

## サッカー協調プレイの強化学習のための状態縮約

The contraction state for reinforcement learning in soccer cooperation play

鹿田 和之進<sup>1</sup> 西野 順二  
Kazunoshin Shikada Junji Nishino

電気通信大学 システム工学専攻  
Dept. of Systems engineering  
The University of Electro-Communications

**Abstract:** In this paper, we show an empirical result on effects of state contraction in cooperative multi agent behavior rule using reinforcement learning method for 'centering tactics' in RoboCup soccer. We examine a difference between finely divided state space and roughly one. Learning result of contracted state with 30 games perform as well as fine-divided state with 60 games. The result of fine-contract combined method also perform well.

### 1 はじめに

サッカーのような複数のエージェントが、協調、競合して点を取り合うマルチエージェント環境における協調行動の実現は、非常に興味のある話題であり、これまで様々な研究がされてきた。マルチエージェント環境では、学習者が1人であったとしても、他のエージェントの行動政策が未知の場合、容易に状態空間を構成できない。内部 [1] らは、学習者の観測と行動を通して、学習者和其他者の行動の関係を局所予測モデルとして推定し、その結果をもとに強化学習を行った。また、複数エージェントが同時にゼロから学習する困難な問題に対して、浅田 [2] は、進化的手法を導入し協調行動の実現に成功している。

本研究で扱う RoboCup Soccer シミュレーションリーグでは、ドリブルやシュートなどといった個人レベルの基礎スキルや、フォーメーションなどが主となって研究されている。現実のサッカーではさらに上のレベル、戦術や戦略といった協調プレイが必要とされる。しかし、協調プレイを Q-learning や強化学習によって

実現するうえで、状態、行動の増加による計算量が問題となり、学習が成り立たないことが多い。

本研究では、代表的な強化学習手法である Q-learning において、その学習効率を上げる為の状態縮約について検討する。そして、戦術の中で特に重要な協調的プレイを獲得することを目的とする。

### 2 強化学習における状態の縮約

強化学習とは、環境から与えられる報酬 (reward) をもとに、期待報酬を最大化するような行動を徐々に選択するように、行動を修正していく教師なしの学習法である。教師付き学習と異なり、状態入力に対する正しい行動出力を明示的に示す教師が存在しない。また、報酬にはノイズや遅れがあり、行動を実行した直後の報酬をみるだけでは、エージェントはその行動が正しかったかどうかを判断できないという困難を伴う。

#### 2.1 強化学習の枠組み

強化学習の枠組を図1に示す。エージェントは環境の状態を表す状態集合  $S$  を観測し、行動集合  $A$  の中の一つの行動を選択する。このと

<sup>1</sup>鹿田 和之進 電気通信大学  
〒182-0021 東京都調布市調布ヶ丘 1-5-1  
Tel:0424-43-5800, Fax:0424-43-8020  
E-mail:shika@fs.se.uec.ac.jp  
本研究の一部は科学研究費補助金、課題番号 14750179 によって行われた。

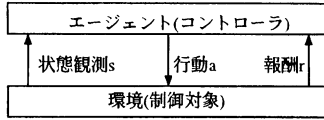


図 1: 強化学習の枠組

き、環境はマルコフ過程としてモデル化され、現在のエージェントがとった行動により確率的な状態遷移を行う。

現在の状態を  $s$ 、エージェントがとった行動を  $a$ 、次状態を  $s'$  とすると、 $T(s, a, s')$  は、そのときの状態遷移確率を表す。ここで、状態から次に選択すべき行動への写像が決まる。この写像を方策 (policy) という。また、それぞれの状態・行動のペア  $(s, a)$  に対し、報酬  $r(s, a)$  が定義される。

エージェントの目的は、環境中を行動しながらタスクを完遂するまでに、できるだけ多くの報酬を獲得することである。ここで、報酬の合計である積算報酬 (retrun)  $R_t$  は、

$$R_t = \sum_{n=0}^{\infty} \gamma^n r_{t+n} \quad (1)$$

で定義される。ここで、 $\gamma$  は割引報酬率を表し、将来の報酬がどの程度行動価値に影響を与えるかを定める。

## 2.2 Q-learning

実際の環境では、状態遷移確率と報酬の分布が完全に既知である場合は少なく、試行錯誤的に探索しながら、最適な行動を学習することが考えられる。以下では、このような探索的な学習法である Q-learning について説明する。

状態  $s$  で行動  $a$  をとり、それ以降最適政策をとったときの積算報酬の期待値、もしくは最適行動価値関数を  $Q^*(s, a)$  をとすると、

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s' \in S} T(s, a, s') \max_{a' \in A} Q^*(s', a') \quad (2)$$

と定義される。最初は、遷移確率や報酬が未知なので、オンラインで逐次的に行動価値  $Q$  (以下  $Q$  値) を更新する。 $Q$  値は初期値 0 から始め、行動がとられるごとに、 $Q$  値を更新する。

本研究では、以下の Q-learning アルゴリズム [4] を用いた。

1. エージェントは環境の状態  $s_t$  を観測する。
2. エージェントは任意の行動選択方法 (探索戦略) に従って行動  $a_t$  を実行する。
3. 環境から報酬  $r_t$  を受け取る。
4. 状態遷移後の状態  $s_{t+1}$  を観測する。
5. 以下の更新式により  $Q$  値を更新：

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha[r_t + \gamma \max_a Q(s_{t+1}, a)] \quad (3)$$

6. 時間ステップ  $t$  を  $t+1$  へ進めて手順 1 へ戻る。

## 2.3 サッカーでの状態とその縮約

サッカーでは、 $108[m] \times 50[m]$  の連続二次元空間のフィールドで試合が行われる。フィールドには、2 チーム 22 人のプレイヤーとボールが存在し、計 23 の状態変数がある。試合の状態数は、 $1[m]$  単位に離散化したとしても、 $(108 \times 50)^{23}$  状態あることになる。この状態では、学習を行う環境に適していない。本研究では、大雑把に空間を分割することにより状態縮約し、このような環境で Q-learning を用い学習を行う。

## 3 縮約の効果の検証

縮約の効果について実験的に検証する。

### 3.1 RoboCup サッカーサーバシステム

RoboCup サッカーサーバは仮想的なフィールドを提供し、ボールとプレイヤーの全ての動

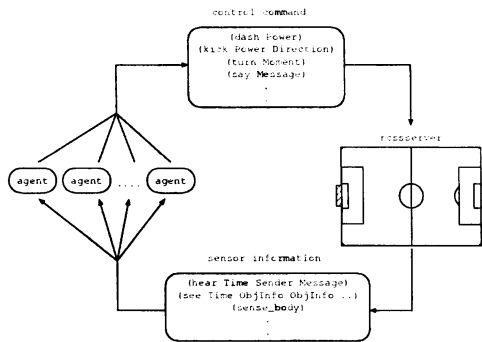


図 2: サッカーサーバシステム

きを2次元でシミュレートする。また、サッカーサーバはプレイヤーの視野である視覚情報や、プレイヤーの位置や体力などが含まれる身体情報などを各プレイヤーに送り、各プレイヤーはこれらの情報から周囲の状況や自らの状態を分析し、走る (dash)、ボールを蹴る (kick) といった行動をサッカーサーバに送る。この一回の送受信を1ステップ (0.1秒) として、試合は6000ステップ (10分) 繰り返行なわれる。システム構成を図2に示す。また、プレイヤー間では情報共有ができないが、かけ声 (say) を用いることによりコミュニケーションをとれるが、本研究では用いないこととする。

### 3.2 サッカーエージェント

プレイヤーはRoboCup2003 (世界大会) で優勝したUvA\_trilearn2003[3]を味方 shooter、敵ディフェンダー (以下DF)、敵ゴールキーパー (以下GK) で使用した。また、学習者である passer プレイヤーは、UvA\_base2003 を使用し作成する。

このUvA\_trilearn2003はアムステルダム大学のJelle Kokが作成したプレイヤーで、coordination graph を用いてマルチエージェント環境での高度な意志決定を行っている。また、得点力が強くGKと1対1では高い確率で得点する事ができる。

UvA\_base2003はUvA\_trilearn2003の基礎ライブラリ (ドリブルやパスなど) を公開し

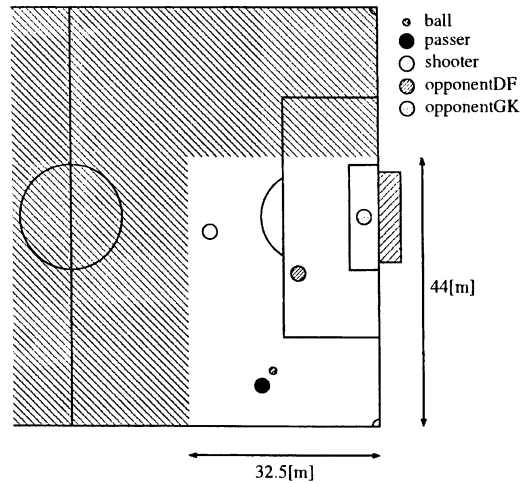


図 3: 実験環境

ている。本研究では、このライブラリを使用し高度な意志決定の学習を試みる。

### 3.3 実験設定

本実験では、図3のような44[m]×32.5[m]のフィールドで行う。この実験では、オフENSEの学習者 passer、味方 shooter、ディフェンスの敵DF、敵GKの計4エージェントのみとする。passerはQ-learningにより意志決定方策を学習していく。ここでは、passerがサイドからshooterにセンタリングをあげ、ゴールを狙う戦術を行う。

passerはボールを保持しているとき、ドリブルをしてゴールに近づくか、あるいはセンタリングをするかの意志決定を学習する。ボールを保持していないときは、shooterよりボールに近いときはボールを取りに行き、それ以外ときは、決められたポジションに移動する。

passerの行動として図4(a)に示した以下の行動を用意した。

1. 上ヘドリブルする。(Dribble1)
2. 右上ヘドリブルする。(Dribble2)
3. 右ヘドリブルする。(Dribble3)

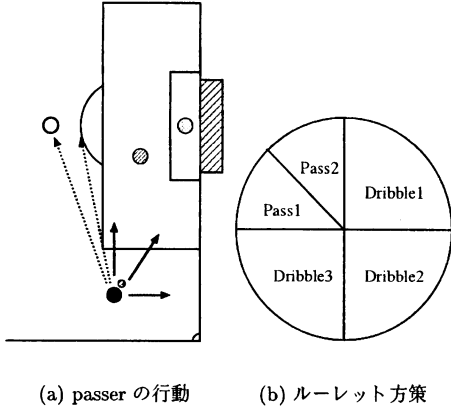


図 4: passer の行動

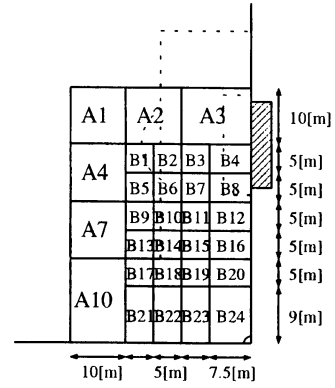
4. shooter の方向へセンタリングをあげる。  
(Pass1)
5. shooter の 5m 前方へセンタリングをあげる。  
(Pass2)

また、学習の過程で用いる方策は図 4(b) のようなルーレット方策を用いた。これは、パスをするとその時点でエピソードが終了し、学習の効率が著しく低くなることを避ける為である。ルーレット方策は、パスをする確率を低くすることで早い終了を回避している。

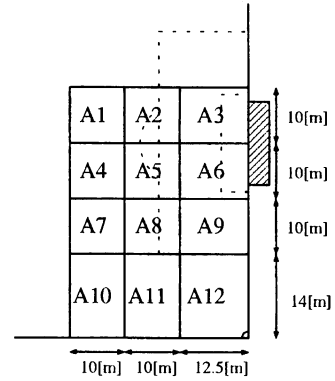
実験 1 では、空間を細かくして精度は良いが、収束速度の遅い学習を行う。実験 2 では状態空間を縮約し、空間を粗く分割し収束速度が早い学習を行う。実験 3 では、まず実験 2 で得られた Q 値を利用して、実験 1 で使用できるように Q 値を拡張し、さらに、実験 1 と同様の状態空間で学習を行う。これより、実験 3 のような状態を改めて変更することで、収束速度が速く精度の良い学習できたかを検証する。

### 3.4 実験 1 (細分)

実験 1 では、図 5(a) のように空間を細かく分割し学習を行った。passer の位置を  $X_p$ 、shooter の位置を  $X_s$ 、敵 DF の位置を  $X_d$  とした。ここで、 $X_p, X_s, X_d$  は以下のように定義する。



(a) 細かい空間分割



(b) 縮約した空間分割

図 5: 空間分割

$$X_p \in \{A7, A10, B9, B10, \dots, B24, o\} \quad (4)$$

$$X_s \in \{A1, A2, A3, A4, A5, A6, o\} \quad (5)$$

$$X_d \in \{A2, A3, A11, A12, B1, B2, \dots, B16, o\} \quad (6)$$

ここで、 $o$  はそれ以外の場所である。

敵 GK はゴール前にいるものとして、状態変数としていない。よって、状態空間  $(X_p, X_s, X_d)$  は、 $19 \times 7 \times 21 = 2793$  通りある。また、Q 値の状態数は  $2793 \times 5 = 13965$  通りとなる。

状態の初期値 ( $X_{p_0}, X_{s_0}, X_{d_0}$ ) は、「その他」のエリアを除き、

$$X_{p_0} \in \{A7, A8, A9, A10, A11, A12\} \quad (7)$$

$$X_{s_0} \in \{A1, A2, A4, A5\} \quad (8)$$

$$X_{d_0} \in \{A2, A3, A5, A6, A8, A9, A11, A12\} \quad (9)$$

で定義され、この  $6 \times 4 \times 8 = 192$  通りからランダムで決められる。シュートが入るか、敵 GK にボールを取られるか、フィールド外にでたら 1 エピソード終了とする。報酬はシュートが入った時に 100 を与える。

サッカーサーバは 1 試合 6000[step] (10 分) であり、本実験では、試合回数を 60 回に設定し学習を行った。1 試合における Q 値の更新回数は約 600 回、エピソード数は約 130 回であり、60 試合の Q 値の更新回数は約 36000 回、エピソード数は約 7800 回であった。また、学習率  $\alpha$  を 0.5 とし、割引報酬率  $\gamma$  を 0.7 に設定した。

学習が終了した後、学習の効果測定して得られた Q 値を使用としてさらに 30 回試合を行った。passer の方策に Q 値の最大値を選ぶグリーディ方策を用いた。最大の Q 値が 0 の場合、ランダムで行動を決定した。

### 3.5 実験 2 (縮約)

実験 1 では図 5(b) のような  $A_1 \sim A_{12}$  の空間分割を行い、passer の位置を  $X'_p$ 、shooter の位置を  $X'_s$ 、敵 DF の位置を  $X'_d$  とした。ここで、 $X'_p, X'_s, X'_d$  は以下のように定義する。

$$X'_p \in \{A7, A8, A9, A10, A11, A12, o\} \quad (10)$$

$$X'_s \in \{A1, A2, A3, A4, A5, A6, o\} \quad (11)$$

$$X'_d \in \{A2, A3, A5, A6, A8, A9, A11, A12, o\} \quad (12)$$

ここで、 $o$  はそれ以外の場所である。

状態空間 ( $X'_p, X'_s, X'_d$ ) は  $7 \times 7 \times 9 = 441$  通りある。また、Q 値の状態数は  $441 \times 5 = 2205$  通りとなる。状態空間の初期値、学習率、割引報酬、学習者の方策は実験 1 と同様である。但し、試合回数は実験 1 の半分の 30 回である。

表 1: 各 30 試合の結果

	得点の平均	得点の標準偏差
学習前	26.7	3.28
実験 1 (細分)	35.4	4.75
実験 2 (縮約)	33.4	4.61
実験 3 (段階)	34.7	3.72

学習が終了した後、実験 1 と同様に学習の効果測定として得られた Q 値を利用してさらに 30 回試合を行った。

### 3.6 実験 3 (段階的縮約)

実験 2 の Q 値を用いて実験 1 の状態空間に適応できるように拡張した。例えば、 $B1, B2, B5, B6$  に  $A5$  のデータを使う。そして、実験 1 と同様な条件で学習させる。但し、学習回数は半分のは 30 回とする。

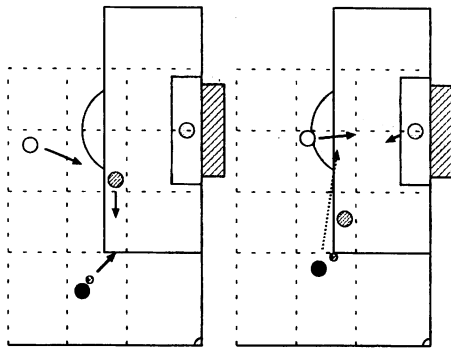
学習が終了した後、実験 1、2 と同様に学習の効果測定として得られた Q 値を利用してさらに 30 回試合を行った。

## 4 縮約した状態の効果

表 1 は実験 1, 2, 3 の各 30 試合の結果である。平均得点、得点の標準偏差を記す。

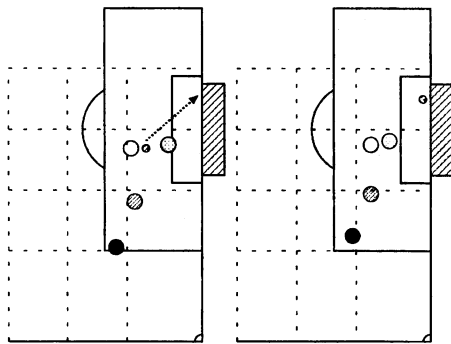
学習前と実験 1 の得点を比較すると、学習者のセンタリングにおける、行動の学習に成功していることがわかる。試合を詳細にみると、図 6 のような、協調行動がみられる。passer は敵 DF を引きつける為に、ドリブルで敵 DF に接近し、近づいて味方にセンタリングをあげるとい、高度な意志決定を行っており、戦術な協調行動を獲得することができた。しかし、成功率が低い理由として、空間の分割が粗く、大雑把な意志決定しかできないこと、行動を 5 個に制限したため、パス精度が悪いことも多かったためと考えらる。

また、shooter である UvA\_trilearn が強すぎ、DF と GK の 2 者がいてさえも、シュートを入れることがあり、最適解に Q 値が収束して



(a)  
(A11, A4, A5) → D2  
敵DFを誘い出す

(b)  
(A11, A5, A8) → P2  
味方ヘセントラリング



(c) シュート

(d) ゴール

図 6: 獲得した協調行動

いない可能性もあると考えられる。

また、実験 1, 2 の結果を比較すると、改善はみられたが小さいものであった。しかし、縮約した状態によって 30 回の試行でもある程度学習することが示され、問題に適応して分割を変更することより、より良い学習ができると考えられる。

## 5 おわりに

本研究により、Q-learning によりマルチエージェント間での、センタリングにおける戦術的な協調行動を獲得した。また、状態縮約を用い

ることにより、短期間である程度学習できることを示した。一方、段階的に細分化することで、学習の効果を上げることは、今回の実験では明らかではない。今後は、初期における状態縮約の方法、段階的な分割を行うための指針について検討が必要である。

## 参考文献

- [1] 内部英治、浅田稔、細田耕,  
複数の学習するロボットの存在する環境  
における協調行動獲得のための状態空間  
の構成, 日本ロボット学会誌 Vol.20  
No.3, pp.281-289, 2002
- [2] 浅田稔, ロボットの行動学習・発達・進化,  
共立出版, 2002
- [3] Jelle Kok, UvA\_trilearn2003  
[http://carol.wins.uva.nl/jellekok/  
robocup/index\\_en.html](http://carol.wins.uva.nl/jellekok/robocup/index_en.html)
- [4] Richard S.Sutton, Andrew G.Barto  
強化学習, 森北出版, 2000