# A framework for reinforcement learning with autocorrelated actions

Marcin Szulc[1], Jakub Łyskawa[1], and Paweł Wawrzyński[1][0000−0002−1154−0470]

Warsaw University of Technology, Institute of Computer Science, Warsaw, Poland
{marcin.szulc,jakub.lyskawa}.stud@pw.edu.pl, pawel.wawrzynski@pw.edu.pl

**Abstract.** The subject of this paper is reinforcement learning. Policies are considered here that produce actions based on states and random elements autocorrelated in subsequent time instants. Consequently, an agent learns from experiments that are distributed over time and potentially give better clues to policy improvement. Also, physical implementation of such policies, e.g. in robotics, is less problematic, as it avoids making robots shake. This is in opposition to most RL algorithms which add white noise to control causing unwanted shaking of the robots. An algorithm is introduced here that approximately optimizes the aforementioned policy. Its efficiency is verified for four simulated learning control problems (Ant, HalfCheetah, Hopper, and Walker2D) against three other methods (PPO, SAC, ACER). The algorithm outperforms others in three of these problems.

**Keywords:** Reinforcement learning · Actor-Critic · Experience replay · Fine time discretization.

## 1 Introduction

The usual goal of Reinforcement Learning (RL) to optimize a policy that samples an action on the basis of a current state of a learning agent. The only stochastic dependence between subsequent actions is through state transition: The action moves the agent to another state which determines the distribution of another action. Main analytical tools in RL are based on this lack of other dependence between actions. E.g., for a given policy, its value function expresses the expected sum of discounted rewards the agent may expect starting from a given state. The sum of rewards does not depend on actions taken before the given state was reached. Hence, only the given state and the policy matter.

Lack of dependence between actions beyond state transition leads to several difficulties. In physical implementation of RL, e.g. in robotics, it usually means that white noise is added to control actions. However, that makes control discontinuous and rapidly changing all the time. This is often impossible to implement since electric motors that are to execute these actions can not operate this way. Even if it is possible, it requires a lot of energy, makes the controlled system shake, and exposes it to damages.

It is also questionable if the lack of dependence between actions beyond state transition does not reduce efficiency of learning. Each action is an experiment that leads to policy improvement. However, due to limited accuracy of (action-)value function approximation, consequences of a single action may be difficult to recognize. The finer the time discretization, the more serious this problem becomes. Consequences of a random experiment distributed over several time instants could be more tangible thus easier to recognize.

The contribution of this paper may be summarized in the following points:

– A framework is introduced in which a policy produces actions on the basis of states and values of a stochastic process. That enables relation between actions that is beyond state transition.
– An algorithm is introduced that approximately optimizes the aforementioned policy.
– The above algorithm is tested on four benchmark learning control problems: Ant, Half-Cheetah, Hopper, and Walker2D.

The rest of the paper is organized as follows. Section 2 overviews related literature. Sec. 3 introduces a policy that produces autocorrelated actions along with tools for its analysis. Sec. 4 introduces an algorithm that approximately optimizes that policy. Sec. 5 presents simulations that compare the presented algorithm with state-of-the-art reinforcement learning methods. The last section concludes the paper.

## 2   Related Work

### 2.1   Stochastic dependence between actions

The idea of introducing stochastic dependence between actions was analyzed in [16] as a remedy to problems with application of RL in fine time discretization. The control process was divided there into "non-Markov periods" in which actions were stochastically dependent. A policy with autocorrelated actions was analyzed in [18] with a standard RL algorithm applied to its optimization that did not account for the dependence of actions.

In [5] a policy was analyzed whose parameters were incremented by the autoregressive stochastic process. Essentially, this resulted in autocorrelated random components of actions. In [8] a policy was analyzed that produced an action being a sum of the autoregressive noise and a deterministic function of state. However, no learning algorithm was presented in this paper that accounted for specific properties of this policy.

### 2.2   Reinforcement learning with experience replay

The Actor-Critic architecture of reinforcement learning was introduced in [1]. Approximators were applied to this structure for the first time in [7]. In order to boost efficiency of these algorithms, they were combined with experience replay for the first time in [17].

Application of experience replay to Actor-Critic encounters the following problem. The learning algorithm needs to estimate quality of a given policy on the basis of consequences of actions that were registered when a different policy was in use. Importance sampling estimators are designed to do that, but they can have arbitrarily large variance. In [17] that problem was addressed with truncating density ratios present in those estimators. In [15] specific correction terms were introduced for that purpose.

Another approach to the aforementioned problem is to prevent the algorithm from inducing a policy that differs too much from the one tried. That idea was first applied in Conservative Policy Iteration [6]. It was further extended in Trust Region Policy Optimization [12]. This algorithm optimizes a policy with the constraint that the Kullback-Leibler divergence between that policy and the tried one should not exceed a given threshold. The K-L divergence becomes an additive penalty in Proximal Policy Optimization algorithms, namely PPO-Penalty and PPO-Clip [13].

A way to avoid the problem of estimating quality of a given policy on the basis of the tried one is to approximate the action-value function instead of estimating the value function. Algorithms based on this approach are Deep Q-Network (DQN) [11], Deep Deterministic Policy Gradient (DDPG) [10], and Soft Actor-Critic (SAC) [4]. In the original version of DDPG the time-correlated OU noise was added to action. However, this algorithm was not adapted to this fact in any specific way. SAC uses white noise in actions and it is considered one of the most efficient in this family of algorithms.

## 3  Policy with autocorrelated actions

Let an action, $a_t$, be computed as

$$a_t = \pi(s_t, \xi_t; \theta) \tag{1}$$

where $\pi$ is a deterministic transformation, $s_t$ is a current state, $\theta$ is a vector of trained parameters, and $(\xi_t, t = 1, 2, \dots)$ is a stochastic process. We require this process to have the following properties:

- Stationarity: The distribution of $\xi_t$ is the same for each $t$.
- Zero mean: $E\xi_t = 0$ for each $t$.
- Autocorrelation decreasing with growing lag:

$$E\xi_t^T \xi_{t+k} > E\xi_t^T \xi_{t+k+1} \geq 0 \text{ for } k \geq 0. \tag{2}$$

  Essentially that means that values of the process are close to each other when they are in close time instants.
- Markov property: For any $t$ and $k, l \geq 0$, the conditional distributions

$$(\xi_t, \dots, \xi_{t+k} | \xi_{t-1}, \dots, \xi_{t-1-l}) \quad \text{and} \quad (\xi_t, \dots, \xi_{t+k} | \xi_{t-1}) \tag{3}$$

  are the same. In words, dependence of future values of $(\xi_t)$ on its past is entirely carried over by $\xi_{t-1}$.

Consequently, if only $\pi$ (1) is continuous for all its arguments, and subsequent states $s_t$ are close to each other, then the corresponding actions are close, although random. In words, they create a consistent, distributed in time experiment that can lead to policy improvement.

*Example: Auto-Regressive* $(\xi_t)$. Let $\alpha \in [0, 1)$ and

$$\epsilon_t \sim N(0, C), \quad t = 1, 2, \ldots$$
$$\xi_1 = \epsilon_1$$
$$\xi_t = \alpha \xi_{t-1} + \sqrt{1 - \alpha^2} \epsilon_t, \quad t = 2, 3, \ldots \tag{4}$$

Fig. 1 demonstrates a realization of both the white noise $(\epsilon_t)$ and $(\xi_t)$. Let us analyze if $(\xi_t)$ has the required properties.

Both $\epsilon_t$ and $\xi_t$ have the same distribution $N(0, C)$. Therefore $(\xi_t)$ is stationary and zero-mean. A simple derivation reveals that



Fig. 1: A realization of the normal white noise $(\epsilon_t)$, and the auto-regressive process $(\xi_t)$ (4).

$$E\xi_t\xi_{t+k}^T = \alpha^{|k|}C \quad \text{and} \quad E\xi_t^T\xi_{t+k} = \alpha^{|k|}\text{tr}(C)$$

for any $t, k$. Therefore, $(\xi_t)$ is autocorrelated, and this autocorrelation decreases with growing lag. Consequently, the values of $\xi_t$ are closer to one another for subsequent $t$ than the values of $\epsilon_t$, namely

$$E\|\epsilon_t - \epsilon_{t-1}\|^2 = E(\epsilon_t - \epsilon_{t-1})^T(\epsilon_t - \epsilon_{t-1}) = 2\text{tr}(C)$$
$$E\|\xi_t - \xi_{t-1}\|^2 = E\left((\alpha-1)\xi_{t-1} + \sqrt{1-\alpha^2}\epsilon_t\right)^T\left((\alpha-1)\xi_{t-1} + \sqrt{1-\alpha^2}\epsilon_t\right)$$
$$= (\alpha - 1)^2\text{tr}(C) + (1 - \alpha^2)\text{tr}(C) = (1 - \alpha)2\text{tr}(C).$$

The Markov property of $(\xi_t)$ directly results from how $\xi_t$ (4) is computed.

In fact, marginal distributions of the process $(\xi_t)$, as well as its conditional marginal distributions are normal, and their parameters have compact forms. We shall not present derivations of these parameters due to lack of space, but we shall denote them for further use. Namely, let as consider

$$\bar{\xi}_t^n = [\xi_t^T, \ldots, \xi_{t+n-1}^T]^T. \tag{5}$$

The distribution of $\bar{\xi}_t^n$ is normal

$$N(0, \Omega_0^n), \tag{6}$$

where $\Omega_0^n$ is a matrix dependent on $n$, $\alpha$, and $C$. The conditional distribution $(\bar{\xi}_t^n|\xi_{t-1})$ is also normal,

$$N(B^n\xi_{t-1}, \Omega_1^n), \tag{7}$$

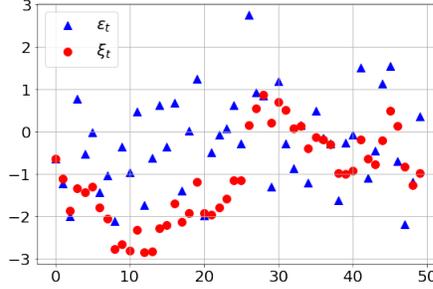where both $B^n$ and $\Omega_1^n$ are matrices dependent on $n$, $\alpha$, and $C$.

*The neural-normal policy.* A simple and practical way to implement $\pi$ (1) is as follows. A feedforward neural network,

$$A(s; \theta), \tag{8}$$

has input $s$ and weights $\theta$. An action is computed as

$$a_t = \pi(s_t, \xi_t; \theta) = A(s_t; \theta) + \xi_t, \tag{9}$$

for $\xi_t$ in the form (4). While the discussion below can be extended to the general formulation (1), in order to make it simpler we will further assume that a policy is of the form (9).

Let us consider

$$\bar{s}_t^n = [s_t^T, \ldots, s_{t+n-1}^T]^T,$$
$$\bar{a}_i^n = [a_t^T, \ldots, a_{t+n-1}^T]^T,$$
$$\bar{A}(\bar{s}_i^n; \theta) = [A(s_t; \theta)^T, \ldots, A(s_{t+n-1}; \theta)^T]^T,$$

and fixed $\theta$. With (9) the distributions $(\bar{a}_t^n | \bar{s}_t^n)$ and $(\bar{a}_t^n | \bar{s}_t^n, \xi_{t-1})$ are both normal, namely $N(\bar{A}(\bar{s}_t^n; \theta), \Omega_0^n)$, and $N(\bar{A}(\bar{s}_t^n; \theta) + B^n \xi_{t-1}, \Omega_1^n)$, respectively (see (6) and (7)). The algorithm defined in the next section updates $\theta$ to manipulate the above distributions. Density of the normal distribution with mean $\mu$ and covariance matrix $\Omega$ will be denoted by

$$\varphi(\cdot; \mu, \Omega). \tag{10}$$

*Noise-value function.* In policy (1) there is a stochastic dependence between actions beyond the dependence resulting from state transition. Therefore, the traditional understanding of policy as distribution of actions conditioned on state does not hold here. Each action depends on the current state, but also previous states and actions. Analytical usefulness of the traditional value function and action-value function is thus limited.

As a valid analytical tool we propose *noise-value function* defined as

$$W^\pi(\xi, s) = E_\pi \left( \sum_{i \geq 0} \gamma^i r_{t+i} \Big| \xi_{t-1} = \xi, s_t = s \right). \tag{11}$$

The course of events starting in time $t$ depends on the current state $s_t$ and the value $\xi_{t-1}$. Because of Markov property of $\xi_t$ (3), the pair $(\xi_{t-1}, s_t)$ is a proper condition for the expected value of future rewards.

The *value function* $V^\pi : \mathcal{S} \mapsto \mathbb{R}$ is slightly redefined, namely

$$V^\pi(s) = E\big(W(\xi_{t-1}, s_t) | s_t = s\big). \tag{12}$$

The random value in the above expectation is $\xi_{t-1}$ and its distribution is conditional with the condition $s_t = s$. The distribution of $\xi_{t-1}$ may differ for different $s_t$. However, being in the state $s_t$ and not knowing $\xi_{t-1}$ the agent may expect the sum of future rewards equal to $V^\pi(s_t)$.

## 4  ACERAC: Actor-Critic with Experience Replay and Autocorrelated aCtions

The algorithm presented here has Actor-Critic structure. It optimizes a policy of the form (9) and uses Critic,

$$V(s; \nu),$$

which is an approximator of the value function (12) parametrized by a vector, $\nu$.

For each time instant of the agent-environment interaction the policy (9) is applied and a tuple, $\langle s_t, A_t, a_t, r_t, s_{t+1} \rangle$, is registered, where $A_t = A(s_t; \theta)$.

The general goal of training is to maximize $W^\pi(\xi_{i-1}, s_i)$ for each state $s_i$ registered during the agent-environment interaction. In this order previous time instants are sampled, and sequences of actions that follow these instants are made more/less probable depending on their return. More specifically, $i$ is sampled from $\{1, \ldots, t-1\}$ and the conditional density of the sequence of actions $(a_i, \ldots, a_{i+n-1})$ is being increased/decreased depending on the return

$$r_i + \cdots + \gamma^{n-1} r_{i+n-1} + \gamma^n V(s_{i+n}; \nu)$$

this sequence of actions yields. At the same time adjustments of the same form are performed for several sequences of actions starting from $a_i$, namely for $n = 1, \ldots, \tau$, where $\tau \in \mathbb{N}$ is a parameter.

### 4.1  Actor & Critic training

The following procedure is repeated several times at each $t$-th instant of agent–environment interaction:

1. A random $i$ is sampled from the uniform distribution over $\{1, \ldots, t-1\}$.
2. If $i$ is the initial instant of a trial, then consider for $n = 1, \ldots, \tau$

$$\mu_{i+j} = E(\xi_{i+j}) = 0, \ j = 0, \ldots, n-1$$
$$\eta_{i+j} = E(\xi_{i+j}) = 0, \ j = 0, \ldots, n-1$$
$$\Omega_2^n = \Omega_0^n.$$

Otherwise, consider

$$\mu_{i+j} = E(\xi_{i+j} | \xi_{i-1} = a_{i-1} - A_{i-1}), \ j = 0, \ldots, n-1$$
$$\eta_{i+j} = E(\xi_{i+j} | \xi_{i-1} = a_{i-1} - A(s_{i-1}; \theta)), \ j = 0, \ldots, n-1$$
$$\Omega_2^n = \Omega_1^n.$$

3. Consider the following vectors for $n = 1, \ldots, \tau$

$$\bar{\mu}_i^n = [\mu_i^T, \ldots, \mu_{i+n-1}^T]^T,$$
$$\bar{\eta}_i^n = [\eta_i^T, \ldots, \eta_{i+n-1}^T]^T,$$
$$\bar{s}_i^n = [s_i^T, \ldots, s_{i+n-1}^T]^T,$$
$$\bar{a}_i^n = [a_i^T, \ldots, a_{i+n-1}^T]^T,$$
$$\bar{A}_i^n = [A_i^T, \ldots, A_{i+n-1}^T]^T,$$
$$\bar{A}(\bar{s}_i^n; \theta) = [A(s_i; \theta)^T, \ldots, A(s_{i+n-1}; \theta)^T]^T.$$

4. Temporal differences are computed for $n = 1, \ldots, \tau$

$$
\begin{aligned}
d_i^n(\theta, \nu) = \big( r_i + \cdots + \gamma^{n-1} r_{i+n-1} + \gamma^n V(s_{i+n}; \nu) - V(s_i; \nu) \big) \times \\
\times \psi_b \left( \frac{\varphi(\bar{a}_i^n; \bar{A}(\bar{s}_i^n; \theta) + \bar{\eta}_i^n, \Omega_2^n)}{\varphi(\bar{a}_i^n; \bar{A}_i^n + \bar{\mu}_i^n, \Omega_2^n)} \right),
\end{aligned}
\tag{13}
$$

where $\psi_b$ is a soft-truncating function, e.g. $\psi_b(x) = b \tanh(x/b)$, for a certain $b > 1$.

5. Actor and Critic are updated. The improvement directions for Actor and Critic are

$$
\Delta\theta = \frac{1}{\tau} \sum_{n=1}^{\tau} \nabla_\theta \ln \varphi(\bar{a}_i^n; \bar{A}(\bar{s}_i^n; \theta) + \bar{\eta}_i^n, \Omega_2^n) d_i^n(\theta, \nu) - \nabla_\theta L(s_i, \theta)
\tag{14}
$$

$$
\Delta\nu = \frac{1}{\tau} \sum_{n=1}^{\tau} \nabla_\nu V(s_i; \nu) d_i^n(\theta, \nu),
\tag{15}
$$

where $L(s, \theta)$ is a loss function that penalizes Actor for producing actions that do not satisfy conditions e.g., they exceed their boundaries. $\Delta\theta$ is designed do increase/decrease the likelihood of the sequence of actions $\bar{a}_i^n$ proportionally to $d_i^n(\theta, \nu)$. $\Delta\nu$ is designed to make $V(\cdot; \nu)$ approximate the value function (12) better. The improvement directions $\Delta\theta$ and $\Delta\nu$ are applied to update $\theta$ and $\nu$, respectively, with the use of either ADAM, SGD, or other method of stochastic optimization.

In Point 1 the algorithm selects an experienced event to replay. In Points 2 and 3 it determines the parameters the distribution of the sequence of subsequent actions, $\bar{a}_i^n$. In Point 4 it determines the relative quality of $\bar{a}_i^n$. The temporal difference (13) implements two ideas. Firstly, $\theta$ is changing due to being optimized, thus the conditional distribution $(\bar{a}_i^n | \xi_{i-1})$ is now different than it was at the time when the actions $\bar{a}_i^n$ were happening. The density ratio in (13) accounts for this discrepancy of distributions. Secondly, in order to limit variance of the density ratio, the soft-truncating function $\psi_b$ is applied. In Point 5 the parameters of Actor, $\theta$, and Critic, $\nu$, are being updated.

## 5 Empirical study

This section presents simulations whose purpose has been to compare the algorithm introduced in Sec. 4 to state-of-the-art reinforcement learning methods. We compared the new algorithm (ACERAC) to ACER [17], SAC [4] and PPO [13]. We used the rllib implementation [9] of SAC and PPO in the simulations. Our implementation of ACERAC is available at `https://github.com/mszulc913/acerac`.

We used four control tasks, namely Ant, Hopper, HalfCheetah, and Walker2D (see Fig. 2) from PyBullet physics simulator [2] to compare the algorithms.

A simulator that is more popular in the RL community is MuJoCo [14].[1] Hyperparameters that assure optimal performance of ACER, SAC, and PPO applied to the considered environments in MuJoCo are well known. However, PyBullet environments introduce several changes to MuJoCo tasks, which make them more realistic, thus more difficult. Additionally, physics in MuJoCo and PyBullets slightly differ [3], hence we needed to tune the hyperparameters. Their value can be found in appendix A.

For each learning algorithm we used Actor and Critic structures as described in [4]. That is, both structures had the form of neural networks with two hidden layers of 256 units each.
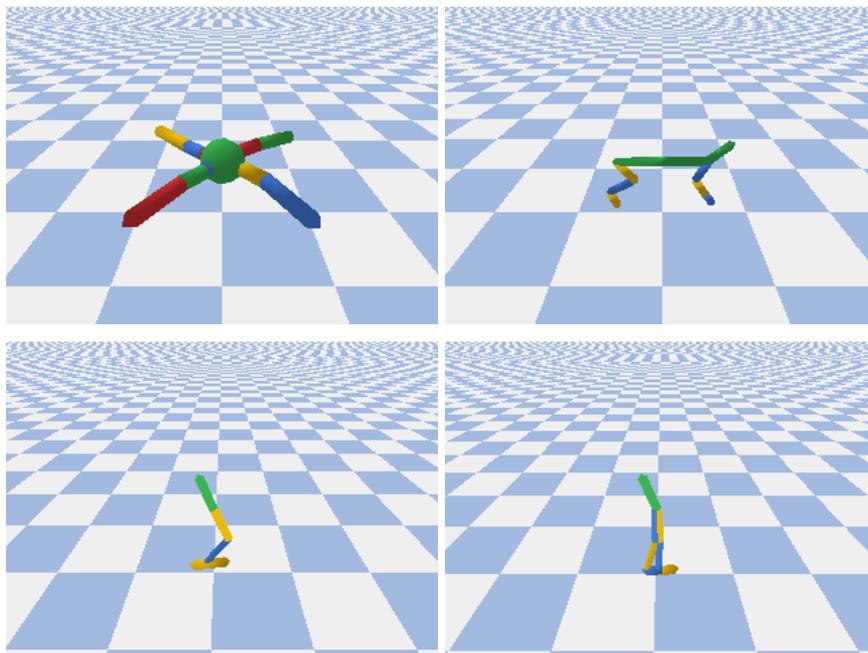


Fig. 2: Environments used in simulations: Ant (left upper), HalfCheetah (right upper), Hopper (left lower), Walker2D (right lower).

### 5.1   Experimental setting

Each learning run lasted for 3 million timesteps. Every 30000 timesteps of a simulation was made with frozen weights and without exploration for 5 test episodes. An average sum of rewards within a test episode was registered. Each run was repeated 5 times.

---

[1] We chose PyBullet because it is a freeware, while MuJoCo is a commercial software.
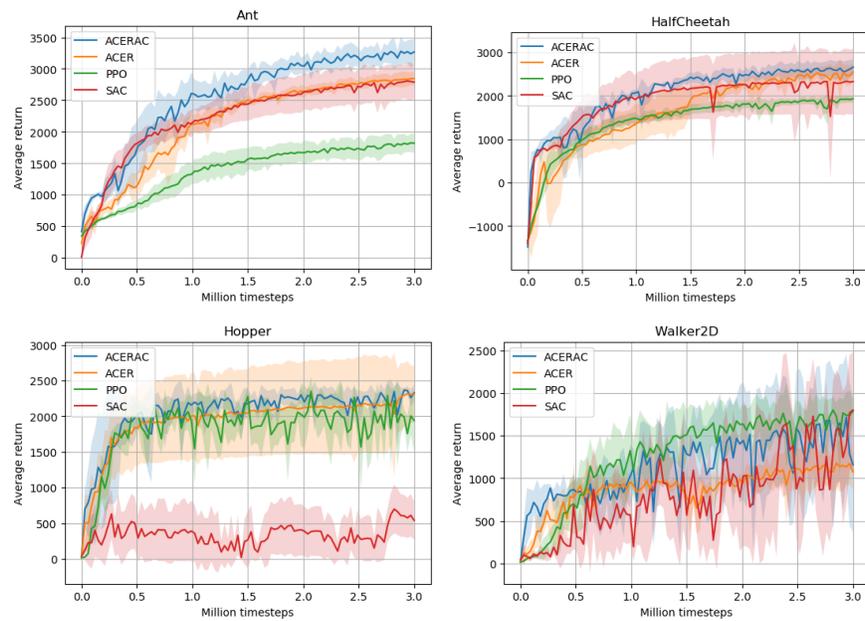
Fig. 3: Learning curves for Ant (left upper), HalfCheetah (right upper), Hopper (left lower) and Walker2D (right lower) environments: Average sums of rewards in test trials.

## 5.2   Results

Figure 3 presents learning curves for all four environments and all four compared algorithms. Each graph reports how a sum of rewards in test episodes evolves within learning. Solid lines represent the average sums of rewards and shaded areas represent their standard deviations.

It is seen that for Ant the algorithm that achieve the best performance is ACERAC, then ACER and SAC, then PPO. For HalfCHeetah, the best performance is achieved by ACERAC which is slightly better than ACER, then SAC, then PPO. For Hopper the algorithms to win are ACERAC *ex aequo* with ACER, then PPO, then SAC; actually SAC fails in this task. Eventually, for Walker2D, PPO achieves the best performance, then ACERAC and SAC, and then ACER.

## 5.3   Discussion

It is seen in Fig. 3 that ACERAC is the best performing algorithm for three environments out of four (in one ACER preforms equally well). ACERAC extends ACER in two directions. Firstly, it admits autocrrelated actions. This enables exploration distributed in many actions instead in one. Secondly, in order to mimic learning with eligibility traces [7], ACER estimates improvement directions with the use of a sum whose limit is random. This increases variance of these estimates. Instead, for each state ACERAC computes an improvement direction as an average of increments similar to those ACER selects on random. Hence smaller variance of improvement direction estimates in ACERAC which enables larger step-sizes and faster learning.

It is important to note that the algorithm introduced here, ACERAC, has been designed for fine time discretization and real life control problems. However, here it has been tested on simulated benchmarks in which time discretization was not particularly fine and control could be arbitrarily discontinuous. Its relatively good performance is a desirable result. It allows to expect that this algorithm will perform relatively even better in real life control problems. That remains to be confirmed experimentally in further studies.

## 6   Conclusions and future work

In this paper a framework was introduced to apply reinforcement learning to policies that admit stochastic dependence between subsequent actions beyond state transition. This dependence is a tool to enable reinforcement learning in physical systems and fine time discretization. It can also yield better exploration thus faster learning.

An algorithm based on this framework, Actor-Critic with Experience Replay and Autocorrelated aCtions (ACERAC), was introduced. Its efficiency was verified in simulations of four learning control problems: Ant, HalfCheetah, Hopper, and Walker2D. The algorithm was compared with PPO, SAC, and ACER.

ACERAC outperformed the competitors in Ant and HalfCheetah. For Hopper ACERAC was the best *ex aequo* with ACER. For Walker2D the best results was obtained by PPO.

It is desirable to combine the framework proposed here with adaptation of dispersion of actions by introducing reward for the entropy of their distribution, as it is done in PPO. The framework proposed here is specially designed for applications in robotics. An obvious step of our further research is to apply it in this field, obviously much more demanding than simulations.

## Acknowledgement

## References

1. Barto, A.G., Sutton, R.S., Anderson, C.W.: Neuronlike adaptive elements that can learn difficult learning control problems. IEEE Transactions on Systems, Man, and Cybernetics B **13**, 834–846 (1983)
2. Coumans, E., Bai, Y.: Pybullet, a python module for physics simulation for games, robotics and machine learning. `http://pybullet.org` (2016–2019)
3. Erez, T., Tassa, Y., Todorov, E.: Simulation tools for model-based robotics: Comparison of bullet, havok, mujoco, ode and physx. In: 2015 IEEE International Conference on Robotics and Automation (ICRA). pp. 4397–4404 (2015)
4. Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Offpolicy maximum entropy deep reinforcement learning with a stochastic actor (2018), arXiv:1801.01290
5. van Hoof, H., Tanneberg, D., Peters, J.: Generalized exploration in policy search. Machine Learning **106**, 1705–1724 (2017)
6. Kakade, S., Langford, J.: Approximately optimal approximate reinforcement learning. In: Proceedings of the Nineteenth International Conference on Machine Learning, ICML02. pp. 267–274 (2002)
7. Kimura, H., Kobayashi, S.: An analysis of actor/critic algorithms using eligibility traces: Reinforcement learning with imperfect value function. In: ICML (1998)
8. Korenkevych, D., Mahmood, A.R., Vasan, G., Bergstra, J.: Autoregressive policies for continuous control deep reinforcement learning. In: Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI-19). pp. 2754–2762 (2019)
9. Liang, E., Liaw, R., Nishihara, R., Moritz, P., Fox, R., Goldberg, K., Gonzalez, J., Jordan, M., Stoica, I.: RLlib: Abstractions for distributed reinforcement learning. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 3053–3062. PMLR, Stockholmsmssan, Stockholm Sweden (10–15 Jul 2018)
10. Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning (2016), arXiv:1509.02971

11. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.: Playing atari with deep reinforcement learning (2013), arXiv:1312.5602
12. Schulman, J., Levine, S., Moritz, P., Jordan, M.I., Abbeel, P.: Trust region policy optimization (2015), arXiv:1502.05477
13. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms (2017), arXiv:1707.06347
14. Todorov, E., Erez, T., Tassa, Y.: Mujoco: A physics engine for model-based control. In: 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems. pp. 5026–5033. IEEE (2012)
15. Wang, Z., Bapst, V., Heess, N., Mnih, V., Munos, R., Kavukcuoglu, K., de Freitas, N.: Sample efficient actor-critic with experience replay (2016), arXiv:1611.01224
16. Wawrzyski, P.: Learning to control a 6-degree-of-freedom walking robot. In: Proceedings of EUROCON 2007 The International Conference on Computer as a Tool. pp. 698–705 (2007)
17. Wawrzyski, P.: Real-time reinforcement learning by sequential actorcritics and experience replay. Neural Networks **22**(10), 1484–1497 (2009)
18. Wawrzyski, P.: Control policy with autocorrelated noise in reinforcement learning for robotics. International Journal of Machine Learning and Computing **5**(2), 91–95 (2015)

# A Algorithms' hyperparameters

This section presents hyperparameters used in simulations reported in Sec. 5. All algorithms used the discount factor equal to 0.99. The rest of hyperparameters for ACERAC, ACER, SAC, and PPO, are depicted in Tab. 1, 2, 3, and 4, respectively.

| Parameter | Value |
|-----------|-------|
| Action std. dev. for Hopper | 0.3 |
| Action std. dev. for other envs. | 0.4 |
| $\alpha$ | 0.5 |
| Critic step-size for Walker2D | $10^{-4}$ |
| Critic step-size for other envs. | $6 \cdot 10^{-5}$ |
| Actor step-size for Walker2D | $5 \cdot 10^{-5}$ |
| Actor step-size for other envs. | $3 \cdot 10^{-5}$ |
| $\tau$ | 4 |
| $b$ | 2 |
| Memory size | $10^6$ |
| Minibatch size | 256 |
| Target update interval | 1 |
| Gradient steps | 1 |
| Learning start | $10^3$ |

Table 1: ACERAC hyperparameters

| Parameter | Value |
|-----------|-------|
| Step-size for Hopper | 0.0001 |
| Step-size for other envs | 0.0003 |
| Replay buffer size | $10^6$ |
| Minibatch size | 256 |
| Target smoothing coef. $\tau$ | 0.005 |
| Target update interval | 1 |
| Gradient steps | 1 |
| Learning start for Ant | $10^4$ |
| Learning start for HalfCheetah | $10^4$ |
| Learning start for Hopper | $10^3$ |
| Learning start for Walker2D | $10^3$ |
| Reward scale for Ant | 1 |
| Reward scale for HalfCheetah | 0.1 |
| Reward scale for Hopper | 30 |
| Reward scale for Walker2D | 30 |

Table 3: SAC hyperparameters

| Parameter | Value |
|-----------|-------|
| Action std. dev. | 0.3 |
| Critic step-size | $10^{-5}$ |
| Actor step-size | $10^{-5}$ |
| $\lambda$ | 0.9 |
| $b$ | 2 |
| Memory size | $10^6$ |
| Minibatch size | 256 |
| Target update interval | 1 |
| Gradient steps | 1 |
| Learning start | $10^3$ |

Table 2: ACER hyperparameters

| Parameter | Value |
|-----------|-------|
| GAE parameter ($\lambda$) | 0.95 |
| Minibatch size | 64 |
| Step-size | 0.0003 |
| Horizon | 2048 |
| Number of epochs | 10 |
| Policy clipping coef. | 0.2 |
| Value function clipping coef. | 10 |
| Target KL | 0.01 |

Table 4: PPO hyperparameters