## マルチエージェントを用いたチーム AI と個別 AI について

日大生産工(院) ○藤田 真広 日大生産工 齋藤 敏雄

# 1. まえがき

近年の人工知能技術の発展はめまぐるしい ものがあり、様々な分野に応用されてきた.

ゲーム業界においては、特にグラフィックに偏り発展してきた歴史と、メモリ容量問題や人工知能技術に計算時間が必要なものが多いなどの理由により、広く利用されているとはいいがたく、ごく一部の分野において利用されていたに過ぎない。しかし、ハードウェアの進化により、シミュレーション能力が向上し、ゲーム基本システムに対する余剰リソースが発生した。これにより、グラフィックのみならず、人工知能技術をゲームに組み込むことが可能となった。

本研究では、"スコットランドヤード"というボードゲームを題材として、AIを備えたエージェントをプレイヤーとして定式化し、チームとしての振る舞いと個人としての振る舞いを比較検討する。結果として、ゲームに人工知能技術を応用することで、ゲーム製作に対し新しい幅を提示する事を目的とする。

### 2. スコットランドヤード

スコットランドヤードは、ドイツのランベスバーガー社より発売された代表的なボードゲームである。ロンドン市内をモチーフにしており、4~5人(本研究では5人を想定)の刑事と1人の怪盗にわかれ、盤上の1~199の地点をそれぞれ、タクシー・バス・地下鉄・船(怪盗のみ)を使い移動する。また、移動手段はそれぞれ使用回数を制限されている。

刑事側は怪盗を捕まえるか、移動できないように追い詰めることが勝利条件で、怪盗側は24時間逃げ切れば勝ちとなる.

## 3. マルチェージェントシステム

## 3.1 マルチエージェントシステム 1)2)

マルチエージェントシステムは,自立的に 行動する多数のエージェントから構成される. それぞれのエージェントは自分が置かれてい る環境を知覚し,自分の目標を達成するよう に行動を選択する.エージェントは互いに影 響を及ぼしあい,それが,各エージェントの 行動選択基準を変化させるキッカケにもなる. マルチエージェントシステムには,協力型マ ルチエージェントシステムと競争型マルチエ ージェントシステムがあるが,本研究ではそ の両方を用いるものとする.

## **3.2 強化学習** 3) 4)

強化学習は教師なし学習のひとつであり、 環境の状態sに対して行動aをとったときに 環境から得られる報酬rをもとに、初期状態 からゴール状態に渡って受け取る報酬が最大 になるような行動戦略を学習する. (Fig. 1) 環境に関する正しい知識をあらかじめ準備す る必要がなく、行動とその行動の評価を繰り 返しながら学習していくため、環境が変化す る動的な環境にも対応することができる.

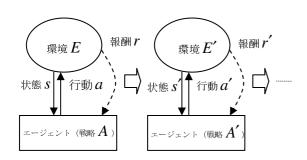


Fig. 1 強化学習における相互作用

### 3.3 マルコフ性

マルコフ性とは、確率論における確率過程の持つ特性の一種で、その過程の将来状態の条件付き確率分布が、現在状態のみに依存し、過去のいかなる状態にも依存しない特性を持つことをいう。本研究では、エージェントの現在位置と残りの移動手段をマルコフ状態として扱う。よって、t+1における環境の応答はtにおける状態と行動表現のみに依存することになり、このときには全てのs'、r、 $s_t$  とa, に対して

$$\Pr\{s_{t+1} = s', r_{t+1} = r | s_t, a_t\}$$
 (1)

のみを指定することで環境のダイナミクスを 定義することができる.

## 3.4 有限マルコフ決定過程(有限 MDP)

マルコフ性を満たす強化学習タスクはマルコフ決定過程 (MDP) と呼ばれる. また,本研究では状態と行動の空間が有限であるので,有限 MDP であるといえる. 有限 MDP は状態と行動の集合と,環境の1ステップダイナミクスから定義される. 次に可能な各状態 s'の確率は,

$$P_{ss'}^{a} = \Pr\{s_{t+1} = s' | s_t = s, a_t = a\}$$
 (2)

これらの量は遷移確率と呼ばれている. 同様

にして,次の報酬の期待値は,

$$R_{ss'}^{a} = E\{r_{t+1} | s_t = s, a_t = a, s_{t+1} = s'\}$$
 (3)

## 3.5 価値関数

価値関数は状態の関数で、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を評価する.ここでは、どれだけ良いのかという概念を将来において期待される報酬に関して定義する.エージェントが将来受け取ることを期待できる報酬は、エージェントがどのような行動を取るかに依存する.したがって、価値関数は特定の方策に関して定義される.

方策 $\pi$  が各状態 $s \in S$  と行動 $a \in A(s)$  から、状態s で行動a を取る確立 $\pi(s,a)$  への写像 であるといえる、その時、MDP に対する

$$V^{\pi}(s) = E_{\pi}\{R_{t} | s_{t} = s\}\}$$

$$= E_{\pi} \{ \sum_{t=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s$$
 (4)

 $E_{\pi}$ {}は、エージェントが $\pi$ に従うとしたときの期待値を表す.関数 $V^{\pi}$ を方策 $\pi$ に対する状態価値関数と呼ぶ.

同様に、方策 $\pi$ のもとで、状態sにおいて行動aを取る事の価値を $Q^{\pi}(s,a)$ で表し、状態sで行動aを取り、その後に方策 $\pi$ に従った期待報酬として定義する.

$$Q^{\pi}(s, a) = E_{\pi}\{R_{t} | s_{t} = s, a_{t} = a\}$$

$$= E_{\pi} \{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} | s_{t} = s, a_{t} = a \}$$
 (5)

 $Q^{\pi}$ を方策 $\pi$ に対する行動価値関数と呼ぶ.

### 3.6 最適価値関数

有限 MDP に対しては、以下のようにして最適方策を定義することができる。価値関数は方策に関して反順序を定義する。すべての状態に対して、方策 $\pi$ の期待収益が $\pi'$ よりも良いか同じであるなら、 $\pi$ は $\pi'$ よりも良いか、

同じであると定義される. 言い換えるなら、 すべての $s \in S$ に対して、 $V^{\pi}(s) \ge V^{\pi'}(s)$ で あるなら、そのときに限り $\pi \ge \pi'$ である. 他 の方策よりも良いか、それに等しい方策が常 に少なくとも1つ以上存在し、これが1つの 最適方策である. 最適方策は1つ以上存在す るかもしれないが、全ての方策を $\pi^*$ と記す. 最適方策群は、最適状態価値関数と呼ばれる、 同じ状態価値関数を共有する. 最適状態関数 は、すべての $s \in S$ に対して

$$V^*(s) = \max_{\pi} V^{\pi}(s) \tag{6}$$

と定義される.

# スコットランドヤードにおけるマルチエージェントの実装

### 4.1 移動方法と移動可能範囲

怪盗は3・8・13・18・24時間目以外は姿が見えないが、毎時間どのように移動したかは刑事にもわかる.移動にはチケットを利用し、最初に刑事には TAXI チケット 10 枚、バスチケット 8 枚、地下鉄チケット 4 枚を渡され、これを用いて移動する.怪盗は TAXI・バス・地下鉄は無制限に移動でき、それとは別に、ブラックチケット 5 枚とダブルムーブ 2 枚が渡される.ブラックチケットは TAXI・バス・地下鉄の他に、船を使って移動する事ができる.また、ブラックチケットを使用したときは、刑事には移動手段がわからない.ダブルムーブはその名の通り、2 時間連続で移動することが可能になる.

199 の移動地点 (Fig. 2) は、すべてある別の移動地点より TAXI で移動可能であり、バスを利用可能な地点は 59 地点、地下鉄が利用可能な地点は 13 地点である。また、船で移動できる地点は 4 地点である。



Fig. 2 スコットランドヤードの地図とチケット 5)

## 4.2 刑事の行動

刑事の移動は5人のエージェントによる協力型マルチエージェントシステムであるとして、怪盗を逮捕するように動く.

### 4.2.1 協力型マルチエージェントシステムの利用

刑事は、怪盗を逮捕するという問題を、直接逮捕に向かう刑事、地下鉄や船を怪盗に利用させないために駅などを押さえる刑事、全体的なゾーンで移動範囲を狭める刑事という副問題へと割り当て、分割する。各エージェントはマルコフ性を持ち合わせているために、この副問題の割り当ては永続的なものではなく、毎時間変化するものとする.

### 4.2.2 移動手段の選択

3時間目までは、怪盗の位置がわからないために地下鉄や船を利用させないように動く.3時間目以降は、それぞれに副問題を割り当て移動させる、怪盗までの距離・怪盗の移動手段による存在しない地域への可能性の除去・残っている移動手段・現在地などを評価することで副問題をどのエージェントが担当するかを割り当てる、その割り当てと残りの移動手段をもとに、移動手段を選択する。また、思考の範囲は怪盗が姿を現した時間から、次に姿を現す時間までの範囲で思考する.

### 4.3 怪盗の行動

全体のシステムとして、怪盗と刑事は競争型マルチエージェントであるとして、刑事から逃げる.

現在から 5 時間後までの刑事の移動可能範囲を考える. 1 時間で移動可能な場所を報酬r=-1 として、そこから、1 時間ごとにr を割引きする. また、複数の刑事が到達可能な場所は、r 同士を加算する. また、移動可能な手段の種類と量により一定のr を加算す

る. そうして、得られた報酬のうち高いものを選択し、移動する. 得られる最大の報酬がある一定以下になった場合には、ブラックチケットやダブルムーブを使用する. また、姿を現した直後からブラックチケットを使うまでは一定の報酬を減算する. これは、刑事側に現在位置を推測されにくくするためである.

## 5. 結果と考察

本研究では、100回の試行を1セットとし、100セット行う.1セット毎に勝敗により報酬の値を変化させ、刑事・怪盗ともに学習させていく.1セット毎の勝率の推移と、刑事が逮捕した場合にかかった時間の推移を比較し、正しく学習できたか確認する.正しく学習できている場合は、刑事・怪盗ともに最適方策が収束していくために、勝率は一定の確率に収束していき、逮捕した場合もかかる時間は増加していくはずである.また、怪盗よりも刑事の方が報酬が曖昧であるために、収束に時間が掛かることが予想される.

また、刑事を協力させずに、個別の行動に対して、報酬を与え学習させる.この場合でも、100回の試行を100セット行い、1セット毎の勝敗により報酬の値を変化させる.協力した場合と、しなかった場合の勝率の推移と、逮捕するまでにかかった時間を比較する.

### 6. まとめ

本研究では、ボードゲームという、エージェントが非常に限られた行動群の中からしか行動を選択しなかったが、行動の幅の広いアクションやシューティングにおいても行動群を一般化することで、マルチエージェントシステムを導入する事が可能であると考える.

また、プランナーやデザイナーの経験則によって学習させるのではなく、繰り返し試行する事によって最適方策を模索するものなので、製作の負担を軽減することが可能であると考える。しかし、経験則による学習ではないために、AI が思考のループに陥ったり、デバッグ時にデバッグ項目を挙げにくいなど、いくつかの問題点を内包している。

今後の展望としては、まだ未完成であるプログラムを完成させ、これらの理論を実証し検証すべきである。また、刑事が、自分たちのいる場所から遠ざかるであろうという推測を報酬に追加した場合と、怪盗が、刑事のいる場所から遠ざかるであろうという推測を一定確率で裏切るという可能性を追加した場合についても比較・検証する必要がある。

#### 参考文献

- 1) 大内東, 山本雅人, 川村秀憲, マルチエー ジェントシステムの基礎と応用, コロナ社, (2002), pp. 1-90.
- 2) 高玉圭樹, マルチエージェント学習, コロナ社, (2003), pp. 21-64.
- 3) 三上貞芳,皆川雅章,強化学習,森北出版, (2000), pp. 2-170.
- 4) 電気学会,学習とそのアルゴリズム,森北 出版,(2002),pp155-179.
- 5) Rabensburger,

http://www.ravensburger.de/web/

6) 三宅陽一郎,人工知能が拓くオンラインゲームの可能性,A0GC2007,(2007).