

# 履歴を用いたQ学習による交渉問題へのアプローチ

水野将史

松本達明

伊藤昭

寺田和憲

岐阜大学工学部

## 1 はじめに

交渉問題は、それぞれのエージェントが独立に最善を尽くすよりも、話し合って協力して行動した方が有利なゲーム理論的状况である。この問題は、経済学・ゲーム理論では Nash による提案以来 [1]、妥当な協力の満たすべき条件をめぐって、様々に議論され続けてきた。

我々は、交渉問題にこれまでのような規範的の接近をとるのではなく、2 個のエージェントがインタラクトする中で、双方が自己にとって最良の妥協点を求める動力学 (dynamic) の結果として解が求まるものとする。この動力学を記述する方法として、各エージェントは相手の行動を観察し、得られた情報に基づき自己の行動を調節する学習エージェントであると考え、マルチエージェント強化学習の理論を適用する。

## 2 交渉問題

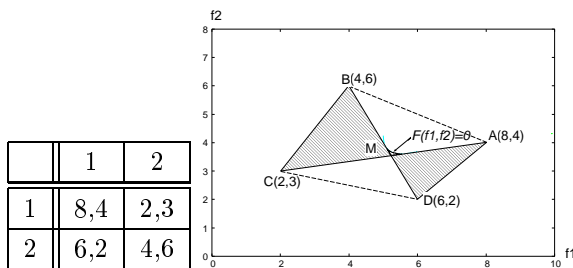


図 1: 交渉問題の利得行列と現実可能領域

我々が取り上げるのは、図 1 にある利得表の非零和二人ゲームであり、「協力すればお互いに利益があるが、その妥協点をめぐって双方が争わなければならない」問題である [2]。この問題の Nash 均衡点は図 1 の A, B 二つであり、その時の利得は (8,4)(4,6)、ともにパレート (Pareto) 最適である。また、それぞれが独立に  $p, q$  の確率で 1, 2 を選択する混合戦略をとると、その実現可能領域は図 1 のようになる。

しかし交渉が可能であれば、それぞれが独立に最適な混合戦略をとるよりも、協力して図の A または B を一定の割合でプレーする方が双方にとって有益であり、その時の報酬  $(r_1, r_2)$  は  $r_1 + 2r_2 = 16$  (線分 AB 上) を満

たす。しかしながら、協調が成立するためには、妥協点  $X(A$  をプレーする割合  $x$ ) を交渉によって合意する必要がある。Nash が提案する交渉問題の妥協点 = 協力の解は、基準点の利得  $(c_1, c_2)$  から  $r_1 + 2r_2 = 16$  の条件の下で  $f_u = (r_1 - c_1)(r_2 - c_2)$  を最大化する  $(r_1, r_2)$  である。この外にも基準点には様々な提案があるが、どれも無矛盾な解以上に深い意味を持たない。

## 3 履歴を使った Q 学習

我々のアプローチは、実際に強化学習エージェントを対戦させることで、どのようにして交渉が行われるのかを観察し、そこから解の意味を考え直そうというものである。学習エージェントは、これまでの対戦履歴を使って報酬の期待値を最大化する自己の政策 (戦略) を学習する。

まず状態  $S_t$  を、自己および相手の過去  $h$  手の情報を用いて  $S_t = \{a^m(\tau), a^o(\tau)\}_{\tau=t-h}^{t-1}$  で定義する。この時の自己の政策  $Q(S, a^m, a^o)$  は、状態  $S$  で自分が手  $a_m$ 、相手が  $a_o$  を選んだ時の得られる報酬の期待値となる。ここで期待値が最大となる政策が最強である。すなわち、

$$Q^\pi(S, a^m, a^o) = r(a^m, a^o) + \gamma \max_{a^m} \bar{Q}^\pi(S', a^m)$$

$$\bar{Q}^\pi(S', a^m) = \sum_{a^o} p(a^o | S', a^m) Q^\pi(S', a^m, a^o)$$

となる。 $(S'$  は、 $S$  から行動  $(a^m, a^o)$  によって遷移した状態である。) また、Q 学習に倣って逐次近似法で  $Q$  の値を求めることにすると、次の式を得る。

$$Q(S, a^m, a^o) \leftarrow (1 - \alpha) Q(S, a^m, a^o) + \alpha (r(a^m, a^o) + \gamma \max_{a^m} \bar{Q}(S', a^m))$$

$$\bar{Q}(S', a^m) = \sum_{a^o} p(a^o | S', a^m) Q(S', a^m, a^o)$$

また行動選択確率は

$$p(a^o | S_t, a^m) = p(a^m, a^o | S_t) / (\sum_{a^o} p(a^m, a^o | S_t))$$

で更新される。

## 4 実験/結果

我々は履歴を用いた確率予測 Q 学習エージェント PQh を履歴長を様々に変えて対戦させてみた。プレイヤー 1, 2 の平均利得を  $r_1, r_2$  とすると  $s = r_1 + 2r_2$  が 16 に近い程、交渉の精度が高いと言える。また、 $x = BX/AB$  で妥協点  $X$  の線分 AB 上の位置を表す。

まずは、同じ履歴長同士 (PQh-PQh) の対戦を行った。実験は利得表を用いた対戦を  $10^8$  回行うのを 1 ラウンドとし、乱数を変えて 10 ラウンド行った。そのときの  $s, x$  の時間的振舞い (10 ラウンドのうちの 1 つのラウンド)

<sup>1</sup>An approach to bargaining problem using Q-Learnig with history

Masafumi Mizuno, Tatsuaki Matumoto,  
Akira Ito, Kazunori Terada  
Faculty of Engineering, Gifu University

を図2に示す。図からわかるようにように  $s$  は急速に上限近くに収束し、多くの場合  $x$  も一定値に収束する。そこで各ラウンド  $10^8$  の内、後ろの半分についての平均を採って収束値とした。用いたパラメタは、 $\alpha = 0.1$ ,  $\beta = 0.9$ ,  $\gamma = 0.1$ ,  $\delta = 0.01$ ,  $T = 0.4$ (P1),  $T = 0.2$ (P2) である。 $T$  の値が P1 と P2 で異なっているのは、交渉のポイントでは双方の利得が  $r_1 + 2r_2 = 16$  を満たしながら争われるためである。

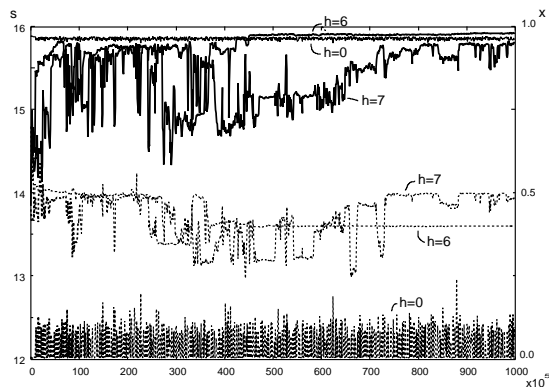


図 2: PQh-PQh 対戦における  $s, x$  の時間変化

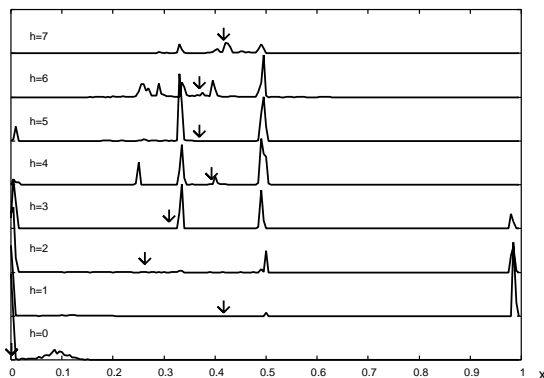


図 3: 履歴長別  $x$  の分布

履歴長を変えたときの、10 ラウンドの  $x$  の分布を図3に示す。図を見ると多くの場合  $x$  の収束値は簡単な分数  $n/m$  になっており、A,B が  $n:m$  の割合で採られる政策であることがわかる。履歴長を  $h$  とすると周期  $(h+1)$  を越えるパターンを学習することはできないため、可能な収束値の分母  $m$  は  $m \leq (h+1)$  に制限される。また、履歴長が小さい時ほど  $x$  の収束は端に偏り、 $h$  の増大に伴い中心近くに移動する。

さらに、同じデータを用いて  $s, x$  の収束値の平均をもとめた結果を表1に示す。履歴長が大きくなるにしたがって  $x$  の平均値が大きくなり、分散が小さくなっている。履歴長  $h = 6$  では  $s$  は上限値に近く、協調的行動がうまく行われているといえる。しかし、 $h = 7$  では逆に

習できなかったことを意味している。

	PQ0	PQ1	PQ3	PQ5	PQ6	PQ7
$s$	15.85	15.87	15.88	15.85	15.72	15.55
$x$	0.039	0.417	0.305	0.372	0.373	0.412
$\sigma(x)$	0.045	0.473	0.299	0.131	0.095	0.057

表 1: PQh-PQh 対戦における  $s, x, \sigma(x)$

つぎに、収束が見られた最大履歴長  $h = 6$  と様々な履歴長とを対戦させた結果を含め  $x$  の平均値を図3に示す。PQ6-PQh は P1 の履歴長を 6 に、PQh-PQ6 は P2 の履歴長を 6 に固定したものである。履歴長の大きいものと小さいものが対戦した場合、小さいものの方が優位な利得を挙げている。これは、履歴長の長い方が「適応能力が大きい」ので、どうしても先に譲歩してしまうと考えられる。

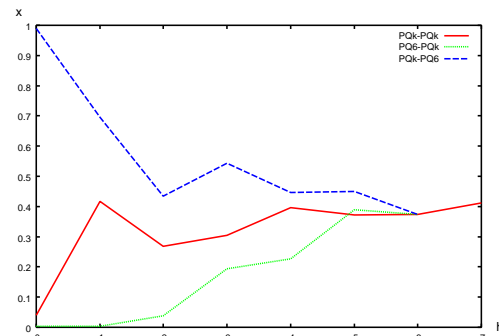


図 4: PQh-PQh, PQ6-PQh, PQh-PQ6 対戦における妥協点  $x$

## 5 まとめ

実験結果は、これまでの理論と少し異なっていてこれまでに提案された様々な解は  $x \geq 0.5$  を提案しているのに対して、実験結果は  $x < 0.5$  を示唆している。また、この収束値  $x$  の決定には、 $(s_1, s_2) = (2, 1)$  での利得が大きく関与していることが別実験からわかっている。

人がこのような問題に直面したときには、交渉が成立しなかった場合の損得情報を踏まえ過去の経験を下に、心の中でシミュレーションし妥協点を計算する。今回我々が行ったマルチエージェント強化学習による交渉問題へのアプローチは、ある意味で人が行っているような交渉過程を再現したものである。実際、人と人でこの問題に取り組んだ場合、今回のシミュレーションに近い結果が得られるかどうか、興味深い研究課題である。

## 参考文献

- [1] Nash, J.: "The bargaining problem," *Econometrica*, 18, pp.155-162, 1950.
- [2] 鈴木光男:「新ゲーム理論」, 勁草書房, 1994.