

接続の変異の場合は、要素ごとにそれぞれ P_{code} の確率で対立遺伝子に変異させた。

3.3. 選択・淘汰

選択・淘汰は親個体群 P と子個体群 R を合わせた全個体群に対して行った。全個体群に対して選択・淘汰を行う利点は、交叉や突然変異によって生まれた子個体の適応率が低い場合の解の悪化を防ぐことができるからである。また、選択・淘汰の基準である適応度はタスクの成功率とし、適応度が一番高い個体をエリート個体として選択し、その後、エリート個体を除く全個体群に対してルーレット選択を行い、局所解に陥ることを防いだ。

4. 実験タスク

時間情報が必要なタスクとして、エージェントがスタート位置から正しいルートを通りゴールまで到達することを成功とするタスクを行った。まず、図2に示すマップを用意し、エージェントをスタート位置に配置する。エージェントは、[上, 下, 右, 左]の4行動をとることができ、学習過程において正しいルートを通った場合は各 step ごとに報酬を与え、それ以外のルートを通った場合は、罰を与えエピソードを終了した。また、ニューラルネットワークの入力としてエージェントに与えられる情報は、周囲8マスの状態であり、障害物がある場合は1を、それ以外は0をそれぞれ入力値として与えた。そのため、例えば、図2の状態1, 2, 3のときに、それぞれエージェントに与えられる情報は全く同じであるため、正しい過去の情報を保持していない限り、状態1, 2, 3で間違った行動をとり、ルートから外れてしまいタスクを達成できないように設定した。

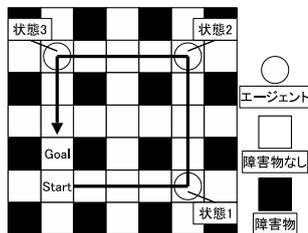


図2: 時間情報を必要とするタスク

5. 実験結果

比較対象として、過去の入力情報を与える Time Dependent Neural Network (TDNN), 過去の中間層の情報を与える Elman 型 Neural Network (ENN), Multi-Context Recurrent Neural Network (MCRNN), 過去の出力情報を与える Jordan 型 Neural Network (JNN) の4種類とした。それぞれのネットワークに対して中間層ニューロン数を {5, 10, 20, 30, 40}, 過去層数を {1, 2, 3, 4, 5, 10} とし合計 90 パターンの構造を用意した。TDNN と JNN はそれぞれ 30 パターン。ENN は過去中間層を1層利用するものなので、5パターン。MCRNN は過去中間層を複数用意するネットワークであり、過去中間層数は2~10層であるため、25パターン。これが90パターンの内訳である。1試行ごとのエピソード数を10000回とし、100エピソードごとにテストを行い、タスク成功率が100%になったら試行を終了し、各パターンをそれぞれ100試行を行った。提案手法のパラメータを表1に、実験結果を表2に示す。表2の比較対象の結果は、各構造で一番良い結果を示したニューロン数、過去層数の組み合わせの結果である。また、表3は、

提案手法によって獲得された個体において、タスクを成功した構造 (94 個体) の過去層の使用割合である。

表1: 実験条件

パラメータ	値
世代数	100
エピソード回数/1世代	100
親個体群 P	30
子個体群 R	40
P_h, M_h	0.25, 5
P_p, M_p	0.25, 2
P_b, M_b	0.25, 2
P_d, M_d	0.25, 2
P_{code}	0.01
報酬	$0.9 \times \text{step 数} / \text{Goal までの step 数}$
罰	-0.9

表2: タスク成功回数とタスク成功率の平均

構造	タスク成功回数	タスク成功率 [%]
提案手法	94	99.1
TDNN	0	32.7
ENN	0	28.2
MCRNN	14	62.2
JNN	61	91.0

表3: 成功した構造の過去層の使用割合

	過去入力層	過去中間層	過去出力層
割合 (個体数) [%]	13(12)	68(64)	85(80)

タスクの成功回数と成功率から任意に設定した一般的な RNN では獲得できなかった知識を同じ学習条件において、提案手法では獲得できている点から、タスクに適した構造を獲得できたといえる。また、JNN のタスクの成功回数が他比較対象と比べて良い結果を示していることからわかる通り、今回のタスクは過去の行動 (出力) が重要な情報であり、提案手法では、過去出力層の使用割合が高いことから、適切な構造を獲得できているといえる。

6. まとめ

本研究では、GAに基づき、DVB 強化学習を用いて、RNN に時系列学習を要するタスクを自律的に学習させることを目的とし、時間情報を必要とするタスクに対する知識の獲得、獲得された構造の検証を行った。実験結果より、リカレントニューラルネットワークを用いることによって、時系列による環境の変化を記憶し、遺伝的アルゴリズムによって、タスクに適したリカレントニューラルネットワーク構造を自律的に獲得、Direct-Vision-Based 強化学習によって自律的なタスクの学習を行うことを示した。

参考文献

- [1] 柴田 克成, 岡部 洋一, 伊藤 広司, “ニューラルネットワークを用いた Direct-Vision-Based 強化学習-センサからモータまで-,” 計測自動制御学会論文集, Vol.37, No.2, pp168-177, 2001.
- [2] M.Delgado, M.P.Cuellar, M.C.Pegalajar, “Multi-objective Hybrid Optimization and Training of Recurrent Neural Networks,” *IEEE Transactions on Systems, Man, and Cybernetics-Part B: Cybernetics*, vol.38, No.2, pp.381-403, 2008.