

Using Model-Free Deep Reinforcement Learning Algorithms with the MuJoCo Physics Simulator

Curtis Brinker
cjbzfd@mst.edu

Tanner May
tmay@mst.edu

Abstract—Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

I. INTRODUCTION

Recent developments in machine learning have made huge strides in approaching problems that were previously unsolvable with traditional programming methods. In general, machine learning techniques are grouped into three categories: supervised learning, unsupervised learning, and reinforcement learning (RL) [1]. While the first two categories require the whole dataset to train, RL is capable of learning on data as it comes in. This data is formulated as a set of interactions with the agent's environment where each interaction consists of the state of the environment, the action the agent performed, and the resulting state. Thus, agents that use RL learn with a circular process that is a result of the agent's previous actions, known as the agent-environment interaction loop. This complex relationship has given rise to multitudes of algorithms, from remembering a mapping between states and actions to using multiple neural networks to determine the next action.

Like many other areas of machine learning, RL has been revolutionized by the introduction of deep learning. Reinforcement learning algorithms can now learn how to behave in complex, high dimensional environments, such as learning to play videogames with only the pixels as input or how to control a robot. Because of the new found power of RL, we believe that it will play an influential role in the future of AI.

To get experience in this exciting field, we implemented some model-free RL algorithms and taught them how to walk with the MuJoCo physics simulator.

II. TESTING ENVIRONMENT

Each of the models were trained using Open AI's Gym library and DeepMind's physics simulator MuJoCo. The Gym toolkit implements numerous environments for testing RL,

such as the MuJoCo physics simulator. MuJoCo stands for **M**ulti-**J**oint dynamics with **C**ontact and was made for applications that require fast, accurate physics simulations such as robotics and machine learning [2]. Using MuJoCo, machine learning algorithms can learn how to control the movement of robotic models, such as a snake, ant, or humanoid.

Each of the Gym's environments adhere to a strict structure that makes it easy to switch between environments. One of the most important parts to keep consistent are the values returned after every step of training. Each environment returns the following four values [3]:

- **Observation:** Represents the current state of the environment, as the agent sees it. In MuJoCo's case, the environment is fully observable and consists of the model's position, velocity, and various forces [4]
- **Reward:** The reward result of the previous action. The reward for the MuJoCo environment depends on a number of things and is unique for each model.
- **Done:** A flag denoting if the environment needs to be reset. Each of the models in MuJoCo define a safe height range that the model is allowed to exist within; if the model leaves that range, the done flag is raised and the model is reset back to the start.
- **Info:** Some information useful for debugging.

III. BACKGROUND

The study of reinforcement learning is about training an *agent* to interact with an *environment*. An action by an agent can influence the environment, which may later affect the actions of the agent. For this reason, reinforcement learning algorithms must be able to handle a changing environment that is causally influenced by the agent itself.

Formalizing the interdependency between the environment and the agent is done with the *agent-environment interaction loop*. In which, at time t the environment is fully described by its *state*, s_t . Then, an agent makes an *observation* of the environment, o_t , where $o_t \subseteq s_t$.¹ The agent responds to this observation with an *action* a_t and is given a *reward*, r_t . After the action is taken, the state of the environment changes with the new state denoted as s_{t+1} .

An environment state can be represented by a tensor of values describing individual aspects of the environment. The

¹While the observation and the state are not necessarily equal, reinforcement learning literature frequently refers to an observation as the state itself.

domain of possible observations of the agent is referred to as the *observation space*, which can be either continuous or discrete. Similarly, the domain of all actions that an agent can take is called the *action space* which can also be continuous or discrete. An action, a_t , is selected from the action space by agent's *policy*. The policy can select actions either deterministically or stochastically. Typically, deterministic policies are denoted with $a_t = \mu(s_t)$ and stochastic policies are denoted as $a_t \sim \pi(\cdot|s_t)$. [5]

The reward given for the transition from s_t to s_{t+1} by action a_t is given by the reward function, where $r_t = R(s_t, a_t, s_{t+1})$. Though, readers should be aware that some literature refers to this value as r_{t+1} . Typically, this reward function is specified by a programmer in a way that rewards a specific goal. With this in mind, the process of training an agent can be formulated as maximizing the cumulative reward, called the *return*. To calculate the return, the *trajectory* is used which contains the sequence of states and actions taken by the agent. The trajectory is given by:

$$\tau = (s_0, a_0, s_1, a_1, \dots)$$

Therefore, as the return is defined in terms of cumulative reward, the trajectory can be used to calculate the rewards as the agent interacts with the environment. Then, in its most simple form, the return is the summation of all rewards. However this introduces two problems: First the return for an action t is defined in terms of actions taken before it. Intuitively, we would to only consider rewards that will happen as a result of the current action. Second, The summation of all rewards weights all rewards equally, however, this means that very long term rewards may eventually dominate the summation. To encourage shorter term rewards we can introduce a discount factor, γ , that discounts the value of future rewards.

Commonly, the reward function R is co-opted to take a trajectory and produce the rewards at each timestep t . With these considerations in mind, the return can be defined as:

$$R_t(\tau) = \sum_{t'=t}^T \gamma^{t'} r_{t'}, \quad \gamma \in (0, 1)$$

Where T is the length of the trajectory and t is the current time for which the return is being calculated. If T is finite, then the return is said to be a *finite-horizon return*. Similarly, if T is infinite, then the return is said to an *infinite-horizon return*. The inclusion of γ makes this return a discounted return. This places more focus on near-term rewards and also helps with convergence in infinite-horizon returns. Using t' makes the return only consider future rewards, which [6] calls the *reward-to-go* return.

However, an agent cannot act to directly maximize its return as it dependent on future actions and states. Instead, an approximation must be used to calculate the expected return, $\mathcal{J}(\pi)$. This this done with *value functions*. Formally, $\mathcal{J}(\pi)$ is defined as:

$$\mathcal{J}(\pi) = \mathbb{E}_{\tau \sim \pi} [R(\tau)]$$

The *state value function* for policy π , denoted as V^π , gives the expected value for policy π given state s . Alternatively, the *state-action value function* for policy π , denoted as Q^π , which can be used estimate the expected value for policy π in state s and immediately taking action a . Formally, these are defined as:

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s]$$

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a]$$

However, these functions are still dependent on the return. To avoid this, we can use the Bellman Equations, which defines both V^π and Q^π in terms of expected values of future states. [5], [7], [8] These equations are given as:

$$V^\pi(s) = \mathbb{E}_{\substack{a \sim \pi \\ s' \sim P}} [R(s, a, s') + \gamma V^\pi(s')]$$

$$Q^\pi(s, a) = \mathbb{E}_{s' \sim P} \left[R(s, a, s') + \gamma \mathbb{E}_{a' \sim \pi} [Q^\pi(s', a')] \right]$$

$$= \mathbb{E}_{s' \sim P} [R(s, a, s') + \gamma V^\pi(s')]$$

Where P is the environment state transition function and s' is the resulting state.

For the agent to maximize $\mathcal{J}(\pi_\theta)$, the agent must learn to improve the policy π_θ , where θ are the parameters for the policy and is normally updated by gradient ascent. Normally this involves using V^π or Q^π to better approximate the optimal value functions. Commonly, these approximation are done via deep neural networks. The exact details of how the policy and value function(s) are updated is dependent on the reinforcement learning algorithm used. Following is a brief discussion on different types of reinforcement learning algorithms followed by discussions on details of specific algorithms studied in this paper.

A simplified taxonomy of deep reinforcement learning algorithms is given by [9]. At the highest level, reinforcement learning algorithms can be split into two categories, *model-based* or *model-free*. Model-based algorithms learn the state transition function, allowing it to predict future states of the environment. Learning the world model has many benefits that can drastically increase the performance of an agent, however, doing so significantly increases the complexity of the algorithms and is outside the scope of this paper. Model-free algorithms do not learn to predict the environment and instead learn other functions to increase the agent's return. Model-free algorithms can be further categorized as using *policy optimization* or *Q-learning*.²

Policy optimization algorithms train by optimizing the parameters of the policy π_θ . This is normally done learning a value approximator V_ϕ and sampling actions taken by the current policy to update parameters accordingly. Because this process requires actions from the current policy, these techniques are considered to be *on-policy*.

²Policy optimization and Q-learning are not mutually exclusive categories.

Another approach, though not necessarily mutually exclusive, is using Q learning. These algorithms learn an approximator Q_θ to determine the value of state action pairs. An advantage to this technique is that it can learn from past experiences that were not from the current policy. Algorithms that can learn from experiences of different policies are called *off-policy*. Actions by these algorithms typically are chosen by selecting the action that maximizes $Q_\theta(s, a)$ for the current state s .

A. The REINFORCE Algorithm

One of the simplest forms of policy optimization algorithms is a class of algorithms presented by Williams as REINFORCE. [10] The core idea this algorithm is to update the policy by sampling trajectories from the policy. Once a sample has been taken from the current policy, the policy π_θ and value approximator V_ϕ can be updated. There is some flexibility in how the parameters are updated, though [6] provides some examples of frequently used techniques. The policy is updated using gradient ascent where the gradient is given by:

$$\nabla_\theta \mathcal{J}(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \Phi_t \right]$$

Where Φ_t is any function that is independent of θ .

To update the value approximator parameters ϕ , it is common to use gradient descent on a loss such as mean squared error on the value approximator and the reward-to-go return.

$$\mathcal{L}(\phi) = \mathbb{E}_{\tau \sim \pi_\theta} \left[(V_\phi(s_t) - R_t(\tau))^2 \right]$$

B. Advantage Actor Critic (A2C)

The advantage actor critic is an on-policy algorithm based on the REINFORCE algorithm. A2C was first proposed by [11] along with an asynchronous variant called A3C. The details about the asynchronous version of the algorithm is not pertinent to this paper and will not be discussed. As a whole, there are two main techniques used to improve the performance of this algorithm.

The first technique is the use of the advantage function, $A^{\pi_\theta}(s_t, a_t)$ as Φ_t . The advantage function is defined as:

$$A^{\pi_\theta}(s_t, a_t) = Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t, a_t)$$

It should be noted that Φ_t must be independent of θ and by consequence a_t . However, [6] provides a proof that $A^{\pi_\theta}(s_t, a_t)$ can be written in terms of s_t only. The advantage can be thought of as quantifying how much better an action in a given state is than average. Negative values signify that the action was less than expected value for the given state, and positive signifies that it was better. Mathematically, the advantage function useful as subtracting by $V^{\pi_\theta}(s_t)$ reduces the variance in the policy gradient leading to more stable learning.

The second technique is the use of an actor-critic framework. The idea of this framework is to separate the concerns of learning different aspects of training. The actor is tasked with learning the policy that will interact with the environment.

The critic is tasked with learning the value function that it will use to criticize the actor's actions. [12]

C. Deep Deterministic Policy Gradient (DDPG)

DDPG is an off-policy learning algorithm for continuous action spaces that was proposed by [13]. It learns a Q function and a deterministic policy to select actions and uses a variety of techniques to improve performance and stability of the algorithm. An overview of unique methods used by this algorithm is provided by [14]. One technique that is commonly used by off-policy algorithms is the use a replay buffer. Because off-policy algorithms can learn from actions that did not originate from the current policy, past experiences can be stored in a buffer. Then, during training, a larger sample size can be used to improve training.

To update the function Q_ϕ , gradient descent is used on ϕ to minimize the mean squared error of $Q_\phi(s, a)$ and the *target*, where the target is defined as:

$$y = r + \gamma(1 - d)Q_\phi(s', \mu_\theta(s'))$$

Where (s, a, r, s', d) is a transition in the replay buffer \mathcal{D} . The resulting state from the action is s' and d is the *done* flag where $d = 1$ if s' is a terminal state, $d = 0$ otherwise. Additionally, μ_θ is the deterministic policy used to select the next action.

However, using gradient descent the mean squared error of Q_ϕ and y is problematic as y is also dependent on ϕ . As a result, both Q_ϕ and y change when ϕ is updated, leading to instability. To resolve this, DDPG uses *target networks*, which uses polyak averaging to remove this direct dependency. The target networks, $Q_{\phi'}$ and $\mu_{\theta'}$ have the same structure as their original networks but are updated as follows:

$$\begin{aligned} \phi' &\leftarrow \tau \phi' + (1 - \tau) \phi \\ \theta' &\leftarrow \tau \theta' + (1 - \tau) \theta \end{aligned}$$

For $\tau \in (0, 1)$. This provides stability to the loss function, therefore Q_ϕ is updated using gradient descent on

$$\mathcal{L}(\phi, \mathcal{D}) = \mathbb{E}_{(s, a, r, s', d) \sim \mathcal{D}} \left[(Q_\phi(s, a) - y)^2 \right]$$

This loss function is called the Mean Square Bellman Error (MSBE) where y is the target, which is calculated as:

$$y = r + \gamma(1 - d)Q_{\phi'}(s', \mu_{\theta'}(s'))$$

To update the policy, we want to maximize the value function Q_ϕ using μ_θ . To do this, the following loss function is used:

$$\mathcal{L}(\theta, \mathcal{D}) = \mathbb{E}_{s \sim \mathcal{D}} [-Q_\phi(s, \mu_\theta(s))]$$

Finally, to gain additional benefit from the algorithm off-policy, DDPG uses another technique to encourage exploration during train time. This is done by adding random noise to actions to stray away from the deterministic policy. The original authors use OU noise, however [14] recommends a zero centered Gaussian distribution which is simpler to implement and still effective.

D. Soft Actor Critic (SAC)

The Soft Actor Critic Algorithm is an stochastic off-policy algorithm that uses an actor-critic framework and entropy regularization. This algorithm was originally proposed by [15], [16], additionally, a guided explanation is given by [17]. SAC can be implemented for either discrete or continuous action spaces, however, the details that follow will be for the continuous version.

SAC works by learning a stochastic policy π_θ and two Q functions, Q_{ϕ_1} and Q_{ϕ_2} . When evaluating the value, the minimum value is taken between $Q_{\phi_1}(s, a)$ and $Q_{\phi_2}(s, a)$. This is called the *clipped double-Q method*. Using this method helps prevent overestimating the value of state action pairs which improves training results.

Similar to DDPG, SAC uses target networks to stabilize the updates using the loss function. Unlike DDPG, SAC only uses target networks for value functions, not for policy. The parameters for the target networks are updated using polyak averaging as shown below:

$$\begin{aligned}\phi'_1 &\leftarrow \tau \phi'_1 + (1 - \tau)\phi_1 \\ \phi'_2 &\leftarrow \tau \phi'_2 + (1 - \tau)\phi_2\end{aligned}$$

For $\tau \in (0, 1)$.

To control exploration of the policy, SAC uses entropy regularization, where entropy is defined as:

$$H(P) = \mathbb{E}_{x \sim P} [-\log P(x)]$$

For any distribution P . Roughly speaking, entropy measures how random a distribution is. An entropy of zero is given by a policy where where an action has probability 1. Larger values of entropy indicate more evenly spread probability distribution. The choice of entropy regularization is controlled with the entropy hyperparameter, α , where $\alpha \in (0, 1)$. The greater the value of α , the more the policy will explore.

Entropy is introduced in several aspects of the SAC algorithm. The first is in the *entropy regularized double-Q clipped target*, which is calculated as:

$$y = r + \gamma(1 - d) \left(\min_{k=1,2} Q_{\phi'_k}(s', a') - \alpha \log \pi_\theta(a' | s') \right)$$

Where $a' \sim \pi_\theta$.

This is used to calculate the entropy regularized MSBE which is used as the loss function for Q_{ϕ_1} and Q_{ϕ_2} , given by:

$$\mathcal{L}(\phi_k, \mathcal{D}) = \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[(Q_{\phi_k}(s, a) - y)^2 \right], \quad k = 1, 2$$

Where \mathcal{D} is the experience replay buffer.

The policy π_θ is updated in a way similar to DDPG. The policy is updated by using gradient ascent to maximize the value function. Therefore the loss function for π_θ becomes:

$$\mathcal{L}(\theta, \mathcal{D}) = \mathbb{E}_{\substack{s \sim \mathcal{D} \\ a \sim \pi_\theta}} \left[- \left(\min_{k=1,2} Q_{\phi'_k}(s, a) - \alpha \log \pi_\theta(a | s) \right) \right]$$

There are two important details regarding sampling a from π_θ in the SAC algorithm. First, sampling is done from a normal distribution using the *reparameterization method*, this the process of sampling a random number differentiable, allowing for actions to backpropagate the the policy's parameters. The second detail is that the SAC policy squashes the action to finite action space, with *tanh*. This is done as most action spaces have finite bounds while normal distributions have infinite support. By limiting the domain of the action to the action space, the policy is learned more efficiently.

IV. METHODOLOGY

Each of the algorithms were implemented on Ubuntu 20.04 with Python 3.8 using the PyTorch library. We used PyTorch because we had previous experience with it and liked how flexible it is compared to other options. To train the models, we used the previously mentioned Gym library. To verify that the implementation of each of the algorithms worked correctly, they were trained in the CartPole and Pendulum environments. For training towards the goal, the Ant-v3 and Humanoid-v3 MuJoCo models were used.

The hyperparameters for each algorithm were as follows:

- A2C: Learning rate of $3e^{-4}$ and $\gamma = 0.99$.
- DDPG: Update freq. of 64 steps, update threshold of 4096 steps, batch size of 128, learning rate of $1e^{-3}$, $\gamma = 0.99$, $\tau = 0.995$, and the noise distribution was Gaussian with a standard deviation of 0.1 (this differs from [13] which used OU noise due to a recommendation by [14] to use the easier to implement Gaussian noise).
- SAC: Update freq. of 64 steps, 64 updates per update step, an update threshold of 4096 steps, batch size of 128, $\alpha = 0.5$, learning rate of $5e^{-4}$, $\gamma = 0.99$, $\tau = 0.995$, and an α decay of 1.

Each of those parameters were chosen because they resulted in the best performing model.

Since the goal is to teach Deep RL algorithms how to walk, the networks for each of the algorithms were defined to have two layers, each with 256 neurons. Originally the layers had 128 neurons each but tests using the 256 neuron configuration resulted in better performance. More layers and neurons can, of course, be used but would have greatly increased the computational cost of training.

Since every environment has its own quirks, each is described in its own subsection below [4].

A. CartPole-v0

The objective of the CartPole environment is to teach an agent to balance a pole on a cart by moving the cart left and right. The environment has a continuous observation space consisting of the cart's position, the cart's velocity, the pole's angle, and the pole's angular velocity. The action space is discrete with the options of moving the cart to the left or to the right. The environment rewards the agent with a value of one, simply for staying healthy.

Initially, the values of both the cart and pole are randomly sampled from a uniform distribution from $[-0.5, 0.5]$. The

environment is terminated when the angle of the pole is $\pm 12^\circ$, the cart position hits the edge of the display, or when the length of the episode surpasses 200 steps. The environment is considered solved when the reward is > 195 for 100 consecutive episodes.

This environment was used as a proof of concept of each algorithm's implementation.

B. Pendulum-v1

The pendulum environment's objective is to balance a frictionless pendulum straight up. The observation space is continuous and describes the angle of the pendulum as well as the angular velocity. The action space is also continuous and details the amount of left or right force to apply to the pendulum. The reward incentivizes keeping the pole at an angle of 0° with as little velocity and force as possible. The reward is formulated as

$$R = -(angle^2 + 0.1 * angular_velocity^2 + 0.001 * action^2)$$

Since the angle is normalized between $[-\pi, \pi]$ before calculating the reward, the reward has a range of $[-16.3, 0]$.

The environment's initial state has the pole at a random angle between $[-\pi, \pi]$ radians with a random velocity between $[-1, 1]$. The environment does not specify a termination state, so a limit of 150 was imposed. Similarly, the environment does not specify when it is solved, so it was allowed to run until it hit the step limit.

Like the CartPole environment, this environment was used as another proof of concept of the implementation for each algorithm.

C. Ant-v3

The ant is a sphere with four legs, each with two joints. The continuous observation space is significantly more complex than the previous environments: it consists of the model's position, velocity, and the forces between the legs and the ground (contact force). The action space is also continuous and describes where and how quickly to move each joint. The reward encourages the model to move as quickly as possible while moving as few joints as little and softly as possible. It is formulated as

$$R = (v + h) - (a + f)$$

where v is the velocity in the x-axis, h is the healthy reward configured to be 1, a is the control_cost calculated by $0.5 * \sum actions^2$, and f is the contact_cost calculated by $5e^{-4} * \sum contact_forces^2$.

The environment starts with each joint in a random position with a random velocity. The episode is terminated when the model exits the safe height range of $[0.3, 2.0]$. The environment did not define a solve condition so the episode was continued until the model made a mistake to trigger the termination condition.

The ant was chosen because it is the simplest provided model that meets the "three dimensional walker" condition. Each of the algorithms had 500,000 episodes to train using

this model. After training was stopped, the most recently saved model was used for analysis. A training length of 500,000 episodes was used because it consistently produced an agent that could walk a non-trivial distance.

D. Humanoid-v3

The humanoid model is a humanoid with two legs, two arms, a head and a torso. The legs, arms, and torso are made up of multiple joints, each individually controlled. Like the ant, the observation space is continuous and describes the model's position, the model's velocity, the forces for each of the joints, and the contact forces. The environment also defines a continuous action space that describes the new position and velocity of each joint. The intuition for the reward is the same as for the Ant environment, but it is calculated slightly differently where v is equal to $1.25 * x_velocity$, h is configured to be 5, a is calculated as $0.1 * \sum actions^2$, and f is defined as $5e^{-7} * \sum contact_forces^2$.

Like the Ant environment, the model starts with each joint in a random position and random velocity. The episode is terminated when the model exits the safe height range of $[1.0, 2.0]$. The environment did not define a solve condition so the episode was continued until the model made a mistake to trigger the termination condition.

The humanoid was chosen because it is the most complex 3D walker and showcased the learning power of the tested algorithms. Since this model is significantly more complex than the Ant, 1,250,000 episodes were used for training. But due to the cost of training, only SAC was trained with this model.

V. RESULTS

Why does loss not matter; only need to see that the reward is trending upward;

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

A. CartPole-v0

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing

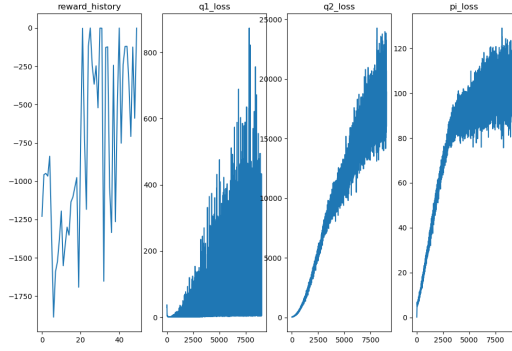


Fig. 1: SAC in the Pendulum-v1 Environment

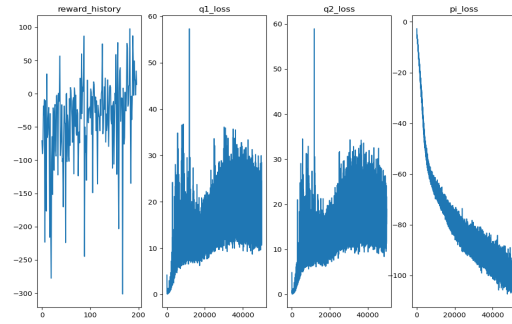


Fig. 2: SAC in the Ant-v3 environment

elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

B. Pendulum-v1

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

C. Ant-v3

¡SAC jumped, DDPG walked¿ ¡SAC reward increase is steeper than DDPG's¿

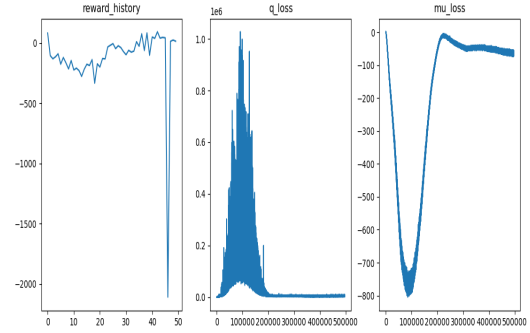


Fig. 3: DDPG in the Ant-v3 environment

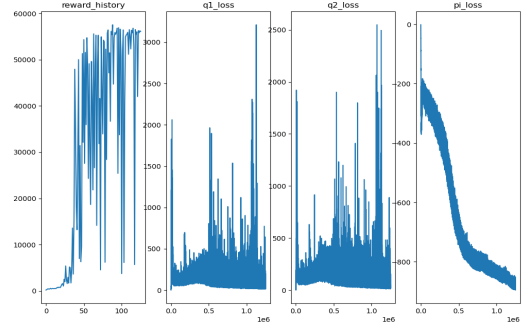


Fig. 4: SAC in the Humanoid-v3 environment

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

D. Humanoid-v3

¡model shuffled along instead of taking steps¿

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit

blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

VI. FUTURE WORK AND CONCLUSIONS

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

REFERENCES

- [1] W. Qiang and Z. Zhongli, "Reinforcement learning model, algorithms and its application," in *2011 International Conference on Mechatronic Science, Electric Engineering and Computer (MEC)*, 2011, pp. 1143–1146.
- [2] DeepMind, "Mujoco documentation." [Online]. Available: <https://mujoco.readthedocs.io/en/latest/overview.html>
- [3] OpenAI, "Gym documentation." [Online]. Available: <https://gym.openai.com/docs/>
- [4] —, Gym library. v0.21.0. [Online]. Available: <https://github.com/openai/gym>
- [5] (2018) Spinning up: Key concepts in rl. [Online]. Available: https://spinningup.openai.com/en/latest/spinningup/rl_intro.html
- [6] (2018) Spinning up: Intro to policy optimization. [Online]. Available: https://spinningup.openai.com/en/latest/spinningup/rl_intro3.html
- [7] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [8] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [9] (2018) Spinning up: Kinds of rl algorithms. [Online]. Available: https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html
- [10] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Mach. Learn.*, vol. 8, no. 3–4, p. 229–256, may 1992. [Online]. Available: <https://doi.org/10.1007/BF00992696>
- [11] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," *CoRR*, vol. abs/1602.01783, 2016. [Online]. Available: <http://arxiv.org/abs/1602.01783>
- [12] H. van Hasselt. (2021) Reinforcement learning lecture 9: Policy gradients and actor critics. [Online]. Available: <https://storage.googleapis.com/deepmind-media/UCL%20x%20DeepMind%202021/Lecture%209-%20Policy%20gradients%20and%20actor%20critics.pdf>
- [13] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [14] (2018) Spinning up: Deep deterministic policy gradient. [Online]. Available: <https://spinningup.openai.com/en/latest/algorithms/ddpg.html#>
- [15] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," *CoRR*, vol. abs/1801.01290, 2018. [Online]. Available: <http://arxiv.org/abs/1801.01290>
- [16] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft actor-critic algorithms and applications," *CoRR*, vol. abs/1812.05905, 2018. [Online]. Available: <http://arxiv.org/abs/1812.05905>
- [17] (2018) Spinning up: Soft actor-critic. [Online]. Available: <https://spinningup.openai.com/en/latest/algorithms/sac.html>