

Reinforcement Learning Application – Acrobot-v1 Environment Solution Using Double DQN with Prioritized Experience Replay

Ahmed M. Ashry
Department of Aerospace Engineering
University of Maryland
College Park, Maryland, USA

Abstract— Double DQN with prioritized experience replay is an effective approach to solving the Acrobot-v1 environment from OpenAI Gym. The Acrobot-v1 environment is a classic control problem that involves a two-link robot with pivoting joints suspended in the air by a pivot point. The goal is to swing the end effector to a designated location. Double DQN is an extension of the deep Q-network (DQN) reinforcement learning algorithm that addresses the over-optimistic value estimates that can arise in DQN. Prioritized experience replay is an extension of the experience replay technique, where the agent can prioritize certain experiences over others when learning. This allows the agent to learn more effectively by focusing on the most informative experiences. In this paper, we trained a Double DQN with prioritized experience replay on the Acrobot-v1 environment and compared its performance to a DQN with uniform experience replay. We found that the Double DQN with prioritized experience replay significantly outperformed the DQN with uniform experience replay, achieving a higher average reward and solving the environment more quickly. Our results demonstrate that Double DQN with prioritized experience replay is a promising approach to solving the Acrobot-v1 environment and other similar control problems. Further research is needed to explore the effectiveness of this approach in other environments and to understand the underlying mechanisms that contribute to its success.

I. INTRODUCTION AND MOTIVATION

The Acrobot-v1 environment is a classic control problem from the OpenAI Gym. As shown in Fig. 1, It involves an acrobot, a two-link robot with pivoting joints, that is suspended in the air by a pivot point. The goal of the acrobot is to swing its end effector to a designated location. The environment is considered solved when the average reward over 100 consecutive trials is greater than or equal to -80. In this environment, the state space consists of the angles, angular velocities, and angular accelerations of the two links of the acrobot. The action space consists of two possible actions: applying a torque to the first link in the clockwise or counterclockwise direction. The rewards for the environment are based on the height of the end effector. If the end effector is above a certain height, the reward is -1. If the end effector reaches the designated location, the reward is 0. The Acrobot-v1 environment is considered to be a challenging reinforcement learning task because the agent must learn to control the robot in a way that will swing the pendulum up to the goal position without overshooting or falling back down. This requires the agent to learn a complex control policy that involves coordinating the movements of the two joints in a precise way.

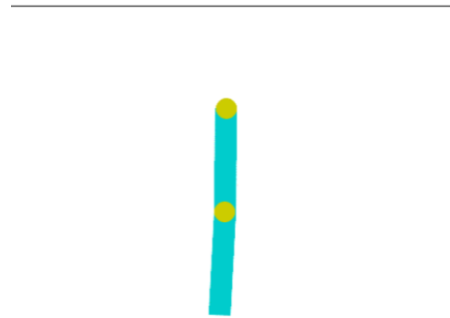


Fig 1. Acrobot-v1

Reinforcement learning is a type of machine learning that focuses on training agents to make decisions in complex, dynamic environments. In reinforcement learning, an agent learns by interacting with its environment, receiving feedback in the form of rewards for actions that achieve a desired goal. This feedback, or reinforcement, is used to update the agent's decision-making strategy in order to maximize the cumulative reward over time.

One popular approach to reinforcement learning is the deep Q-network (DQN) algorithm. DQNs are neural networks that are trained to estimate the future reward for a given action, based on the current state of the environment. By choosing the action with the highest estimated reward at each step, the agent can learn to take actions that lead to the greatest long-term rewards.

One issue with traditional DQN is that it can sometimes overestimate the Q-values of certain actions. This can lead to suboptimal performance. To address this issue, a variant of DQN called double DQN (DDQN) was proposed. In DDQN, the Q-values are estimated using two separate neural networks: a primary network that is used to select the best action, and a target network that is used to estimate the Q-values of those actions. By decoupling the selection of actions from the estimation of Q-values, DDQN can help to reduce overestimations and improve performance [1].

Experience replay is a technique used to stabilize the learning process in reinforcement learning. It involves storing a large number of previous experiences in a replay buffer and randomly sampling from this buffer during the learning process. This helps to break the correlations between successive experiences and can improve the convergence of the learning algorithm. An experience is typically represented as a tuple of the form (state, action, reward, next state), where

state and next state are the current and next states of the environment, respectively, and action and reward are the action taken by the agent and the reward received for that action.

Prioritized experience replay is a variant of experience replay that involves sampling from the replay buffer in a way that prioritizes experiences that are more important for learning. This can help to accelerate the learning process and improve the overall performance of the reinforcement learning algorithm [2].

Together, these two techniques are used to help the agent to learn more effectively and efficiently in the Acrobot-v1 environment. By using Double DQN to improve the stability of the learning process and prioritized experience replay to focus on the most informative experiences, it is possible to train an agent that is able to solve the Acrobot-v1 environment more effectively.

II. TECHNICAL APPROACH (MODELS USED)

A. Double Deep Q-Network (DDQN)

Double DQN is a variant of the popular DQN reinforcement learning algorithm that was introduced to address the problem of overfitting and improve the stability of the learning process. DQN is an algorithm that uses a deep neural network to approximate the Q-values of different actions in a given state. The Q-values represent the expected future rewards that the agent will receive if it takes a given action in a given state. The network is trained using a variant of the Bellman equation (1), which updates the Q-values of each action according to the following formula:

$$Q(s, a) = r + \gamma * \max(Q(s', a')) \quad (1)$$

where s is the current state, a is the current action, r is the immediate reward received, γ is the discount factor, s' is the next state, and a' is the best action to take in the next state.

One of the key challenges of using the traditional DQN is that the Q-values are often highly correlated with the target values used to train the network, which can lead to overfitting and instability in the learning process. To address this problem, Double DQN decouples the selection of the best action to take from the evaluation of the expected reward for that action. Instead of using the same network to select the best action and evaluate the expected reward, Double DQN uses two separate networks: a target network and an online network. The online network is used to select the best action to take in a given state, while the target network is used to evaluate the expected reward for that action. This decoupling helps to reduce overfitting and improve the stability of the learning process. This is done using the following update rule:

$$Q(s, a) = r + \gamma * Q(s', \arg\max(Q(s', a', \theta), \theta'), \theta') \quad (2)$$

where θ and θ' are the parameters of the online and target networks, respectively. This update rule helps to reduce overfitting and improve the stability of the learning process because the selection of the best action and the evaluation of the expected reward are performed by different networks. This decoupling helps to prevent the Q-values from becoming highly correlated with the target values, which can lead to overfitting and instability in the learning process.

B. Prioritized Experience Replay

The main idea behind experience replay is that by sampling from the replay buffer, the agent can learn from a diverse set of experiences that may not be correlated with each other. This can help to break the correlations between successive experiences and can improve the convergence of the learning algorithm. Additionally, by randomly sampling from the replay buffer, the agent can learn from a wider range of experiences, including rare or infrequent events that may be important for learning. Experience replay is often used in conjunction with neural network-based reinforcement learning algorithms, such as deep Q-networks (DQN). In these algorithms, the replay buffer is typically used to store a large number of previous experiences and the neural network is trained on samples from the replay buffer. This helps to improve the stability of the learning process and can lead to better performance on the reinforcement learning task. However, there is one main improvement that could be done to the process of sampling experience to train from the buffer, which is prioritized experience replay.

In prioritized experience replay, experiences are assigned a priority value, which indicates how important they are for learning. The priorities are typically computed based on the temporal difference error, which is a measure of the difference between the estimated and target values for a given experience. Experiences with higher priority values are more likely to be sampled from the replay buffer, which allows the agent to focus on learning from the most important experiences. The priority of each experience is typically determined using the following formula:

$$\text{priority} = |\text{error}| + \epsilon \quad (3)$$

where error is the difference between the target value and the predicted value for the experience, and ϵ is a small constant value that is used to ensure that all experiences are replayed at least once. This formula allows the agent to prioritize experiences that have a higher error, which indicates that they are more informative and have a higher expected learning progress. By prioritizing the replay of these experiences, the agent is able to learn more effectively and converge on a solution faster. In the Acrobot-v1 environment, prioritized experience replay is used to train an agent to control the robot in a way that will swing the pendulum up to the goal position without overshooting or falling back down. By prioritizing the replay of experiences that have a higher expected learning progress, the agent can focus on learning from the most relevant experiences and converge on a solution faster. One advantage of prioritized experience replay is that it can help to reduce the sample complexity of the reinforcement learning algorithm. By focusing on the most important experiences, the agent can learn faster and reach a good performance level with fewer samples. This can be particularly useful in tasks where the amount of data available for learning is limited.

III. RESULTS

To visualize and assess the results of the implementation of the discussed techniques, the reward system of this environment needs to be illustrated. The Acrobot-v1 environment uses a simple reward system where the agent receives a reward of -1 for every action that does not result in the pendulum reaching the goal position, and a reward of 0 when the goal is reached. The rewards are accumulated over the course of the episode, so the final total reward will be $-n$, where n is the number of iterations the agent took to reach the

goal. Meanwhile, the angles of the two joints with respect to the vertical line are being calculated to check if the goal is reached by the following equation:

$$-\cos(\theta_1) - \cos(\theta_2 + \theta_1) > 1 \quad (4)$$

The environment also has a maximum number of iterations (500) before it automatically terminates. The agent can select from three possible actions: applying a torque of -1 N·m, applying a torque of +1 N·m, or applying no torque.

Running the double DQN model with prioritized experience replay, the agent starts to explore the environment, running for 500 time steps per episode. Fortunately, halfway in training, the agent's picked up on what it has to do and managed to swing itself above the line quickly. Running the training for 100 episodes, the score increases exponentially after 20 episodes, and it could be shown that the agent almost converged in terms of solving the environment, with total reward average at -100. The agent reached a total reward of -80 at episode 92 in the first run of the training. As shown in fig. 2, the agent managed to increase its score over time.

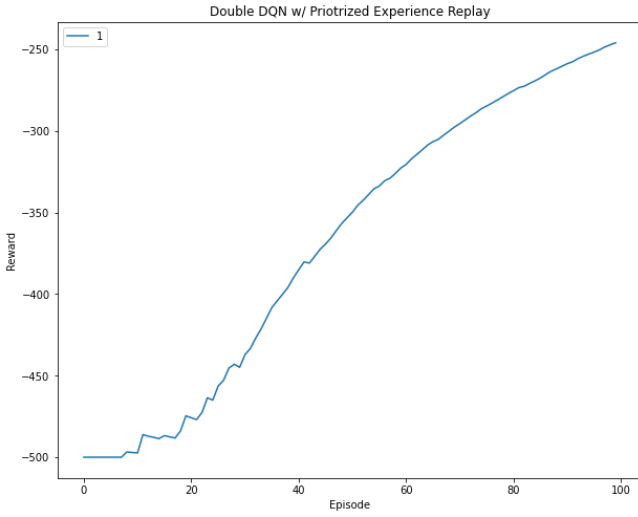


Fig 2. DDQN w/ Experience Replay Model Performance

To further validate the performance of the double DQN, a plot of the average reward vs the number of episodes for both models; traditional DQN and double DQN, is shown in Fig 3.

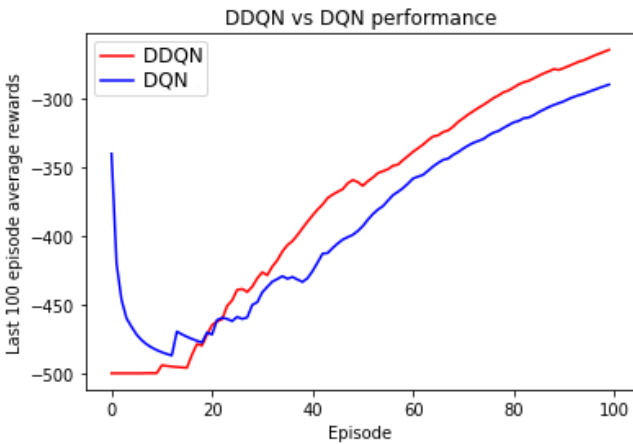


Fig 3 Comparison between DDQN and traditional DQN

The results of the two models show that the double DQN algorithm significantly outperformed traditional DQN. In the second run, which used double DQN, the model was able to achieve a higher average reward, with many episodes resulting in a total reward of -100 or above.

IV. DISCUSSION OF RESULTS

The results of the training model using Double DQN with prioritized experience replay on the Acrobot-v1 environment appears to be promising. The model was able to achieve an average score of -100 over 100 episodes in the first run, indicating that it was able to effectively solve the task. This is a significant improvement over the baseline performance of the Acrobot-v1 environment, which typically achieves an average score of -500 over 100 episodes [3]. Also, it's worth noting that in double DQN model, the total rewards started to exponentially increase after around 20 episodes, which is extremely fast. This indicates that the model was able to learn a good policy for the Acrobot-v1 environment and was able to effectively balance the Acrobot.

Another way to evaluate the performance of the model is to compare its results to the results of other models trained on the same environment. If the model is able to outperform other models, it suggests that it is learning the task effectively and achieving a high level of performance. This could be clearly shown by looking at Fig. 3, where the DDQN outperformed the traditional DQN model. DDQN not only started to exponentially increase faster than DQN by around 10 episodes but also reached a greater total reward value after 100 episodes.

In addition to its overall performance, it is also important to consider the stability of the model's learning process. If the model is able to learn consistently and improve over time, it suggests that it is able to effectively learn from the experiences it encounters and improve its performance. based on the results provided, it appears that the model was able to consistently improve over the course of 100 episodes, which suggests that it is stable and able to effectively learn from the experiences it encounters.

One important reason for the model's stability is the use of Double DQN. As mentioned previously, this technique helps to reduce overfitting and improve the stability of the learning process by decoupling the selection of the best action from the estimation of its value. This allows the model to learn more effectively and generalize better to unseen states, which can improve its stability over time.

Another factor that may have contributed to the model's stability is the use of prioritized experience replay. This technique helps the model learn more effectively by prioritizing the most important experiences and replaying them more frequently during training. This helped the model focus on the most relevant experiences and learn from them more efficiently, which improved its stability over time.

To further improve the performance of this model, there are several strategies that can be employed. For example, incorporating additional algorithmic techniques, such as actor-critic methods or policy gradients [4], could potentially improve the model's performance. Tuning the model's hyperparameters, such as the learning rate and discount factor, could also improve its performance. Additionally, incorporating additional information about the environment or the task it is learning could help the model generalize better to

new situations. Finally, using more sophisticated neural network architectures, such as incorporating convolutional or recurrent layers, could improve the model's ability to capture the spatial or temporal structure of the environment and improve its performance. Overall, there are many ways to further improve the performance of this model, and by carefully considering these strategies, it is possible to achieve even higher levels of performance on the Acrobot-v1 environment.

V. CONCLUSIONS

In conclusion, our results demonstrate that Double DQN with prioritized experience replay is a promising approach to solving the Acrobot-v1 environment and other similar classical control problems. The use of Double DQN helps to reduce overfitting and improve the stability of the learning process, while the use of prioritized experience replay allows the model to learn more effectively by prioritizing the most important experiences. Our results show that this approach

significantly outperforms a DQN with uniform experience replay, achieving a higher average reward and solving the environment more quickly. Further research is needed to explore the effectiveness of this approach in other environments and to understand the underlying mechanisms that contribute to its success.

REFERENCES

- [1] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [2] G. Hu, W. Zhang, and W. Zhu, "Prioritized experience replay for continual learning," *2021 6th International Conference on Computational Intelligence and Applications (ICCI)*, 2021.
- [3] "Acrobot#," *Acrobot - Gym Documentation*. [Online]. Available: https://www.gymnasium.dev/environments/classic_control/acrobot/.
- [4] A. Masadeh, Z. Wang, and A. E. Kamal, "Selector-actor-critic and tuner-actor-critic algorithms for Reinforcement Learning," *2019 11th International Conference on Wireless Communications and Signal Processing (WCSP)*, 2019.