進化的計算と方策勾配法による学習を用いた 3次元制御タスクにおけるマルチタスク深層強化学習

Multi-task Deep Reinforcement Learning with Evolutionary Algorithm and Policy Gradient Method in 3D Control Tasks

今井 翔太	清 雄一	田原 康之	大須賀 昭彦
Shota IMAI	Yuichi SEI	Yasuyuki TAHARA	Akihiko OHSUGA

電気通信大学大学院情報理工学研究科

Graduate School of Informatics and Engineering. The University of Electro-Communications

In deep rerinforcement learning, it is difficult to converge when the exploration is insufficient or a reward is sparce. Besides, in a specific tasks, the number of exploration may be limited. Therefore, it is considered effective to learn in source tasks previously to promote learning in the target tasks. In this research, we propose a method to train a model that can work well on variety of target tasks with evolutionary algorithm and policy gradient method. In this method, agents explore multiple environments with diverce set of neural networks to train a general model with evolutionary algorithm and policy gradient methid. In the experiments, we assume multiple 3D control source tasks. After the model training with our method in the source tasks, we shows how effective the model is for the 3D Control tasks of the target tasks.

1. はじめに

近年,深層学習と強化学習を組み合わせた深層強化学習が ゲーム AI[Mnih 15] や,実世界における機械制御 [Levine 16] 等の分野で大きな成果をあげている.深層強化学習は,深層 ニューラルネットワークを関数近似器として用い,Q 値や方 策の出力を学習する手法である.

深層強化学習においては、学習のために膨大な探索によるサ ンプルを必要とする.しかし、探索する環境の行動・状態空間 が大きすぎる場合は、十分なサンプルを得るのに大きな時間が かかり、得られたサンプルによっては収束しないこともある. 実世界において実機を用いた探索を行う場合には、物理的な条 件による制約から、複数の探索自体が困難である.また、学習 が完了していない方策で探索を行うことで、機器が危険な行動 を行う可能性がある.この問題を解決するため、解きたいタス ク(ターゲットタスク)とは別のタスク(ソースタスク)で事 前に学習を行うことでニューラルネットワークの転移可能なパ ラメータを獲得し、解きたいタスクにおいて少数のサンプルで 高速に学習が可能な汎用的学習モデルを作成することが望ま しい.

深層強化学習で学習を行いたいタスク(ターゲットタスク) があるとき,最終的に解きたいタスクと性質が似ており,なん らかの理由(単にタスクが単純,シミュレータで学習可能)で 学習が容易なタスク(ソースタスク)が他に存在するのであれ ば、ソースタスクで学習を行うことで,効率的に両タスクに共 通するパラメータを獲得できる可能性がある.また,ソース タスクが複数ある場合,それら全てのタスクに対して高いパ フォーマンスを発揮するようなパラメータの学習を行うこと で、ソースタスク,ターゲットタスクの広い範囲のタスクで共 通する良いパラメータを学習できると考えられる.

従来の研究では、勾配降下を用いた学習を工夫するなどして、 複数のタスクで高速な学習が可能なパラメータを得る事前学習 手法が提案されてきた [Finn 17].しかし、事前学習を行うタス クで報酬がスパースである場合の学習の難しさ [Sutton 98],強 化学習の一般的な課題であるサンプルの偏り [Henderson 17], 初期パラメータへの依存 [Kolen 90] に対応した事前学習手法 の検討は不十分である.

これらの背景を踏まえ、本研究では、進化的計算のランダ ム要素と、勾配降下による学習を組み合わせ、上記の課題を 解決しつつ事前学習可能な手法を提案する.本研究では、解 くべきタスクを3次元制御問題と定め、行動空間が連続であ る場合に用いられる深層強化学習アルゴリズムである Deep Deterministic Policy Gradient(DDPG)[Lillicrap 15]を用い、 ニューラルネットワークに事前学習手法を適用する.

2. 関連研究

本研究では、行動空間が連続な場合に使用される深層強化学 習アルゴリズムとして, DDPG[Lillicrap 15]を用いる. DDPG は,状態の入力に対して,決定論的に行動を出力する方策を パラメータで表し, 期待報酬和が最大になるような勾配を用 いてパラメータを更新する深層強化学習アルゴリズムである. 行動空間が連続な場合における Q 学習 [Watkins 92] では, -般的に特定の状態において最も Q 値が高い行動を求めること が困難である.一方 DDPG は入力に対して,それぞれの行動 の出力で、決定論的に一つの値を出力するため、主に行動空間 が連続である場合のタスクで使用される. DDPG のアーキテ クチャは, 行動の数と同数の出力を持ち, 観測した状態を入力 とすることで行動の値を出力する Actor と、観測した状態と Actor の出力を入力とし、観測した状態における Actor の出 力を評価する Critic で構成される. Critic は一般的な教師あ り学習によって学習を行い、Actor は Critic の出力を使用し、 決定論的な方策を学習する方策勾配定理 Deterministic Policy Gradient(DPG)[Silver 14] によってパラメータの更新を行う. パラメータ θ^Q を持つ Critic Q は, リプレイバッファからの サンプル (s_i, a_i, r_i, s_{i+1}) を用いて、以下の損失関数を最小化 するようにパラメータを更新する.

$$L = \frac{1}{N} \Sigma_i (y_i - Q(s_i, a_i | \theta^Q))^2 \tag{1}$$

ここで y_i は, パラメータ $\theta^{Q'}$ を持つ Critic のターゲットネットワーク Q', パラメータ $\theta^{\pi'}$ を持つ Actor のターゲットネッ

連絡先: 今井翔太,電気通信大学大学院情報理 工学研究科,東京都調布市調布ヶ丘 1-5-1, Email:imai.shota@ohsuga.lab.uec.ac.jp

トワーク,割引報酬率 γ を用いて,以下の式で与えられる.

$$y_i = r_i + \gamma Q'(s_{i+1}, \pi'(s_{i+1}|\theta^{\pi'})|\theta^{Q'})$$
(2)

Actor のパラメータは上記の Critic の出力を使用し,以下の 式で与えられる勾配によって更新する.

$$\nabla_{\theta^{\pi}} J \approx \frac{1}{N} \Sigma_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=a_i} \nabla_{\theta^{\pi}} \pi(s | \theta^{\pi}) |_{s=s_i} \quad (3)$$

深層強化学習において複数タスクに適用可能な事前学習手法 としては MAML(Model-Agnostic Meta-Learning)[Finn 17] がある.これは、画像認識や強化学習などの深層学習のタスク において、汎用的に使用できる学習モデルのパラメータの初 期値を勾配降下法による学習で得る手法である.この手法は、 学習の仕方そのものを学習するメタラーニングと呼ばれる手法 の一種である.

深層強化学習における効果的な学習として、進化的計算と勾配 降下法による学習を組み合わせた Evolutionary Reinforcement Learning(ERL)[Khadka 18] が提案されている. ERL では、進 化的計算のランダム要素によって学習するニューラルネット ワークの個体群の中に、勾配降下で学習するニューラルネット ワークを混ぜ、定期的に勾配降下で学習するネットワークのパ ラメータを個体群の一部にコピーすることで、一つのタスクに 対する最適化を効率化する.進化的計算では、複数のニューラ ルネットを用いて学習を行うため、学習が安定しているものが 最適化され、また、パラメータの初期値に対して頑健であると いう特性がある.これに加え、複数の個体による探索を行うた め、報酬がスパースなタスクであっても報酬の獲得がしやすい という利点がある.ERL は強化学習における課題を進化的計 算を用いて解決しつつ、一般的な勾配降下による学習を取り入 れることで、タスク特化の学習を効率化したものといえる.

3次元における制御課題のように、行動、状態空間が膨大 である場合は、報酬がスパースであり、事前学習に必要な探 索自体が難しいという問題がある [Sutton 98]. また、単一の ニューラルネットワークを用いた学習は、ネットワークの初期 パラメータによって学習が左右される傾向がある [Kolen 90]. その他,単一のニューラルネットワークを用いた探索では,得 られるサンプルに偏りが発生し、局所解に陥るなどして学習が うまくいかない場合がある [Henderson 17]. これらは,一般 的な事前学習手法における「転移可能なパラメータを得られる かどうか」という問題設定からは独立した課題である.そのた め, 仮に転移可能なパラメータを得る手段があったとしても, タスクによっては、これらの課題により事前学習が阻害される 可能性がある.これらの課題は上述の ERL のように進化計算 的な手法を用いることで解決が可能である.したがって, ERL のように進化計算を行い、ニューラルネットワークの最適化の 対象を複数のタスクに設定することにより、課題を克服しつつ 様々なタスクで利用可能なパラメータを取得することができる と考えられる.本研究では、進化的計算による最適化と、ソー スタスクにおける複数タスクに対する勾配降下の学習を事前 学習として行うことで、3D 制御タスクに共通するニューラル ネットの良いパラメータを獲得し、ターゲットタスクとなる複 数のタスクで高速に学習することを目指す.

3. 提案手法

3.1 問題設定

本手法の問題設定は以下の通りである.



図 1: 進化的計算と方策勾配法を用いた深層強化学習

タスクの質が同じタスクの集合を仮定する.タスクの集合 のうち,モデルの事前学習に用いられるものをソースタスクと 呼ぶ.タスクの集合のうち,事前学習の時点では与えられず, 学習後に解くべき未知のタスクをターゲットタスクと呼ぶ.す なわち,本研究では,ターゲットタスクと似たタスクのソース タスクで事前学習を行うことで,ニューラルネットがタスク間 に共通する良いパラメータを獲得し,複数のターゲットタスク において,高速に学習可能なニューラルネットワークを得るこ とを目指す.

3.2 提案手法の詳細

3.2.1 個体集団による探索

図1に提案手法の概要を示す.最初に,進化的計算の個体 集合となる複数のニューラルネットワークを用意する.本研究 では,状態を入力とし,行動を出力するニューラルネットワー ク(以後 Actor)をこの個体集団のニューラルネットワー クとする.個体集合 pop_{π} の Actor π はそれぞれランダムなパラ メータ θ^{π} で初期化される.ソースタスクの集合を Tとし,各 Actor π は各世代で全てのタスクを探索し得られた報酬 r_{π} を 記録する.また,この探索で得られたサンプル (s_i, a_i, r_i, s_{i+1}) は各タスクのリプレイバッファに記録する.複数の Actorの 探索により,偏りのないサンプルを得ることが期待でき,これ らのサンプルは後の DDPG による Actor の学習に使用する.

3.2.2 適応度によるエリート選択,学習

全ての Actor π による探索終了後,記録された報酬の合計 を元に各 Actor π の適応度 f_{π} を計算し,適応度が最も高い Actor をエリート個体として選択する.ここで選択されたエ リート個体は、後述のノイズ付加をパスする.選択されたエ リート個体の Actor のコピーは、勾配降下による学習で追加 の最適化を行うため、確率的に選ばれたタスクのリプレイバッ ファR のサンプルを用いて DDPG による学習を行う.ここで、 あるタスクのリプレイバッファ R_{T_i} が選択される確率 $P_{R_{T_i}}$ は、 その世代の各 Actor が各ソースタスク T_i で得た報酬 r_{T_i} を元 に以下の式で与えられる.

$$P_{R_{T_i}} = \frac{r_{T_i}}{\sum_i r_{T_i}} \tag{4}$$

3.2.3 適応度による個体の選択、ノイズの付与と最終選択

エリート個体の選別後,次世代の個体集合となる Actor を, 適応度を元にルーレット選択によって選択する.ここで選ばれ たネットワークは,確率的にノイズを付与することにより,突 然変異を行う.エリート個体,コピーされ勾配降下による学習 を行ったエリート個体,エリートを除いて選択された個体が 次世代の個体集合 *pop*[']_π となり,現世代 *pop*_π にコピーされる. Algorithm 1 提案手法の疑似コード

- 1: Initialize actor π and critic Q with weight θ^{π} and θ^{Q} , respectively
- 2: Initialize a population of k actors pop_{π}
- 3: Initialize replay buffers ${\cal R}$
- 4: Define a random number generator $r() \in [0, 1)$
- 5: for generation = 1, ∞ do
- 6: for actor $\pi \in pop_{\pi}$ do
- 7: for all source tasks T_i do
- 8: Explore T_i using θ^{π}
- 9: Append transition to replay buffer R respectively
- 10: end for
- 11: end for
- 12: Select the elite actor π based on fitness score f_{π}
- 13: Select the replay buffer R based on all fitness scores f_{π}
- 14: Sample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from R
- 15: Update Q by minimizing the loss
- 16: $L = \frac{1}{N} \sum_{i} (y_i Q(s_i, a_i | \theta^Q))^2$
- 17: Update copied elite actor π using the sampled policy gradient
- 18: $\nabla_{\theta^{\pi}} J \approx \frac{1}{N} \Sigma_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=a_i} \nabla_{\theta^{\pi}} \pi(s | \theta^{\pi}) |_{s=s_i}$
- 19: Select the (k-2) actors based on fitness scores f_{π} and insert selected actors into next generation's population pop'_{π}
- 20: for (k-2) actors $\in pop'_{\pi}$ do
- 21: **if** $r() < mut_{prob}$ **then**
- 22: Add noise to θ^{π}
- 23: end if
- 24: end for
- 25: Insert the elite actor into pop'_{π}
- 26: Insert the copied elite actor pop'_{π}
- 27: $pop_{\pi} \leftarrow pop'_{\pi}$
- 28: end for

以降,上記の手順を最終世代まで繰り返し,最終世代で得られ たエリート個体の Actor が最終的に得られるニューラルネッ トワークである.

4. 実験

実験では、物理演算エンジン Pybullet [Coumans 19] で提 供される OpenAI Gym[Brockman 16] の 3 次元物体制御タス クを用いて、事前学習手法の有効性の評価を行う.

4.1 実験設定

PyBullet で提供されるタスクの中から, Minitaur, Hopper, Walker をソースタスクとし, HalfCheetah, Ant, Humanoid をターゲットタスクとする(図2).各タスクはそれぞれ,物体 の部位の角度や回転,速度などを状態とし,関節など,物体の 各部位の調節を出力とする.出力を調節することで,物体が適 切に行動(立ち上がる,進む)すると報酬が与えられる.ソー スタスクにおいて,提案手法を用いた複数世代の探索,選択, ノイズ付与,学習を行い,最終世代におけるエリート個体とな る1つのActorを全てのターゲットタスクに用いる学習モデ ルとする.各世代におけるActorの個体数は10とし,1つの 環境における最大探索数は1000回,ノイズ付与のフェーズで



図 2: 実験で用いるソースタスクとターゲットタスク

は、エリートを除いた次世代個体にそれぞれ 1/4 の確率でラン ダムに発生させた平均 0 のガウシアンノイズを付与する.パラ メータの過度な変動を防ぐため、このノイズの中で-0.5 以下、 0.5 以上の値は 0 として扱う.本実験における事前学習の世代 数は 100 とする.Actor として用いるニューラルネットワー クの設定は、入力と出力はタスクごとに可変とし、隠れ層は 3 層、各層のノード数は 128 とする.Critic として用いるニュー ラルネットワークの設定は、状態を入力する側のネットワーク をノード数 200 の隠れ層 1 層とし、Actor の入力と、状態を入 力した側からの出力を受け取るニューラルネットワークは 300 ノードの隠れ層 1 層とし、出力は 1 である.活性化関数として は ReLU を使用し、最適化手法として Adam[Kingma 14] を 用いた.

4.2 実験結果と考察

図3に各ターゲットタスクにおける,提案手法によって得 られた事前学習モデルと、ランダムに初期化したモデルの学習 過程を示す. 横軸はエピソード数, 縦軸は各エピソードで得ら れた報酬を指す. Ant, HalfCheetah においては, 事前学習モ デルを使用した場合のエピソードの初期の段階で獲得する報 酬量が、ランダムに初期化したモデルの報酬量を上回っている ことが確認できる、両タスク共にエピソードによって得られる 獲得報酬量に大きな幅があるものの,全体的に事前学習を行っ たモデルの方が大きな報酬を得ており,少ないサンプルで高速 に学習できているといえる.一方, Humanoid においては事 前学習モデルが、ランダムな初期の場合よりも高いパフォーマ ンスを発揮できていないことが確認できる.これについては, Humanoid タスクは、他のタスクと異なる多くの部位の観測、 行動が要求されるため、ソースタスクでそれらと紐づいた良い パラメータが学習できなかったこと,また今回実験で使用した ニューラルネットワークの表現力が足りていないことなどが原 因として考えられる.

5. まとめ

本稿では、深層強化学習における既存の事前学習手法にお いて解決されていない課題に対応する手法として、進化的計算 と勾配降下を用い、特に3次元物体制御タスクにおいて複数 のタスクで高速学習可能なパラメータを得る事前学習手法を 提案した.実験では、物演算エンジン PyBullet で提供される Open AI Gym の3次元制御タスクに提案手法を適用し、複 数タスクで高速学習可能なパラメータを持ったニューラルネッ トが得られることを示した.



図 3: ターゲットタスクにおける事前学習モデルとランダムに 初期化したモデルの学習

今後は、さらに探索空間が大きく、報酬がスパースである課 題に対して本手法を適用し、有効性を検証したい.

6. 謝辞

本研究は,JSPS 科研費 16K00419,16K12411,17H04705, 18H03229,18H03340の助成を受けたものです.

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)
- [Coumans 19] Coumans, E. and Bai, Y.: PyBullet, a Python module for physics simulation for games, robotics and machine learning, http://pybullet.org (2016-2019)
- [Finn 17] Finn, C., Abbeel, P., and Levine, S.: Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks, in Precup, D. and Teh, Y. W. eds., Proceedings of the 34th International Conference on Machine Learning, Vol. 70 of Proceedings of Machine Learning Research, pp. 1126–1135, International Convention Centre, Sydney, Australia (2017), PMLR
- [Henderson 17] Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D.: Deep Reinforcement Learning that Matters, *CoRR*, Vol. abs/1709.06560, (2017)
- [Khadka 18] Khadka, S. and Tumer, K.: Evolution-Guided Policy Gradient in Reinforcement Learning, in Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. eds., Advances in Neural Information Processing Systems 31, pp. 1196–1208, Curran Associates, Inc. (2018)
- [Kingma 14] Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, *CoRR*, Vol. abs/1412.6980, (2014)
- [Kolen 90] Kolen, J. F. and Pollack, J. B.: Back Propagation is Sensitive to Initial Conditions, in *Proceedings of*

the 1990 Conference on Advances in Neural Information Processing Systems 3, NIPS-3, pp. 860–867, San Francisco, CA, USA (1990), Morgan Kaufmann Publishers Inc.

- [Levine 16] Levine, S., Finn, C., Darrell, T., and Abbeel, P.: End-to-End Training of Deep Visuomotor Policies, *Journal of Machine Learning Research*, Vol. 17, No. 39, pp. 1–40 (2016)
- [Lillicrap 15] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D.: Continuous control with deep reinforcement learning, *CoRR*, Vol. abs/1509.02971, (2015)
- [Mnih 15] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D.: Human-level control through deep reinforcement learning, *Nature*, Vol. 518, No. 7540, pp. 529–533 (2015)
- [Silver 14] Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M.: Deterministic Policy Gradient Algorithms, in Xing, E. P. and Jebara, T. eds., Proceedings of the 31st International Conference on Machine Learning, Vol. 32 of Proceedings of Machine Learning Research, pp. 387–395, Bejing, China (2014), PMLR
- [Sutton 98] Sutton, R. S. and Barto, A. G.: Introduction to Reinforcement Learning, MIT Press, Cambridge, MA, USA, 1st edition (1998)
- [Watkins 92] Watkins, C. J. C. H. and Dayan, P.: Qlearning, *Machine Learning*, Vol. 8, No. 3, pp. 279–292 (1992)