

Title	プレイヤーの意図や価値観を学習し行動選択するチーム プレイAIの構成
Author(s)	吉谷, 慧
Citation	
Issue Date	2013-03
Type	Thesis or Dissertation
Text version	author
URL	<a href="http://hdl.handle.net/10119/11300">http://hdl.handle.net/10119/11300</a>
Rights	
Description	Supervisor:池田 心, 情報科学研究科, 修士

修 士 論 文

プレイヤーの意図や価値観を学習し行動選択する  
チームプレイ AIの構成

北陸先端科学技術大学院大学  
情報科学研究科情報科学専攻

吉谷 慧

2013年3月

修士論文

プレイヤーの意図や価値観を学習し行動選択する  
チームプレイ AI の構成

指導教官 池田 心 准教授

審査委員主査 池田 心 准教授  
審査委員 飯田 弘之 教授  
審査委員 東条 敏 教授

北陸先端科学技術大学院大学  
情報科学研究科情報科学専攻

1110067 吉谷 慧

提出年月: 2013年2月

## 概要

近年までのゲーム AI 研究は，“強い” AI を作るという目的のものが多く，そのための様々な手法が提案されてきた．チェスや将棋などのボードゲームにおいてゲーム AI の強さは既に人間のプロレベルに達しており，1997 年にはチェスで Deep Blue が世界チャンピオンのカスパロフ氏に勝利している．ゲーム AI 研究の対象はボードゲームだけに留まらず，現在では一般的になってきているコンピュータゲームも，研究の対象となっている．コンピュータゲームを対象としている研究では，FPS ゲームや RTS ゲームなどのゲーム AI が題材として扱われており，今までに多くの“強い”ゲーム AI を作る方法が提案されてきた．これらは人間プレイヤーの“敵”として，人間プレイヤーを楽しませるための研究と言える．しかし一方で，人間プレイヤーの“仲間”として，“強さ以外”の要素に着目して人間プレイヤーを楽しませるためのゲーム AI 研究の例は少ない．

近年市販されているコンピュータゲームでは，ゲーム AI が操作するキャラクターとチームを組んで遊べるものも多いが，しばしば AI プレイヤーは人間プレイヤーが望まない行動を取ることがあり，これが人間プレイヤーにとってストレスとなる場合がある．そこで本研究では，「人間と一緒にプレイして“楽しい”“仲間”」の AI プレイヤーを生成することを目的とし，この問題の解決にあたった．仲間のゲーム AI が前述のような問題行動を取ってしまう原因を「人間の価値観（効用）を理解できていないため」と推定し，ゲーム AI に人間の効用を表現する何らかのモデルを与え，人間プレイヤーの効用を学習させることで，この問題の解決を図った．

本研究では，「人間はゲーム中の状態の価値を“状態を構成する要素”から総合的に判断してしており，人間プレイヤー自身（もしくは仲間のゲーム AI プレイヤー）の行動によってその価値がどう変化するかを予想して行動を決定している」と仮定した．この仮定のもと，“状態を構成する要素”を入力・価値を出力とする入出力モデルを効用のモデルとして設定し，これを調整することで人間プレイヤーの効用を表現した．モデルの調整には，人間プレイヤーから AI プレイヤーへ出された指示や，人間プレイヤー自身の行動を使い，効用モデルを近似させていく．本研究では，特に日本での人気の高いコマンド形式 RPG ゲームを独自に設定し，提案手法の性能を評価し有望であることを示した．

# 目次

第1章	はじめに	1
第2章	本研究の目的	2
第3章	関連研究	3
第4章	人間の思考について	5
4.1	人間の状態評価に関する仮説	5
4.2	人間の行動決定に関する仮説	5
4.3	人間のストレスとなりうる AI プレイヤの行動選択	6
第5章	本研究で扱うゲームの概要	8
5.1	キャラクタの持つパラメータ	8
5.2	行動とその効果	8
5.3	状態遷移	9
5.4	遷移例	9
5.5	ゲームの要点	10
第6章	提案手法	11
6.1	効用モデルと期待効用値	11
6.2	学習部分	11
6.3	行動決定部分	13
第7章	実験	14
7.1	戦闘参加キャラクタの設定	14
7.2	提案手法を実装したチームメイト AI の詳細	15
7.2.1	効用モデル	15
7.2.2	学習部分	16
7.2.3	行動選択部分	20
7.3	その他の AI の概要	20
7.4	実験の設定	21
7.5	実験結果	21
7.5.1	効用モデルの結合強度による評価	21

7.5.2	「人間が不快に感じ得る行動選択」の出現頻度による評価 . . . . .	22
<b>第 8 章</b>	<b>まとめ</b>	<b>25</b>
<b>第 9 章</b>	<b>付録</b>	<b>26</b>
9.1	実験 1 . . . . .	26
9.2	実験 2 . . . . .	27

# 第1章 はじめに

本論文では、ゲームにおいて敵や味方などのキャラクタを知的にふるまわせる技術、およびプログラムのことをAI、ゲームAIと呼ぶことにする。ゲームAI研究の目標のひとつとして「人間を楽しませるゲームAI」を作ることが挙げられる。近年までのゲームAIの研究は“強さ”に焦点をあてたものが多く、既にオセロやチェス、将棋といったボードゲームにおいて、ゲームAIの強さは人間のプロレベルに達している。チームとしての“強さ”に焦点をあてた研究としては、Sander,B.らの研究[1]で、QuakeIIIのゲームAIを動的に敵チームに適応させることで強いチームを生成する手法が提案されている。これら、ゲームAIの“強さ”に焦点をあてた研究は、人間プレイヤーの「敵」として人間を楽しませるための研究といえる。

Aswin,T.A.らの研究[2]では変化していくゲームの難易度に合わせて、チームメイトであるNon Player Characterのサポート行動をプレイヤーに適応させる試みがなされている。また、Kenneth O.S.らの研究[3]では、ゲームプレイ中のようなリアルタイムでも、プレイヤーとのやりとりを通じてゲームキャラクタが弱点を克服することで、ゲームの面白さを保つ試みがなされており、人間がゲームAIを訓練するという新たなゲームジャンルを提唱している。これらは人間の「仲間」として、もしくは“強さ”だけでなく別の要素にも着目し、人間を楽しませるための研究といえる。

しかし、「仲間」として人間を楽しませるAIの研究は「敵」としての研究と比べるとまだ少なく、また“強さ”以外の要素に焦点をあてた研究も少ない。特に日本での人気の高いコマンド形式RPG（ドラゴンクエスト等）ではこのような技術の需要は大きい。

## 第2章 本研究の目的

本研究の目的は「人間と一緒にプレイして“楽しい（ストレスや不満が少ない）”，“仲間”」のAIプレイヤを生成することである。近年市販されているコンピュータゲームではコンピュータ AI が操作するキャラクタとチームを組んで遊べるものも多いが、しばしばAIプレイヤは人間プレイヤが意図しない（望んでいない）行動を取ることがあり、これが人間プレイヤにとってストレスとなる場合がある。

人間プレイヤがAIプレイヤに望む行動と、AIプレイヤが実際に選択する行動が異なってしまう原因を本論文では以下4つに分類する。

1. 人間プレイヤとAIプレイヤの効用（何を目標して行動しているか、またはその重み）が異なるため、それぞれにとって最適な行動が異なるから
2. それぞれの効用が同じだとしても、人間プレイヤの情報処理能力が不十分であり最適な行動が分かっておらず、“間違っただけ”不満を持っているため
3. それぞれの効用が同じだとしても、AIプレイヤの探索などが不十分で最適な行動が求まっていないため
4. 本質的に混合戦略が必要なゲームであり、「一定確率で期待と違ってもしょうがない」ため

本研究では1.の理由に注目し、AIプレイヤに人間プレイヤの効用を学習させることで、行動選択の差異による人間プレイヤへのストレスを軽減したい。

## 第3章 関連研究

協力ゲーム・非協力ゲームにおける戦略はゲーム理論の枠組で論じられることが多い。コマンド形式RPGの多くは同時着手性をもち、仲間との間には利得を共有する協力ゲームの関係がある。

ただし、こういったゲームの利得は確定的なものではなくランダム性に影響され、リスクをどの程度重く見るかなどを扱うためには効用の考え方を導入する必要がある。

- 効用

効用とは個々人が持つ、物事への価値観のことである [4]。人間の価値観を定量化することを目的とした効用理論の歴史は、D. ベルヌーイにさかのぼる。人々の行動原理として、ベルヌーイは「人々は期待値を最大化して行動しているのではなく、期待効用を最大化している」と説明し、お金に対する効用関数として対数関数を提案した。これが効用関数の発祥である。図 3.1 は単属性効用関数の一例である。効用を表す関数として、他にも様々な関数が提案されている。

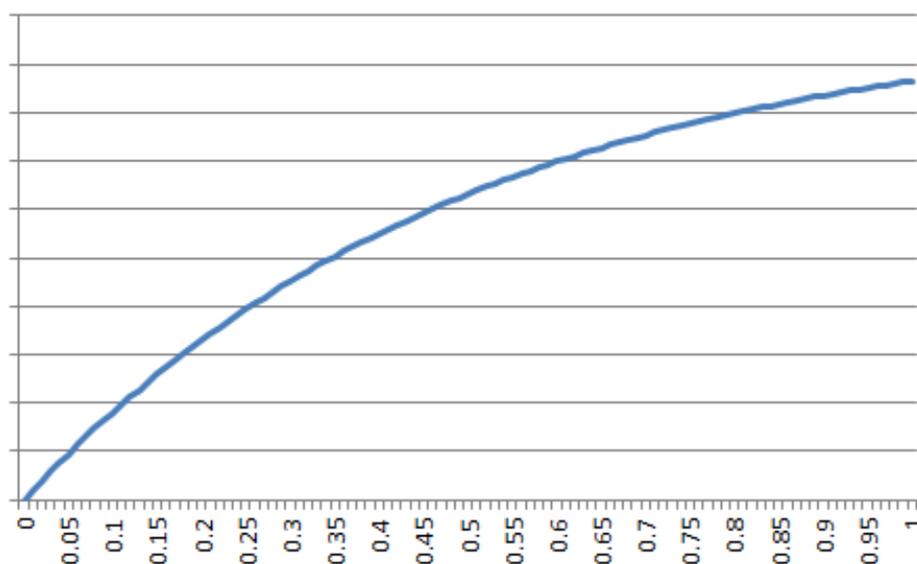


図 3.1: 効用関数の一例

- プロスペクト理論

Kahneman は、期待効用に基づく選択の矛盾を数多くの実験結果によって指摘したうえで、期待効用値理論に代わる理論として、プロスペクト理論を提案した [5]. 人は、ほぼ確実な収益に対しては、実際の確率よりも低い確率を考えがちであり、これを確実性効果と呼んでいる. 人は収益の領域ではリスク回避的な確実性効果が強く表れるが、損失の領域では不確実なものを好むようになる. 簡単な例を次に示す.

「つぎの二つのうちのどちらを選ぶだろうか？」

1. A : 80% の確率で 4,000 円,     B : 100% の確率で 3,000 円
2. C : 20% の確率で 4,000 円,     D : 25% の確率で 3,000 円
3. E : 45% の確率で 6,000 円,     F : 90% の確率で 3,000 円
4. G : 0.1% の確率で 6,000 円,     H : 0.2% の確率で 3,000 円

1. の質問に対して多くの人が B と答え、2. の質問に対しては C と答える. この違いは統計的にも 1% の水準で有意であった. また同じように 3. と 4. の質問を比較してみると、3. では F が、4. では G が多くの人に選ばれた. 1. と 2. は、確率が 100% から 80% に落ちる方が、25% から 20% に落ちるより満足度に対する影響が大きいことを示しており、人間はほぼ確実な収益に対しては、実際の確率よりも低い確率を考えがちであるとみることができる. また 3. と 4. からは、確率が極端に低い収益に関して、その生起確率を実際よりも高い確率として考えがちであることがわかる.

## 第4章 人間の思考について

### 4.1 人間の状態評価に関する仮説

人間は“状態を構成する要素”それぞれに対して好ましさを持っていて、各要素を総合的に評価した結果がその状態の良さとなる、と仮定する。

“状態を構成する要素”とは、人間プレイヤーや仲間のAIプレイヤー、もしくは敵対プレイヤーなどの行動選択によって影響されるゲーム中のパラメータのことである。ここで“状態を構成する要素”について、一般的なターン制のRPGゲームにおける戦闘シーンを例に説明する。この種のゲームでは戦闘に参加している敵・味方キャラクターが持つパラメータとしてHP（キャラクターの体力のこと。この数値が0になると行動できなくなる）やMP（キャラクターが一部の特殊な行動を取るために必要となる数値）などが設定されている場合が多く、これらが“状態を構成する要素”となりうる。また他にも、戦闘に参加しているキャラクターそれぞれの生死や、戦闘の経過時間（経過ターン数）なども、“状態を構成する要素”となりうる。

ひとつ気を付けなければならないのは、人間の持つ効用は個人個人で違うという点である。同じ状態を評価したとしても個人個人の評価は違ってくる場合がある。例えば図4.1のような状態を複数の人が評価した場合、ある人は2人共HPが半分以下という理由で状態を「悪い」と判断するかもしれないが、別の人はMPが残っているので「普通」と判断する可能性がある。

### 4.2 人間の行動決定に関する仮説

人間は複数の行動選択肢がある場合、それぞれの選択肢について行動した後の“未来の状態”を予想・評価して行動選択をしている。

人間が予想する“未来の状態”がどの程度先の状態かは、個人差やゲームの種類によって様々である。ターン制のゲームならば、1ターン先の状態や数ターン先の状態、もしくは戦闘終了時の状態などが、人間が予想する“未来の状態”の候補として挙げられる。 $\alpha\beta$ 法は前者の、モンテカルロ方は後者の代表的な探索方法である。

図4.2は二人同時着手確定ゲームの状態遷移図と、遷移後の各状態に対する人間の評価を表にしたものである。人間プレイヤーの行動選択肢が $a_1$ と $a_2$ 、AIプレイヤーの行動選択肢が $b_1$ と $b_2$ である。なお、評価表はあくまで人間プレイヤーが自身の持つ効用に従って得た表であり、AIプレイヤーにとっての評価表が同じとは限らない。状態は両プレイヤーの選択

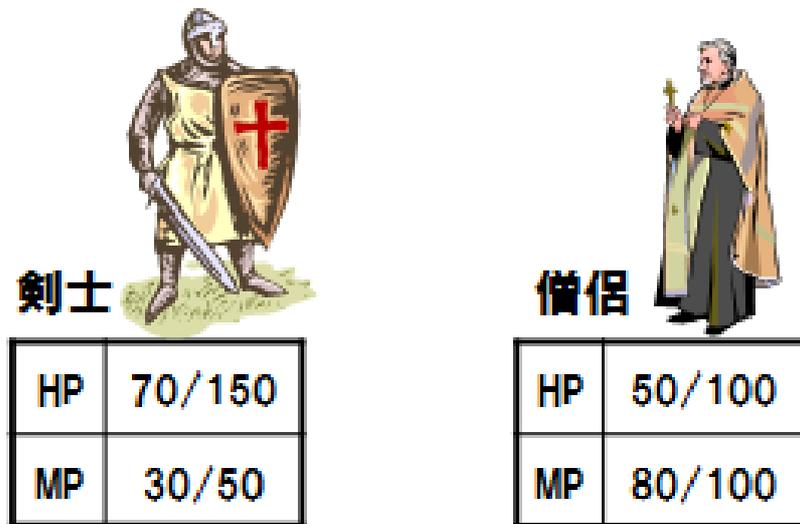


図 4.1: 一般的な RPG ゲームにおける状態の一例

の組合せによって確率的に遷移する。

人間プレイヤーが図 4.2 のように評価した場合、実際に選択する行動は AI プレイヤの立場や個人差で違って来るであろうが (AI プレイヤが仲間の場合、AI プレイヤが賢いと思っている人は  $a_1$  を選ぶかもしれないし、AI プレイヤを信頼していない人はリスク回避的に  $a_2$  を選ぶかもしれない) いずれにせよ “未来の状態” を予想・評価して行動選択をしている。

### 4.3 人間のストレスとなりうる AI プレイヤの行動選択

人間が大きなストレスを感じるパターンのひとつとして、「合理的と思えない」「選択理由が分からない」行動を仲間の AI プレイヤが選択したときが挙げられる。AI プレイヤの選択によって評価の低い状態に遷移した場合も、確かに人間プレイヤーはストレスを感じるであろうが、それ以上に「選択理由が分からない」行動を AI プレイヤが選択した場合の方が大きなストレスを与える可能性がある。表 4.1 は、ある状態で人間プレイヤーが “未来の状態” を評価し、その良さを数値化して表にしたものである。人間プレイヤーの行動選択肢が  $a_1 \sim a_3$ 、AI プレイヤの行動選択肢が  $b_1 \sim b_3$  である。

人間プレイヤーの評価では  $(a_1, b_3)$  の組合せの評価が一番高いので、AI プレイヤが  $b_3$  の行動を選択する理由を人間プレイヤーは理解できる。また、 $(a_2, b_2)$  の組合せも  $(a_1, b_3)$  程ではないが評価は高いので、AI プレイヤが  $b_2$  を選んだ場合も人間プレイヤーはその選択に理解を示せるはずである (2 章 4. を参照)。人間プレイヤーの選択次第では  $(a_2, b_3)$  や  $(a_1, b_2)$  といった評価の低い状態に遷移する可能性もあるが、AI プレイヤの選択に理由付けができるので、ある程度ストレスが緩和されることが期待できる。これらは、効用は個々で異

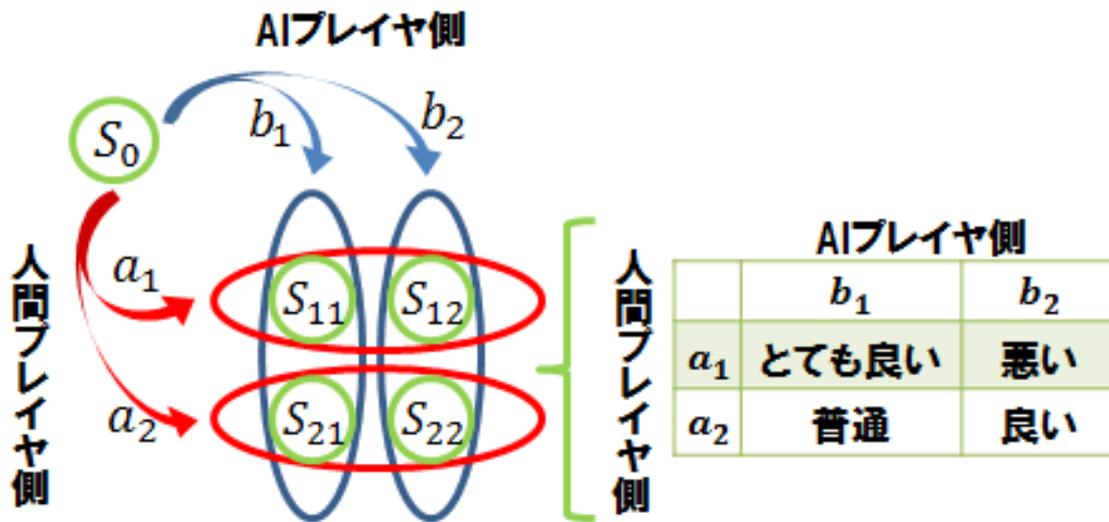


図 4.2: 同時着手確定ゲームの状態遷移図と人間の評価例

なるものであり自分と相手の効用が違う可能性や、未来の状態の予想に伴うノイズの存在が影響している。

ところが一方で、AIプレイヤーが  $b_1$  の行動を選択した場合、人間プレイヤーにとって大きなストレスとなる可能性がある。これは、 $b_1$  の選択肢は  $b_2$  や  $b_3$  の選択肢に対して優位性が少なく選択する理由として弱いことや、人間プレイヤーが  $a_1$  の選択肢を選んだ場合に唯一評価が勝ることになる  $(a_3, b_1)$  の組合せも、人間プレイヤーにとって  $a_3$  は  $a_2$  の完全な被支配戦略であり、そもそも  $a_3$  の選択肢を考慮に入れること自体がありえないと思えるからである。これらは、効用の差異やノイズの影響を鑑みても「ありえないだろう」と人間プレイヤーが思ってしまう場合である。

このように納得できる評価の低い状態に遷移したうえ、AIプレイヤーの行動選択に理由が見つからない場合、人間プレイヤーは大きなストレスを感じる可能性が高い（2章1.を参照）。

人間プレイヤーが不快に思わない行動をAIプレイヤーが選択するためには、人間の効用を理解し、明らかにその効用に反するような選択肢を選ばないことが大切である。

表 4.1: 人間による“未来の状態”の評価表の一例

	$b_1$	$b_2$	$b_3$
$a_1$	0.2	0.5	0.8
$a_2$	0.4	0.7	0.4
$a_3$	0.3	0.2	0.2

## 第5章 本研究で扱うゲームの概要

本研究では独自にコマンド形式のターン制ゲームを設定し，提案手法の実装や実験を行う．このゲームは Markov Decision Process で記述できる多人数同時着手不確定ゲームである．ゲームの各種設定は以降で説明する．

### 5.1 キャラクタの持つパラメータ

ゲームに参加するキャラクターには以下の通り 4 つのパラメータを設定する．キャラクターを操作するプレイヤーは，自身が操作しているキャラクターのパラメータは勿論，敵・味方全てのキャラクターのパラメータを知ることができる．

- 体力 (HP)  
体力が 0 になると，そのキャラクターは行動不能となる．チーム全員の体力が 0 になった場合，そのチームは戦闘に敗北する．可変値．
- 精神力 (MP)  
一部の術技を選択するために必要な数値で，対象となる術技を使用すると値が消費される．可変値．
- 攻撃力  
攻撃をした際，相手に与えるダメージの指標．固定値．
- 使用可能な術技  
そのキャラクターが選択できる行動の種類．固定．

### 5.2 行動とその効果

各キャラクターの行動は，術技（行動の種類）と対象キャラクターの組合せで表現される．術技それぞれの効果は以下の通りである．

- 攻撃  
対象キャラクターに攻撃力と同じ値のダメージを与える。
- 回復  
精神力を一定量消費して、対象キャラクターの体力を一定量回復する。この術技を使うキャラクターに消費される分の精神力が残っていない場合は選択できない。

よって、取りうる行動の数は生存する敵と味方の数と等しい。

## 5.3 状態遷移

このゲームでは複数のキャラクターが敵チーム・味方チームに分かれて、チーム同士で戦闘を行う。戦闘は戦闘終了条件を満たすまでターンを繰り返す。1つのターンは以下3つのステップに分けられる。

1. 行動選択ステップ  
行動可能な全キャラクターが現在の状態に基づき行動を選択する。全員が行動を選択したら次のステップに進む。この際、パーティ内で相談はできない（同時着手性 []）
2. 行動処理ステップ  
ランダムな順番で、各キャラクターが選択した行動を順次処理していく。処理順の関係で、あるキャラクターの行動が処理される前にその行動主体のキャラクター、もしくは行動対象のキャラクターが行動不能になってしまった場合、そのキャラクターの行動は処理されない。全員分の行動を処理したら次のステップに進む。
3. 戦闘終了判断ステップ  
どちらかのチーム全員が行動不能になった場合に戦闘を終了する。両チームに行動可能なキャラクターが1人でもいる場合は1. のステップへ戻る。

## 5.4 遷移例

図 5.1 はこのゲームにおける状態遷移の一例で、左側が遷移前の状態、右側が遷移後の状態を表している。表 5.1 は図 5.1 における、状態遷移によるパラメータの変化と、行動選択ステップで各キャラクターが選択した行動の表である。

行動処理ステップでの処理順番が、敵・僧侶・剣士の順番であった場合、最初に敵の攻撃で僧侶の HP は 0 となり行動不能となる。次に僧侶の行動順が来るが、先の攻撃で行動不能となっているので、僧侶の行動は処理されない。最後に剣士の行動が処理され敵にダメージを与えて、図 5.1 の右側の状態に遷移する。

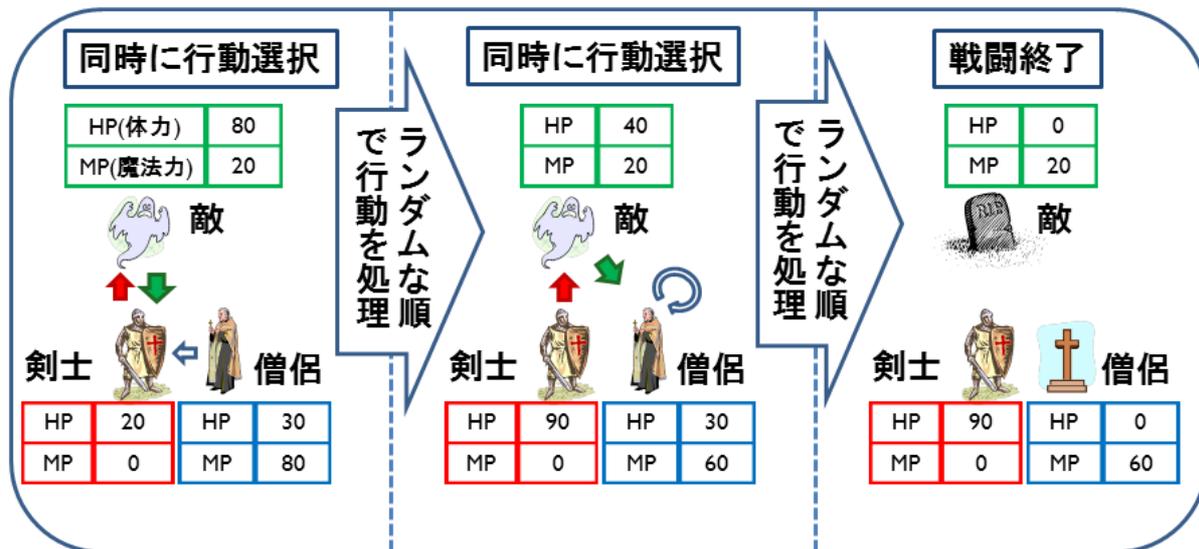


図 5.1: ゲームにおける状態遷移の例

表 5.1: 図 5.1 の状態での各キャラクターのパラメータなど

	剣士	僧侶	敵
HP	50 → 50	20 → 0	60 → 20
MP	50 → 50	80 → 80	10 → 10
攻撃力	40	20	20
選択行動	攻撃	回復	攻撃

## 5.5 ゲームの要点

今回のゲームは、市販されている RPG を意識し設定を行った。一般的な RPG では、様々な敵が出現し、それらと繰り返し何度も戦闘を行う場合が多い。このことから、一回一回の戦闘ではただ単に勝利すれば良いというわけではなく、どう勝利するかも重要となってくる。

## 第6章 提案手法

本研究では、人間プレイヤーの選択した行動からそのプレイヤーの効用を推定し、自身の行動選択に活用するチームメイトAIを構成する手法を提案する。

提案するチームメイトAIの全体像は図6.1の通りである。このチームメイトAIは人間プレイヤーの効用を表現するための効用モデルを持っている。学習部分では人間プレイヤーの効用を近似するように効用モデルを調整し、行動選択部分では効用モデルを使い自身の行動を決定する。

### 6.1 効用モデルと期待効用値

効用モデルには、ニューラルネットワークなどに代表される入出力モデルを使用する。この入出力モデルは4章で説明した“状態を構成する要素”を入力とし、出力が状態の効用値となる。状態  $s$  で行動  $a$  を取った場合の期待効用値  $Eu$  は6.1式で求められる。

$$Eu(s, u, a) = \sum_{s'} \{u(s')Tr(s, s', a)\} \quad (6.1)$$

ここで、 $u$  は状態から効用値を算出する効用関数（効用モデル）、 $Tr(s, s', a)$  は状態  $s$  で行動  $a$  を取ったとき“未来の状態”  $s'$  に遷移する確率である。5章で説明したゲームのように、同時着手性があるゲームの場合は6.1式に他者の行動  $b$  を加えた6.2式で期待効用値を計算することができる。

$$Eu(s, u, a, b) = \sum_{s'} \{u(s')Tr(s, s', a, b)\} \quad (6.2)$$

### 6.2 学習部分

学習段階では人間プレイヤーの効用を近似するために効用モデルを調整する。学習を行うため事前に学習用のデータとして、効用を近似させたい人間プレイヤーの行動データ  $d_i = (s_i, a_i)$  を複数用意し、これを学習用データの集合  $D(\ni d_i)$  とする（図6.1の①）。 $d_i$  は「人間プレイヤーが状態  $s_i$  で行動  $a_i$  を選択した」という情報である。学習段階は次の4つのステップから構成される。

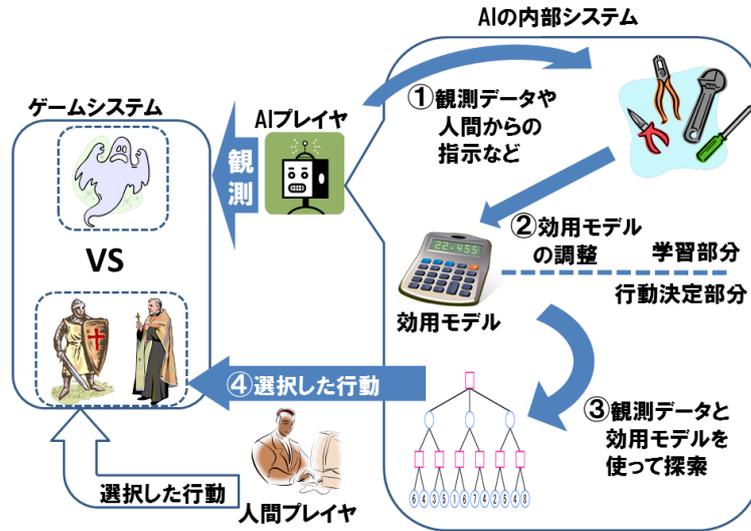


図 6.1: 提案するチームメイト AI の全体像

### 1. 期待効用値の計算

状態  $s$  で人間プレイヤーが選択できる行動の集合を  $A$  とし、任意の要素を  $a_n \in A$  とする。同時着手ゲームの場合はさらに、他者が選択できる行動の集合を  $B$  とし、任意の要素を  $b_m \in B$  とする。

同時着手ゲームの場合は  $(a_n, b_m)$  の組合せ全てについて 6.2 式で期待効用値を計算する。この様にして得られた表を期待効用値表  $EuTable$  とし、6.3 式のように表す。

$$EuTable(u, s) = \{Eu(s, u, a_n, b_m)\}_{n,m} \quad (6.3)$$

表 4.1 は期待効用値表  $EuTable(u, s)$  の一例であり、図 4.2 が  $EuTable(u, s)$  を計算する際のイメージとなっている。用意した学習用データに含まれる  $s_i$  全てについて  $EuTable(u, s_i)$  の計算を行い、 $EuTable(u, s_i)$  の集合である

$$EuTables(u, D) = \{EuTable(u, s_i)\}_i \quad (6.4)$$

を得る。

### 2. 選択確率の計算

前ステップで計算した期待効用値表  $EuTable(u, s_i)$  から人間プレイヤーが状態  $s_i$  で行動  $a_i$  を選択する確率  $p(s_i, u, a_i, EuTable)$  を計算する。  $p(s_i, u, a_i, EuTable)$  の計算方法として様々な方法があるが、気を付けなければならないのは、計算方法を「人間らしい」方法にすることである。この算出方法が「人間らしくない」場合、そのチームメイト AI の性能は低いものになってしまう。本論文での  $p$  の計算の詳細は 7.2.2 の 2. で示す。

### 3. 平均対数尤度の計算

前ステップで計算した選択確率から、効用モデル  $u$  の平均対数尤度を計算する。計算式は 6.5 式の通りである。

$$\logLikelihood(u, D) = \frac{\sum_i \log\{\epsilon + p(s_i, u, a_i, EuTable)\}}{i} \quad (6.5)$$

ここで、 $\epsilon (> 0)$  は  $p(s_i, u, a_i, EuTable)$  が 0 になってしまう場合を考慮した調整用パラメータである。

### 4. 効用モデルの調整

効用モデルの調整方法としては、ローカルサーチや遺伝的アルゴリズムといった生成検査法、バックプロパゲーションなどの勾配法が挙げられる。2つの効用モデル  $u$  と  $u'$  について比べたとき、 $\logLikelihood(u', D) > \logLikelihood(u, D)$  ならば、 $u'$  の方が学習用データの集合  $D$  を  $u$  より尤もらしく説明している。つまり、よりサンプルとなった人間の効用に近いということである。このステップでは、対数尤度を最大化するように効用モデルを最適化していく（図 6.1 の ②）。

このように、平均対数尤度を用いてデータの再現度を表し、モデルパラメータを最適化する手法はしばしば用いられる。その典型は Remi らによる BTMM 法 [1]、保木らによるボナンザ法 [2] である。

## 6.3 行動決定部分

行動決定部分では、自身が持つ効用モデル  $u_t$  を使用して行動選択を行う（図 6.1 の ③と ④）。ここではまず現在の状態  $s_t$  から期待効用値表  $EuTable(u_t, s_t)$  を計算し、人間プレイヤーがストレスを感じないように行動を選択しなければならない。

$EuTable(u_t, s_t)$  から自身の行動を決定する方法として様々な方法があるが、4章で説明した「人間のストレスとなりうる AI プレイヤの行動選択」については気を付けなければならない。詳細は 7.2.3 節で述べる。

## 第7章 実験

提案手法のチームメイト AI を実装し、仲間プレイヤーの効用を近似する実験を行った。効用モデルを近似させる対象として、本論文では人間ではなく、戦闘終了時の状態に対して特定の、典型的な効用を持つエージェントモンテカルロ AI を用意した。これは異なる効用を持つ人間プレイヤーそれぞれに対して提案手法が有用であるかを確かめるためである。

提案手法の AI にモンテカルロ AI の効用を学習させるため、提案手法の AI とモンテカルロ AI とでチームを組ませて 100 戦分の戦闘を行い、調整された効用モデルの結合強度のパラメータを評価した。

また、「人間が不快に感じ得る行動選択のパターン」を設定し、その出現頻度から提案手法の AI を評価した。モンテカルロ AI は自身と味方の行動の組合せについてモンテカルロシミュレーションを行い、自身の行動を決定する AI である。つまり、モンテカルロ AI は自身と味方の行動についての期待効用値表を作成している。この期待効用値表から、人間が納得できないであろう仲間 AI の行動選択を定義した。

### 7.1 戦闘参加キャラクターの設定

今回、実験を行うために設定した敵・味方キャラクターのパラメータは表 7.1 と表 7.2 の通りである。敵キャラクターは環境の一部に過ぎず、着目したいのは味方キャラクター 1 がキャラクター 2 の効用を学習して協調できるか、という点である。

表 7.1: 味方チームに属するキャラクターのパラメータ

	味方キャラクター 1	味方キャラクター 2
最大 HP	150	100
最大 MP	50	100
攻撃力	40	20
使用可能術技	攻撃・回復	攻撃・回復
操作 AI	提案手法 AI	モンテカルロ AI

表 7.2: 敵チームに属するキャラクターのパラメータ

	敵キャラクター 1	敵キャラクター 2	敵キャラクター 3
最大 HP	60	60	80
最大 MP	0	0	60
攻撃力	40	40	10
使用可能術技	攻撃	攻撃	攻撃・回復
操作 AI	ランダム AI	ヒューリスティック AI	原始モンテカルロ AI

## 7.2 提案手法を実装したチームメイト AI の詳細

### 7.2.1 効用モデル

6章で説明した提案手法に則り、チームメイト AI を実装した。今回実装したチームメイト AI は“未来の状態”として「戦闘終了時の状態」を評価し、状態の良さを効用値として数値化する AI とした。“状態を構成する要素”として以下 4 つの要素を使用した。

1. 戦闘の勝敗

自チームが戦闘に勝利した場合は 1 を、戦闘に敗北した場合は 0 を入力値とする。

2. 自チームキャラクターの生死

戦闘終了時、自分を含めた自チームのキャラクターの生死によって入力値が変わる。今回の実験の設定では自チームは 2 人なので、2 人共生きている場合は 1 を、1 人だけ生きている場合は 0.5 を、2 人共行動不能な場合は 0 を入力値とする。

3. チームの残り MP

自チーム全体でみた MP の残りパーセンテージを入力値とする。

4. 経過ターン数

入力値は 7.1 式のシグモイド関数で計算する。

$$Sigmoid(x) = \frac{1}{1 + \text{Exp}(-(x - 8))} \quad (7.1)$$

図 7.1 は実装した効用モデルのイメージである。効用モデルの調整にはローカルサーチ法を用い、各入力ノードと出力ノードを繋ぐコネクシヨンの結合強度  $w_j$  を変化させることで行った。

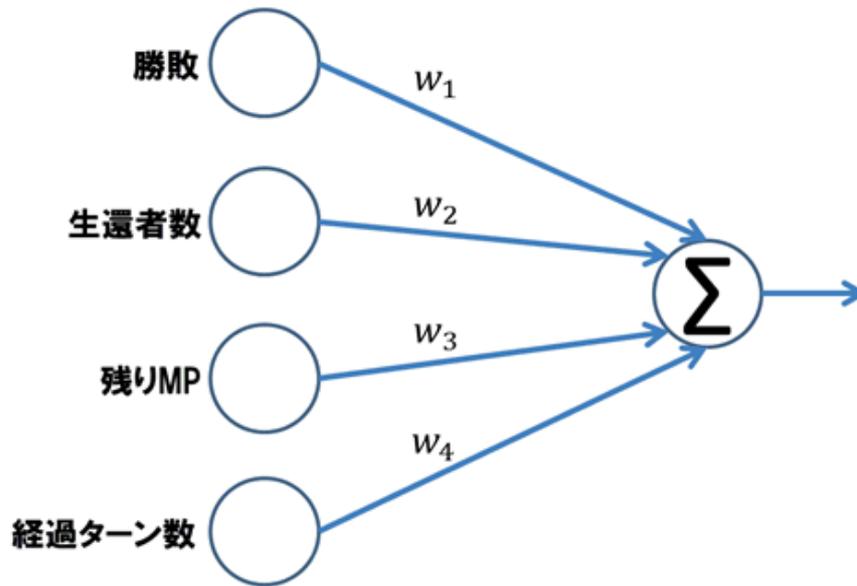


図 7.1: 入出力モデルのイメージ

## 7.2.2 学習部分

今回の実験では「味方キャラクタ 2」を人間プレイヤーに見立てて効用の近似を行うこととする。提案手法の AI は「味方キャラクタ 2」とチームを組み、100 戦分戦う中で学習用データ  $D$  を収集し学習を行う。図 7.2 は学習部分全体のイメージである。詳細は以降で説明する。

### 1. 期待効用値の計算

「味方キャラクタ 2」が選択可能な行動の集合を  $A$ 、自身が選択可能な行動の集合を  $B$  として  $EuTables(u, D)$  を計算する。今回は状態遷移における分岐数の多さ・ランダム性の強さなどの理由から、モンテカルロシミュレーションを使って各行動の組合せ  $(a_n, b_m)$  の期待効用値を近似し、期待効用値表  $EuTable(u, s_i)$  を得る。モンテカルロシミュレーションの処理手順は以下の通りである。シミュレーションでは「良い」方策を高い確率で選ぶ方が正しい良さの推定ができることが知られているが [9]、ここではランダムな選択とした。

#### (a) 最初のターンの行動選択

状態  $s_i$  から戦闘のシミュレーションを開始する。最初のターン、自チームのキャラクタは期待効用値を調べたい行動の組合せ  $(a_n, b_m)$  の行動を選択する。敵チームのキャラクタはランダムに行動選択を行う。

#### (b) 次ターン以降の行動選択

次ターン以降は、戦闘に参加する全キャラクタがランダムに行動を選択する。これを戦闘終了まで繰り返す。

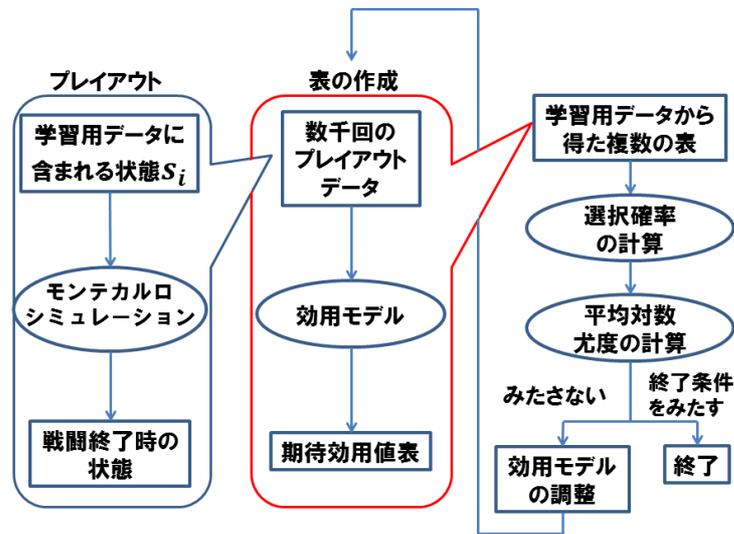


図 7.2: 実装した学習部分全体のイメージ

(c) 戦闘が終了したら

戦闘が終了したら効用モデルを使って戦闘終了時の状態を評価し効用値を計算する.

## 2. 選択確率の計算

ここでは前段階で得られた期待効用値表  $EuTable(u, s_i)$  それぞれに対して、人間側の選択枝の選択確率を計算する. 今回の実験では「人間らしい」選択確率の計算方法として、以下のような方法で状態  $s_i$  において人間プレイヤーが行動  $a_i$  を選択する確率  $p(s_i, u, a_i, EuTable)$  を計算した. これらの式は、数人分の被験者実験の結果から予想したものに過ぎず、実際に人間がどのような確率で行動するのかは明らかではない.

(a) 人間側の選択枝の削減

効用値にはその人の主観が入っていること、また、人間は完全に合理的ではないことなどから、ゲーム理論などで提唱されている考え方を適用すれば人間の選択確率を正しく計算できるとは限らない. また、今回の AI は期待効用値の計算にモンテカルロシミュレーションを使用しているため、シミュレーションによる誤差も考慮しなければならない. そこで、この段階では「人間らしい」考え方として「明らかな被支配戦略」と「仮の弱被支配戦略」を定義し、これらに該当する選択枝を以降の計算の対象から外すことで、選択確率を計算するための前処理を行う.

「明らかな被支配戦略」を以下の通り定義する.

“ 相手がどの戦略を取ったとしても，ある戦略の期待効用値の有意水準 5 % の下限以下の期待効用値しか得られない戦略 ”

表 7.3 は「明らかな被支配戦略」の一例である．各行動の組合せは 100 回ずつモンテカルロシミュレートされているものとする．この場合， $(a_1, b_1)$  の有意水準 5 % の下限は 0.4， $(a_1, b_2)$  の有意水準 5 % の下限は 0.216 となり，いずれも  $(a_2, b_1)$  と  $(a_2, b_2)$  の期待効用値より大きくなる．よって  $a_2$  は「明らかな被支配戦略」となる．

表 7.3: 「明らかな支配戦略」の一例

	$b_1$	$b_2$
$a_1$	0.5	0.3
$a_2$	0.3	0.0

「仮の弱被支配戦略」以下のように定義する．

“ 自分側の 2 つの戦略について見たとき，ある相手の戦略に関して低い方の期待効用値が高い方の期待効用値の有意水準 5 % の下限より高い場合，それらを同程度と見なす．このとき同程度と見なした値以外に注目し，相手がどの戦略を取ったとしても，ある戦略の期待効用値の有意水準 5 % の下限以下の期待効用値しか得られず，かつ平均期待効用値が 90 % 以下の戦略 ”

表 7.4 は「仮の弱被支配戦略」の一例である．各行動の組合せは 100 回ずつモンテカルロシミュレーションされているものとする．この場合， $(a_2, b_1)$  と  $(a_1, b_2)$  共に有意水準 5 % の下限は 0.4 となる．相手が  $b_2$  の戦略を取った場合， $(a_2, b_2)$  の期待効用値 0.45 は  $(a_1, b_2)$  の有意水準 5 % の下限 0.4 より大きいので，これらを同程度とみなす． $(a_2, b_1)$  の有意水準 5 % の下限 0.4 は  $(a_1, b_1)$  の期待効用値 0.1 より大きい．各戦略の平均期待効用値も  $0.3 \leq 0.475 \times 0.9$  を満たしている．よって  $a_1$  は「仮の弱被支配戦略」として以降の計算の対象から外れる．このようにして削減した人間側の選択肢の集合を  $A_{cut} (\subseteq A)$  とし，任意の要素を  $a_c (\in A_{cut})$  と表す．

表 7.4: 「仮の弱被支配戦略」の一例

	$b_1$	$b_2$
$a_1$	0.1	0.5
$a_2$	0.5	0.45

(b) 人間側の選択肢の評価値計算

前段階で削減されなかった選択肢について、それぞれの評価値を計算する。この評価値は、人間がその選択肢に対して抱く期待度を数値化したものである。計算式は7.2式の通りである。

$$Em(a_c, s_i) = max_c - ((1 - belief) * (risk_c - average_c)) \quad (7.2)$$

ここで、 $max_c$  は  $a_c$  を選んだ場合に得られる可能性がある最大期待効用値、 $average_c$  は  $a_c$  を選んだ場合の平均期待効用値、 $risk_c$  は  $a_c$  を選んだ場合の最大期待効用値と最低期待効用値の差の絶対値、 $belief$  は人間プレイヤーからAIプレイヤーに対する信頼度を表す調整用パラメータである。本研究では、ある程度の期間共にプレイしてきた人間プレイヤーとAIプレイヤーを想定しているので、 $belief$  の値を0.5とした。 $belief$  値が高ければ最大値、低ければリスク回避的な選択を行う。

(c) 人間側の選択肢の確率を計算

状態  $s_i$  において人間プレイヤーが行動  $a_i$  を選択する確率を7.4式で計算する。 $Em$  値が1違うと選択確率が  $e$  倍になるような式である。

$$sumEm = \sum_c Exp(Em(a_c, s_i)) \quad (7.3)$$

$$p(a_i, s_i) = \begin{cases} \frac{Exp(Em(a_i, s_i))}{sumEm} & (\exists a_i \in A_{cut}) \\ 0 & (otherwise) \end{cases} \quad (7.4)$$

3. 対数尤度の計算

対数尤度は、 $\epsilon$  を0.1として6.5式の通りに計算する。

4. 効用モデルの調整

今回はローカルサーチにより効用モデルを調整した。現在の効用モデルを  $u_t$  (図7.1を参照) とする。

ローカルサーチでは最初に、 $u_t$  の入力ノードと出力ノードを繋ぐコネクションの中からランダムにひとつのコネクションを選択肢し、その結合強度を0.8~1.2倍の範囲でランダムに変化させ、その後最大強度の絶対値が1になるよう正規化する。

次に、全ての結合強度から絶対値が最大となるもの ( $maxAbs$ ) に注目し、この絶対値を使い全ての結合強度を7.5式のように正規化する。 $w'_i$  は正規化後の各結合強度である。

$$w'_i = \frac{w_i}{maxAbs} \quad (7.5)$$

上記のようにして、より尤度の高い効用モデルの候補  $u'$  を生成する。  $u'$  を使用して学習用データ集合  $D$  の対数尤度を計算し、この値が現在の効用モデル  $u_t$  の対数尤度よりも大きかった場合、  $u_t$  を  $u'$  と置き換える。これを繰り返すことで最適化を行う。

### 7.2.3 行動選択部分

現在の状態を  $s_t$ 、現在の効用モデルを  $u_t$  とする。行動選択部分では、最初にこの AI が持つ  $u_t$  を使い、人間プレイヤーが選択可能な行動と、自身（「味方キャラクタ 1」）が選択可能な行動の組合せについての期待効用値表  $EuTable(u_t, s_t)$  を求める。このとき、7.2.2 で紹介したモンテカルロシミュレーションを使い、期待効用値の近似を求めている。 $EuTable(u_t, s_t)$  が求まったら、最大期待効用値となった組合せに含まれている自身の行動を、このターンの行動として選択する。

## 7.3 その他の AI の概要

提案手法以外の、実験で使用する AI の概要を説明する。

### 1. モンテカルロ AI (MCAI)

モンテカルロ AI は戦闘結果について簡単な効用を持っており、自身と味方（この実験の場合は「味方キャラクタ 1」）が選択できる行動の組合せについてモンテカルロシミュレーションを行う。得られた期待効用値表から期待効用値が最大となる行動の組合せに含まれる自身の行動を選択する。モンテカルロ AI が持つ効用の内容は次の通りである。

#### (a) 勝てばよい派

「戦闘終了時、自チームのキャラクタ 2 人中 1 人でも生きていれば 1 の効用を持ち、2 人共行動不能の場合は 0 の効用を持つ。」

#### (b) 死んだら負け同然派

「戦闘終了時、自チームのキャラクタ 2 人中 1 人でも行動不能ならば 0 の効用を持ち、2 人共生きている場合は 1 の効用を持つ。」

#### (c) MP 重視派

「戦闘終了時、自チームが勝っている場合はチーム全体の残り MP のパーセン

テージに比例した効用を持ち、負けている場合は0の効用を持つ。」

## 2. ランダム AI

ランダム AIはその時選択可能な行動の中からランダムにひとつの行動を選択する AIである。

## 3. ヒューリスティック AI

このヒューリスティック AIは、敵チームの中で一番 HP の数値が低い敵を狙って攻撃する AIである。

## 4. 原始モンテカルロ AI

原始モンテカルロ AIは自身の行動選択肢のみについてモンテカルロシミュレーションを行い、一番期待勝率が高くなった行動を選択する AIである。

## 7.4 実験の設定

提案手法の AI の設定は以下の通りである。

1. 一回の行動選択で5000回（原始モンテカルロの場合はひとつの選択肢につき500回）のシミュレーションを行う
2. シミュレーションの対象となる行動の組合せはUCB値を使って決定する
3. 効用モデルを使用し、シミュレートした戦闘終了時の状態から効用値を算出する
4. 期待効用値表中、一番期待効用値が高い組合せに含まれる自身の手を選択する

## 7.5 実験結果

提案手法の AI には、各モンテカルロ AI の効用を100回分の戦闘を通じて学習させた。提案手法が保持する学習用データの最大数は50とし、収集したデータ数が50を超える場合は古いものから破棄していくようにした。

### 7.5.1 効用モデルの結合強度による評価

学習終了時の、提案手法の AI が調整した効用モデルの結合強度パラメータは表7.5のようになった。また、各種モンテカルロ AI に対する理想的な結合強度のパラメータを表しているのが表7.6である。学習後の効用モデルの結合強度は理想値の通りの値にはなら

なかったが、モンテカルロ AI(a) とモンテカルロ AI(b) に対して理想値に近い偏りとなっていることが分かる。また、結合強度のパラメータは学習の対象となったモンテカルロ AI が持つ効用の違いによって、その数値の偏りが異なっていることも分かる。

結合強度のパラメータが理想値通りにならなかった理由として以下のような可能性が挙げられる。

1. 局所解に捕まってしまった可能性
2. 学習用データの数が有限なので、その有限個数のデータに対して高い尤度のパラメータに収束してしまった可能性

表 7.5: 学習後の結合強度のパラメータ

	MCAI(a) の結合強度	MCAI(b) の結合強度	MCAI(c) の結合強度
戦闘の勝敗	1	-0.63	1
生還者数	0.158	1	0.847
チームの残り MP	-0.576	0.376	1
経過ターン数	0.25	0.025	-0.847

表 7.6: 理想的な学習後の結合強度のパラメータ

	MCAI(a) の結合強度	MCAI(b) の結合強度	MCAI(c) の結合強度
戦闘の勝敗	1	-0.5	1
生還者数	0	1	0
チームの残り MP	0	0	1
経過ターン数	0	0	0

### 7.5.2 「人間が不快に感じ得る行動選択」の出現頻度による評価

「人間が不快に感じ得る仲間 AI の行動選択」のパターンの抽出方法を以下のように定義し、モンテカルロ AI のシミュレーションの記録と戦闘の記録からその出現頻度を算出した。

1. モンテカルロ AI の作成した期待効用値表中の、最大期待効用値に注目する

2. 最大期待効用値との差が閾値以内の行動の組合せをリストアップする
3. リストアップした行動の組合せに含まれる仲間の行動を列挙していく
4. 列挙した行動の中に、実際に仲間が選択した行動が無ければ該当パターンとする

このパターンの抽出方法は、仲間の AI の行動選択において、その行動を選択した理由が理解できない場合に人間プレイヤーは大きな不満感を感じるはずであるという考えに基づいて定義した。

例えば人間プレイヤーが表 7.7 のような期待効用値表を得たとする。  $a_1 \sim a_2$  が人間プレイヤー側の選択肢で、  $b_1 \sim b_4$  が仲間の AI プレイヤーの選択肢である。この場合、仲間の AI プレイヤーが  $b_1$  もしくは  $b_3$  を選択をしたならば、人間プレイヤーは不満感を感じることはないだろう。何故ならば、  $b_1$  の選択肢には  $(a_1, b_1)$  の組合せが、  $b_3$  には  $(a_2, b_3)$  の組合せが、十分に期待効用値が高い組合せとして存在し、その行動選択の理由付けができるからである。しかし、  $b_4$  の組合せは  $(a_2, b_4)$  の組合せが 0.7 ではあるが、先の 2 つと比べると得られる最大の期待効用値は低いので、人間プレイヤーが不満に感じてしまう可能性がある。また、  $b_2$  については  $b_1$  の被支配戦略となっており、この選択をする理由が人間プレイヤーからは理解できない。このような選択をした場合に人間プレイヤーは大きな不満を感じる可能性が高い。「人間が不快に感じ得る仲間 AI の行動選択」は、表 7.7 の  $b_2$  や  $b_4$  のような行動を仲間 AI が選択してしまった場合を認識するように定義されている。

表 7.7: 期待効用値表の一例

	$b_1$	$b_2$	$b_3$	$b_4$
$a_1$	0.9	0.6	0.3	0.4
$a_2$	0.6	0.3	0.8	0.7

以上のように「人間が不満に感じ得る行動選択」を定義し、その出現頻度を調査した。表 7.8 は各モンテカルロ AI と提案手法の AI がチームを組んだ場合の、問題パターンの出現頻度である。同じ効用を持つモンテカルロ AI 同士がチームを組んだ場合には、問題となるパターンの出現頻度がかなり低いことが分かる。また、仲間の AI プレイヤーの効用が違う場合には、問題パターンの出現頻度が高くなっていることも分かる。提案手法の AI は、同じ効用を持つモンテカルロ AI 同士でチームを組んだ場合には及ばないが、違う効用を持つモンテカルロ AI 同士で組んだ場合よりも上手く相手の効用を近似できていることが分かる。

表 7.8: パターンの出現頻度

	MCAI(a)	MCAI(b)	MCAI(c)
MCAI(a)	0.000%	6.236%	16.457%
MCAI(b)	1.930%	0.303%	23.759%
MCAI(c)	1.589%	24.157%	0.213%
提案手法 AI	0.304%	4.175%	4.868%

## 第8章 まとめ

本研究では人間プレイヤーの効用を推定・学習し，その効用を使って行動選択を行うことで，人間プレイヤーの満足度が高いAIプレイヤーの構成手法を提案し，独自に設定したコマンド形式RPGゲームに適用した．提案手法のAIはシステムの内部に人間の効用を近似するための入出力モデルを持ち，学習部分ではそのパラメータを調整することで人間プレイヤーの効用を近似し，行動選択部分では調整した入出力モデルを使うことで，人間が不満に思わないような行動を選択するようにした．

今後の研究予定として，提案手法は実際に人間の不満度を改善できるのかを調べるために被験者実験を行うことや，効用モデルを多層パーセプトロンなどの複雑なモデルに変えて性能を調査することなどを考えている．

## 第9章 付録

### 9.1 実験 1

第7章の実験で使用した戦闘の設定で人間プレイヤーがゲームをプレイした場合に、仲間のAIプレイヤーの行動選択に不満を示す可能性があるかどうかを調査した。以下に不満を持った状況と、その際に仲間プレイヤーに選択して欲しかった行動を例示する。

表 9.1: 味方キャラクター 2 が MCAI(a) の際に、人間プレイヤーが不満を持った状況の一例

	残り HP	残り MP
味方キャラクター 1 (被験者)	150	50
味方キャラクター 2 (MCAI)	50	100
敵キャラクター 1	0	0
敵キャラクター 2	0	0
敵キャラクター 3	80	40

表 9.2: 味方キャラクター 2 が MCAI(b) の際に、人間プレイヤーが不満を持った状況の一例

	残り HP	残り MP
味方キャラクター 1 (被験者)	140	50
味方キャラクター 2 (MCAI)	100	80
敵キャラクター 1	0	0
敵キャラクター 2	0	0
敵キャラクター 3	40	40

表 9.1 の状況で、人間プレイヤーは味方プレイヤー 2 に「敵キャラクター 3 に攻撃」する行動を望んでいたが、実際に味方プレイヤー 2 が選択したのは「味方プレイヤー 1 を回復」する行動であった。この状況では、シミュレーションの結果はどの行動を選択しても勝率が 100% になってしまうと予想される。MCAI(a) は勝敗にのみ効用を持っているので、全ての行動の期待効用値が等しくなってしまうために、このような人間が不満を感じる行動でも選択肢に入ってしまうと考察する。

表 9.3: 味方キャラクタ 2 が MCAI(c) の際に、人間プレイヤーが不満を持った状況の一例

	残り HP	残り MP
味方キャラクタ 1 (被験者)	110	50
味方キャラクタ 2 (MCAI)	50	100
敵キャラクタ 1	0	0
敵キャラクタ 2	60	0
敵キャラクタ 3	80	60

表 9.2 の状況で、人間プレイヤーは味方プレイヤー 2 に「敵キャラクタ 3 に攻撃」する行動を望んでいたが、実際に味方プレイヤー 2 が選択したのは「味方プレイヤー 1 を回復」する行動であった。この状況でも、シミュレーションの結果はどの行動を選択しても二人共生きて勝つ確率が 100% になってしまうと予想される。MCAI(b) は 2 人共に生きて勝利することに効用を持っているので、全ての行動の期待効用値が等しくなってしまうため、このような人間が不満を感じる行動でも選択肢に入ってしまうと考察する。

表 9.3 の状況で、人間プレイヤーは味方プレイヤー 2 に「味方キャラクタ 2 を回復」する行動を望んでいたが、実際に味方プレイヤー 2 が選択したのは「敵プレイヤー 2 を攻撃」する行動であった。MCAI(c) はより MP を多く残して勝利することに効用を持っているので、「味方キャラクタ 2 を回復」してより安全に勝利するよりも、一斉に攻撃して敵を倒してしまう方が期待効用値が高くなってしまい、このような人間が不満を感じる行動を選択したと考察する。

## 9.2 実験 2

第 7.2.2 章で定義した「仮の弱支配戦略」の削減が実際に行われるのか調査した。また、同章で触れた信頼度によって手の選択確率が変化するかも調査した。調査では利得表を使った簡単なゲームを被験者に行ってもらい、アンケート形式で回答してもらった。ゲームの概要は以下の通りである。

1. ゲームにはゲームマスターがおり、被験者は仮想 AI プレイヤーと一緒にチームを組んでゲームをプレイする。
2. 被験者と仮想 AI プレイヤーは、最初にゲームマスターから 5 枚の利得表を配られる。ゲームは、この配られた利得表を一枚ずつ使いながら進む。
  - ゲームマスターは、あらかじめ 2 種類の利得表のグループを用意している。

- 被験者は、仮想A I プレイヤに配られた利得表のグループが何なのか知ることができない。
- 仮想A I プレイヤも、被験者に配られた利得表のグループが何なのか知ることができない。
- 被験者と仮想A I プレイヤには、同じ利得表のグループが配られる可能性もあれば、違う利得表のグループが配られる可能性もある。

### 3. ゲームの進行は次の通り。

- (a) 被験者と仮想A I プレイヤは、配られた利得表を見て自身の行動（戦略）を決定し、ゲームマスターに伝える。
- (b) 2人共に行動（戦略）を伝え終わったら、ゲームマスターは被験者と仮想A I プレイヤがそれぞれどの行動（戦略）を選んだか発表し、それぞれに利得を与える。
- (c) 利得表は“一回戦目用の利得表”，“二回戦目用の利得表”……“五回戦目用の利得表”のようにになっているので、(a) と (b) を5回繰り返したら終わり。

### 4. 最終的に獲得した利得がなるべく多くなるようにゲームを進め。

以上を被験者に伝えてゲームを行った。

実験2の前半では人間プレイヤからAIプレイヤへの信頼度が低くなるようにゲームを進行し、最後にアンケートを行った。アンケートでは、第1回戦から第4回戦までを踏まえて100回行動を選択できるならば、各選択肢を何回ずつ選択するかを聞いた。

アンケートの結果として表9.4の利得表では平均値として、 $a_1$ が20.625回・ $a_2$ が43.75回・ $a_3$ が6.875回・ $a_4$ が28.75回選ばれる結果となった。 $a_1$ は $a_2$ の、 $a_3a_4$ の「仮の弱支配戦略」となっており、それぞれ選択される回数が低めになっている。

また、アンケート結果ではこの時点で「相手に配られた利得表が違う」と感じているプレイヤが多く、そのためか選択肢毎の最低値がより高い $a_2$ の選択肢がより選択されやすい傾向が出ているように見えた。

実験2の後半では人間プレイヤからAIプレイヤへの信頼度が高くなるようにゲームを進行し、最後にアンケートを行った。アンケートでは前半と同様に、第1回戦から第4回戦までを踏まえて100回行動を選択できるならば、各選択肢を何回ずつ選択するかを聞いた。

アンケートの結果として表9.5の利得表では平均値として、 $a_1$ が26.25回・ $a_2$ が48.125回・ $a_3$ が13.125回・ $a_4$ が12.5回選ばれる結果となった。 $a_2$ は $a_1$ の、 $a_4a_3$ の「仮の弱支配戦略」となっているが、前半とは違う傾向が見てとれた。

アンケート結果ではこの時点で「相手に配られた利得表が同じ」と感じているプレイヤが多く、 $a_2$ は $a_1$ の弱支配戦略ではあるがより高い効用値を求めて $a_2$ が選択されやすいとい

う傾向の結果となった。これは、 $b_3$  が仮想 AI プレイヤにとって支配戦略である影響もあると推測した。

表 9.4: 実験 2 の前半で使用した利得表 1

	$b_1$	$b_2$	$b_3$	$b_4$
$a_1$	0.4	0.3	0.95	0.3
$a_2$	0.45	0.5	0.9	0.45
$a_3$	0.6	0.35	0.75	0.5
$a_4$	0.7	0.3	0.7	0.7

表 9.5: 実験 2 の後半で使用した利得表 2

	$b_1$	$b_2$	$b_3$	$b_4$
$a_1$	0.45	0.45	0.9	0.5
$a_2$	0.4	0.3	0.95	0.3
$a_3$	0.7	0.7	0.7	0.3
$a_4$	0.6	0.5	0.75	0.35

# 謝辞

本研究を進めるにあたり，池田心准教授には研究の構想や機械学習に関する要点での確かなアドバイスをいただきました。本論文の添削など，様々な場面でご迷惑をおかけしましたが，丁寧にご指導をいただき深く感謝いたします。

また，東条敏教授には研究の方向性などを相談したところ，有益な助言を数多くいただきました。ありがとうございました。

研究室の大町洋君には，選択確率の計算について相談したところ親身に考察をしていただきました。ありがとうございました。

最後に，被験者実験などに協力してくださった研究室の皆様に感謝いたします。

## 参考文献

- [1] Sander Bakkes, Pieter Spronck and Eric Postma TEAM : The Team-oriented Evolutionary Adaptability Mechanism, Entertainment Computing - ICEC 2004, pp.273-282, 2004.
- [2] Aswin Thomas Abraham and Kevin McGee, AI for Dynamic Team-mate Adaptation in Games, Computational Intelligence and Games(CIG), 2010 IEEE Symposium on, pp.419-426, 2010.
- [3] Kenneth O.Stanley, Bobby D. Bryant, Student Member, IEEE, and Risto Mikkulainen, Real-Time Neuroevolution int the NERO Video Game, Evolutionary Computation , IEEE Transactions on, Vlo.9, No.6, pp.653-668, 2005.
- [4] 田村坦之, 中村豊, 藤田眞一 効用分析の数理と応用, コロナ社, pp.1-28, 1997.
- [5] 加藤英明, 行動ファイナンス - 理論と実証 - 朝倉書店, pp.61-66 2003.
- [6] 大町洋, 池田心 同時進行ゲームのためのモンテカルロ木探索 The 17th Game Programming Workshop 2012, IPSJ Symposium Series Vol.2012, No.6, pp.197-204, 2012.
- [7] Remi Coulom Computing Elo Ratings of Move Patterns in the Game of Go ICGA journal, Vol.30, No.4, 2008.
- [8] 保木邦仁, 渡辺明 ボナンザ VS 勝負脳 - 最強将棋ソフトは人間を超えるか 角川書店, 2007.
- [9] Gelly, G., Wang, W., Munos, R. and Teytaud, O. Modification of UCT with Patterns in Monte-Carlo Go Technical Report, No.6062, INRIA, 2006.