

複素価値関数を用いた強化学習に関する基礎的検討

A fundamental study of Q-learning with complex value function

澁谷 長史†

Takeshi SHIBUYA

濱上 知樹†

Tomoki HAMAGAMI

1 はじめに

強化学習において、エージェントは環境から与えられる報酬に基づいて行動政策の獲得を行う。強化学習は、環境との相互作用をもとに自律的に学習を行うため、自律移動ロボットなどへの応用に向いている。

一方で、不完全知覚問題や次元の呪いなどをはじめとする多くの本質的問題も存在し [1], 重要な課題となっている。これら本質的問題を解決するような研究もさまざまに行われている [2][3]。しかし、多くの方法は環境に対してサブゴールなど一定の仮定を導入することで環境への依存性が高まってしまう。

本研究では、従来法とは異なり、仮定を用いずにこれらの問題の解決を図り、汎用性のある手法を検討している。特に本稿では価値関数に複素数を用いる手法について基礎的な検討を行った。複素価値関数による強化学習を複素強化学習と呼ぶことにする。具体例として本稿では Q-learning における Q 値を複素数として扱う。これを \hat{Q} -learning と呼ぶことにする。計算機実験によって複雑なアルゴリズムを用いることなく不完全知覚問題のある環境において提案手法が有効であることを示す。

2 提案手法

2.1 方針

Q 値が複素数であることを明示的に \hat{Q} 値と表す。 \hat{Q} 値の更新式において遷移先の状態に関連する \hat{Q} 値との荷重平均を取る際に、位相回転を加えて時系列の情報を含ませる。

本稿では、直前の行動の \hat{Q} 値との位相差を考慮して、次のステップで選択されるであろう \hat{Q} 値を予測する。図1は予測された \hat{Q} 値 (\hat{R}) と選択可能な行動に対応する \hat{Q} 値 (\hat{Q}_1, \hat{Q}_2) の関係を示したものである。 \hat{R} が異なると、原点と各 \hat{Q} 値からの垂線の足を結ぶ線分の長さも異なる。このことを利用すると、選択確率を文脈に基づいて動的に変化させることができる。 \hat{Q} 値を実数として扱う場合には単純に大きさの比較を行うことしかできないが、複素数である \hat{Q} 値を用いることで位相差を含んだ比較が可能となる。これによって複雑なアルゴリズムを使わなくても文脈に応じた行動選択をすることができるようになる。

2.2 更新アルゴリズムの定式化

状態 s_i から行動 a_i をとって状態 s_{i+1} へと遷移し報酬 r を受け取ったときの、 \hat{Q} 値の更新則を以下のように定義する。

$$\hat{Q}(s_{i-k}, a_{i-k}) \leftarrow (1 - \alpha)\hat{Q}(s_{i-k}, a_{i-k}) + \alpha(r + \gamma \hat{Q}_{\max}^{s_i \rightarrow s_{i+1}}) \hat{u}(k) \quad (1)$$

k ステップ前の状態, 行動をそれぞれ s_{i-k}, a_{i-k} とする。

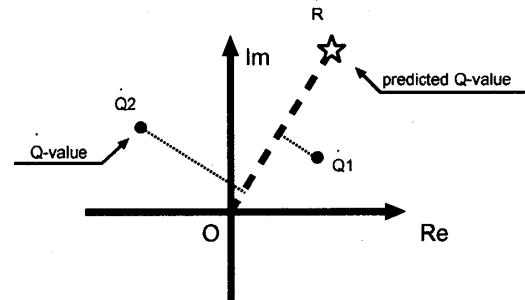


図1 ある状態における \hat{Q} 値を複素平面上に示した図

ここで、 k ステップ前の状態, 行動をそれぞれ s_{i-k}, a_{i-k} とする。 $\hat{u}(k)$ は形式上の適格度トレースであり、式2のように定義する。

$$\hat{u}(k) = \beta^k \quad (2)$$

式1の適用は、予め定めた整数 N_e を用いて、 $0 \leq k \leq N_e$ の範囲で行う。ただし、 β は絶対値が1以下の複素数である。 $\hat{Q}_{\max}^{s_i \rightarrow s_{i+1}}$ は、式3のように定義する。

$$\hat{Q}_{\max}^{s_i \rightarrow s_{i+1}} = \max_a \hat{Q}(s_{i+1}, a) \overline{\hat{R}_i} \quad (3)$$

ここで、予想される \hat{Q} 値 (\hat{R}_i) は次のように定義する。

$$\hat{R}_i = \hat{Q}(s_i, a_i) / \beta \quad (4)$$

2.3 行動選択アルゴリズムの定式化

本稿では、Max-Boltzmann 選択を用いる。すなわち、状態 s_i に居るエージェントは、確率 $1 - P_{\max}$ で Boltzmann 選択を行い、確率 P_{\max} で greedy 方策を行うことにする。

状態 s_i , 行動 a_i に対応する \hat{Q} 値を $\hat{Q}(s_i, a_i)$ とする。また、状態 s_i における行動 a の Boltzmann の選択確率を $Prob(s_i, a)$ とする。状態 s_i における行動集合を $A(s_i)$, 直前の状態と行動に対応する \hat{Q} 値を $\hat{Q}(s_{i-1}, a_{i-1})$, Boltzmann 選択の温度パラメータを T とするとき、 $Prob(s_i, a)$ を次式のように定める。

$$Prob(s_i, a) = \frac{\exp(\text{Re}[\hat{Q}(s_i, a) \overline{\hat{R}_{i-1}}] / T)}{\sum_{a' \in A(s_i)} \exp(\text{Re}[\hat{Q}(s_i, a') \overline{\hat{R}_{i-1}}] / T)} \quad (5)$$

ただし、 $\text{Re}[\cdot]$ は複素数の実部を表す。

greedy 方策は $\arg \max_a Prob(s_i, a)$ を選択することにする。

† 横浜国立大学大学院工学府

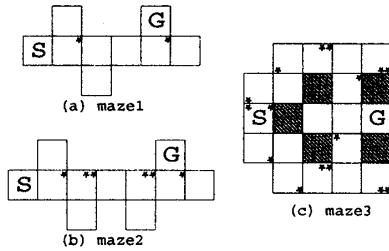


図2 実験環境。Sはスタートを、Gはゴールをそれぞれ表す。*や**は不完全知覚の影響のある状態を表す。

表1 パラメータ設定

	trials	α	$1 - P_{\max}$
phase 1	1 to 20	0.05	0.1
phase 2	21 to 80	(100-try)/400	(100-try)/1600
phase 3	81 to 100	0	0

3 計算機実験

図2のような簡単なグリッドワールドにおける迷路問題を対象として計算機実験を行い、提案手法の有効性を確認する。

3.1 状態空間と行動集合

エージェントが観測可能な情報は、東西南北周囲4マスの壁の有無のみとし、この情報を直接状態として割り当てることにする。すなわち観測可能な状態数は $2^4 = 16$ となる。

これらの環境において不完全知覚の影響のある状態が存在する。例えば、*においてはそれぞれにおいて選択すべき行動が異なり、**においては同じ行動をとらなければならない。

エージェントが任意の状態において選択することができる行動は、壁のない方向に進むのみとする。すなわち、行動集合 $A = \{\text{東, 西, 南, 北}\}$ の空集合でない部分集合とする。

3.2 パラメータ設定と実験結果

エージェントは、ゴールにたどり着くと環境から報酬 $r = 100$ を受け取り、初期状態であるスタートに再配置されるものとした。エージェントの行動1ステップごとに負の報酬を与えることや、ゴールにたどり着くのかかったステップ数に応じて報酬を変えることなど、早くゴールにたどり着く学習を助長するような報酬の与え方はしない。

試行数 100 を3つのフェーズに分け、それぞれについてパラメータの設定を行った。ステップごとに変化するパラメータについては表1のように設定し、それ以外のパラメータについては各フェーズにおいて共通とし $\beta = 0.9 \exp(j\pi/6)$, $\gamma = 0.999$, $T = 3000$, $N_e = 1$ とした。ただし、 $j^2 = -1$ である。

計算機実験の結果を図3に示す。この結果は100試行を1学習として100学習行い、収束したものに關しての平均である。

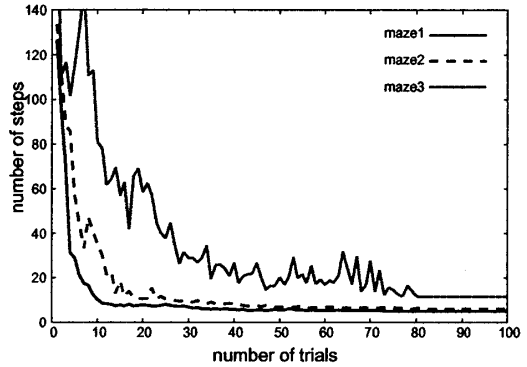


図3 実験結果

4 考察

maze1, maze2 においては100%が収束し、maze3 においては95%が収束した。本稿では maze1, maze2 において100%が最短経路を実現するような方策を獲得した。

maze1 では獲得した一連の行動について、ある Q 値は次の行動に対応する Q 値と β の偏角だけずれるような位相をもつよう学習がされた。

maze2 においては maze1 のように一連の行動に対する Q 値は単純な関係にはならなかった。しかし、 Q 値の位相を自律的に調整することで、不完全知覚問題を解決している様子が観察された。

maze3 では、最短経路を実現するような方策を学習しなかったが、環境中を一部往復する行動をとることで自律的に環境を多重化して不完全知覚問題に対処している様子が観察された。

5 おわりに

強化学習の価値関数に複素数を用いる手法を提案した。本稿では、その具体例として Q -learning の Q 値を複素数とする Q -learning を検討しその有効性を示した。今後は、同じ行動に対して複数の Q 値を割り当てることを行うなど、さまざまな観点からさらなる性能の向上を目指す。

参考文献

- [1] Kaelbling and Littman and Moore. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, pp. 237-285, 1996.
- [2] T.Hamagami, S.Koakutsu, and H.Hirata. Reinforcement learning to compensate for perceptual aliasing using dynamic additional parameter: Motivational value. In *Proc. of IEEE Int' conf. on SYSTEMS (SMC2002)*.
- [3] 山城, 上野, 武田. 遅れ報酬に基づく遺伝的アルゴリズムによる部分観測マルコフ決定問題の解決手法. *信学論誌 (D-I)*, Vol. J84-D-I, No. 12, pp. 1635-1647, 2001.