

迷路探索問題Frozen Lakeにおける最適な行動選択ポリシー及び学習エージェントの提案

1.はじめに

本実験では,迷路探索問題Frozen Lakeにおける最適な行動選択ポリシー及び学習エージェントの決定を行った.性能評価には未学習の状態から1500エピソードの学習を1シミュレーションとし,30回のシミュレーションを行ったときに得られたデータを使用した.また,性能評価の指標として以下の2つを用いた.

-**第一指標** : 最後の100エピソードにおける平均獲得報酬.ここで,計算方法として,30回のシミュレーションを行うため,30個の平均獲得報酬が得られるが,さらにその平均を取ることにする.また,最低基準を0.30とする.

-**第二指標** : 平均獲得報酬が0.4以上になるまでのエピソード数.

今回の実験では,特にこの第一指標に焦点を当て,最後の100エピソードにおける平均獲得報酬がより大きくなるような学習方法を模索した.第一指標の値の大きい学習では,必然的に第二指標の値も小さくなると考えたため,第二指標については最後のまとめで確認することとする.また,今回想定しうる学習エージェントとして,以下の学習則をもつようなものを考える.

s:状態 a:行動 r:報酬 t:時刻 Q:行動価値関数 V:状態価値関数 α :学習係数 γ :割引率

$$\text{Qlearner} \quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \{r_{t+1} + \gamma \max_a Q(s_{t+1}, a') - Q(s_t, a_t)\} \quad (1)$$

$$\text{SARSAlearner} \quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \{r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)\} \quad (2)$$

$$\text{ActorCritic} \quad \text{TD error} : \delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \quad (3)$$

$$\text{Critic} : V(s_t) \leftarrow V(s_t) + \alpha \delta_t \quad (4)$$

$$\text{Actor} : Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \delta_t \quad (5)$$

さらに,今回想定しうる行動選択ポリシーとして,以下のようなものを考える.

①**greedy法**...常に行動価値関数Qが最大となるような行動を選択する.

② **ϵ -greedy法**... ϵ の確率でランダムな行動をとり,それ以外は行動価値関数Qが最大となるような行動を選択する.

③**softmax法**...行動価値関数を $Q(s,a)$ としたとき,状態sとしたときの行動確率 $P(a_i|s)$ を以下のように定める手法.

$$P(a_i|s) = \frac{e^{\frac{Q(s,a_i)}{T}}}{\sum_{j=1}^n e^{\frac{Q(s,a_j)}{T}}} \quad (n: \text{状態数}, T: \text{正の定数}) \quad (6)$$

以上の条件のもと,実験1で行動選択ポリシー及び学習エージェントの決定,実験2,3で最適なパラメータの決定を行った.

2.実験1 : 行動選択ポリシー及び学習エージェントの決定

【実験の目的】 行動選択ポリシー及び学習エージェントの決定

【実験の手法】 行動選択ポリシー①**greedy法**② **ϵ -greedy法**③**softmax法**に対して,それぞれ学習エージェントQ,SARSA,ActorCriticにおける第一指標の平均獲得報酬の値を調べた.学習エージェントのパラメータとして学習係数 α ,割引率 γ ,行動価値関数の初期値 p_0 ,状態価値関数の初期値 v_0 があるが,私は特に α, γ が重要なのではないかと考えたため,ひとまず $p_0 = 0, v_0 = 0$ で固定し, $\alpha(0 < \alpha < 1), \gamma(0 < \gamma < 1)$ 対して0.1区切りで数値を与えて,横軸 α 縦軸 γ として平均獲得報酬の値をヒートマップを用いて図示を行う.なお,② **ϵ -greedy法**では行動選択ポリシーのパラメータ $\epsilon(0 < \epsilon \leq 1)$ があるので,0.2区切りで数値を与え,一つの学習エージェントにつき5個のヒートマップを作製した.同様に,③**softmax法**についても行動選択ポリシーのパラメータT(正の定数)があるので,T=0.25,0.50,1.0,2.0,4.0で数値を与え,一つの学習エージェントにつき5個のヒートマップを作成した.ただし,②③に関しては紙面の都合上,すべてのヒートマップを示すことはできないので,各ヒートマップの最大値を表にまとめ,一番大きい値をとるパラメータのヒートマップのみを図として示すこととする.

【実験結果】

①**greedy法**...Q学習,SARSA学習,ActorCriticの結果を,それぞれ図1,2,3に示す.greedy法ではどの学習エージェントにおいても第一指標の最低基準(0.30)を下回っており,良い結果が得られなかった.また,パラメータの値による法則性も全く見受けられなかった.

② **ϵ -greedy法**...パラメータ ϵ に対する,Q学習,SARSA学習,ActorCriticの各ヒートマップにおける最大平均獲得報酬を表1に示す.表より, $\epsilon=0.2$ の時,最大平均獲得報酬が一番大きいものとなったので, $\epsilon=0.2$ の時のQ学習,SARSA学習,ActorCriticのヒートマップをそれぞれ図4,5,6に示す.結果として全ての場合において,第一指標の最低基準を下回っていたが,共通して γ を大きくしたり, ϵ を小さくすることで平均獲得報酬が大きくなることや,Q学習及びSARSA学習では α を小さくしたり,ActorCriticでは α を中間あたりの値にすると平均獲得報酬が大きくなるといったパラメータの値による法則性が見られた.

③**softmax法**...パラメータTに対する,Q学習,SARSA学習,ActorCriticの各ヒートマップにおける最大平均獲得報酬を表2に示す.表より,T=0.25の時,最大平均獲得報酬が一番大きいものとなったので,T=0.25の時のQ学習,SARSA学習,ActorCriticのヒートマップをそれぞれ図7,8,9に示す.結果としてActorCritic(T=0.25,0.5,1.0)のとき最低基準を上回るものが現れた.

また,全てに共通して γ を大きくしたり,Tを小さくすることで平均獲得報酬が大きくなることや,ActorCriticでは α を中間あたりの値にすると平均獲得報酬が大きくなるといったパラメータの値による法則性が見られた.

⇒今回は,第一指標の最低基準を上回ることがあるActorCritic+softmax法を採用する.

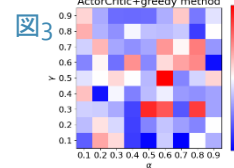
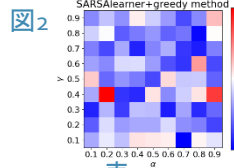
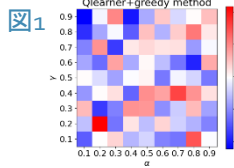


表1

平均獲得報酬 (最大)	$\epsilon=0.2$	$\epsilon=0.4$	$\epsilon=0.6$	$\epsilon=0.8$	$\epsilon=1.0$
Qlearner	0.18($\alpha=0.1, \gamma=0.9$)	0.08($\alpha=0.1, \gamma=0.9$)	0.04($\alpha=0.1, \gamma=0.9$)	0.03($\alpha=0.3, \gamma=0.9$)	0.02($\alpha=0.1, \gamma=0.8$)
SARSAlearner	0.14($\alpha=0.1, \gamma=0.9$)	0.07($\alpha=0.1, \gamma=0.9$)	0.04($\alpha=0.2, \gamma=0.9$)	0.02($\alpha=0.2, \gamma=0.8$)	0.02($\alpha=0.5, \gamma=0.9$)
ActorCritic	0.18($\alpha=0.4, \gamma=0.9$)	0.09($\alpha=0.7, \gamma=0.9$)	0.05($\alpha=0.3, \gamma=0.8$)	0.03($\alpha=0.4, \gamma=0.8$)	0.02($\alpha=0.8, \gamma=0.8$)

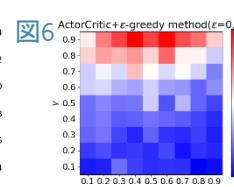
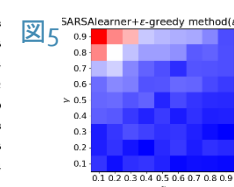
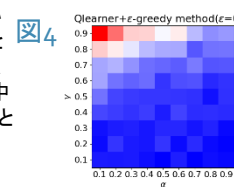
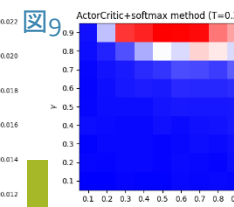
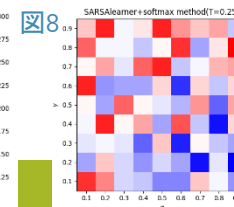
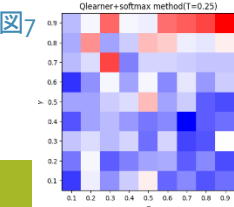


表2

平均獲得報酬 (最大)	T=0.25	T=0.5	T=1.0	T=2.0	T=4.0
Qlearner	0.03($\alpha=0.9, \gamma=0.9$)	0.02($\alpha=0.6, \gamma=0.9$)	0.02($\alpha=0.8, \gamma=0.9$)	0.02($\alpha=0.2, \gamma=0.7$)	0.02($\alpha=0.1, \gamma=0.7$)
SARSAlearner	0.02($\alpha=0.2, \gamma=0.9$)	0.02($\alpha=0.6, \gamma=0.7$)	0.02($\alpha=0.1, \gamma=0.5$)	0.02($\alpha=0.5, \gamma=0.7$)	0.02($\alpha=0.7, \gamma=0.9$)
ActorCritic	0.62($\alpha=0.5, \gamma=0.9$)	0.59($\alpha=0.6, \gamma=0.9$)	0.38($\alpha=0.8, \gamma=0.9$)	0.11($\alpha=0.9, \gamma=0.9$)	0.04($\alpha=0.9, \gamma=0.9$)



3.実験2 : 最適なパラメータの決定(alpha, gamma, T)

【実験の目的】 ActorCritic+softmax法におけるパラメータ (α, γ, T) の決定

【実験の手法】 実験2では,先ほど採用したActorCritic+softmax法における最適なパラメータを一部決定する.先ほどと同様にひとまず $p_0 = 0, v_0 = 0$ とし,それ以外のパラメータ α, γ, T を決定する.まず, γ に関しては,実験1の図9のヒートマップの傾向から γ を大きくすると平均獲得報酬が大きくなることがわかったので, $\gamma=1.0$ で確定させた.残り, α とTについて考えればよいので, α とTの値を変えて第一指標の平均獲得報酬を調べれば良いということとなる.ということで, $\alpha(0 < \alpha < 1)$ は0.1区切りで数値を与え,Tは実験1の表2の結果から

ActorCriticでは,1以下で平均獲得報酬が最低基準を上回ることがわかったので,Tの範囲を0<T<1に絞り,0.1区切りで数値を与え,第一指標の平均獲得報酬を調べた.その結果を横軸 α 縦軸平均獲得報酬scoreとして,折れ線グラフにて示す.折れ線グラフでは,平均獲得報酬の値と共に標準偏差も表示し,さらに図10ではT=0.1~0.3,図11ではT=0.4~0.6,図12ではT=0.7~0.9のグラフを同時に表示することとする.

【実験結果】 前述の図10~12を以下に示す. α については中間あたりの値をとると平均獲得報酬が大きくなることがわかった.また,Tについては α が中間あたりの値を取る場合,平均獲得報酬は多少の違いはあれど,ほとんど違いが無いことがわかった.

⇒今回は, $\gamma=1.0$,Tは α の値によって平均獲得報酬の値の変動が比較的小さいT=0.4, α はT=0.4の中で標準偏差の小さい $\alpha=0.4$ を採用する.

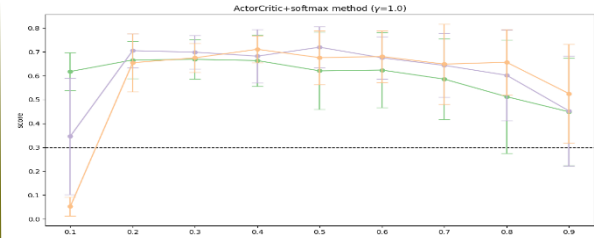


図10

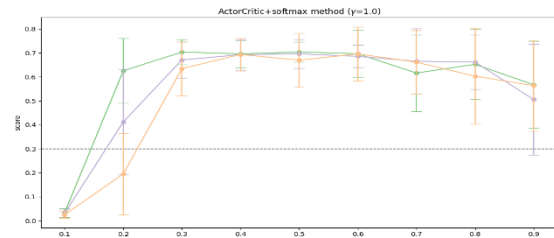


図11

4.実験3 : 最適なパラメータの決定(q0, v0)

【実験の目的】 ActorCritic+softmax法におけるパラメータ (q_0, v_0) の決定

【実験の手法】 実験3では,ActorCritic+softmax法における最適なパラメータの最終決定を行う.先ほど最適なパラメータとして $\gamma=1.0, T=0.4, \alpha=0.4$ を決定したので,最後にこのパラメータの下で最適なパラメータ q_0, v_0 の決定を行う. q_0, v_0 にそれぞれ-2.0,-1.0,0.0,+1.0,+2.0の数値を与えて,第一指標の平均獲得報酬の値を調べた.その結果を横軸 v_0 縦軸平均獲得報酬scoreとして,折れ線グラフにて示す.折れ線グラフでは,平均獲得報酬の値と共に標準偏差も表示し,さらに $q_0=-2.0, -1.0, 0.0, +1.0, +2.0$ の時のグラフを同時に表示することとする.

【実験結果】 実験結果を図13に示す. v_0 が負だと極端に報酬が小さくなり,非負ならばほとんど報酬の違いが生じなかった.また, q_0 の値による違いは多少標準偏差の違いはあれどほとんど見受けられなかった.⇒今回は, v_0 が一番平均獲得報酬の高くなった $v_0=1.0$, q_0 が一番標準偏差の小さくなった $q_0=2.0$ を採用する.

5.考察

これまでの実験結果をもとに,考察を行っていく.大きく①**行動選択ポリシー**及び**学習エージェント**②**パラメータ**に分けて考察を行う.

①**行動選択ポリシー及び学習エージェント**...最適なものとしてActorCritic+softmax法になった理由について考察をしてみる.

まず,学習エージェントActorCriticとなった理由として,学習則のTD誤差(3式)によるものだと考える.Frozen Lakeではスリップという必ずしも自分の思った行動がとれるとは限らない要素があるため,行動を含まない学習則の方が良いのではないかと考えた.したがってそのような学習則を唯一もつActorCriticで平均獲得報酬が高くなり,行動を含む学習則をもつQ学習及びSARSA学習では平均獲得報酬がそこまで高くならなかったのではないかと考える.

また,行動選択ポリシーがsoftmax法となった理由として,学習エージェントと行動選択ポリシーの相性によるものだと考える.実験1の結果から,同じ学習エージェントをもつ場合でも,行動選択ポリシーによって平均獲得報酬が大きく異なり,Q学習とSARSA学習は ϵ -greedy法,ActorCriticはsoftmax法と組み合わせると報酬が高くなった.したがって,学習エージェントとしてActorCriticが選ばれたことで,相性の良いsoftmax法が選ばれたのではないかと考える.

②**パラメータ**...最適なパラメータとして $\alpha=0.4, \gamma=1.0, T=0.4, v_0=1.0, q_0=2.0$ を採用したが,この値になった理由について考察してみる. α →学習の速度を表すパラメータであり, $\alpha=0.4$ と α が中間あたりの値を取るとき平均獲得報酬が高くなった.これはActorCriticの適正学習係数が中間あたりの値であるからだと考えられる.図6,9から α を中間あたりに設定すると高い報酬を得られることがわかるので,ActorCriticでは学習速度を速くもなく遅くもなくすることで良い結果が得られると考えられる.ちなみに図3を考慮しなかった理由としては,図1,2よりgreedy法ではそもそも学習が上手くいっていないと判断したからである.

γ →目先の報酬をどれくらい重視するかを表すパラメータであり,小さければ小さいほど目先の報酬を重視してしまう.したがって, $\gamma=1.0$ と大きな値にし,将来の報酬を最大限に重視させることにより,高い平均獲得報酬が得られたのだと考える.T→softmax法で用いられるパラメータであり,Tが大きくなるほどrandomに近い行動選択をしてしまう.したがって,T=0.4と比較的小さい値にし,行動確率をより行動価値関数Qの値に依存させたことにより,高い平均獲得報酬が得られたのだと考える.

q_0, v_0 →実験3では q_0 による報酬の変化はほぼなかったで, v_0 についてのみ考察を行う.実験3では v_0 の値が負になると極端に獲得報酬が低くなることが,これはFrozenLakeでは報酬が0~1であるので,初期値を負にすることで,学習がうまく進まなかったことが原因だと考えられる.したがって, $v_0=1.0$ と非負にすることで学習がうまく進み高い平均獲得報酬が得られたのだと考える.

6.まとめ

以上より,私はActorCritic+softmax法($\alpha=0.4, \gamma=1.0, T=0.4, v_0=1.0, q_0=2.0$)を提案する.最後にこのときの第一指標と第二指標がどうなったか確認を行う.第一指標は,最後の100エピソードにおける平均獲得報酬30シミュレーション分の箱ひげ図である図14,第二指標は,各エピソードで30シミュレーション分の獲得報酬の平均をとったものの推移を横軸episode縦軸平均獲得報酬として表した図15にて確認を行う.

まず,図14より,第一指標の平均獲得報酬が0.702となり,一番悪いシミュレーションで最小値0.63,一番良いシミュレーションで最大値0.78をとることがわかった.

次に図15では,最低基準である0.30を赤線,第二指標の基準である0.40を青線で表示しているが,約200エピソードを超えたあたりで報酬0.40を超えることがわかった.

図14

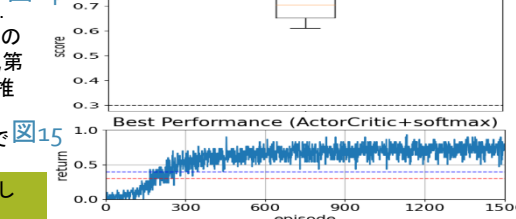


図15

