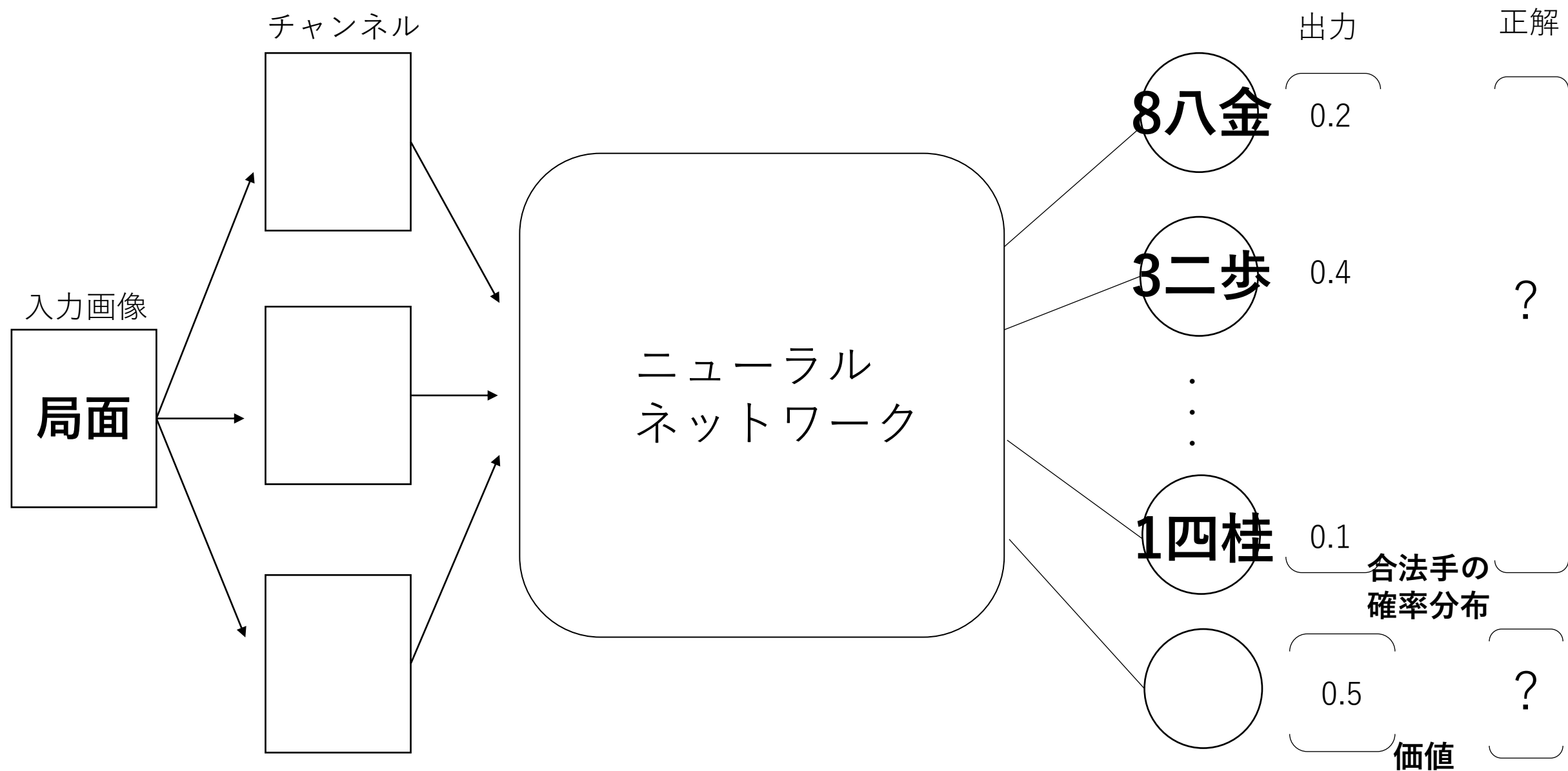
ニューラルネットワークの将棋への活用



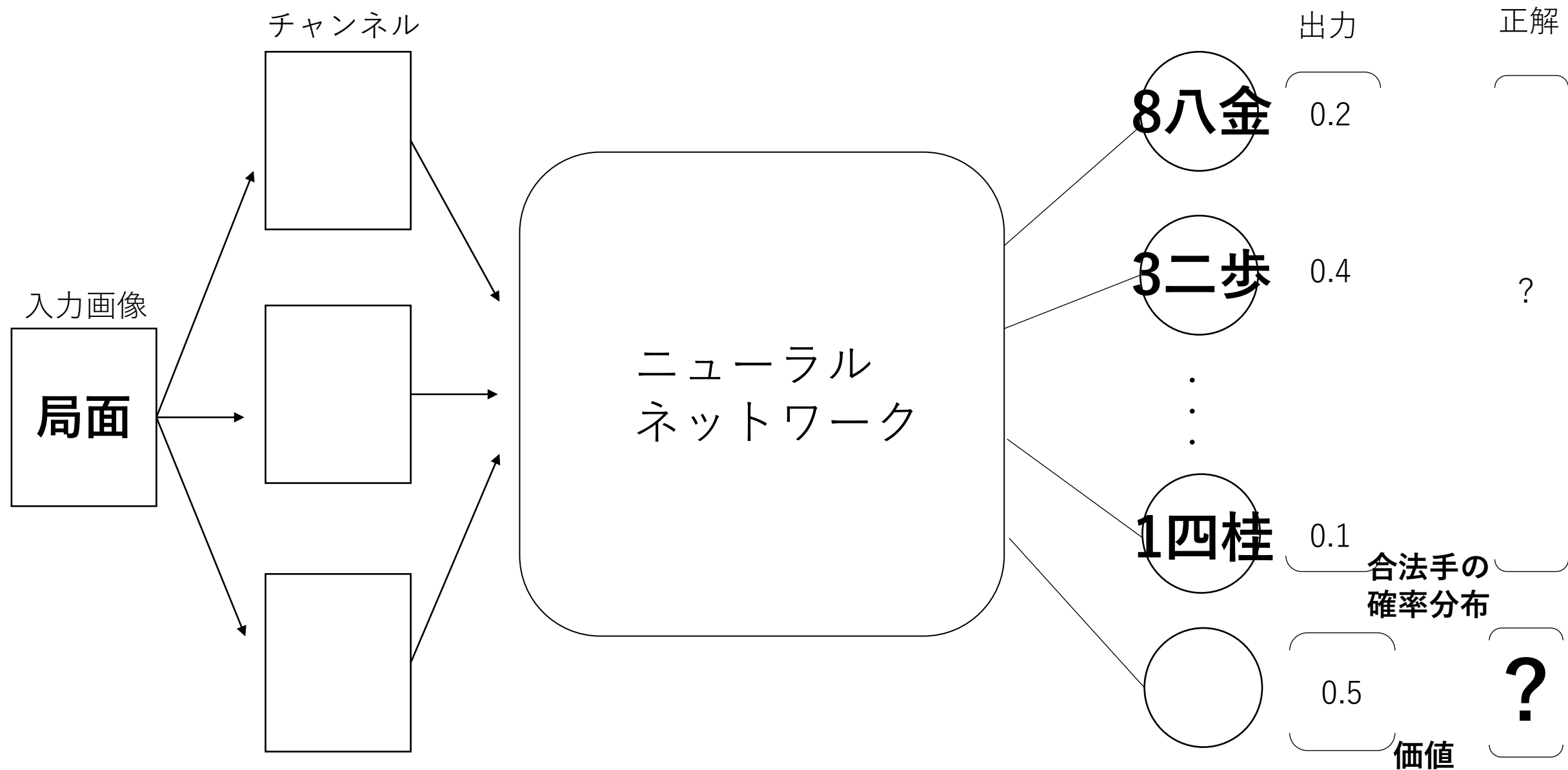
ニューラルネットワークの将棋への活用



AlphaZeroのニューラルネットワーク



AlphaZeroのニューラルネットワーク



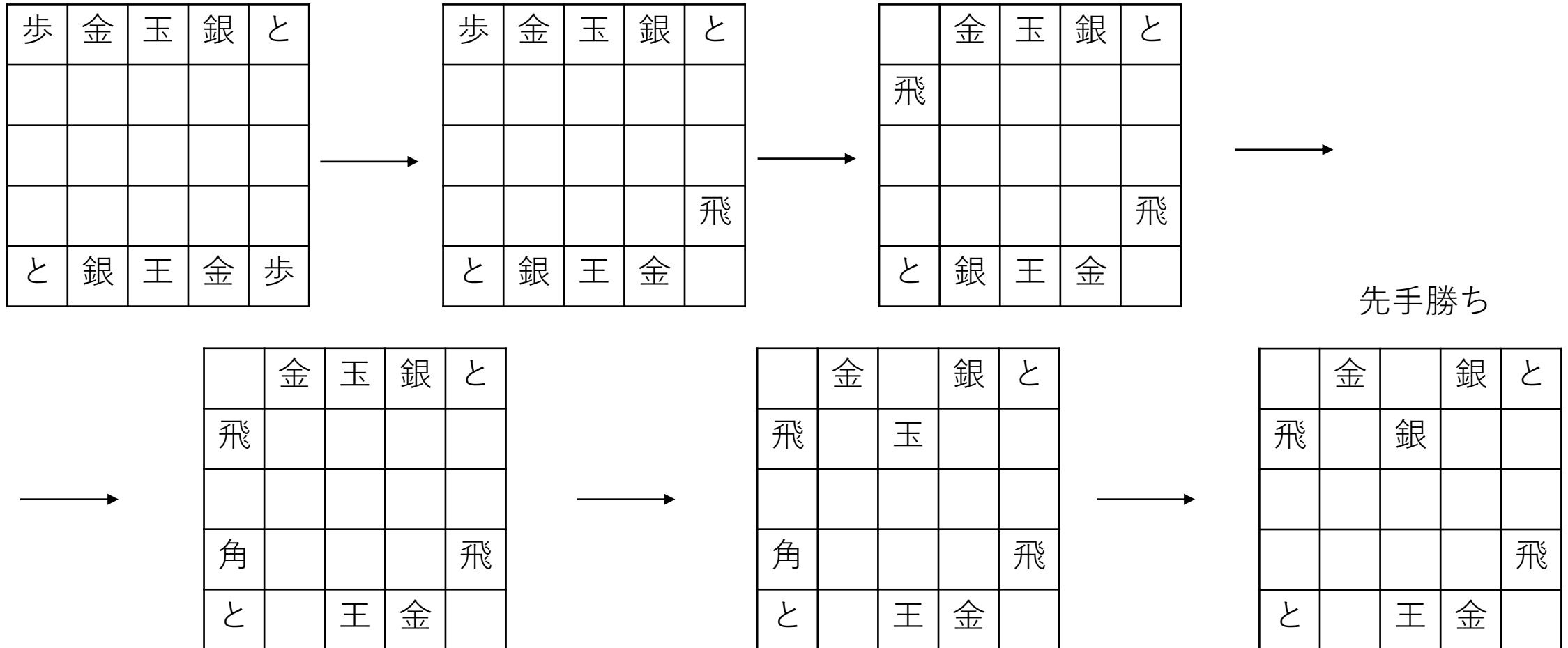
AlphaZeroのニューラルネットワークの「価値」の正解データ

勝った場合はその勝ちに至った局面の「価値」を1、
負けた場合は負けに至った局面を-1として学習させる。

学習後にさらに対戦させてデータをとって、
そのデータでさらに学習させてさらにさらに対戦させて
データを取って.....ってのを繰り返してよいなら、こんなに
単純なやり方でも多くの場合でうまくいく（らしい）

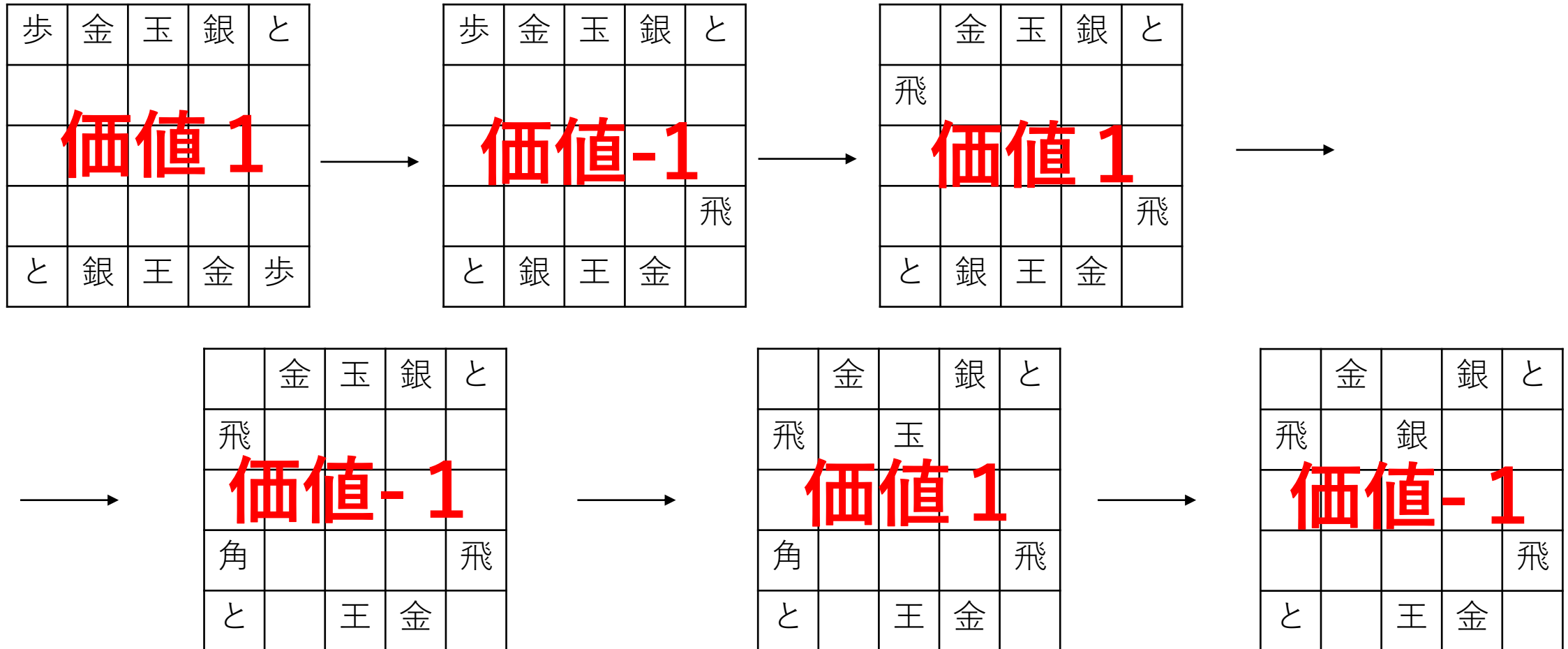
AlphaZeroのニューラルネットワークの「価値」の正解データ

勝った場合はその勝ちに至る局面の「価値」を1、
負けた場合は-1を正解とする。

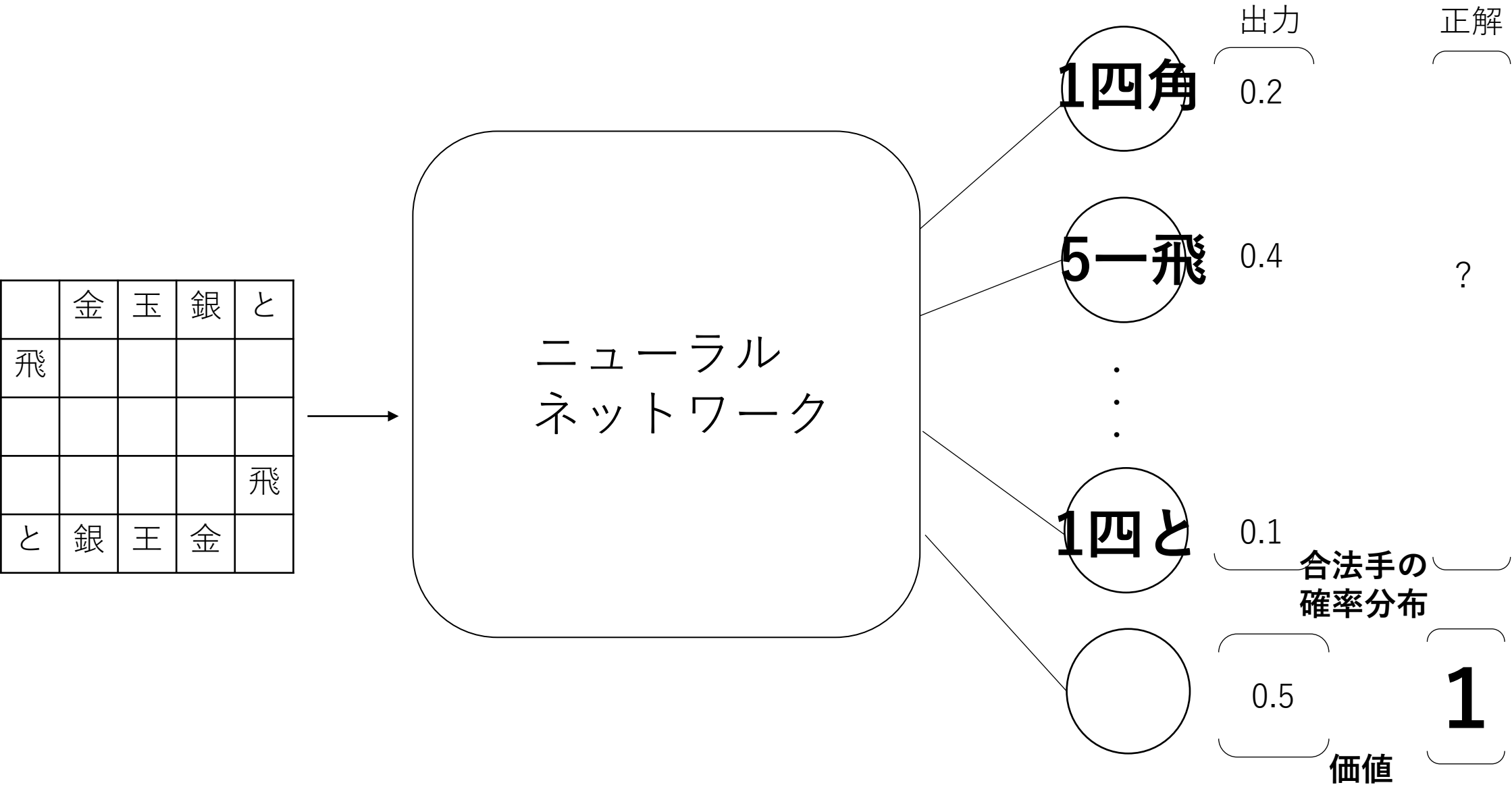


AlphaZeroのニューラルネットワークの「価値」の正解データ

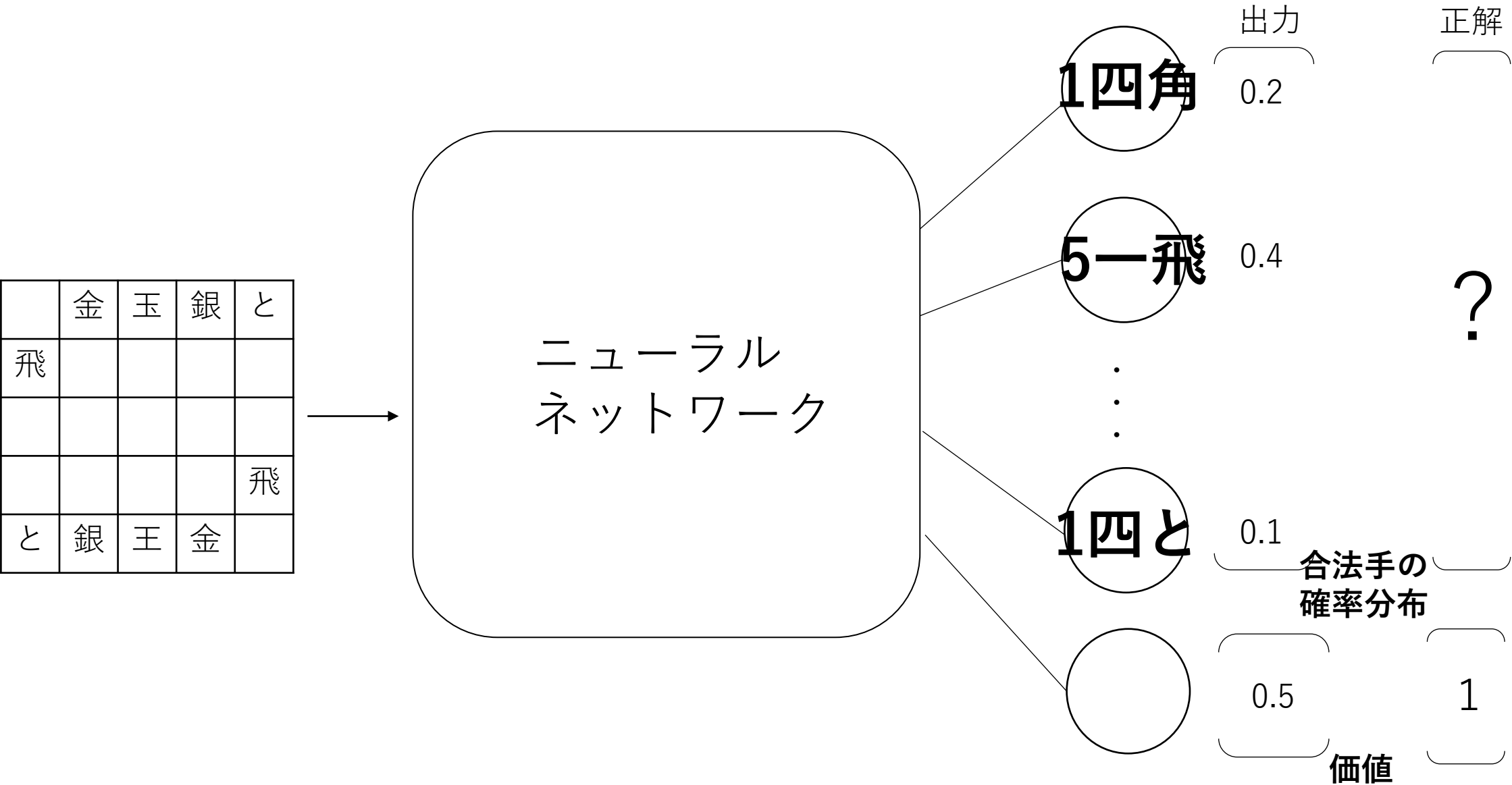
勝った場合はその勝ちに至る局面の「価値」を1、
負けた場合は-1を正解とする。



AlphaZeroのニューラルネットワーク



AlphaZeroのニューラルネットワーク



AlphaZeroのニューラルネットワークの「方策」の正解データ

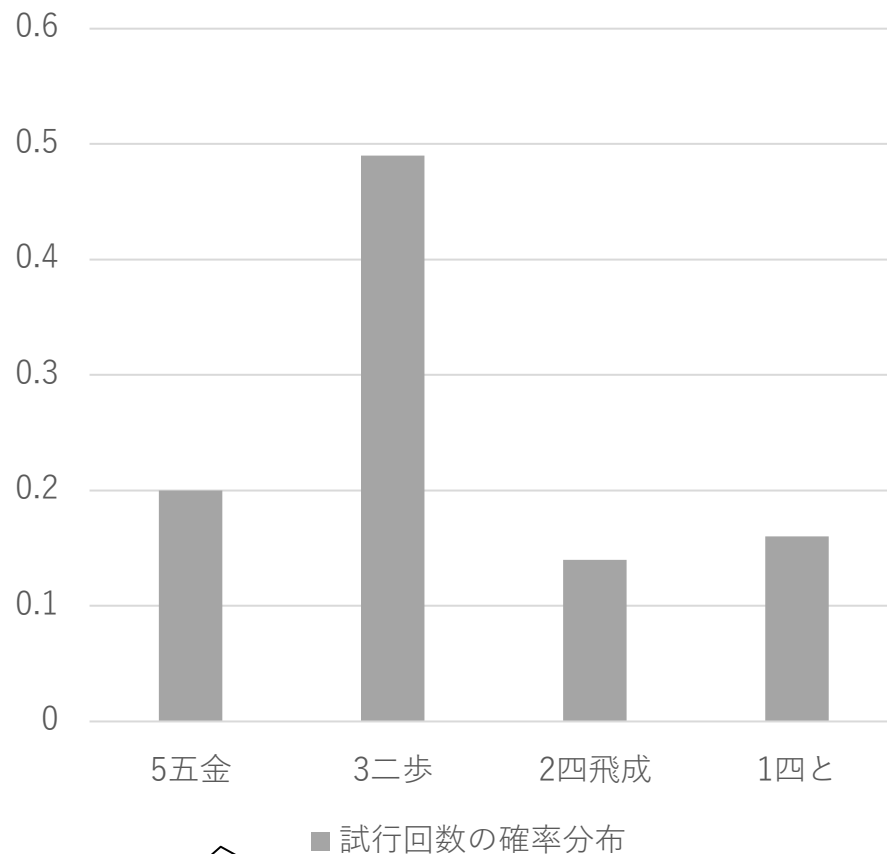
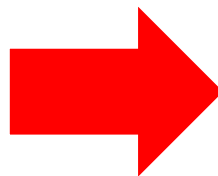
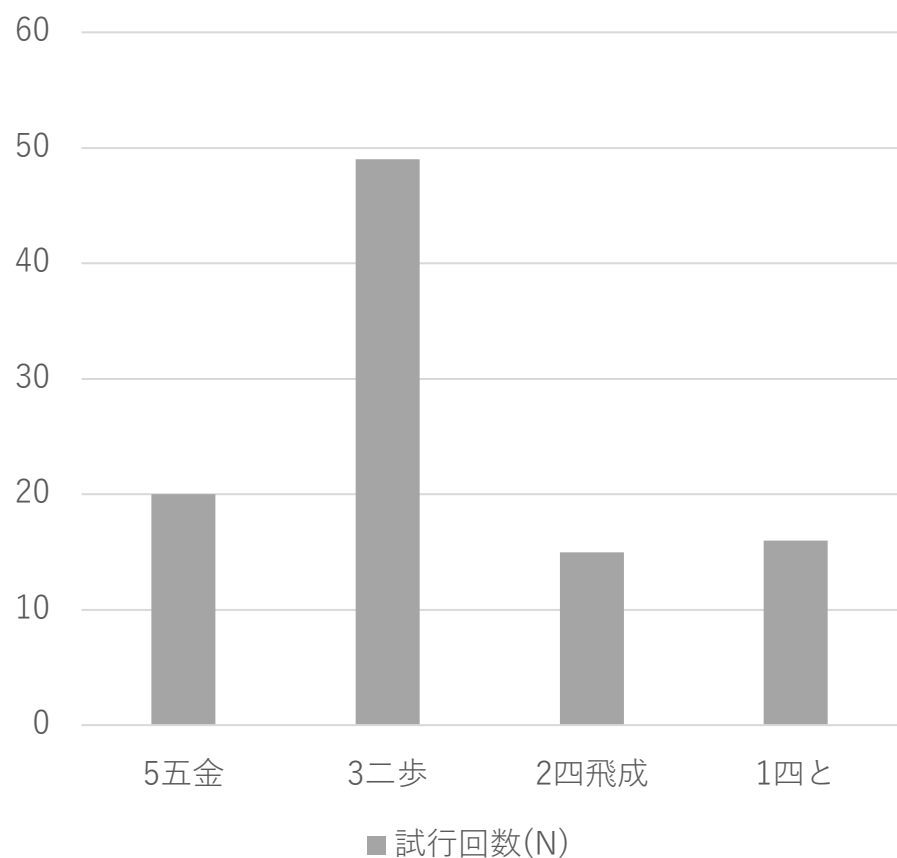
モンテカルロ木探索の試行回数の確率分布の値を正解にする。

勝ったにせよ負けたにせよ試行回数の値はモンテカルロ木探索を行ったため、深く読むべき局面と浅く読むべき局面を表現している、正解に近い値だと考える。

また、価値が正しく設定されるなら、その価値に沿う形に次回の対戦での試行回数が変わって、次は今よりもさらに正解に近い値となる。

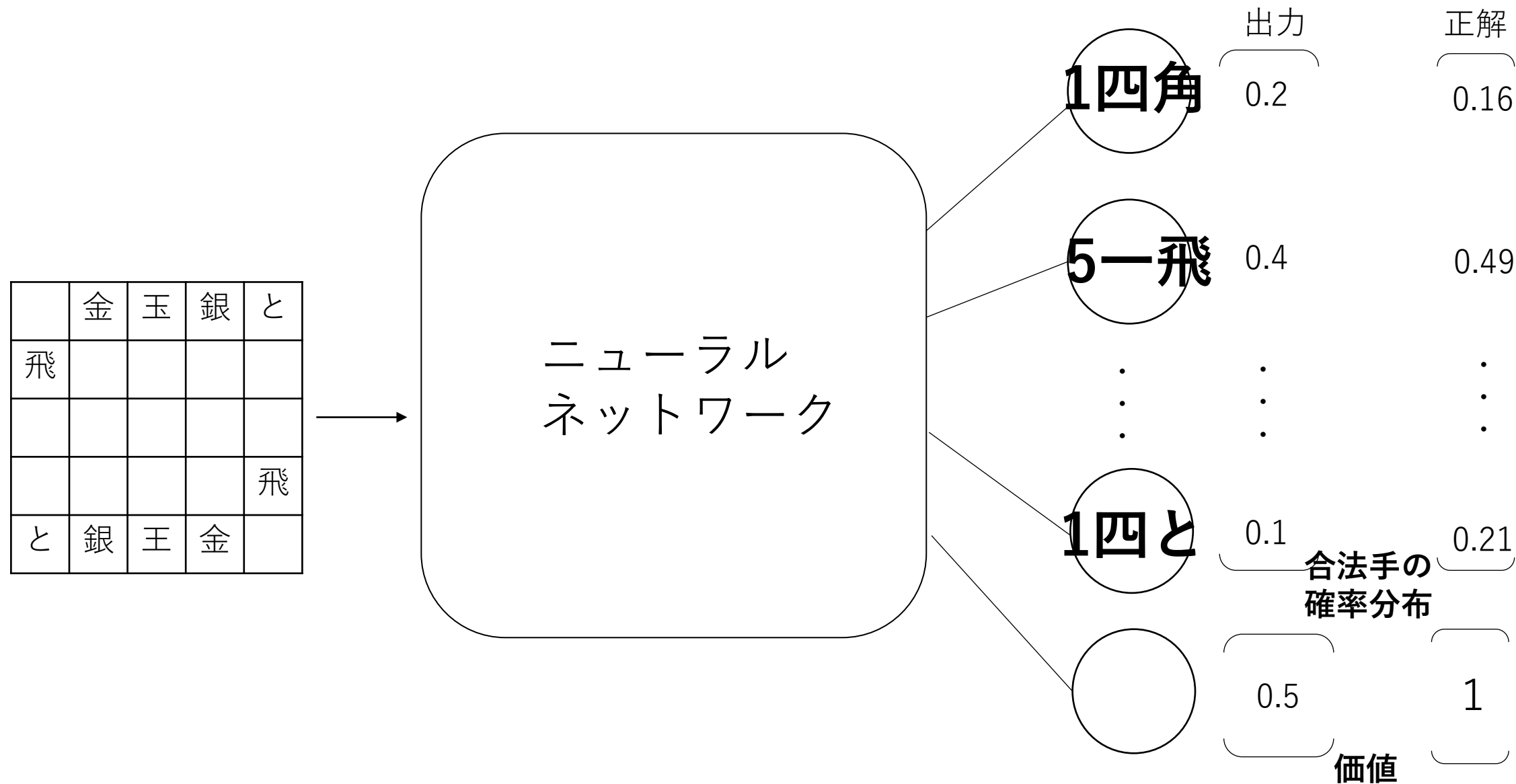
AlphaZeroのモンテカルロ木探索一行動選択

(学習時) (再掲)



この確率分布に従って1手を選ぶ

AlphaZeroのニューラルネットワーク



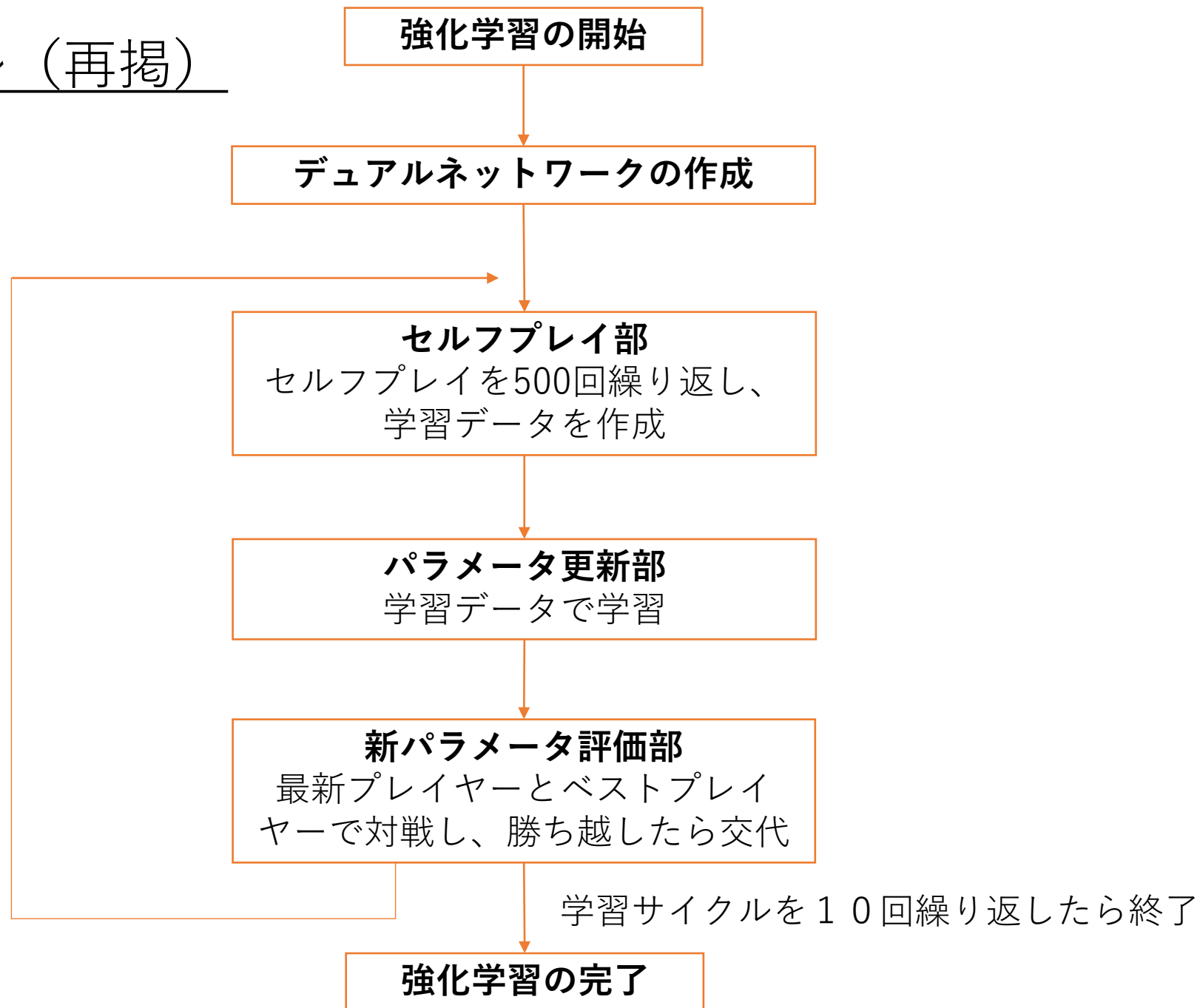
目次

- 京都将棋について
 - 京都将棋の概要
 - 京都将棋のルール
- 準備
 - 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
 - モンテカルロ木探索
 - ニューラルネットワーク
- AlphaZeroの仕組み
- 結果
- 参考文献

結果

- 10 サイクル
VS モンテカルロ木探索 0.78
VS AlphaBeta 0.20
- 20 サイクル
VS モンテカルロ木探索 0.86
VS AlphaBeta 0.22
- 30 サイクル
VS モンテカルロ木探索 0.94
VS AlphaBeta 0.21
- 40 サイクル
VS モンテカルロ木探索 0.95
VS AlphaBeta 0.18

学習サイクル（再掲）



目次

- 京都将棋について
 - 京都将棋の概要
 - 京都将棋のルール
- 準備
 - 「AlphaGo」、「AlphaGo Zero」、「AlphaZero」の違い
 - モンテカルロ木探索
 - ニューラルネットワーク
- AlphaZeroの仕組み
 - 探索
 - 学習
- 結果
- 参考文献

参考文献

- 布留川栄一,『AlphaZero 深層学習・強化学習・探索・人工知能プログラミング実践入門』,株式会社ボーンデジタル,2019.
- 曾我部東馬,『強化学習アルゴリズム入門ー「平均」からはじめる基礎と応用ー』,オーム社,2019.
- 『今年49歳になるおっさんでも作れたAlphaZero』,閲覧日 2020.-1-30,
<https://tail-island.github.io/programming/2018/06/20/alpha-zero.html>