

強化学習を用いた
拡張型七五三ゲームの自動生成

執筆者：田邊 繁

指導教官：新納 浩幸

平成14年3月1日

目次

第1章 序論	5
1.1 概要.....	5
1.2 本論文の構成.....	6
第2章 強化学習	7
2.1 強化学習.....	7
2.2 動的計画法.....	9
2.2.1 動的計画法.....	9
2.2.2 方策評価.....	10
2.2.3 方策改善.....	11
2.2.4 反復方策.....	12
2.2.5 価値反復.....	13
2.3 モンテカルロ法.....	16
2.4 TD学習.....	17
2.4.1 TD予測.....	17
2.4.2 方策オン型TD制御.....	18
2.4.3 方策オフ型TD制御.....	19
第3章 ニューラルネットワーク	21
3.1 ニューラルネットワーク.....	21
3.2 バックプロパゲーション法.....	22

第4章	七五三ゲーム	27
4.1	七五三ゲーム	27
4.2	最適方策の取得	27
4.2.1	七五三ゲームのモデル化	28
4.2.2	DPによる強化学習	28
4.2.3	VBによるインターフェイスの作成	29
4.3	評価	30
第5章	拡張型七五三ゲーム	31
5.1	拡張型七五三ゲーム	31
5.2	モデル化	33
5.3	学習アルゴリズム	33
5.4	評価結果	34
第6章	考察	37
第7章	結論	39
	謝辞	40
	参考文献	41
付録A	プログラムリスト (七五三ゲーム)	42
付録B	プログラムリスト (拡張型七五三ゲーム)	48

目次

2. 1	強化学習におけるエージェントと環境間の相互作用.....	8
2. 2	反復方策評価アルゴリズム.....	11
2. 3	反復方策アルゴリズム.....	13
2. 4	価値反復アルゴリズム.....	14
2. 5	動的計画法の学習モデル.....	15
2. 6	モンテカルロ法アルゴリズム.....	16
2. 7	TD を用いた V^{π} の予測アルゴリズム.....	18
2. 8	Q 学習 : 方策オフ型 TD 制御アルゴリズム.....	20
2. 9	TD 学習のモデル.....	20
3. 1	BP 法における出力のモデル.....	25
3. 2	BP 法における ω の更新のモデル.....	26
4. 1	VB による七五三ゲームの流れ.....	29
4. 2	VB による七五三ゲームインターフェイス.....	30
5. 1	ニューラルネットワークを用いたモデル化.....	32
5. 2	TD と BP 法を用いた学習アルゴリズム.....	35
5. 3	VB による拡張型七五三ゲームインターフェイス.....	36

表目次

4. 1	七五三ゲームにおける状態数と行動.....	28
5. 1	ニューラルネットワークの設定.....	33

第1章

序論

1. 1 概要

現在，機械学習において様々な問題の解決が試みられているが，その多くは分類問題である．しかし，現実の問題では分類問題として定式化できない問題も数多く存在する．

その中の一つに相互作用型の問題というものがある．これはある目的を持ったエージェントと提供される環境が相互に影響を与えあう問題で，強化学習を用いることが問題解決に有効である．

本研究では強化学習を実際に相互作用問題に適用できるようにすることを目的としてその例を解くものである．

これには状態数が非常に多数である場合の汎化手法が重要な問題となっている．そこで，今回は七五三ゲーム及び七五三ゲームを拡張したゲームの最適な行動を検討することとする．次の流れで実験を行う．

1 七五三ゲームでの最適解の取得

これにより実際に相互作用問題に対して強化学習を利用することができることを確認する．また強化学習の手法を学ぶ．

2 拡張型七五三ゲームでの最適解の取得

最も重要な問題である状態数が非常に多い問題の場合の汎化を検討し、最適方策の取得を試みる。

本研究の目的は状態数が非常に多い問題である拡張型七五三ゲームの最適な方策を、強化学習を用いて取得することである。

1. 2 本論文の構成

第2章において教科学習について述べる。強化学習にはいくつかの手法がありそれぞれについての説明をする。本研究では最終的にはTD学習を使用する。

第3章では汎化の手法として今回利用したニューラルネットワークについて述べる。基本的な内容と実際に用いるアルゴリズムである逆誤差伝播法について説明する。

第4章では七五三ゲームにおける最適方策を得るための実験について述べる。

第5章では拡張型七五三ゲームの最適方策の取得を試みる実験について述べる。概要、工夫点、また実験に使うアルゴリズムも同時に説明する。また実験結果とその評価を述べる。

第6章では実験の結果の考察、さらに本研究の目的である状態数が多数である問題の最適方策を得るための方針を述べる。

第2章

強化学習

2.1 強化学習

ここではまず強化学習について説明する。

まず、強化学習とは相互作用から学習して目標を達成する問題の枠組みそのものである。そこで最初に問題の枠組みから説明する。学習と意思決定を行うものは「エージェント」と呼ばれる。このエージェント外部から構成され、エージェントが相互作用を行う対象は「環境」と呼ばれる。両者は継続的に相互作用を行い、エージェントは行動を選択し、環境はその行動に応答し、エージェントに新しい状況を提示する。環境は報酬の発生源でもあり、この報酬はエージェントが時間の経過の中で最大化しようとする特別な数値である。エージェントの目標は最終的に受け取る報酬を最大化することである。

もっと詳細に言うならば、エージェントと環境は離散的な時間ステップ、 $t = 0, 1, 2, 3, \dots$ の各々において相互作用を行う。各時間ステップ t において、エージェントは何らかの環境の状態(state)の表現 $s_t \in S$ (S は可能な状態の集合)を受け取り、これに基づいて行動 $a_t \in A(s_t)$ を選択する($A(s_t)$ は状態 $s_t \in S$ において選択することのできる全ての行動の集合である)。1時間ステップ後に、エージェントはその結果として数値化された報酬 $r_{t+1} \in \mathcal{R}$ を受け取り、新しい状態 s_{t+1} にいることを知る。図にエージェントと環境との相互作用を示す。

この枠組みは抽象的かつ柔軟で、種々の問題に対して種々の方法で適用することができる。例えば、時間ステップを固定間隔の実時間で参照する必要はない。時間ステップは意思決定と行動の連続的な任意の段階とすることもできる。

このように考えることによって今回用いる七五三ゲーム，もしくはその他のゲームに対応できる。

次に学習について述べる。

強化学習問題を解くことは大雑把に言えば，最終的に多くの報酬獲得を達成するような方策(最適方策)を見つけることを意味する。これを価値関数に基づく評価により，行っている。この関数は状態(あるいは状態行動対)の関数で，エージェントがある状態にいることがどれだけ良いのか(もしくはある状態において，ある行動を行うことがどれだけ良いか)を評価する。全ての状態に対する最適な価値関数が得られたなら最適方策を得ることは容易である。1ステップ探索によって最適な行動が得られるのだ。

よって，強化学習とは最適な価値関数を得ることによって最適方策を得るものであるということが出来る。

次に最適価値関数を得るための学習アルゴリズムについて述べる。

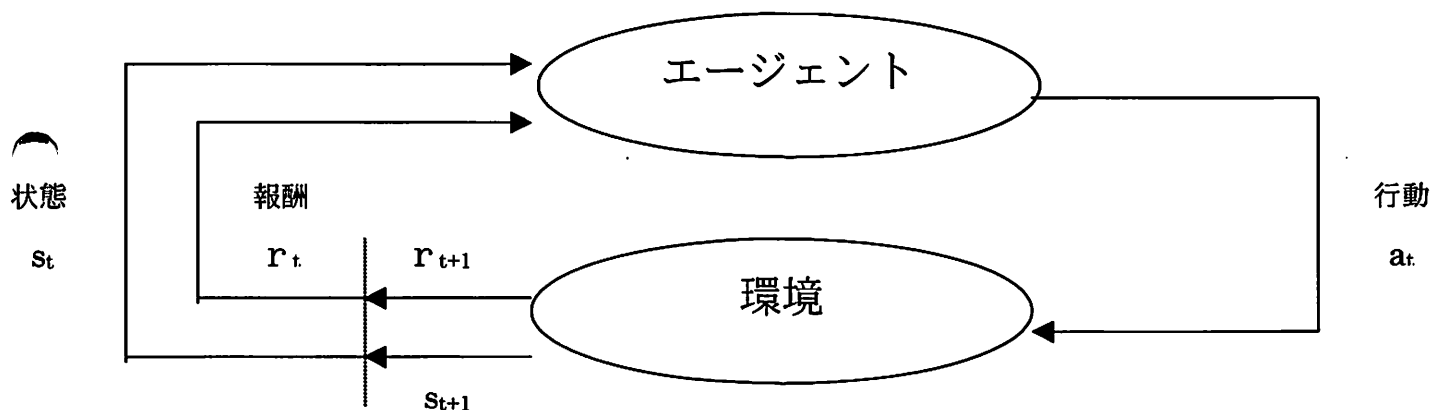


図2. 1 強化学習におけるエージェントと環境間の相互作用

2.2 動的計画法

2.2.1 動的計画法

動的計画法(DP)とは、最適方策を計算するためのアルゴリズムであり、環境の完全なモデルがマルコフ決定過程として与えられる場合に適用できる。DPアルゴリズムは、完全なモデルを前提としていて計算量も多くなる古典的なものだが理論的にはこの後の方法につながる重要なものである。

まずマルコフ決定過程について述べる。

状態信号が過去の知覚をコンパクトに集約しさらに過去の関連情報を全て保持していることをマルコフ性を持つという。このマルコフ性を満たす強化学習タスクをマルコフ決定過程という。また状態と行動空間が有限であれば有限マルコフ決定過程となる。DPは主に方策評価、方策改善、方策反復の三つからなる。

手法としてはさきも述べたように最適価値関数を用いたものである。価値関数には状態価値関数と行動価値関数があり状態価値関数は状態そのものの価値、行動価値関数はその行動がどの程度有効かを表す。下に式に示す。

$$V^*(s) = \max_a E\{r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a\}$$

$$= \max_a \sum_s P_{ss'}^a [R_{ss'}^a + \gamma V^*(s')]$$

$$Q^*(s, a) = E\{r_{t+1} + \gamma \max_a Q^*(s_{t+1}, a') | s_t = s, a_t = a\}$$

$$= \sum_s P_{ss'}^a [R_{ss'}^a + \gamma \max_a Q^*(s', a')]$$

2. 2. 2 方策評価

方策評価とは任意の方策 π に対する状態価値関数 V^π を計算するものである。これは予測問題とも呼ばれる。

まず、状態価値は以下の式で示す。

$$\begin{aligned} V^\pi(s) &= E_\pi\{r_{t+1} + \gamma r_{t+2} + \dots | s_t = s\} \\ &= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')] \end{aligned}$$

$\pi(s, a)$ は方策 π の下で、状態 s で行動 a を選択する確率を表す。添え字 π は π に関する条件付であることを示す。割引率 $\gamma < 1$ であるか、または方策 π の下で最終的にすべての状態から終端状態に至ることが保証されるなら、 V^π の存在と一意性は保証される。

環境のダイナミクスが完全に知っている場合これは簡単に求められる。計算時間のかかる場合もあるが、反復解法を用いるのが適している。 V^π を1ステップ観測で以下のように近似できる。これを全ての状態を対象としてバックアップを行うので完全バックアップと呼ぶ。

$$\begin{aligned} V_{k+1}(s) &= E_\pi\{r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s\} \\ &= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V_k(s')] \end{aligned}$$

これを各状態 s に適用して V_k から続く V_{k+1} への近似を得る。つまり以前の s の価値と即時報酬の期待値を使って新たな s の価値を計算し、それを以前の s の価値と置き換え、現在評価している方策の下で可能な1ステップ遷移のすべてに対してこの操作を行う。このような操作を完全バックアップと呼び、これを方策評価の中でも特に、反復方策評価と呼ぶ。

<p>評価対象の方策πを入力</p> <p>すべての$s \in S$に対して$V(s) = 0$とする</p> <p>繰り返し:</p> <p style="padding-left: 2em;">$\Delta \leftarrow 0$</p> <p style="padding-left: 2em;">各$s \in S$について:</p> <p style="padding-left: 4em;">$v \leftarrow V(s)$</p> <p style="padding-left: 4em;">$V(s) = \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$</p> <p style="padding-left: 4em;">$\Delta \leftarrow \max(\Delta, v - V(s))$</p> <p>$\Delta < \theta$ (正の小さな値) ならば, 繰り返しを終了</p>
--

図 2. 2 反復方策評価アルゴリズム

2. 2. 3 方策改善

方策に対して評価関数を計算することで, さらに良い手がかりが得られる. いま決定論的な方策 π に対して, 価値関数 V^π がきまっていたとする時, ある行動 $a \neq \pi(s)$ を行うことが良いことかを知ることがそのまま方策の改善につながる. そこで状態 s で行動 a を選択し, その後は既存の方策 π に従うことで得られるのが以下の式でこれを行動価値関数という.

$$\begin{aligned}
 Q^\pi(s, a) &= E_\pi \{r_{t+1} + \gamma V_k(s_{t+1}) \mid s_t = s, a_t = a\} \\
 &= \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')]
 \end{aligned}$$

この値を $V^\pi(s)$ と比べた場合の大小が判断基準となる. つまり, これを利用して, すべての状態で, すべての可能な行動に関して, 各状態で $Q^\pi(s, a)$ が最良となる行動を選択する変更が有効である. このようにして決定した方策をグリーディ方策といい次式で表せる.

$$\begin{aligned}
\pi'(a) &= \arg \max_a Q^\pi(s, a) \\
&= \arg \max_a E \{r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s, a_t = a\} \\
&= \arg \max_a \sum_s P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')]
\end{aligned}$$

この方策改善をすることにより最適な方策へ収束していく。

2. 2. 4 方策反復

これまでの流れで方策評価によって得られた V^π を使って方策改善を行い、より良い方策 π' を得できた。これを反復して行うことで最適方策を得るのが方策反復である。以下に詳しいアルゴリズムを載せる。

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \pi_2 \rightarrow \dots \rightarrow \pi^* \xrightarrow{E} V^*$$

のように π が更新されていく。ここで \xrightarrow{E} は方策評価を \xrightarrow{I} は方策改善を表す。

この場合の方策評価の開始点は一つ前に方策評価をして得た値となるので方策評価の収束スピードは 2 回目からは大幅に速くなる。また、この繰り返しで得られた方策は常に最小ステップで終端状態へ行くような最適方策となる。

しかし、大規模な問題においては方策評価と方策改善の繰り返しがたくさん必要となってしまうことがある。

以下に方策反復のアルゴリズムを示す。

1. 初期化
 $s \in S$ に対して $V(s) \in \mathbb{R}$ と $\pi(s) \in A(s)$ を任意に初期化する
2. 方策評価
3. 方策改善
 $policy\text{-}stable \leftarrow true$
 各 $s \in S$ について：
 $b \leftarrow \pi(s)$
 $\pi(s) \leftarrow \arg \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$
 もしも $b \neq \pi(s)$ ならば $policy\text{-}stable \leftarrow false$
 もし $policy\text{-}stable$ ならば、終了; それ以外は、2 から繰り返す

図 2. 3 反復方策アルゴリズム

2. 2. 5 価値反復

先ほど述べたように反復方策は大規模な問題だと手間がかかる。しかも、方策評価の際にスイープ操作を繰り返し行わなければならない。さらに厳密な収束は極限においてのみ得ることができるため、その際の反復の停止条件にも疑問が残る。

そこで考えられたのが価値反復である。これは方策評価と方策改善を同時に行うもので常に最適な行動、つまりグリーディ方策のみを1スイープ操作するものです。厳密に言えば、価値反復とは方策改善と方策評価のステップ打ち切りを組み合わせたものです。これはきわめて単純なバックアップ操作で、次の式で表される。

$$V_{k+1}(s) = \max_a E_{\pi} \{r_{t+1} + \gamma V_k(s_{t+1}) \mid s_t = s, a_t = a\}$$

$$= \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V_k(s')]$$

厳密には方策評価同様に無限回の反復で極限となり V^* へ収束するが、実際の終了条件は、一回のスweep操作による価値関数の変化量が十分に小さくなった時とする。

この価値反復は有限のマルコフ決定過程である強化学習問題にはとても有効である。

また動的計画法のモデル図を載せる。

V を任意の値で初期化
 繰り返し：
 $\Delta \leftarrow 0$
 各 $s \in S$ について：
 $v \leftarrow V(s)$
 $V(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 $\Delta < \theta$ (正の小さな数) ならば、繰り返しを終了

次式の決定論的方策 π を出力

$$\pi(s) = \arg \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$$

図2.4 価値反復アルゴリズム

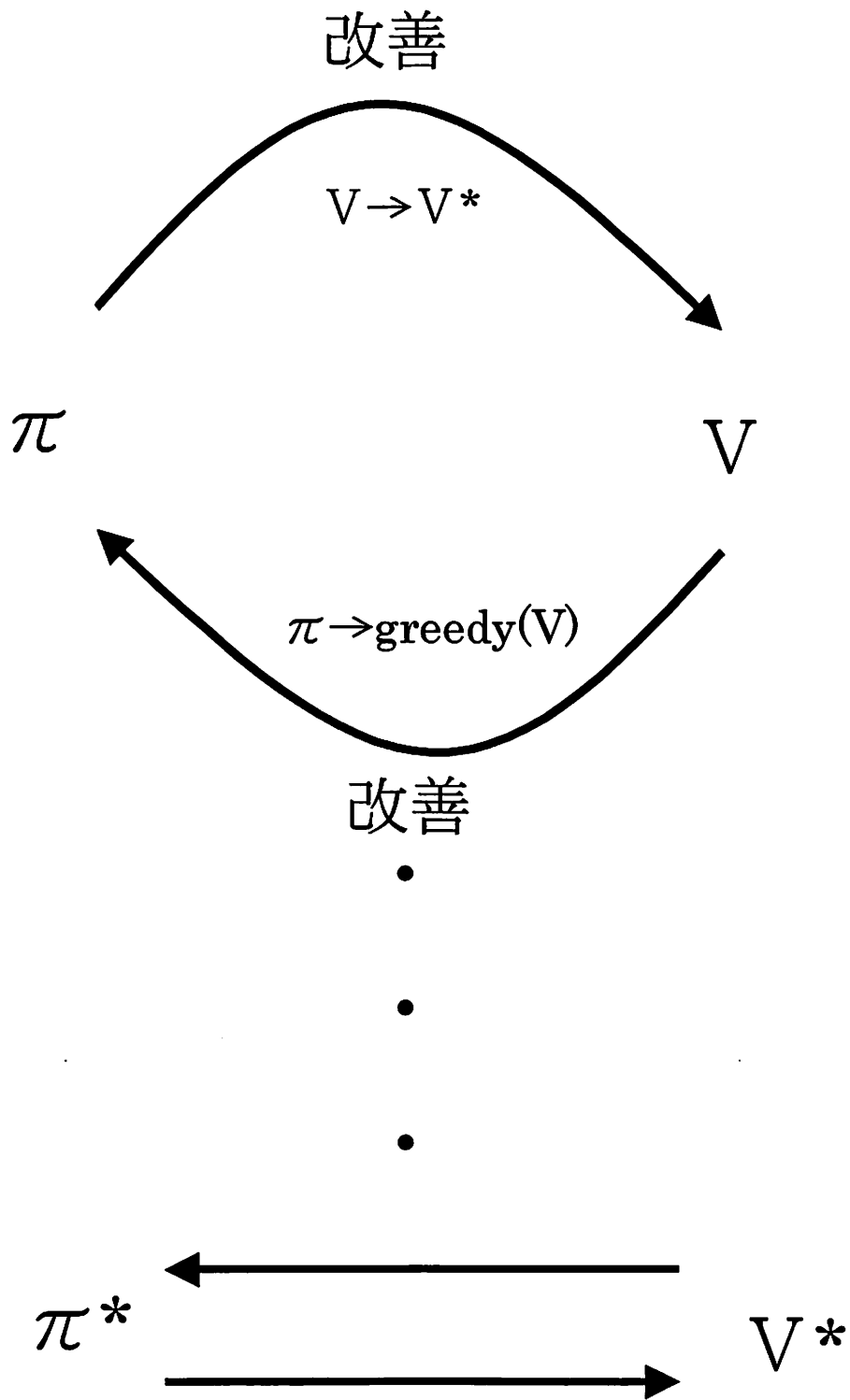


図2. 5 動的計画法の学習モデル

2.3 モンテカルロ法

前述の動的計画法は環境の完全なモデルが必要な学習アルゴリズムであったが、モンテカルロ法は環境のダイナミクスを必要としない、経験のみを必要とするものです。つまり、モンテカルロ法で必要としているのは遷移のサンプル生成だけで、遷移の確率分布がわからない場合にも使うことができます。

モンテカルロ法は、サンプル収益を平均化することに基づいて強化学習を解く方法であり、また、エピソード的なタスクのみ定義できます。経験はエピソードごとで、どのような行動が選択された結果、どのような最終結果を得るかを元にエピソードの終了時に価値の推定と方策の改善がされます。

モンテカルロ法では大数の法則から期待値を推測するのであらゆるすべての状態に無限回訪れることが理想です。しかし、方策によっては訪れることのない状態がある場合があります。そこで ϵ -グリーディ方策が用いられることがあります。これは確率 ϵ でランダムな行動をとるもので、これにより幅広い行動を助ける事ができる。

初期化： $\pi \leftarrow$ 評価すべき方策 $V \leftarrow$ 任意の状態価値関数 $returns(s) \leftarrow$ 空のリスト(すべての $s \in S$ に対して)
繰り返し： (a) π を用いてエピソードを生成する (b)エピソード中に出現する各状態 s について： $R \leftarrow$ s の発生後の収益 R を $returns(s)$ に追加する $V(s) \leftarrow average(returns(s))$

図2.6 モンテカルロ法アルゴリズム

2.4 TD学習

2.4.1 TD予測

動的計画法，モンテカルロ法これらを進化させたのがTD学習法である。TDとモンテカルロ法は，いずれも経験を用いて予測問題を解決し，方策 π に従って経験をいくつか得ることで V^π の推定値 V を更新する。しかし，TDはモンテカルロ法とは違い，エピソードの終わりまで更新を待つ必要がない。動的計画法と同様にそれまでの学習から得た次の状態の推定値を用いて1ステップごとに更新する。これを，ブートストラップという。

つまり動的計画法とモンテカルロ法をあわせたのがTDとすることが出来る。

時刻 t で訪問した非終端状態を s_t として，その後起きる出来事によって $V(s_t)$ を更新する。式に表すと次のようになる。

$$V(s_t) \leftarrow V(s_t) + \alpha[R_t + V(s_t)]$$
$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

α はステップサイズ・パラメータを表す定数で，更新する幅を指定する。
また状態価値関数は以下になる。

$$V^\pi(s) = E_\pi\{R_t | s_t = s\}$$
$$= E_\pi\{r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s\}$$
$$= E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s\}$$

これを利用して V^π の推定値を予測する。

<p>$V(s)$を任意に初期化し, πを評価対象の方策に初期化する</p> <p>各エピソードに対して繰り返し:</p> <p> sを初期化</p> <p> エピソードの各ステップに対して繰り返し:</p> <p> $a \leftarrow s$に対してπで与えられる行動</p> <p> 行動aをとり, 報酬rと次状態s'を観測する</p> <p> $V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$</p> <p> $s \leftarrow s'$</p> <p> sが終端状態ならば繰り返し終了</p>
--

図 2. 7 TDを用いた V^π の予測アルゴリズム

2. 4. 2 方策オン型TD制御

次にTD予測法を制御問題に適用する方法について述べる。評価を予測部分にTD法を利用し、一般化方策反復の考え方を適用していく。この時、調査と知識の利用のトレードオフの問題がある。これに対するアプローチには方策オン型と方策オフ型の二つの手法がある。

方策オン型のTD制御では、現在の方策 π における状態 s と行動 a の全てに対して $Q^\pi(s, a)$ を評価しなければならない。これはTD予測法を行動価値関数に置き換えることによって簡単に定義できる。

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

この更新式は状態が遷移するたびに更新されるが、とるべき行動はエピソードの前にあらかじめ決められた行動となる。つまり、はじめは Q^{π^1} のそった行動を行い、次のエピソードでは Q^{π^2} のそった行動を採る。

ある状態行動対から次の状態行動対への遷移を定義する

$$s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$$

から Sarsa と呼ばれる。

2. 4. 3 方策オフ型 TD 制御

方策オン型はエピソードの開始時点ですべての状態に対する行動を観測しておくものだったが、それに対して方策オフ型はステップごとに行動を決定する。今回は方策オフ型の手法を用いて TD 学習を行うこととした。これを特に Q 学習と呼ぶ。

Q 学習は以下のように定義される。

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a)]$$

ここで学習される行動価値関数は、使われている方策とは独立に、 Q^* を直接近似する。これによって 1 ステップごとに新しい方策を生成、使用することになるので、収束が早まる。今回の実験ではこれを状態価値関数に置き換え使用した。

$Q(s, a)$ を任意に初期化
 各エピソードに対して繰り返し：
 s を初期化
 エピソードの各ステップに対して繰り返し：
 Q から導かれる方策を使い，
 s での行動 a を選択する
 行動 a を取り，報酬 r と次状態 s' を観測する
 $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$
 $s \leftarrow s'$
 s が終端状態ならば繰り返しを終了

図2. 8 Q学習：方策オフ型TD制御アルゴリズム



図2. 9 TD学習のモデル図

第3章

ニューラルネットワーク

3.1 ニューラルネットワークの基礎

今回、状態をニューラルネットワークを用いた関数近似とするため、ニューラルネットワークに基礎と実験で使用したバックプロパゲーション法について述べる。

ニューラルネットワークとは、生物の神経系の特徴的な機能に着目してそのモデル化を行ったものである。簡単に言うとある信号がシナプスへ入力として送られるとシナプス結合強度により変化した信号が出力される、ニューロンをモデル化したものである。

このニューロンモデルを、複数接続することによってニューラルネットワークを作ることができる。このとき、ニューロンの接続のしかたでネットワークの構造と特徴が決まる。ニューラルネットワークは構造によって大きく分けて2つの種類に分けられる。1つは階層型ネットワークといわれるもの。これはニューロンを層状にならべ、前の層から次の層へ方向にのみ信号が伝わるというネットワークである。このタイプのネットワークは入力層のニューロンに加えた信号（入力信号）に対して出力層のニューロンの出力（出力信号）が一意的に定まる。今回使用するバックプロパゲーション法はこれにあたる。

他には相互結合型ネットワーク、フィードバック付き相互結合型ネットワークがある。これらは流れが一方向に定まっていず、また、フィードバックを持つ。

3. 2 バックプロパゲーション法

バックプロパゲーション法は階層型のニューラルネットワークで最急降下法と呼ばれる、誤差Eが小さくなるようにする結合荷重 ω を求める方法を利用したものである。この方法では、関数のある点での勾配を求めて誤差Eが少なくなる方向へ結合荷重 ω を変化させていくこのことになる。

BPで最急降下法の考え方をを使うために、1つ準備しておくことがある。それは、微分可能な連続関数を用意することで、そこで使うものが関数をシグモイド関数である。これを用いると連続な関数になると共に出力が0～1の値になる。

$$f(x) = \frac{1}{1 + \exp(-x)}$$

上の関数がシグモイド関数である。

次にバックプロパゲーション法での出力について述べる。これは

$$x_{ip} = h_i(z_{ip})$$
$$z_{ip} = \sum_j \omega_{ij} y_{jp}$$

で求められる。これは入力側のニューロンの値とそれぞれの層の間の結合強度をかけた値、これをすべて足したものが出力となることを表している。関数hはシグモイド関数である。このように出力は前向きの演算を行うことによって求められる。

次に近似関数の更新。これは先ほども述べたように結合荷重 ω による近似関数なので ω を更新していく。まず、簡単に仕組みを説明する。

ある入力と出力をの時、BP法は教師あり学習なので、本来出力されるべき値、教師信号がある。この教師信号と出力の値の誤差を小さくしていくのが最急降下法という手法である。この誤差(P番目の入力に対するものとする)を誤差関数 $E_p(\omega)$

といい、すべての入力に対する誤差を相誤差関数 $E(\omega)$ と言う。

$$E_p(\omega) = \frac{1}{2} \sum_i (\bar{x}_{ip} - d_{ip})$$

$$E(\omega) = \sum_p E_p(\omega)$$

BP法はこの相誤差関数を最小化するため、各パターンに対する誤差関数の勾配ベクトルを計算し、その逆方向に ω を改良する。 ω はいくつもの成分を持つ多次元関数なのでそれぞれの成分ごとに考える。

$$\frac{\partial E_p(\omega)}{\partial \omega_{ij}} = \frac{\partial E_p(\omega)}{\partial z_{ip}} \frac{\partial z_{ip}}{\partial \omega_{ij}}$$

$$\begin{aligned} &= \frac{\partial E_p(\omega)}{\partial z_{ip}} y_{ip} \\ &= -\delta_{ip} y_{ip} \end{aligned}$$

となる。よって $-\delta_{ip}$ を求めることによって ω のそれぞれの変化させる量がわかることになる。この値は第 i ニューロンが出力層の場合とそうでない場合によって違う値となる。

(i) 第 i ニューロンが出力層の場合

出力 x_{ip} は相誤差関数における \bar{x}_{ip} に相当するので

$$\frac{\partial E_p(\omega)}{\partial \omega_{ij}} = \bar{x}_{ip} - d_{ip}$$

となるので

$$\delta_{ip} = -(\bar{x}_{ip} - d_{ip}) h'_i(z_{ip})$$

となる。

(ii) 第 i ニューロンが出力層でない場合

第 i ニューロンを s 層とすると出力 x_{ip} が変化することによって $(s+1)$ 層の各ニューロンへの入力の総量 z_{kp} が変化し $(s+1)$ 層の各ニューロンからの出力も変化する。よってこれらの変化が $(s+1)$ 層以降のニューロンを伝わっていき $E_p(\omega)$ の値を変動させる。

$$\frac{\partial E_p(\omega)}{\partial x_{ip}} = \sum_k \frac{\partial E_p(\omega)}{\partial z_{kp}} \frac{\partial z_{kp}}{\partial x_{ip}}$$
$$\frac{\partial E_p(\omega)}{\partial z_{kj}} = \frac{\partial E_p(\omega)}{\partial x_{kp}} h'_k(z_{kp}) = -\delta_{kp}$$
$$\frac{\partial z_{kp}}{\partial x_{ip}} = \omega_{ki}$$

以上の式から、

$$\delta_{ip} = h'_i(z_{ip}) \sum_k \delta_{kp} \omega_{ki}$$

となる。

よって、以下のような形で記述される結合強度ベクトル変更則により変更される。 η は固定された値で、更新する幅、ステップ幅を表している。

$$\omega^{(k+1)} = \omega^{(k)} - \eta \left. \frac{\partial E_p(\omega)}{\partial \omega_{ij}} \right|_{\omega = \omega^{(k)}}$$

ただし、BP 法では ω の初期値によっては、 $E(\omega)$ の大域的最小値へ収束するとは限らず、局所的最小値に陥ってしまう場合がある。

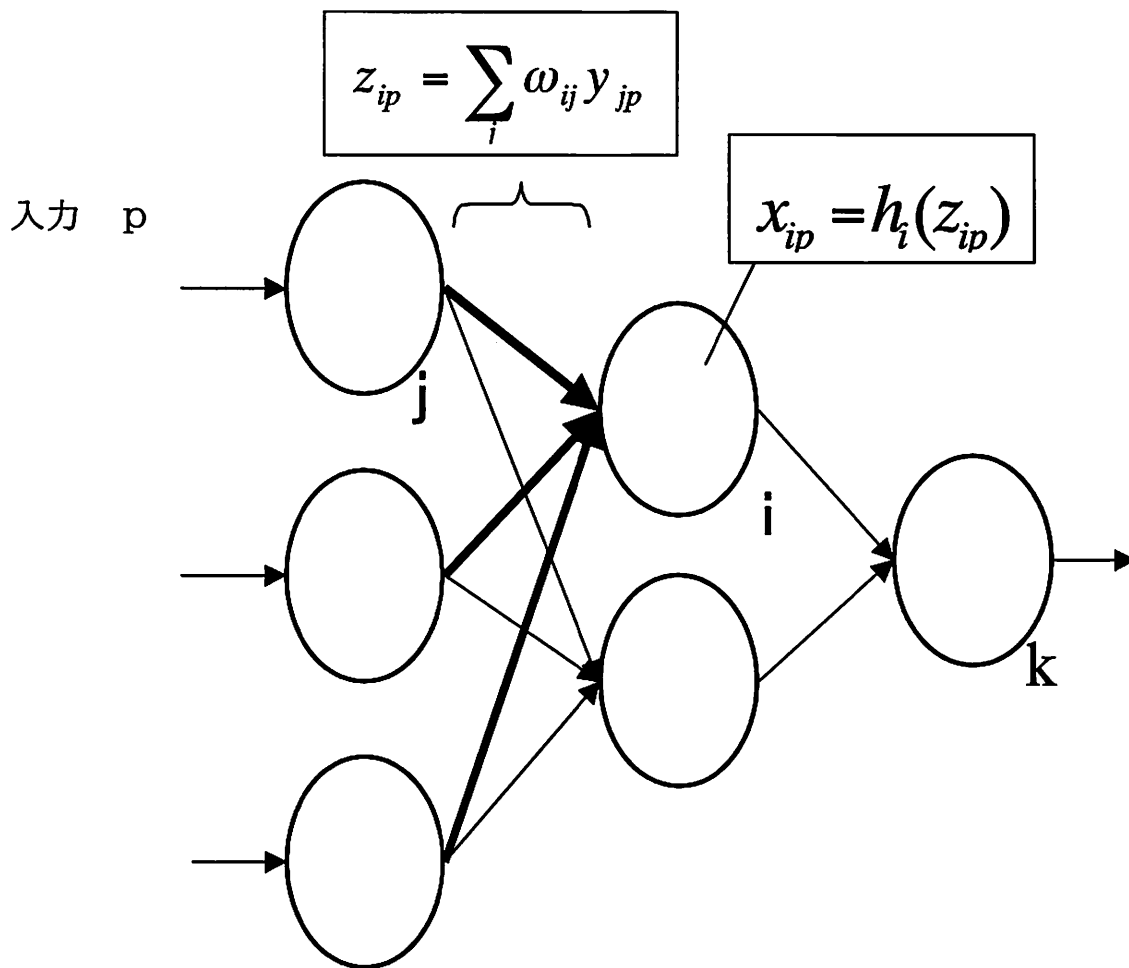


図3. 1 バックプロパゲーションにおける出力のモデル図

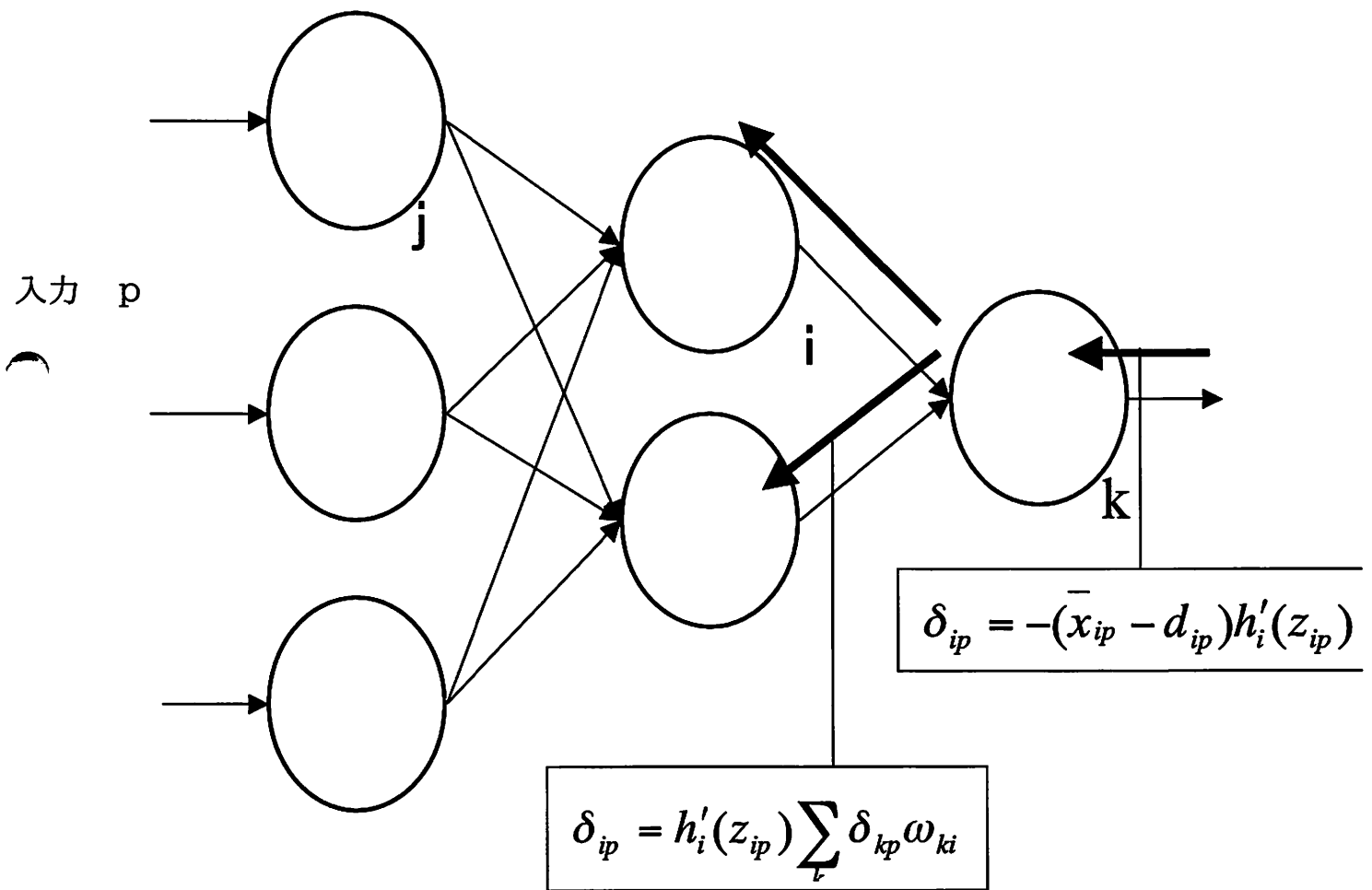


図3.2 BP法における ω の更新のモデル図

第4章

七五三ゲーム

4.1 七五三ゲーム

一つ目の実験として実際に七五三ゲームにおける最適方策の取得を行う。

まず、七五三ゲームについての説明をする。七五三ゲームは3本、5本、7本の棒を二人で交互に消していくゲームである。最後の棒を消した方が負けとなる。細かいルールは、一つで棒の消し方である。これにはいくつかの規則があり一度に消す棒の数が制限される。まず、棒を横に並べた塊(それぞれ3本、5本、7本)がある訳だが、棒を複数消す場合にはそれらの塊の一つからしか消すことが出来ない。また、隣り合う棒がすでに消されている場合、それを飛び越えてさらにもう一つ隣の棒を消すことは出来ない。

このゲームではすでに必勝法、つまり最適方策が存在するので、ここでは強化学習による最適方策の取得が有効であることの確認、検討を行う。

4.2 最適方策の取得

次の流れにそって最適方策取得、取得した方策の評価をする

- (1) 七五三ゲームのモデル化
- (2) DPによる強化学習
- (3) VBによるインターフェイスの作成

4. 2. 1 七五三ゲームのモデル化

七五三ゲームを数学的なモデルに置き換える。七五三ゲームのような盤面によるゲームはすべての盤面の状態を変数として持つことがもっとも簡単である。しかし、それでは状態数が増え計算量が増えることになる。よって、ゲームの特性に合わせて状態を定義することが有効になる。

七五三ゲームでは「何本の棒の集まり(消されていない棒)がいくつあるか？」がもっとも重要になるため、すべての状態を「N本の集まりがM個残っている状態」として定義する。

次にこのモデルにそった状態の抽出と、それぞれの状態における行動とそれによる状態の遷移を抽出する。状態と状態の遷移の抽出には JAVA を用いたプログラムを作成。状態は、すべての状態における N(1~7)に対する M を調べ重複を取り除くことで抽出する。行動は、対象とする N 本の塊をどのように消すか、つまり消す対象と消し方の二つに分けた行動を行い、それぞれの行動後の状態を観測、先に得た状態と対応させた。以下の結果となる。

状態数	2 5 1
行動(すべての状態において共通)	5 1

表 4. 1 七五三ゲームにおける状態数と行動

4. 2. 2 DPによる強化学習

七五三ゲームは状態数 2 5 1 と比較的小さな問題であり 環境(行動後の詳細な遷移の分布)が得られる。また七五三ゲームはすべての状態においてすべての棒が消される終端状態への遷移が約束されたマルコフ決定過程なので DP を用いることができる。

アルゴリズムは反復方策を用いる。方策評価を行う value メソッドと方策改

善を行う **kaizen** メソッドからなる。ゲーム性から $\gamma=0$ で割引はないものとする。また方策評価の際の終了条件は $\theta = 0.001$ とする。

巻末にプログラムのソースを添付する。

4. 2. 3 VBによるインターフェイスの作成

VBを使い強化学習で導き出した方策を実行するゲームを作成する。この方策とはすべての状態の最適勝ち関数を指す。ゲームでは先攻、後攻を任意に選択でき、それぞれの場合での結果を評価する。

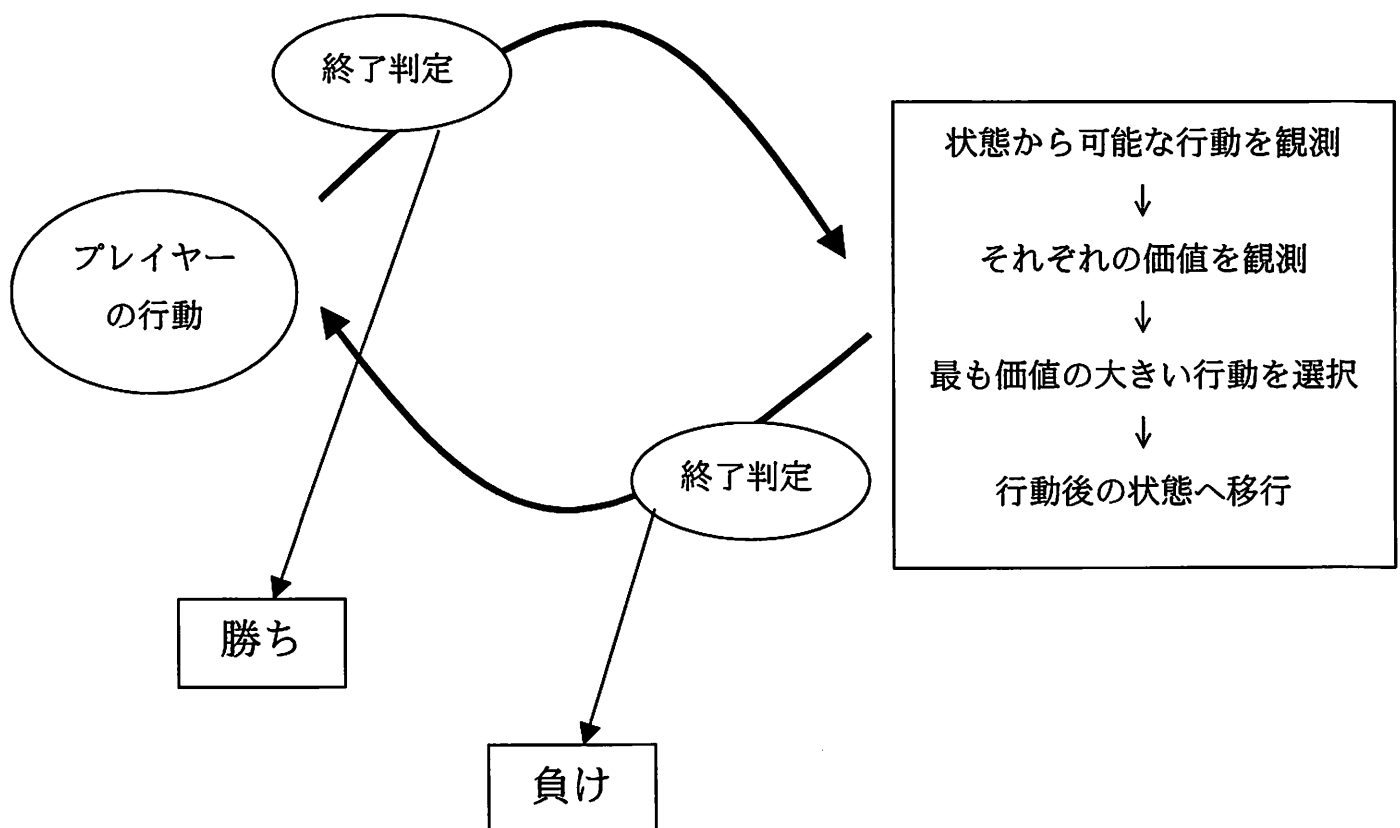


図4. 1 VBによる七五三ゲームの流れ

4.3 評価

実際にVBによる最適方策を搭載したゲームでの試行を行った結果、最適な方策がエラ亭たことが確認された。

後攻(プレイヤー)の場合、このゲームで勝つことは不可能となった。これは既存の七五三ゲームの必勝法とも合致する。先攻(プレイヤー)においても、必勝法を知らない者だと、つまり一度でも手を間違えると勝利することは出来なかった。

このことから強化学習を用いることによって七五三ゲームの最適方策を得ることができることが確認された。

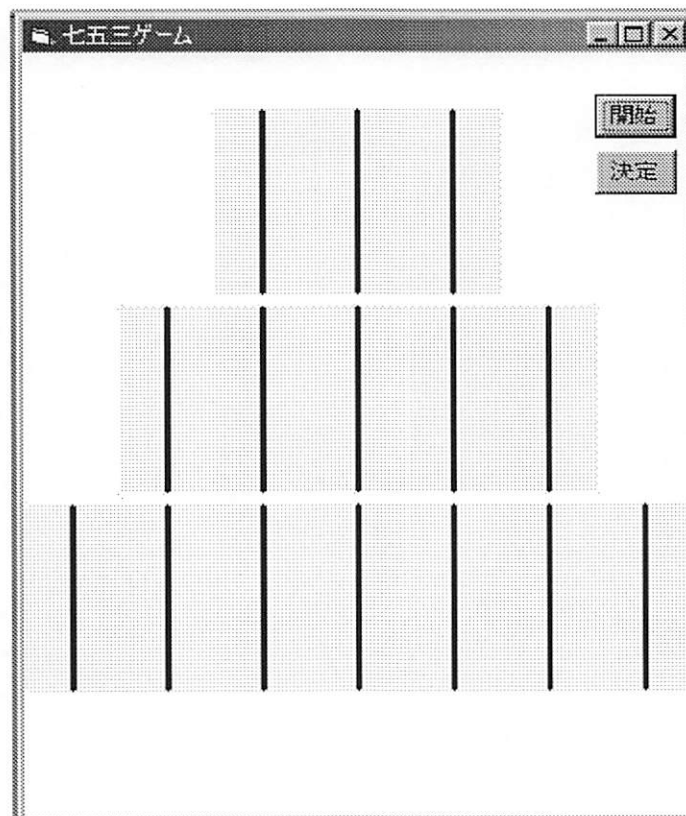


図4.2 VBによる七五三ゲームインターフェイス

第 5 章

拡張型七五三ゲーム

5. 1 拡張型七五三ゲーム

先ほどは七五三ゲームの最適方策の取得を行った。ここで状態数が非常に s 多数である問題として拡張型七五三ゲームを提案する。

拡張型七五三ゲームとは、先ほどの七五三ゲームを拡張したもので七五三ゲームが 3, 5, 7 本の棒の集まりを状態としたのに対して、8 行 8 列の升目状の状態をもつゲームである。

これによって状態数はすべてのマスの状態が、消されたか、消されていないかによって決定するので、

$$2^{(8 \times 8)} = 2^{64}$$

の状態数を持つことになる。さらにルールを一部改良した。升目状にしたことにより一度に複数のマスを横方向だけでなく、縦方向にも消せるように変更した。

その他の基本的なゲームの流れは七五三ゲームと同様とする。

実験は以下の流れで行う。

- (1) 拡張型七五三ゲームのモデル化
- (2) 学習アルゴリズムの模索
- (3) 学習プログラムの作成
- (4) 取得した最適方策の評価

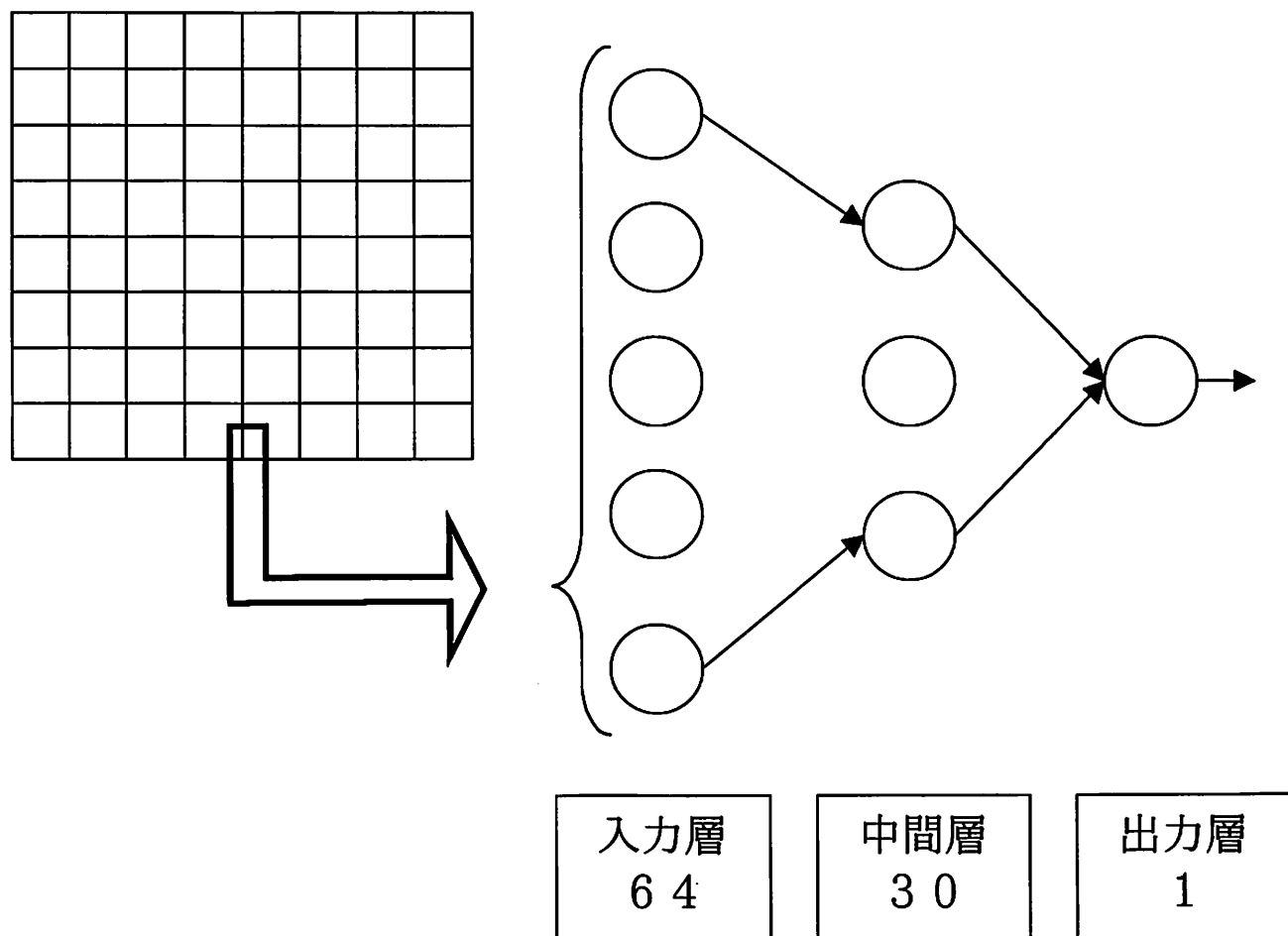


図5.1 ニューラルネットワークを用いたモデル化

5. 2 モデル化

拡張型七五三をモデル化しなければならない訳だが、七五三ゲームと同様に等価値の状態動詞をまとめたとしても、状態数が非常に多くすべての状態を変数として持つことは難しくメモリの確保も出来ない。

また状態の遷移も複雑で完全な状態の遷移モデルを作るのは困難である。

よって、サンプル生成だけを必要とするTD学習とニューラルネットワークを使った近似関数による汎化を用いてモデル化を行う。

ニューラルネットワークでの近似関数の ω はランダムに初期値を決定することにする。ニューラルネットワークの設定は以下の表のとおりとする。これはニューラルネットワークの学習の際のデータと参考資料から決定した。

入力層	64 (8×8)
中間層	30 (入力半分程度)
出力層	1
ステップ幅	0.13
ω の初期値	ランダム

表5. 1 ニューラルネットワークの設定

5. 3 学習アルゴリズム

5. 2のような近似関数を用いてTD学習を行う。よってTD学習を近似関数を用いた形に改良する必要がある。

TDとは簡単に言えば、次状態の価値を観測し、現在の状態の価値をそれに近づけることによって最適状態価値観を見つけることによって最適方策を得るものである。

これに対してニューラルネットワーク、特にBPとはある入力に対する出力

を教師信号と比較し、この誤差から w を改善していくものである。

両者に共通する点は、向かうべき目標値と現在提示されている推定値の誤差から改善を行うことである。そのことからTD学習から今日しい信号を生成するという手法を試みることにする。

ニューラルネットワークのBP法は本来教師あり学習である。しかし、強化学習は教師を必要としない学習なのでBP法における教師を与えることが出来ない。そこで、TD学習の手法を用いて教師信号を生成することにする。

具体的に言うと、ある状態と次状態があった時、次状態を入力した時の出力をそのまま、現在の状態を入力した時の教師信号として利用する。これをエピソードの毎ステップ行う。そして状態が終端状態になった時、終端状態を入力、報酬である0を教師信号としてBPによる関数近似を行う。

するとまず終端状態を入力とした時の出力が0へと改善され、以降はTD学習のように一つずつ前の状態へと伝播されて行く。

この手法を用いて学習プログラムを作成。作成にはJAVAを用いる。

評価方法として、学習プログラムから得られた近似関数を用いて次状態の価値を観測しグリーディ方策をとるものと、ランダムな行動を選ぶプログラムとの対戦を一万回戦行うプログラムを作成。先手後手を順番に行う。同様にJAVAを用いて作成。

またVBによるインターフェイスも作成する。

5. 4 評価結果

一万回の対戦を行う評価プログラムの結果、勝率はわずかに55%しか得ることは出来なかった。また最適方策が得られればどちらかが絶対的に有利になると思われるが、今回の実験では先手、後手どちらもほぼ変わらない結果となった。またVBにより作成したゲームにおいては人間に勝つことはほとんどなかった。

ニューラルネットワークの結合荷重 ω をランダムに初期化

エピソードの対して繰り返し:

状態 s を初期化

各ステップに対して繰り返し:

s で可能な行動を A を検索

確率 ϵ でランダムに行動

(i) ϵ

$a \in A$ をランダムに選択する

行動 a をとり次状態 s' を観測する

(ii) $1 - \epsilon$

すべての A の次状態 s' を観測する

s' を入力としてニューラルネットワークによる近似関数から

状態価値 V を観測する

V から行動 $a \in A$ を選択する

行動 a をとり次状態 s' と次状態の価値 $V(s')$ を観測する

BP法により関数を近似する, ω を更新する

(入力を s , 教師信号を $V(s')$ とする)

$s \leftarrow s'$

s が終端状態なら

BP法により関数を近似する, ω を更新する

(入力を s (終端状態), 教師信号を0とする)

終了

図5.3 TDとBP法を用いた学習アルゴリズム

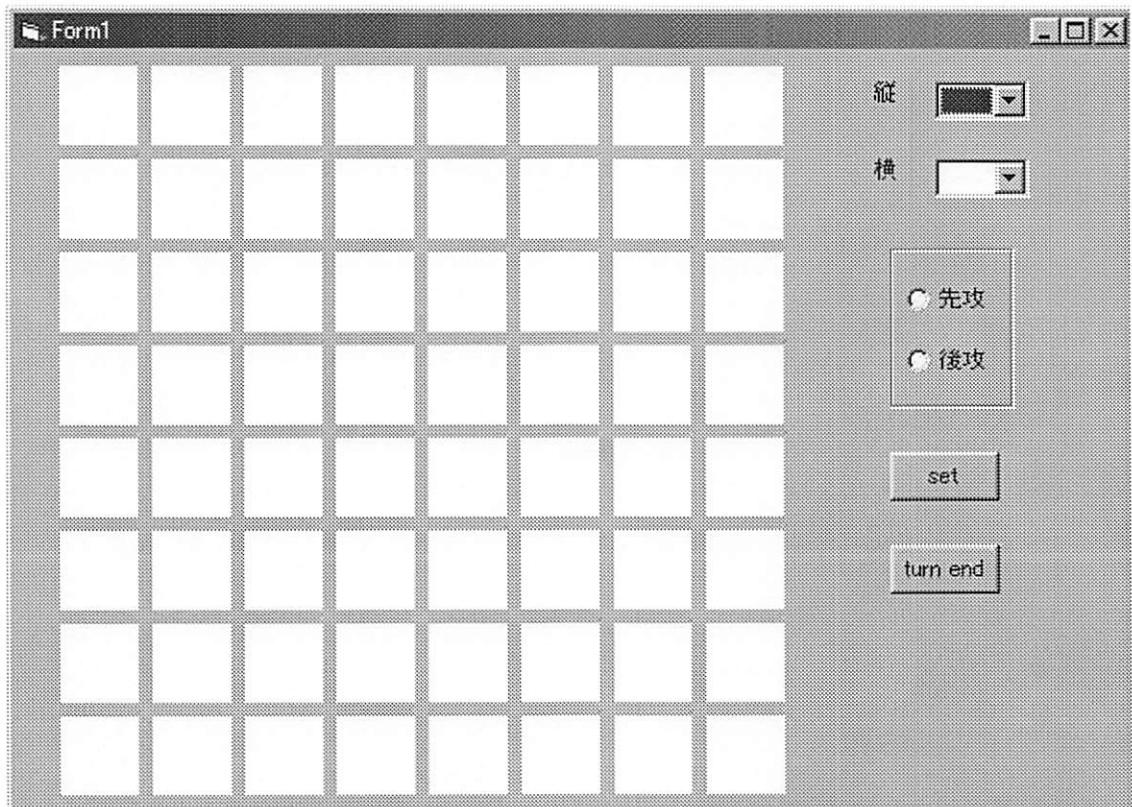


図5.4 VBにより作成した拡張型七五三ゲームのインターフェイス

縦と横のマス of 長さを自由に設定できる, また先攻後攻を選択できる

第6章

考察

実験の結果、最適方策を得るには至らなかった。つまり近似関数を使い最適状態価値関数を得ることが出来なかった。

いくつかの状態においては状態価値が最適な値を指していたが、その他の多くの状態においてはうまく状態価値をいることが出来ていなかった。このためある状態においては最適な行動を選ぶが、その他の状態では全く無意味な行動をしてしまうと言うことがあった。

しかし、いくつかの状態では最適状態価値を得ることが出来たことから手法としての間違いではないと思われる。

原因として考えられるのは以下の2点である。

(1)学習不足

(2)結合強度 ω の初期値によるBPの誤作動

(1)

もともとTD学習も、BP法も幾度もの改善を繰り返し解を得る手法であり、それを併用することによって、単なる2倍以上の時間がかかることになった。

このアルゴリズムを用いた学習だと終端状態に対する近似関数を得る必要がある。そこからTDによって伝播する各ステップ毎に禁じ関数を得なくてはならない。このようなことから膨大な時間がかかることになる。

さらにTDでは全ての状態に無限回訪問するのが理想だが、大きな問題になるにつれそれぞれの状態への訪問に偏りなどがでるため、一つのルートを何度も学習することになる。それで、いくつかの状態では学習が正確に行われたものと思われる。

(2)

第3章でも少し触れた通り、BP法には ω の初期値によっては局所解に陥ってしまう欠点がある。この欠点を補うためにはランダムで ω を与え、これまで学習してきた物との性能の比較を行うなどの対策があった。しかし、これには正確な教師信号が必要となり、教師信号にもそれまで学習から得られた推定値を用いる今回のアルゴリズムには適用できなかった。

現在採るこのできる対策としては ω の初期値をかえて何度も試行することしかない。しかし、この方法だと、一度でさえ時間がかかるので、さらに時間を要するようになり現実的で有効な手法とは言えなくなってしまう。

第7章

結論

本研究で提案したアルゴリズム，作成したプログラムでは膨大な時間がかかりうまくいかなかった．収束に時間のかかる手法を用いたことが原因であり，うまく汎化は行うことが出来なかった．

今後はプログラムの高速化を計り短時間で方策を得ることができる必要になる．また局所解に陥った時の脱出法を考える必要がある．

もしくはニューラルネットワーク以外の関数近似、もしくは別のアプローチによる汎化手法を検討する必要がある．

謝辞

本研究の遂行及び論文の作成に多大な御助言及び御指導を賜った新納 浩幸 教官（茨城大学システム工学科）に深い感謝の意を表します。

また、御指導を頂きましたシステム工学科計算機応用学講座の教官の方々にも深く感謝いたします。

最後に、本研究を進めるにあたり助言，協力を頂きました，同研究室の阿部修也氏（茨城大学大学院理工学研究科システム工学専攻），高橋篤史氏（茨城大学大学院理工学研究科システム工学専攻），進藤邦夫氏（茨城大学システム工学科 4 回生），杉田尚史氏（茨城大学システム工学科 4 回生）に深く感謝いたします。

参考文献

- [1] Ricard S. Sutton and Andrew G. Barto (訳 三上貞芳, 皆川雅章):
‘強化学習’, 森北出版 (2000)
- [2] 馬場規夫, 小島史男, 小澤誠一: ‘ニューラルネットの基礎と応用’,
共立出版 (1994)
- [3] ジョゼフ・オニール: ‘独習 Java’, トップスタジオ訳, 武藤武志監修,
翔泳社 (1999)
- [4] 玉井浩: ‘VisualBasic プログラミング’, サイエンス社 (2000)

付録 A プログラムスリスト

(七五三ゲーム)

動的計画による七五三ゲームの最適方策学習プログラム

```
import t13;
import t14;
//状態と遷移のリスト(詳細省略)

class Keisan{
    int i,j;
    public int s[][]=new int[252][8];
    int k[][]=new int[1500][8];
    double r[][]=new double[252][51];
    double val[]=new double[252];
    double va[][]=new double[252][51];
    double vs[]=new double[252];
```

```
Keisan(){

for(int m1=0;m1<=4;m1++){
for(int m5=0;m5<=10;m5++){
for(int m7=0;m7<=21;m7++){
    i++;

    if(m1==0){k[i][3]++;}
    if(m1==1){k[i][2]++;}
    if(m1==2){k[i][1]++;}
    if(m1==3){k[i][1]+=2;}
    if(m5==0){k[i][5]++;}
    if(m5==1){k[i][4]++;}
```

```
if(m5==2){k[i][3]++;}
if(m5==3){k[i][2]++;}
if(m5==4){k[i][1]++;}
if(m5==5){k[i][3]++;
           k[i][1]++;}
if(m5==6){k[i][2]++;
           k[i][1]++;}
if(m5==7){k[i][1]+=2;}
if(m5==8){k[i][2]+=2;}
if(m5==9){k[i][1]+=3;}
if(m7==0){k[i][7]++;}
if(m7==1){k[i][6]++;}
if(m7==2){k[i][5]++;}
if(m7==3){k[i][4]++;}
if(m7==4){k[i][3]++;}
if(m7==5){k[i][2]++;}
if(m7==6){k[i][1]++;}
if(m7==7){k[i][5]++;
           k[i][1]++;}
if(m7==8){k[i][4]++;
           k[i][1]++;}
if(m7==9){k[i][3]++;
           k[i][1]++;}
if(m7==10){k[i][2]++;
            k[i][1]++;}
if(m7==11){k[i][1]+=2;}
if(m7==12){k[i][4]++;
            k[i][2]++;}
if(m7==13){k[i][3]++;
            k[i][2]++;}
if(m7==14){k[i][2]+=2;}
if(m7==15){k[i][3]+=2;}
if(m7==16){k[i][1]+=4;}
if(m7==17){k[i][1]+=2;
            k[i][3]++;}
if(m7==18){k[i][1]+=2;
```



```
}}}}}
```

```
val[251]=1;
```

```
vs[251]=1;
```

```
t13.def();
```

```
t14.defp();
```

```
//System.out.println("ロード完 "+t13.a[1][1]);
```

```
}
```

```
double ab=0.0;
```

```
double delta=0.0;
```

```
double v=0.0;
```

```
void value(){
```

```
    for(i=1;i<=250;i++){
```

```
        vs[i]=0;
```

```
    }
```

```
    do{
```

```
        delta=0;
```

```
        for(i=1;i<=250;i++){
```

```
            v = val[i];
```

```
            val[i] = atai1(i);
```

```
            ab = v-val[i];
```

```
            if(ab<=0){ab=ab*(-1);} 
```

```
            if(ab>delta){delta=ab;} 
```

```
        }
```

```
        for(i=1;i<=250;i++){
```

```
            vs[i]=val[i];
```

```
        }
```

```
    }while(delta>=0.001);
```

```
}
```

```
double atai1(int n){
```

```
    double gamma=1.0;
```

```

double valu=0.0;
for(int y=1;y<=50;y++){
va[n][y]=r[n][y]-gamma*vs[t13.a[n][y]];
valu=valu+t14.p[n][y]*(r[n][y]-gamma*vs[t13.a[n][y]]);
}

//System.out.println("a "+n);
return(valu);
}

void kaizen(){

for(i=1;i<=250;i++){
double c=0;
double vm=-1;
for(j=1;j<=50;j++){
if(t14.p[i][j]>0){
if(vm<va[i][j]){vm=va[i][j];}
}}

for(j=1;j<=50;j++){
if(t14.p[i][j]>0){
if(va[i][j]==vm){c++;}
}}
for(j=1;j<=50;j++){
if(t14.p[i][j]>0){
if(vm==1){
if(va[i][j]==vm){t14.p[i][j]=1/c;}
else{t14.p[i][j]=0;}}
}

if(va[i][j]==vm){//System.out.println(i+" "+s[i][1]+s[i][2]+s[i][3]+s[i][4]+s[i]
[5]+s[i][6]+s[i][7]+")"+";"+j+"("+t13.a[i][j]+") "+t14.p[i][j]+" "+vm);
System.out.println("s("+i+"="+j);
}
}
}}

```

```
}  
}
```

```
class nana{  
    public static void main(String arg[]){  
  
        Keisan obj =new Keisan();  
        double pai[][]=new double[252][51];  
        int time=0;  
        obj.value();  
        obj.kaizen();  
        for(int i=1;i<=250;i++){  
            for(int j=1;j<=50;j++){  
                if(pai[i][j]==t14.p[i][j]){  
                    }else{time++;  
                        System.out.println(time);  
                        i=0;  
                        j=0;  
                        obj.value();  
                        System.out.println(time);  
                        i=0;  
                        j=0;  
                        for(i=1;i<=250;i++){  
                            for(j=1;j<=50;j++){  
                                pai[i][j]=t14.p[i][j];  
                            }}  
  
                            i=0;  
                            j=0;  
                            obj.kaizen();  
                            System.out.println(time);  
                            i=0;  
                            j=0;  
                            }  
                        }}  
        }}  
    }
```

付録B プログラムリスト

(拡張型七五三ゲーム)

T DとB Pによる拡張型七五三ゲームの最適方策の学習プログラム

```
import java.io.*;
import java.lang.*;
import java.util.Random;

class game{
    /*-----
       ゲームに関する設定
    -----*/
    static int tate=8;
    static int yoko=8;
    static int action;

    static int state[][]=new int[yoko][tate];
    static int next_s[][]=new int[yoko][tate];
    static int koudou[][]=new int[10000][4];

    /*-----
       ニューラルネットワークに関する
    -----*/
    //入力層
    static int n1=tate*yoko;
    //中間層
    static int n2=30;

    //結合強度
    static double w1[][]=new double[n1][n2];
    static double w2[]=new double[n2];
```

```

/*-----
    ゲームに関する設定
-----*/
// ε グリーディ方策の ε
static int ipusiron=1;
static int c[]=new int[n1+1];

//          ×          ソ          ツ          ド
-----

/*-----
    重みを出力する(java)
-----*/

static void output(){
    for(int i=0;i<n1;i++){
        for(int j=0;j<n2;j++){
            System.out.println("w1["+i+"]["+j+"]="+w1[i][j]+");");
        }
    }
    for(int i=0;i<n2;i++){
        System.out.println("w2["+i+"]="+w2[i]+");");
    }
}

/*-----
    重みを出力する (VB)
-----*/

static void outputvb(){
    for(int i=0;i<n1;i++){
        for(int j=0;j<n2;j++){
            System.out.println("w1("+i+", "+j+")="+w1[i][j]);
        }
    }
    for(int i=0;i<n2;i++){

```

```

        System.out.println("w2("+i+")="+w2[i]);
    }
}
/*-----
    エピソードの初期状態にする
-----*/

```

```

static void s0(){
    for(int i=0;i<yoko;i++){
        for(int j=0;j<tate;j++){
            state[i][j]=0;
        }
    }
}

/*-----
    エピソードの終端状態にする
-----*/

```

```

static void s1(){
    for(int i=0;i<yoko;i++){
        for(int j=0;j<tate;j++){
            state[i][j]=1;
        }
    }
}

```

//行動選択に関するメソッド

```

/*-----
    確率  $\epsilon$  でランダムに行動を選択する
-----*/

```

```

static int koudou_1(){
    int i=(int)(Math.random()*100)%10;
    if(i<=ipusiron){
        return(1);
    }
}

```

```

    }else{
        return(0);
    }
}

/*-----
   選択できる行動を選ぶ
-----*/

static void koudou_2(){
    int sita,migi;
    action=0;
    for(int i=0;i<yoko;i++){
        for(int j=0;j<tate;j++){
            if(state[i][j]==0){
                koudou[action][0]=i;
                koudou[action][1]=j;
                koudou[action][2]=0;
                koudou[action][3]=0;
                action++;
                //System.out.println("a"+i+ " "+j);
                migi=1;
                point1: while(i+migi<yoko){
                    if(state[i+migi][j]==0){
                        koudou[action][0]=i;
                        koudou[action][1]=j;
                        koudou[action][2]=1;
                        koudou[action][3]=migi;
                        action++;
                        migi+=1;
                        //System.out.println("yoko"+"migi"+migi);
                    }else{
                        //System.out.println("break migi");
                        break point1;
                    }
                }
            }
        }
    }
}

```

```

        sita=1;
        point2: while(j+sita<tate){
            if(state[i][j+sita]==0){
                koudou[action][0]=i;
                koudou[action][1]=j;
                koudou[action][2]=2;
                koudou[action][3]=sita;
                action++;
                sita+=1;
                //System.out.println("tate"+"sita"+sita);
            }else{
                break point2;
            }
        }
    }
}
}
}
//action--;
//System.out.println("a="+action);
}

```

```

/*-----
    行動 a の選択
-----*/

```

```

static int c_a(){
    double max=0;
    int select=0;
    for(int a=0;a<action;a++){
        for(int i=0;i<yoko;i++){
            for(int j=0;j<tate;j++){
                next_s[i][j]=state[i][j];
            }
        }
        koudou_4(a);
        seiretu_1();
        //System.out.println(a+": "+nyural());
    }
}

```

```

        if(max<nyural()){
            max=nyural();
            select=a;
        }
    }
    return(select);
}

/*-----
    次の状態に移る
-----*/

static void koudou_3(int a){
    if(koudou[a][2]==0){
        state[koudou[a][0]][koudou[a][1]]=1;
    }else{
        if(koudou[a][2]==1){
            for(int i=0;i<=koudou[a][3];i++){
                state[koudou[a][0]+i][koudou[a][1]]=1;
            }
        }else{
            for(int j=0;j<=koudou[a][3];j++){
                state[koudou[a][0]][koudou[a][1]+j]=1;
            }
        }
    }
}

/*-----
    行動後の状態を観測する
-----*/

static void koudou_4(int a){//-----
    if(koudou[a][2]==0){
        next_s[koudou[a][0]][koudou[a][1]]=1;
    }else{

```

```

    if(koudou[a][2]==1){
        for(int i=0;i<=koudou[a][3];i++){
            next_s[koudou[a][0]+i][koudou[a][1]]=1;
        }
    }else{
        for(int j=0;j<=koudou[a][3];j++){
            next_s[koudou[a][0]][koudou[a][1]+j]=1;
        }
    }
}
}
}

```

//ニューラルネットワークに関するメソッド

```

/*-----
    状態をNNで使うために並べ替える

    [8][8]→[64]
-----*/

```

```

static void seiretu(){
    int num=1;
    for(int j=0;j<tate;j++){
        for(int i=0;i<yoko;i++){
            c[num]=state[i][j];
            num++;
        }
    }
}
}

```

```

static void seiretu_1(){
    int num=1;
    for(int j=0;j<tate;j++){
        for(int i=0;i<yoko;i++){
            c[num]=next_s[i][j];
            num++;
        }
    }
}

```

```

    }
}

/*-----
    ニューラルネットワークによる出力
-----*/

static double nyural(){
    int i,j;
    double x,y;
    double x1[]=new double[n2];
    //double y1[]=new double[n3];
    //double z1[]=new double[n4];
    for(i=0;i<n2;i++){
        x = 0;
        for(j=0;j<n1;j++){
            x += c[j+1] * w1[j][i];
        }
        x1[i] = 1.0 / (1.0 + Math.exp(-1.0 * x));
    }
    y = 0;
    for(i=0;i<n2;i++){
        y += x1[i] * w2[i];
    }
    y = 1.0 / (1.0 + Math.exp(-1.0 * y));
    return(y);
}

/*-----
    wの初期化値ランダムをランダムに与える
-----*/

static void syokika(){//-----
    for(int i=0;i<n1;i++){
        for(int j=0;j<n2;j++){
            int a=(int)(Math.random()*10);

```

```

        int zoku=(int)(Math.random()*10)%2;
        double b=Math.random();
        if(zoku==0){
            w1[i][j]=b+a;
        }else{
            w1[i][j]=-(b+a);
        }
    }
}
for(int i=0;i<n2;i++){
    int a=(int)(Math.random()*10);
    int zoku=(int)(Math.random()*10)%2;
    double b=Math.random();
    if(zoku==0){
        w2[i]=b+a;
    }else{
        w2[i]=-(b+a);
    }
}
}

/*-----
   -終端状態かの確認
   -----*/

static int end(){
    for(int j=0;j<tate;j++){
        for(int i=0;i<yoko;i++){
            if(state[i][j] != 1)return(0);
        }
    }
    //System.out.println("end");
    return(1);
}

/*-----

```

現在の状態を出力

-----*/

```
static void s1(){
    for(int j=0;j<tate;j++){
        for(int i=0;i<yoko;i++){
            System.out.print(state[i][j]+" ");
        }System.out.println();
    }
}
```

/*-----

メイン

-----*/

```
public static void main(String arg[]){
```

```
    int act=0;
```

```
    int i,j,k,l;
```

```
    //δの値
```

```
        double delta1;
```

```
        double delta2[]=new double[n2];
```

```
    //出力
```

```
        double ans;
```

```
    //誤差
```

```
    double err;
```

```
    double gosa=0.0;
```

```
    //それぞれのニューロンの値
```

```
    double x,y;
```

```
        double x1[]=new double[n2];
```

```
    //ステップ幅
```

```
        double a = 0.13;
```

```

syokika());

for(l=0;l<1000000;l++){//-----エピソードの開始
  s0();
  do{

    seiretu();
    koudou_2();//-----行動の検索

    if(koudou_1()==1){//-----方策に従わないランダムな行動
      act=(int)(Math.random()*100)%action;
      koudou_3(act);
    }else{//-----方策に従った行動
      act=c_a();
      for(i=0;i<yoko;i++){
        for(j=0;j<tate;j++){
          next_s[i][j]=state[i][j];
        }
      }
      koudou_4(act);
      seiretu_1();
      gosa=(1.0-nyural());//-----教師信号の決定
      seiretu();//-----現在の入力のセット

//-----学習開始

      for(i=0;i<n2;i++){
        x = 0;
        for(j=0;j<n1;j++){
          x += c[j+1] * w1[j][i];
        }
        x1[i] = 1.0 / (1.0 + Math.exp(-1.0 * x));
      }
      y = 0;

```

```

        for(i=0;i<n2;i++){
            y += x1[i] * w2[i];
        }

ans = 1.0 / (1.0 + Math.exp(-1.0 * y));
err = gosa - ans;

delta1 = ans * (1 - ans) * err;
    for(i=0;i<n2;i++){
        w2[i] = w2[i] + a * delta1 * x1[i];
    }
for(i=0;i<n2;i++){
    delta2[i] = x1[i] * (1 - x1[i]) * delta1 * w2[i];
    for(j=0;j<n1;j++){
        w1[j][i] = w1[j][i] + a * delta2[i] * c[j+1];
    }
}
koudou_3(act);//-----1ステップの終了
}
}while(end(<1));//-----1エピソード終了

s1();//-----端末状態を教師信号0
で学習
seiretu();
    for(i=0;i<n2;i++){
        x = 0;
        for(j=0;j<n1;j++){
            x += c[j+1] * w1[j][i];
        }
        x1[i] = 1.0 / (1.0 + Math.exp(-1.0 * x));
    }

    y = 0;
    for(i=0;i<n2;i++){
        y += x1[i] * w2[i];
    }

```

```

ans = 1.0 / (1.0 + Math.exp(-1.0 * y));
err = 0.0 - ans;

delta1 = ans * (1 - ans) * err;
    for(i=0;i<n2;i++){
        w2[i] = w2[i] + a * delta1 * x1[i];
    }
for(i=0;i<n2;i++){
    delta2[i] = x1[i] * (1 - x1[i]) * delta1 * w2[i];
    for(j=0;j<n1;j++){
        w1[j][i] = w1[j][i] + a * delta2[i] * c[j+1];
    }
}
if(end()==1){
    if(1%1000==0){
        System.out.println("learn"+l);
        if(1%100000==0){
            outputvb();
            output();
        }
    }
}
} //-----エピソード終了
}
}
}

```