

学習による協力関係構築に関する 1・2・5じゃんけんの考察

仁科 諒子[†] 竹川 高志[†]
[†]工学院大学情報学部情報デザイン学科

1. はじめに

非ゼロ和ゲームの「1・2・5じゃんけん」では、相互に協調することで高得点を得られることがわかっている[1]. しかし選択された手を見ただけではプレイヤーの明確な意図がわからず、また協調といえる状態が複数ステップに及ぶため協調が難しい.

本研究では「1・2・5じゃんけん」を用いて、強化学習プレイヤーに行動戦略を学習させる. プレイヤーが対戦相手と自分の行動履歴に基づき、対戦を行うことでQ学習によって行動戦略を学習する様子をシミュレーションし、人間社会においてどのように行動すればお互いに協調し合い社会全体の利益を増やすことができるか考察を行った.

2. 1・2・5じゃんけん

1・2・5じゃんけんとは、通常のじゃんけんと同じ「グー(G)」、「チョキ(C)」、「パー(P)」で勝敗を決めるが、勝ったときの手により得点が変わるゲームである.

グー・チョキ	グー選択側は+1点 チョキ選択側は-3点
チョキ・パー	チョキ選択側は+2点 パー選択側は-3点
パー・グー	パー選択側は+5点 グー選択側は-3点
あいこ	各々加点数なし

両者がG, C, Pを等確率で選んだ場合は、お互いに得られる平均利得が-1/9となる. この得点構造の場合、プレイヤー2人が協調し交互にGとPを出し合う戦略は、5点と-3点を交互に得ることで安定して得点を増やすことができる. 両者が始めから協調できた場合、お互いに得られる平均利得は、1となる.

3. 方法

実験では、学習しないプレイヤーとして対戦相手に「歩み寄り戦略(CMP)」「GP 交互戦略(ALT)」「GP 等確率戦略(GPR)」「等確率戦略(EQP)」「常にP戦略(POL)」を用いた.

CMPは「囚人のジレンマ」における「しっぺ返し戦略」に相当し、相手が協調していると思われるなら自分も協

調を続け、逆に相手が協調する様子を見せないならば自分も協調しない(表1). 状態sを過去2手の対戦履歴、行動aを出す手としてQ学習を行った.

表1: 歩み寄り戦略

最初の4行はそれぞれが対戦時に出した手の履歴である. CMP(A), CMP(B)の列は選択する手である. RはGCPの中から1/3の確率で手を選択し、GPはGPの中から1/2の確率で手を選択する.

相手の1手前	G		C		P									
相手の2手前	G	C	P		G		C		P					
自分の1手前					G	C	P	G	C	P	G	C	P	
自分の2手前												G	C	P
CMP(A)	G	P	R	G	R	P	R					GP		
CMP(B)	G	P	R	G	R	P	C	P	R	C	R	GP	C	GP

4. 実験

学習結果は5つの戦略のプレイヤーとの対戦結果で評価した. 単独の単純戦略で学習した場合、全ての場合において相手の戦略の攻略法を見つけた. 例えば、CMP, GPRで学習したQ学習プレイヤーはPOLに対してPを出し続けた. GPRと対戦したQ学習プレイヤーは、CMP(A), ALT, GPR, POLに対してPを出し続けた. 未来報酬の割引率 γ が高いとき、CMP(B)で学習したQ学習プレイヤーは、CMPに対してGPPの順で手を出し続けた. EQPと対戦したQ学習プレイヤーはALTに対して、 γ が低いときはCPPPの順で手を出し続けた.

5. まとめ

今回の条件のように個々の単純な戦略で学習させた場合、Q学習の戦略はパターンを読みきり、協調より絶対的な勝利に推移する. 人間社会のように対戦相手が1人に限らず、またその戦略もさまざまである場合の学習で得た戦略を評価することが今後の課題である.

参考文献

[1] 伊藤昭, 大橋資紀, 寺田和憲:「非零和ゲームの強化学習-相手の行動を読むプログラム」, 情報処理学会研究報告. ICS, [知能と複雑系], 53-60, 2005