SUNAによるシミュレーション上での二足歩行動作学習

Learning to Bipedal Walking on Simulation with SUNA

井上 湧太 ダニロ・ヴァスコンセロス・ヴァルガス Yuta Inoue Danilo Vasconcellos Vargas

九州大学

Kyushu University

SUNA is currently one of the most adaptive neuroevolution methods which is able to tackle different problems efficiently. However, many questions remain unanswered. In this research, we applied SUNA to the bipedal-walking problem and evaluate it general learning properties. The results show that even without any modificiations SUNA is able to learn in this environment. Moreover, contrary to many other methods, it is continuously improving its average rewards showing a near open-ended learning.

1. はじめに

Spectrum-diverse Unified Neuroevolution Architecture (SUNA)[Vargas 17] は、異なる問題を同一の手法で効率的に 学習するためのニューロエボリューションの手法である。しか しながらその性質についての多くは調べられておらず、汎用性 や計算速度での優位性などの性質は明らかになっていない。本 研究では SUNA の汎用性について調べるため、二足歩行シミュ レーションを用いて SUNA による学習を行った。

SUNA は Unified Neural Model (UNM) と呼ばれる特殊な 構造を用いて汎用性を実現しており、また効率性は UNM の各 ニューロンの特徴にしたがって探索を行うことで実現されてい る。UNM では、同じニューロエボリューションの手法である NEAT[Stanley 02] で使用されるニューロンに改良を加えた拡 張ニューロンが使用される。このニューロンでは NEAT で使 用されるニューロンの持つ特徴に加え、それぞれ異なる時定数 を持ち、異なる活性化関数を選択できるという特徴がある。異 なる活性化関数を選択できることは、単一の活性化関数を持つ ニューロンのみで構成されるネットワークと比較して、ある小 さなネットワークを 1 つのニューロンで表すことができる可能 性があるため、構造を簡単にすることに貢献する。

また UNM では、多様な問題に対応できるように、状況に応じ てネットワークの構造を変えるための制御ニューロンとニュー ロモジュレーションを導入している。制御ニューロンは入力に 応じて、制御対象のニューロンごとに重み付けされた制御信号 を出力する。各ニューロンは受け取った制御信号の総和が閾値 以上であればそれぞれのニューロンの入力に従って出力を行い、 そうでなければ出力を行わない。まだニューロモジュレーショ ンはあるニューロンの内部状態に応じて結合強度を変化させる ことでネットワークの構造を変化させる。

以上で上げたように UNM は高い自由度を持つため、その 探索空間は膨大なものになる。これを効率的に探索するため、 Danilo & Murata[Vargas 17] ではネットワークの含むニュー ロンの特徴によって個体を種別し、それぞれの種の中で最も高 い適応度を持つものを変異させたものについて探索を行った。 その世代のすべての個体を Novelty Map を用いて含まれる個



図 1 OpenAI Gym[Brockman 16] の二足歩行シミュレー ション (BipedalWalker-v2)

体同士の距離が最も大きくなるように選抜し、選抜された個体 を中心として下位個体群とする。それぞれの下位個体群のうち 最も適応度の高い個体と、それらを変異させた個体を次の世代 として評価を行う。これを行うことで探索の方向を抽出した特 徴量に応じて分散することができるため、効率的な探索が可能 になると考えられる。

2. 方法

学習環境には OpenAI Gym の 2 次元二足歩行シミュレー ション環境である BipedalWalker-v2(図1)を用いた*¹。入力 は、すべての入力を使用できる場合と、使用できる入力が制限 される場合に分けて実験を行った。すべての入力が使用できる 場合、使用可能な入力は胴体(図1中の紫色の多角形)の進行 方向に対する速度、上下方向に対する速度、角度、角速度、両足 の股関節、膝関節の角度、角速度、両足それぞれが地面と接し ているかどうかなどの24つの入力を与えた。また入力は両足そ れぞれの股関節、膝関節の角度、それぞれの足が地面に付いて いるかどうかの計6つとした。出力は両足の股関節、膝関節の 角速度の計4つであった。報酬はロボットの前進に応じて与え られた。また胴体の傾きや出力の大きさによって負の報酬が与 えられた。到達地点まで達した場合+300、ロボットの胴体(図 1中の紫色の多角形)が地面(図1中の下部の緑色の部分)に 接触した場合-100の報酬を与え、試行を終了した。前述した条 件で終了しなかった場合は、1 試行で最大 300 ステップを行っ た。各個体に対して5回の試行を行い、試行ごとの報酬の平 均を適応度とした。SUNA のパラメータは Danilo & Murata

Contact: 井上 湧太、九州大学、福岡県福岡市西区周船 寺1丁目9番地 19パークホームズ周船寺 901号、 yi1306c12@gmail.com

^{*1} https://gym.openai.com/envs/BipedalWalker-v2



図2 各世代で最も高かった適応度の、30回の平均



図3 すべての入力を使用可能にした場合の各世代で最も高かった適応度の、実験30回のそれぞれの遷移

(2017)[Vargas 17]と同じものを使用した。1世代あたりの個体数は100とし、1000世代分の学習を行った。

3. 結果

実験を 30 回行った場合の、それぞれの学習において各世代で の最も高い適応度の平均を図 2 に示す。また以下に 30 回行っ た実験それぞれでの、各世代での最も高い適応度の世代ごとの 遷移を図 3 と図 4 に示す。

4. 考察

表1に、本実験の結果と同一の環境を用いて実験を行った [Zhang 17] の結果を示す。ただしSUNAの結果は10⁷ステッ プの結果と合わせるため、100世代目前後のものを使用してい る。ステップ数は正確に対応していないため直接の比較はでき ないが、SUNAの結果はNEATと同程度であった。

入力数が24の場合と6の場合では、24の場合のほうが良い 結果を出している(図2)。入力を限定した場合では実験ごとの



図 4 入力を制限した場合に各世代で最も高かった適応度の、 実験 30 回のそれぞれの遷移

表 1 SUNA とその他の手法の結果の比較。Shangtong & Osmar(2017) [Zhang 17] の結果を使用した。ただしすべて の値はグラフ(Shangtong & Osmar(2017), Figure. 2) から読み取った値である。

methods	SUNA	NEAT	NES	CMA	P3O
rewards	10	10	50	-100	300

適応度が世代数の増加とともに安定して上がらず(図4)、入力 数が多いほうが学習はスムーズであった(図3)。SUNAでは突 然変異の際、ニューロン間の結合強度を変化させるかどうか決 めてから、変化させるのであればそれぞれの結合強度を一様分 布で得られた値を用いて変化させるという手段を取る。この方 法では、これまでの進化で得られた有用なネットワークを保持 しながら新しい構造を探索するようなことは難しいため、予め 必要な情報が入力として用意されている場合のほうが学習に有 利であると考える。

参考文献

- [Brockman 16] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W.: OpenAI Gym, *CoRR*, Vol. abs/1606.01540, (2016)
- [Stanley 02] Stanley, K. O. and Miikkulainen, R.: Evolving Neural Networks Through Augmenting Topologies, *Evol. Comput.*, Vol. 10, No. 2, pp. 99–127 (2002)
- [Vargas 17] Vargas, D. V. and Murata, J.: Spectrum-Diverse Neuroevolution With Unified Neural Models, *IEEE Transactions on Neural Networks and Learning* Systems, pp. 1–15 (2017)
- [Zhang 17] Zhang, S. and Zaiane, O. R.: Comparing Deep Reinforcement Learning and Evolutionary Methods in Continuous Control, arXiv (2017)