

CS6700: Reinforcement Learning

Assignment 2



Royyuru Sai Prasanna Gangadhar - ME19B190

Abhigyan C - EE19B146

Course Professor : Prof. Balaraman Ravindran

INDEX

March 29, 2023

Contents

Contents	i
List of Figures	i
1 Abstract	1
2 Environment Description	1
3 DQN	1
3.1 CartPole-v1	1
3.2 Prioritized Sampling from Replay Buffer	1
3.3 Methodology followed:	2
3.4 Acrobot-v1	5
3.5 MountainCar-v0	8
3.6 Discretization of State Space for mountain car	8
4 Actor-Critic	10
4.1 CartPole-v1	12
4.2 Acrobot-v1	15
4.3 MountainCar-v0	17
5 Conclusions	17

List of Figures

1	Summary of the results of hyperparameter tuning for CartPole-v1	4
2	Initial Hyper parameter setting	4
3	Performance of the tutorial DQN agent	5
4	Best Hyper parameter setting for Cartpole DQN	5
5	Summary of the results of hyperparameter tuning for Acrobot-v1	6
6	Initial Hyper parameter Settings for Acrobot	6
7	Performance of Acrobot with respect to tutorial performance	7
8	Best Hyper parameter setting for Acrobot using DQNs	7
9	Initial setting of Acrobot	9
10	Performance with respect to the tutorial hyperparameters	9
11	Success with the use of state space discretization for mountain car	10
12	Actor Critic Method on the OpenAI Cartpole-v1 Environment for 3 different runs	14
13	Best full return for Cartpole	14

14	Actor Critic using N-step return for Cartpole-v1	15
15	Actor Critic Method on the OpenAI Acrobot-v1 Environment for 3 different runs	16
16	Actor Critic using Full-step return for Acrobot-v1	16
17	Actor Critic using N-step return for Acrobot-v1	17
18	Relative importance of hyperparameters and its correlation with the performance metric	18

1 Abstract

In this assignment, we explore Deep Q-Networks and Actor Critic Methods to find the best possible solution for 3 different OpenAI Gym RL environments with discrete actions - Acrobot-v1, CartPole-v1 and MountainCar-v1.

The aim is to find the best hyperparameters and configurations for the three OpenAI gym environments CartPole, Acrobot and MountainCar.

In order to keep the report concise only important plots and tables have been added in this assignment. All the plots that have been obtained during the working of assignment have been stored in a folder and zipped together along with the code files and the report. Code snippets where ever necessary have been added in this report. The performance metric that has been used to compare various models is chosen to be the number of episodes that an agent takes to solve the environment

2 Environment Description

- **Acrobot-v1:** The acrobot system includes two joints and two links, where the joint between the two links is actuated. Initially, the links are hanging downwards, and the goal is to swing the end of the lower link up to a given height.
- **CartPole-v1:** A pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The system is controlled by applying a force of +1 or -1 to the cart. The pendulum starts upright, and the goal is to prevent it from falling over. A reward of +1 is provided for every timestep that the pole remains upright. The episode ends when the pole is more than 15 degrees from vertical, or the cart moves more than 2.4 units from the center.
- **MountainCar-v0:** A car is on a one-dimensional track, positioned between two "mountains". The goal is to drive up the mountain on the right; however, the car's engine is not strong enough to scale the mountain in a single pass. Therefore, the only way to succeed is to drive back and forth to build up momentum.

3 DQN

3.1 CartPole-v1

Now, let us implement our DQN agent on the cart pole environment.

Techniques Used: Prioritized Replay Buffer, Target network and Hyper-parameter tuning

3.2 Prioritized Sampling from Replay Buffer

One way of sampling from the replay buffer is to sample a random set of experiences of a given batch size. Instead we sample using a technique called prioritized sampling. We sample batches in proportion to the magnitude of their TD error. Thus, we preferentially learn the values of states that have been approximated poorly. This helps to speed up the training process. Prioritized sampling drastically decreases the number of training episodes but increases the computational overhead considerably as well.

With prioritized sampling not only the computation increases, but also the number of hyperparameters the bias correction(a) in finding the probability and the weight decay in the importance sampling weights(b). To simplify our analysis the weight decay is considered the same as the decay in epsilon and we keep the value of a as a hyperparameter to get tuned.

3.3 Methodology followed:

Listing 1: DQN Algorithm

```
1 state_shape = env.observation_space.shape[0]
2 action_shape = env.action_space.n
3 env = gym.make('Acrobot-v1') #Replace this with the appropriate ...
   environment for CartPole and MountainCar
4
5 def dqn(agent,n_episodes, max_t, eps_start, eps_end, eps_decay,a):
6     #necessary initializations
7     for i_episode in range(1, n_episodes+1):
8         state = env.reset() #initializes the environment
9         score = 0
10        step = 0
11        for t in range(max_t):
12            action = agent.act(state, eps) #action taken in an epsilon ...
               greedy fashion
13            next_state, reward, done, _ = env.step(action)
14            agent.step(state, action, reward, next_state, done,a = a) ...
               #updates the replay buffer with the current experience and ...
               updates the weights of neural network
15            state = next_state
16            score += reward
17            step += 1
18            if done:
19                break
20            steps.append(step)
21            scores_window.append(score)
22            scores_window_printing.append(score)
23            scores_exit_bad_params.append(score)
24
25            eps = max(eps_end, eps_decay*eps)
26            #necessary print statements are added, please refer to the code ...
               files submitted
27            #Termination criteria
28            if np.mean(scores_window)>=-100.0: #termination criteria
29                print('\nEnvironment solved in {:d} episodes!\tAverage Score: ...
                   {:.2f}'.format(i_episode-100, np.mean(scores_window)))
30                break
31        return [np.array(steps),np.array(scores),i_episode,False]
```

Listing 2: Prioritize replay buffer

```
1
2 def modify_priorities(self,indices,TD_errors,offset=0.001):#
3     #updates the value of the priorities of experiences
4     np.put(self.priorities ,indices, np.abs(TD_errors).flatten()+offset)
5
6 def return_probabilities(self, a):#a = priority scale
7     #returns the probability distribution of the priorities of each ...
       experience, with a as the bias correction factor
8     arr = np.array(self.priorities)**a #modify a
9     probab = arr/arr.sum()
10    return probab
11
12 def return_importance_weights(self,probabilities):
13     #return importacne weights correspnding to the current probability ...
       distribution
```

```

14     importance_weights = 1/(probabilities*len(self.memory))
15     normalized_weights = importance_weights/max(importance_weights)
16     return normalized_weights
17
18     def sample(self, batch_size=None, a=1.0):
19         """Randomly sample a batch of experiences from memory."""
20         num_samples = self.batch_size if batch_size == None else batch_size
21         probab = self.return_probabilities(a)
22         get_sample_idxs = np.random.choice(np.arange(len(self.memory)), ...
23             size = num_samples, replace = False, p = probab)
24         importance_weights = ...
25             self.return_importance_weights(probab[get_sample_idxs])
26         experiences = self.memory[get_sample_idxs]
27         states = torch.from_numpy(np.vstack([e[0] for e in experiences if ...
28             e is not None])).float().to(device)
29         actions = torch.from_numpy(np.vstack([e[1] for e in experiences if ...
30             e is not None])).long().to(device)
31         rewards = torch.from_numpy(np.vstack([e[2] for e in experiences if ...
32             e is not None])).float().to(device)
33         next_states = torch.from_numpy(np.vstack([e[3] for e in ...
34             experiences if e is not None])).float().to(device)
35         dones = torch.from_numpy(np.vstack([e[4] for e in experiences if e ...
36             is not None])).astype(np.uint8).float().to(device)
37         return (states, actions, rewards, next_states, ...
38             dones), importance_weights, get_sample_idxs

```

Listing 3: Target Q network

```

1 class TutorialAgent():
2
3     def __init__(self, state_size, action_size, ...
4         seed, num_layers, neurons_each_layer, lr, gamma, update_every, batch_size
5         , max_truncation):
6
7         ''' Agent Environment Interaction '''
8         self.state_size = state_size
9         self.action_size = action_size
10        self.seed = random.seed(seed)
11        self.batch_size = batch_size
12        self.gamma = gamma
13        self.MAX_TRUNCATION = max_truncation
14        self.update_every = update_every
15        ''' Q-Network '''
16        self.qnetwork_local = QNetwork1(seed=seed, state_size= state_size, ...
17            action_size=action_size ...
18            , num_layers=num_layers, neurons_each_layer=neurons_each_layer).
19            to(device)
20        self.qnetwork_target = QNetwork1(seed=seed, state_size= state_size, ...
21            action_size=action_size ...
22            , num_layers=num_layers, neurons_each_layer=neurons_each_layer).
23            to(device)
24        self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=lr)
25
26        ''' Replay memory '''
27        self.memory = ReplayBuffer(action_size, int(1e5), self.batch_size, ...
28            seed=seed)

```

Initial setting: We start our analysis with the network architecture and network hyper parameters given in the tutorial 5. We compare the performance of our model with respect to

the performance of this model given in tutorial 5. We will initially start with a random set of hyper parameters and then tune them to reach an optimal performance. The following are the hyper parameters that are tuned in training a DQN Agent

- Epsilon decay
- Gamma
- Learning rate
- Number of layers of the neural network
- Number of neurons in each layer
- Update frequency of target network
- α = bias correction in priority sampling

To be concise, only the graphs corresponding to the initial set of hyper parameters and final graphs corresponding to optimal values of hyper parameters have been illustrated in this report. However, a complete summary of Hyper parameter tuning of CartPole-v1 has been displayed in the Figure 1.

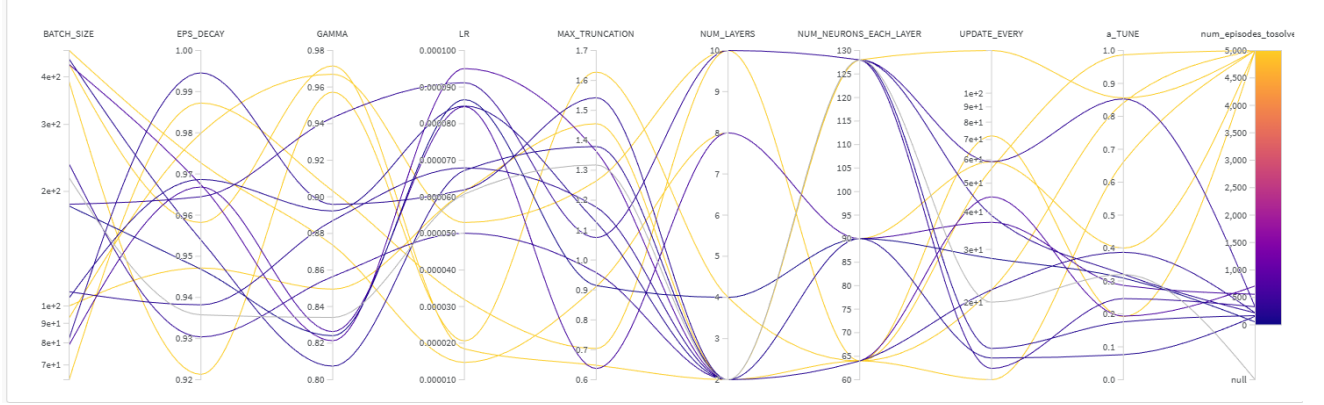


Figure 1: Summary of the results of hyperparameter tuning for CartPole-v1

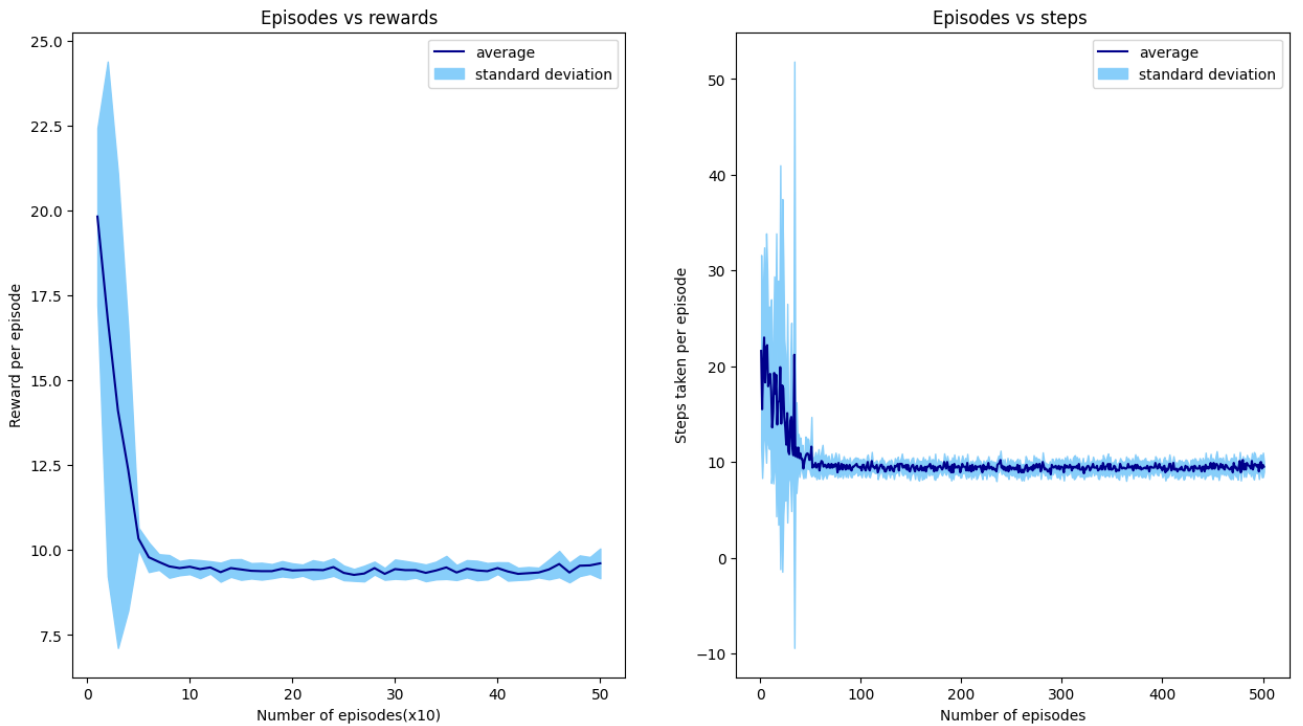


Figure 2: Initial Hyper parameter setting

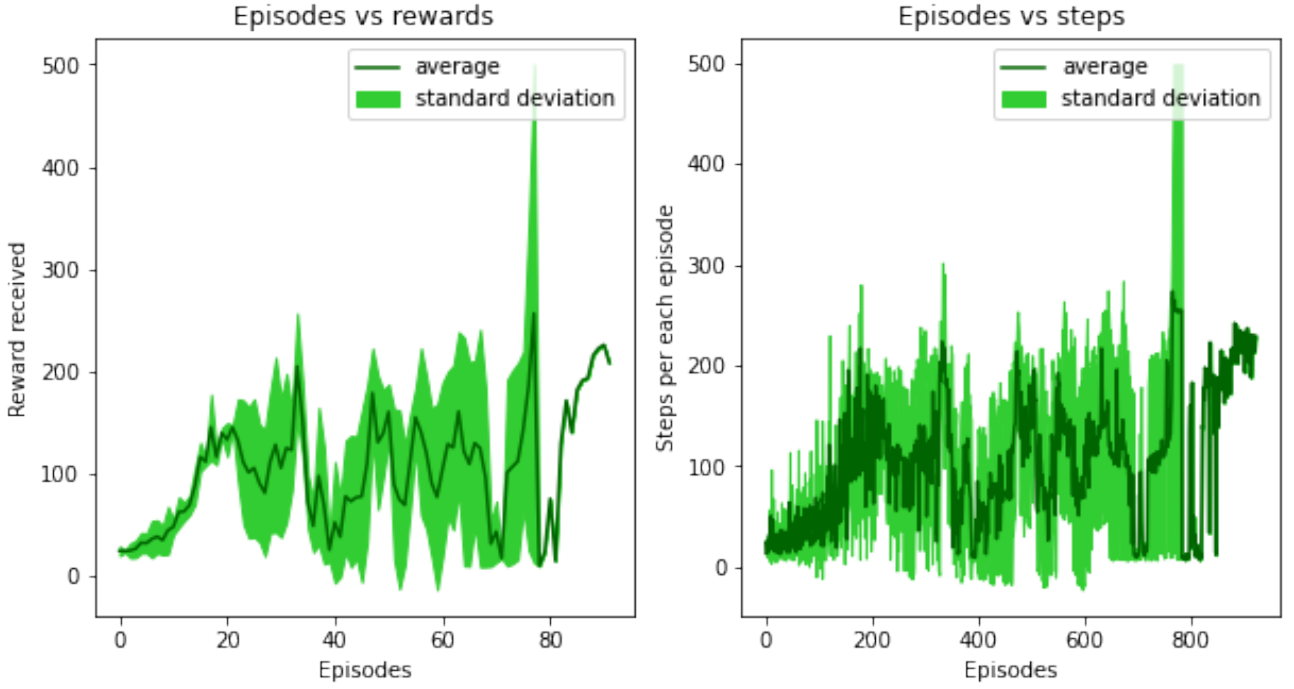


Figure 3: Performance of the tutorial DQN agent

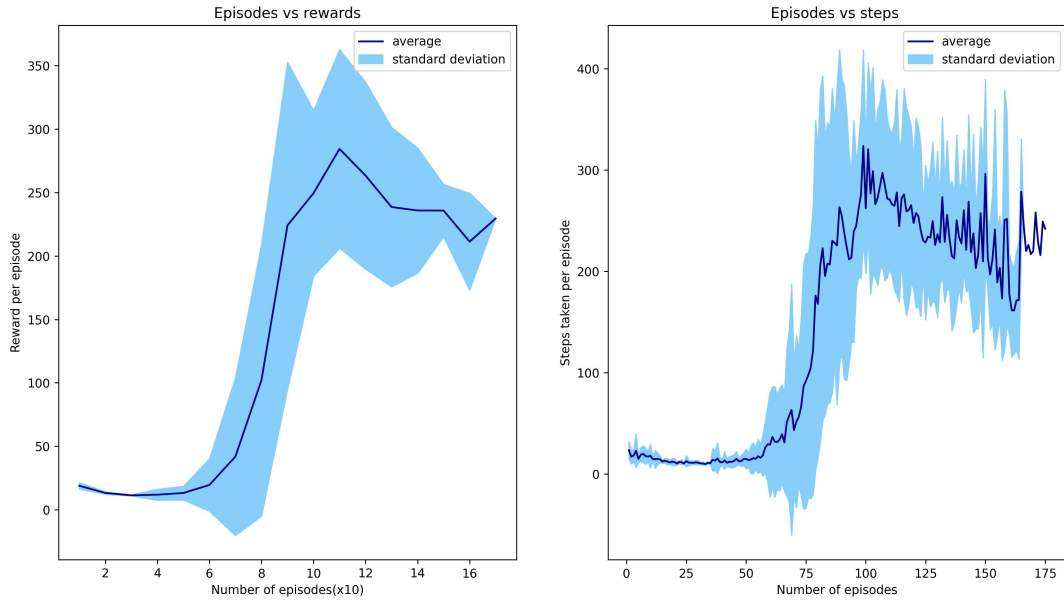


Figure 4: Best Hyper parameter setting for Cartpole DQN

State	Batch size	Epsilon decay	Gamma	Learning rate	Num network layers	Neurons each layer	Update frequency (Target network)	a(bias correction)	Prioritized sampling	Number of episodes to solve the env
Initial Hyperparameter setting	432	0.958	0.971	1.84e-05	2	64	11	0.648	Yes	Did not solve!
Tutorial Agent	64	0.995	0.99	5e-4	2	64	20	None	No	685
Optimal Agent	183	0.946	0.8241	8.6e-5	4	90	28	0.308	Yes	54

3.4 Acrobot-v1

Methodology Followed:

The methodology followed here is the same as that followed in the Cartpole environment, except, we have the following initial hyperparameters:

Following is the summary of the hyperparameters over the course of their tuning:

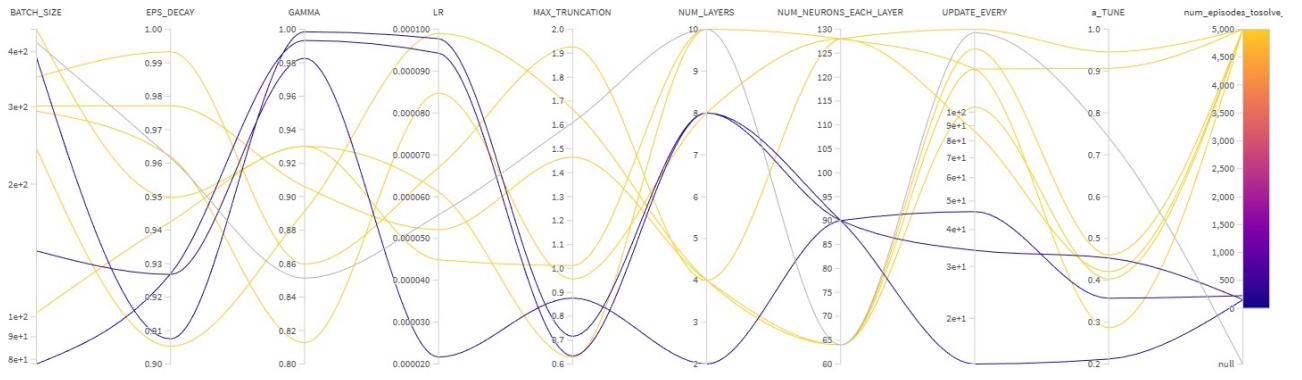


Figure 5: Summary of the results of hyperparameter tuning for Acrobot-v1

State	Batch size	Epsilon decay	Gamma	Learning rate	Num network layers	Neurons each layer	Update frequency (Target network)	a(bias correction)	Prioritized sampling	Number of episodes to solve the env
Initial Hyperparameter setting	450	0.949	0.93	4.49e-05	10	64	140	None	No	Did not solve!
Tutorial Agent	64	0.995	0.99	5e-4	2	64	20	None	No	685
Optimal Agent	388	0.90	0.9985	9.774e-5	8	90	14	None	NO	150

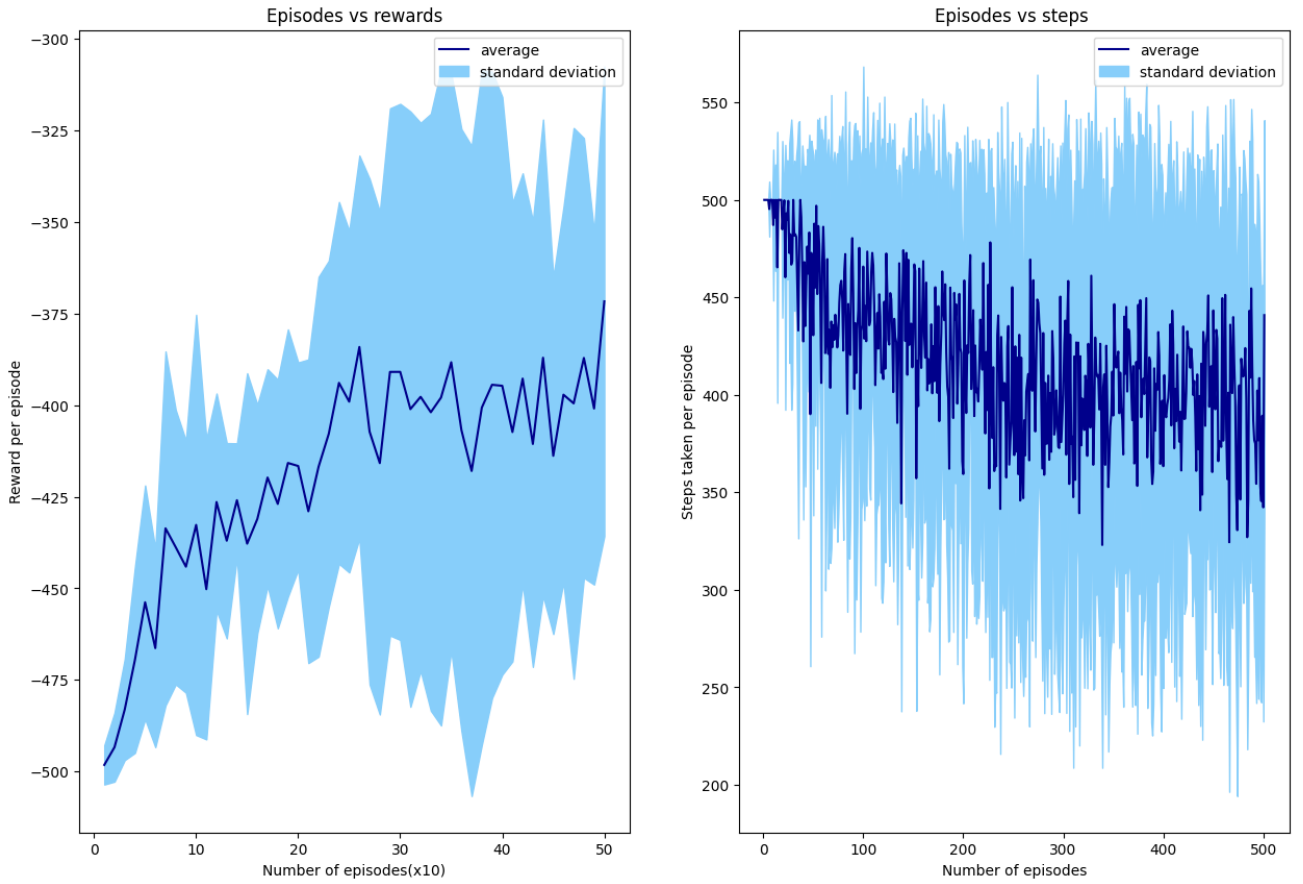


Figure 6: Initial Hyper parameter Settings for Acrobot

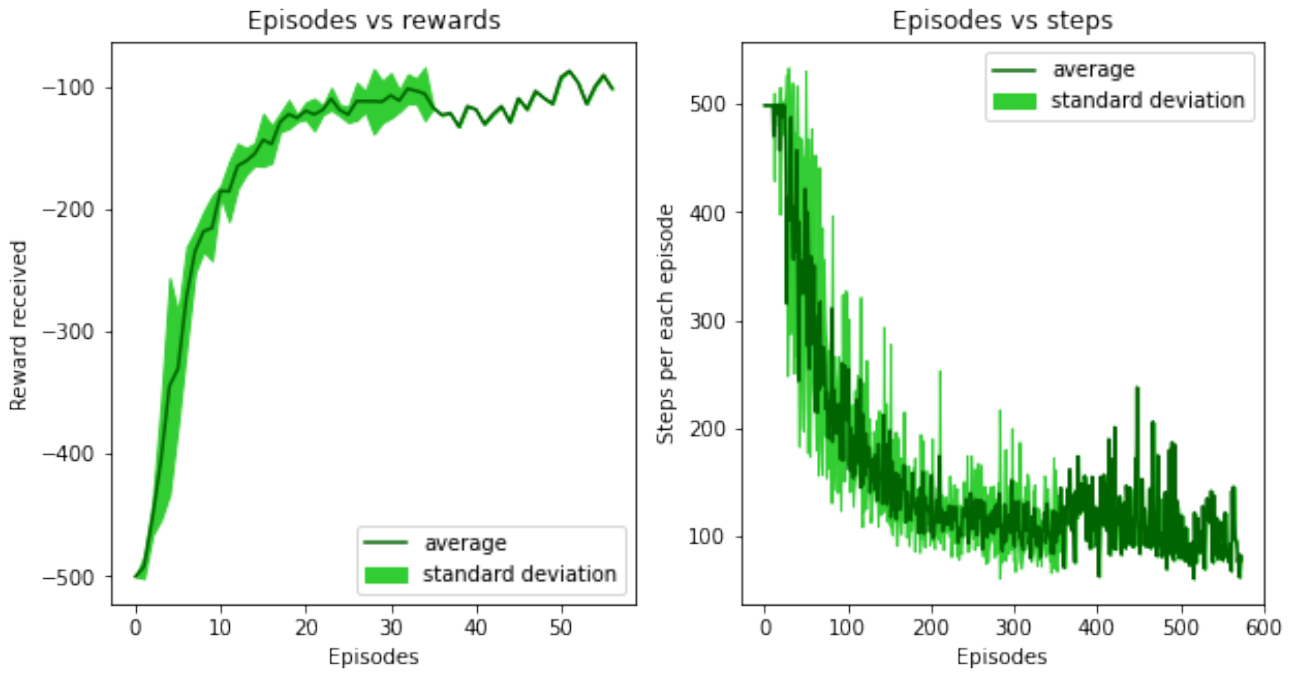


Figure 7: Performance of Acrobot with respect to tutorial performance

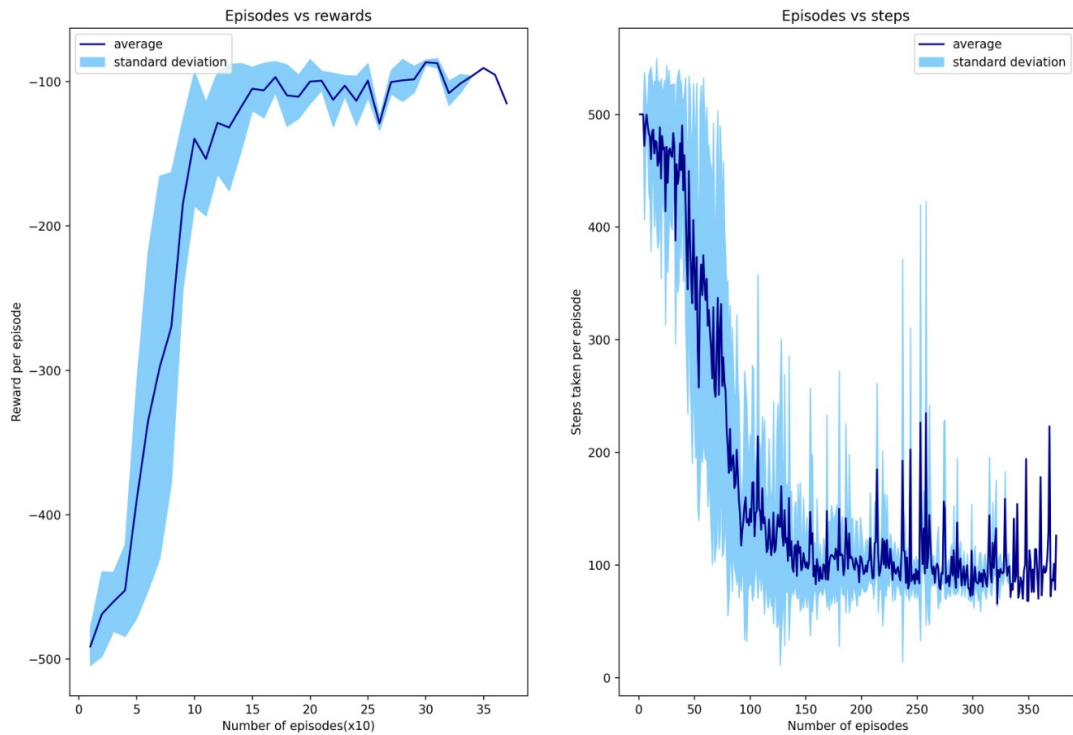


Figure 8: Best Hyper parameter setting for Acrobot using DQNs

We are able to reach a performance, where the agent is able to solve the environment in 154 episodes over an average trial of 154 runs.

3.5 MountainCar-v0

3.6 Discretization of State Space for mountain car

In the MountainCar experiment, due to the low speed of convergence, the model did not converge easily. We therefore chose to discretize the state space (varying from -1.2 to 0.6) into steps of 0.1, leading to 19 discrete states. This might help the agent to learn better in continuous state spaces, as the agent is having very limited exploration (with a maximum time steps of 200) in each episode. We also in the same way multiply the velocity by 50 and round it of to nearest integer to discretize it. We have also increased the truncation steps to 500 instead of 200.

The agent now learns a value function over these states with the help of a Deep-Q Network. Let us compare our results before and after implementation of discrete state spaces.

Listing 4: Discretization of continuous states for MountainCar-v0

```
1 for i_episode in range(1, n_episodes+1):
2     state = env.reset()
3     score = 0
4     step = 0
5     for t in range(max_t):
6         state_adj = (state-env.observation_space.low)*np.array([10,50])
7         state_adj = np.round(state_adj,0).astype(int)
8         action = agent.act(state_adj, eps)
9         next_state, reward, done, _ = env.step(action)
10        next_state_adj = ...
11        (next_state-env.observation_space.low)*np.array([10,50])
12        next_state_adj = np.round(next_state_adj,0).astype(int)
13        agent.step(state_adj, action, reward, next_state_adj, done,a = a)
14        state = next_state
15        score += reward
16        step += 1
17        if done:
18            break
```

State	Batch size	Epsilon decay	Gamma	Learning rate	Num network layers	Neurons each layer	Update frequency (Target network)	Discrete States	Prioritized sampling	Number of episodes
Initial Hyperparameter setting	314	0.90	0.81	0.002	2	64	15	No	No	-500
Tutorial Agent	64	0.995	0.99	5e-4	2	64	20	No	No	-300
Optimal Agent	64	0.995	0.99	5e-4	2	64	20	Yes	NO	-150

Before discrete state space implementation: The following picture illustrates the learning of Mountain car before state space discretization

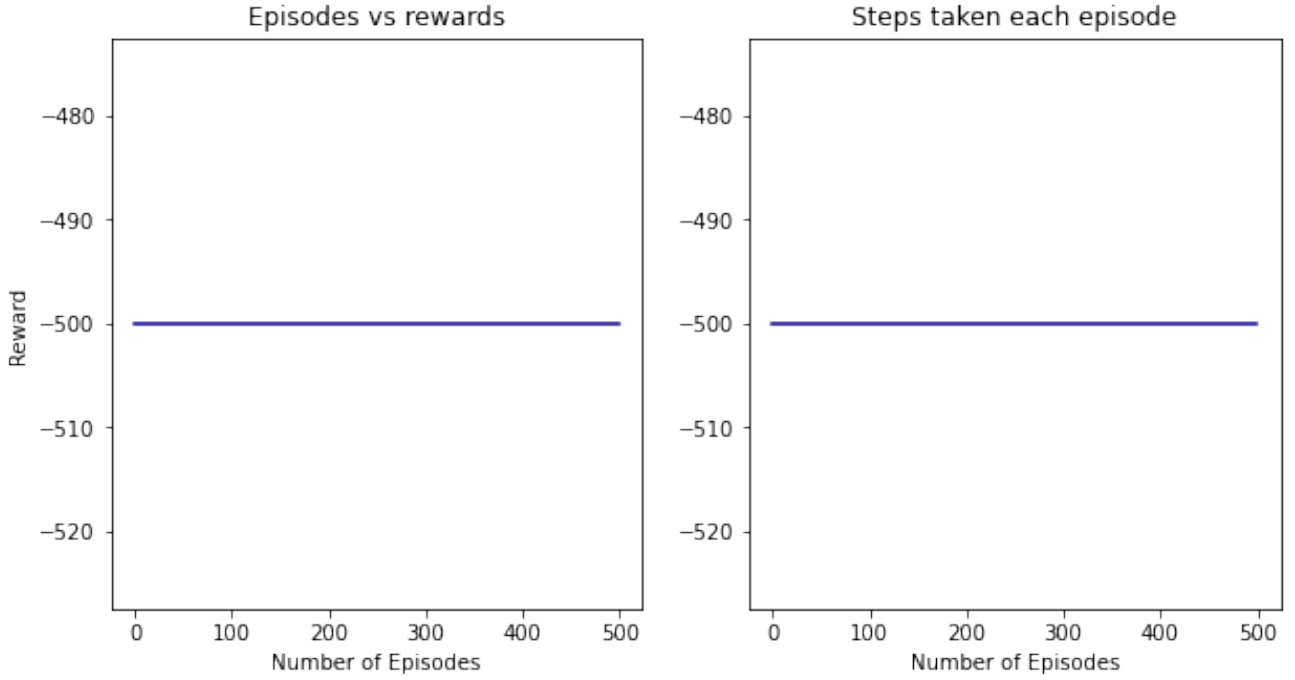


Figure 9: Initial setting of Acrobot

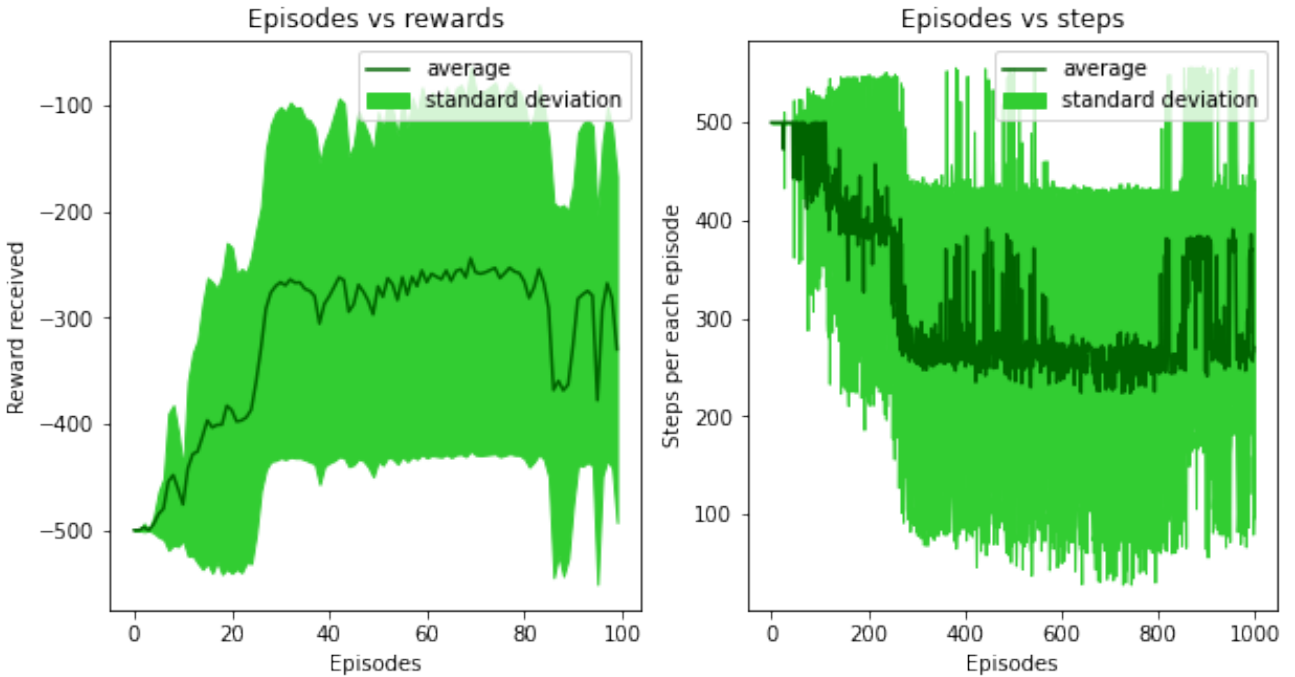


Figure 10: Performance with respect to the tutorial hyperparameters

Results for mountain car after state space discretization: The results are much better when we discretize the state space for mountain car, when the training is done for 5000 episodes. The truncation steps of mountain car has been increased to 500.

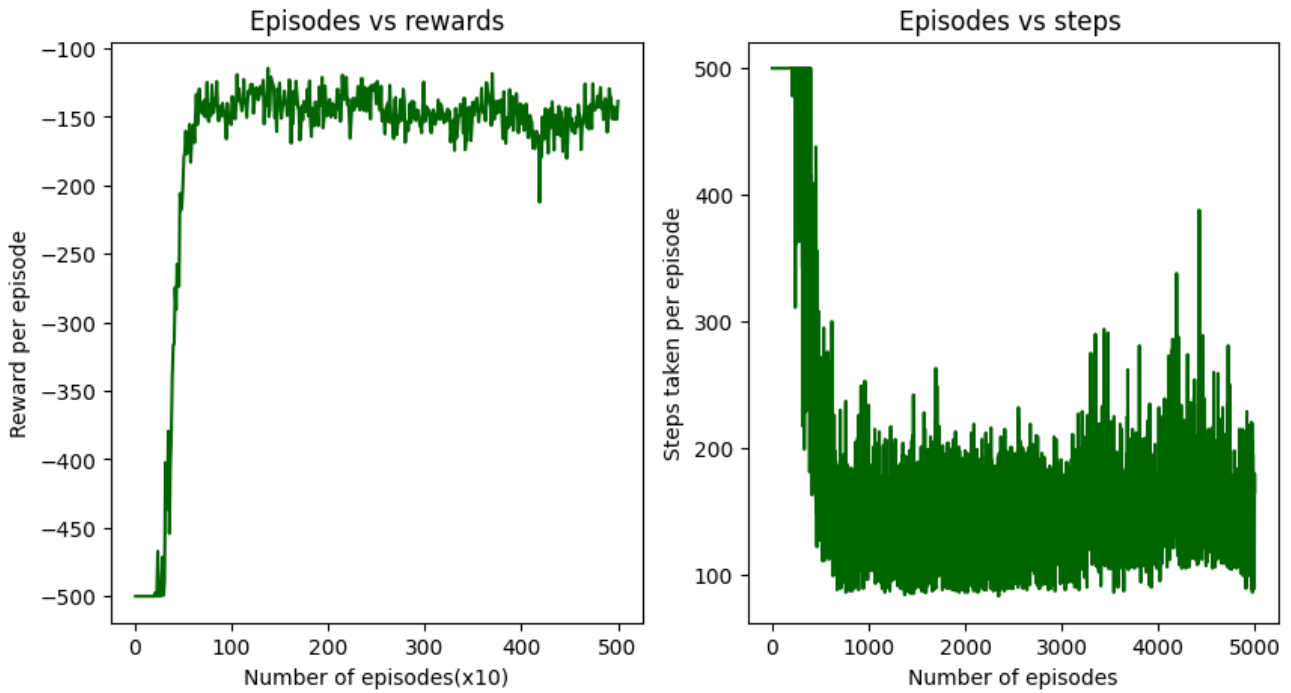


Figure 11: Success with the use of state space discretization for mountain car
Although not illustrated in this report, some success compared to normal training of DQN is observed when we use prioritized sampling for mountain car. The learning through this approach is very very slow and not very significant (reached a average reward value of approximately -400).

4 Actor-Critic

Listing 5: Neural network architecture that acts as "Actor Critic"

```

1 class ActorCriticModel(tf.keras.Model):
2     def __init__(self, ...
      action_size,num_layers=2,num_neurons_each_layer=[1024,512]):
3         super(ActorCriticModel, self).__init__()
4         #defining the hidden layers of the model
5         self.num_layers =num_layers
6         self.neurons_each_layer = num_neurons_each_layer
7         self.linears = [tf.keras.layers.Dense(self.neurons_each_layer[0], ...
          activation=tf.nn.relu)]
8         self.linears.extend([tf.keras.layers.Dense(self.neurons_each_layer[i], ...
          activation=tf.nn.relu) for i in range(1, self.num_layers)])
9         #outputs the policy and the value of the state
10        self.pi_out = tf.keras.layers.Dense(action_size, ...
          activation=tf.nn.softmax)
11        self.v_out = tf.keras.layers.Dense(1)
12
13    def call(self, state):
14        h = None
15        for i in range(self.num_layers+1):
16            if i == 0:
17                h = tf.nn.relu(self.linears[0](state))
18            elif i < self.num_layers:
19                h = tf.nn.relu(self.linears[i](h))
20            else:
21                return self.pi_out(h), self.v_out(h)

```

The one step TD loss and learning the weights from the observed TD loss is defined as below.

Listing 6: Learn function

```

1  def actor_loss(self, action, pi, Δ):  #actor loss
2
3      return -tf.math.log(pi[0,action]) * Δ
4
5  def critic_loss(self,Δ):  #critic loss
6      return Δ**2
7
8  @tf.function #function used to update the weights of the network.
9  def learn(self, state, action, reward, next_state, done):
10
11      with tf.GradientTape(persistent=True) as tape:
12          pi, V_s = self.ac_model(state)
13          _, V_s_next = self.ac_model(next_state)
14
15          V_s = tf.squeeze(V_s) #get the value of current and next ...
16                               states using Actorcritic neural network
17          V_s_next = tf.squeeze(V_s_next)
18          #TD return update
19          Δ = reward+((self.gamma)*V_s_next) - V_s
20          loss_a = self.actor_loss(action, pi, Δ)
21          loss_c = self.critic_loss(Δ)
22          loss_total = loss_a + loss_c
23
24      gradient = tape.gradient(loss_total, ...
25                              self.ac_model.trainable_variables)
26      self.ac_model.optimizer.apply_gradients(zip(gradient, ...
27                                                  self.ac_model.trainable_variables))

```

The full TD return algorithm for actor critic is defined as follows

Listing 7: Full TD return

```

1  def learn_full_return(self, states, actions, rewards, next_states, dones,T):
2  #get the complete trajectory and iteratively loop over each time and ...
3  calculate the complete return and then accumulate it in reward.
4      with tf.GradientTape(persistent=True) as tape:
5          Δ_ts = []
6          pis = []
7          for t in range(0,T):
8              Δ_t = 0
9              for t_ in range(t,T):
10                  Δ_t += (((self.gamma)**(t_-t))*rewards[t_+1] if t_+1< ...
11                      len(rewards) else 0) #calculating the total return
12
13          pi, V_s = self.ac_model(states[t]) #get the value of current state
14          V_s = tf.squeeze(V_s)
15          Δ_t = Δ_t - V_s #Full TD return
16          Δ_ts.append(Δ_t)
17          pis.append(pi)
18          loss_a = 0
19          loss_c = 0
20          for i in range(T): #loop through all time steps and accumulate ...
21              the actor loss and critic loss
22              loss_a += self.actor_loss(actions[i], pis[i], Δ_ts[i])
23              loss_c +=self.critic_loss(Δ_ts[i])
24          loss_total = loss_a + loss_c

```

```

22     gradient = tape.gradient(loss_total, ...
        self.ac_model.trainable_variables, ...
        unconnected_gradients=tf.UnconnectedGradients.ZERO)
23     self.ac_model.optimizer.apply_gradients(zip(gradient, ...
        self.ac_model.trainable_variables))

```

Listing 8: N step TD return

```

1  def learn_nstep_return(self, states, actions, rewards, next_states, done, ...
    n, T):
2  #get the entire trajectory but at each time step, only accumulate the ...
    rewards till t+n and approximate the remaining return with the state ...
    s(t+n) and then calculate the n-step return
3      with tf.GradientTape(persistent=True) as tape:
4          Δts = []
5          pis = []
6          for t in range(0,T):
7              Δt = 0
8              for t_ in range(t,t+n-1):
9                  Δt += (((self.gamma)**(t_-t))*rewards[t_+1] if t_+1< ...
                          len(rewards) else 0)
10             pi, V_s = self.ac_model(states[t])
11             _, V_tn = self.ac_model(states[t+n]) if t+n < T else 0,0
12             Δt += (self.gamma**n*V_tn - V_s)
13             V_s = tf.squeeze(V_s)
14             pi, V_s = self.ac_model(states[t])
15             Δt = Δt - V_s
16             Δts.append(Δt)
17             pis.append(pi)
18             loss_a = 0
19             loss_c = 0
20             for i in range(T):
21                 loss_a += self.actor_loss(actions[i], pis[i], Δts[i])
22                 loss_c += self.critic_loss(Δts[i])
23                 loss_total = loss_a + loss_c
24                 gradient = tape.gradient(loss_total, ...
                    self.ac_model.trainable_variables, ...
                    unconnected_gradients=tf.UnconnectedGradients.ZERO)
25                 self.ac_model.optimizer.apply_gradients(zip(gradient, ...
                    self.ac_model.trainable_variables))

```

Since, we are averaging over multiple runs, we need to be careful of the averaging strategy as not all runs end at the same termination value.

Listing 9: Calculating the mean of various arrays that have different lengths

```

1  def calculate_mean(array):
2      lens = [len(i) for i in array]
3      arr = np.ma.empty((np.max(lens), len(array)))
4      arr.mask = True
5      for idx, l in enumerate(array):
6          arr[:len(l), idx] = l
7      return arr.mean(axis = -1), arr.std(axis=-1)

```

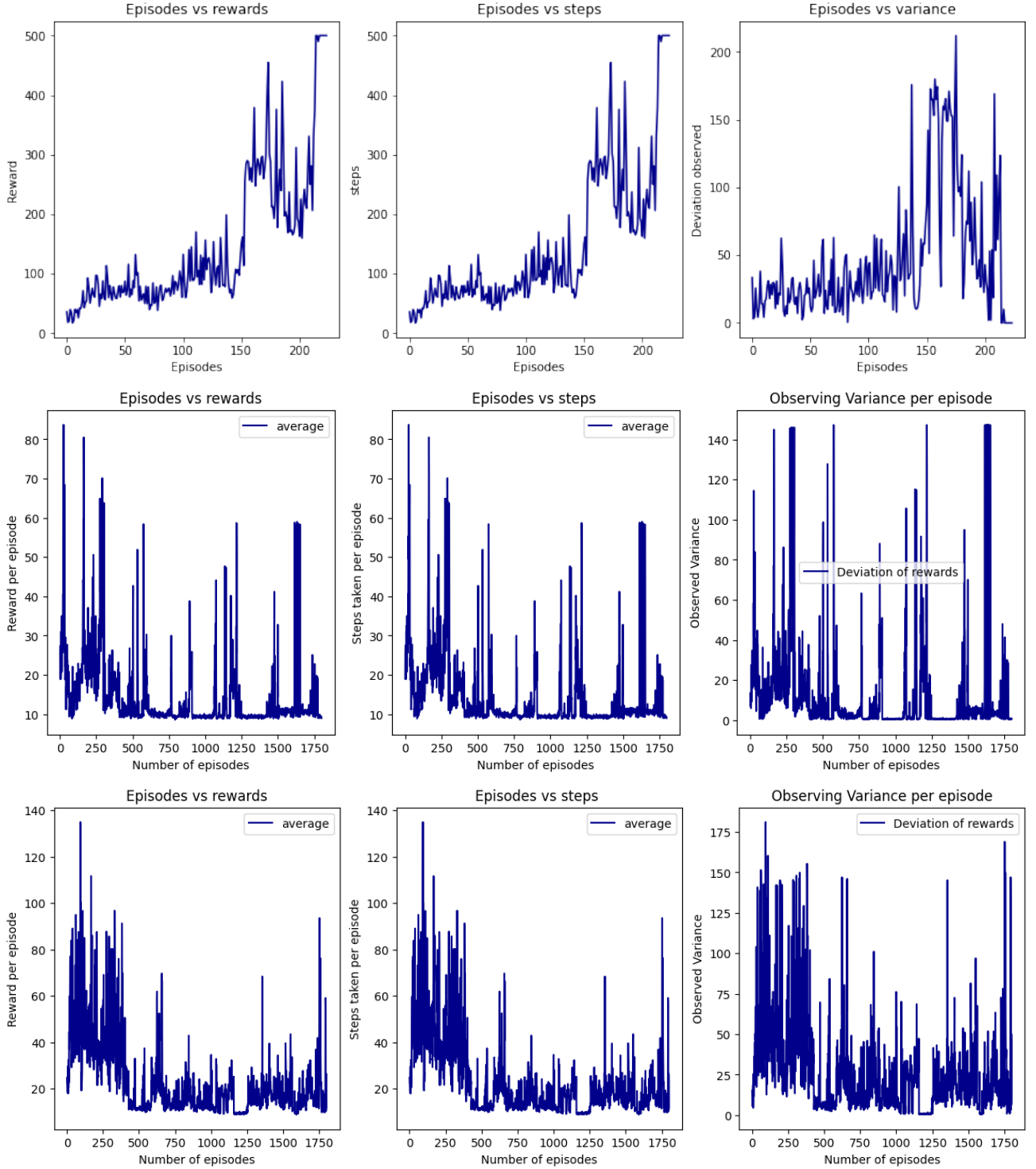
4.1 CartPole-v1

The Cartpole algorithm is run for various actor critic methods, including one step actor critic, n step actor critic and full returns actor critic. By intuition, one step actor critic has a higher

variance compared to n-step actor critic and n-step has a higher variance compared to full returns actor critic, let us compare our results.

The graphs are plotted for the following configurations:

State	Learning rate	Num_episodes	Num_layers	Num_neurons_each_layer
Tutorial agent	1e-05	1800	2	1024,512
Hyperparameter setting 1	8.8e-05	1800	08	724
Hyperparameter setting 2	8.62e-05	2000	02	724
Hyperparameter setting 3	4.66e-05	1800	08	724



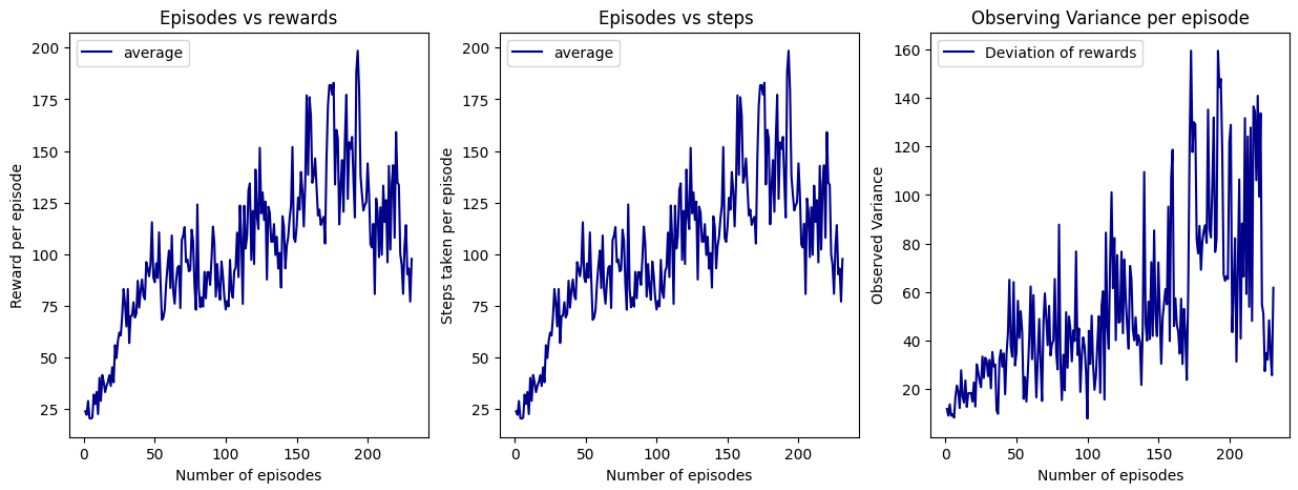


Figure 12: Actor Critic Method on the OpenAI Cartpole-v1 Environment for 3 different runs
Now we plot the best full returns for cartpole:

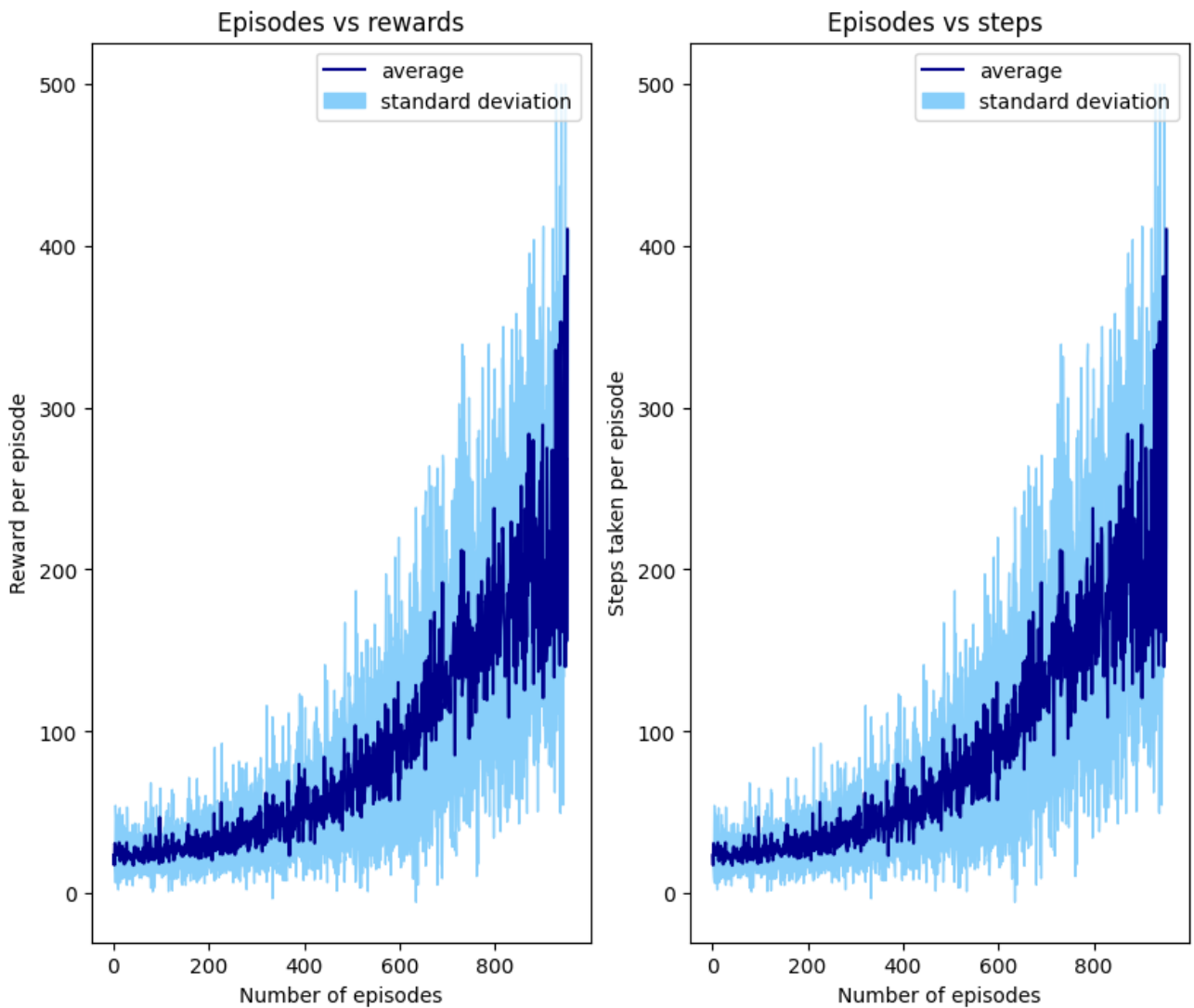


Figure 13: Best full return for Cartpole

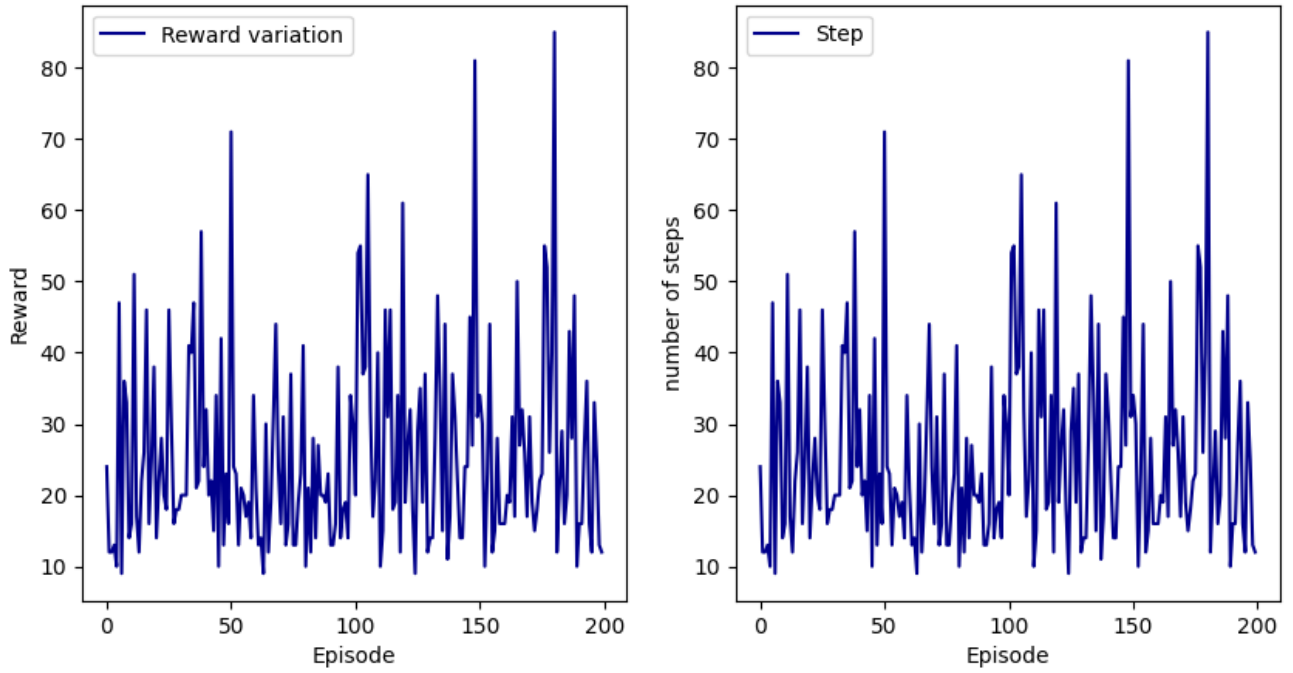


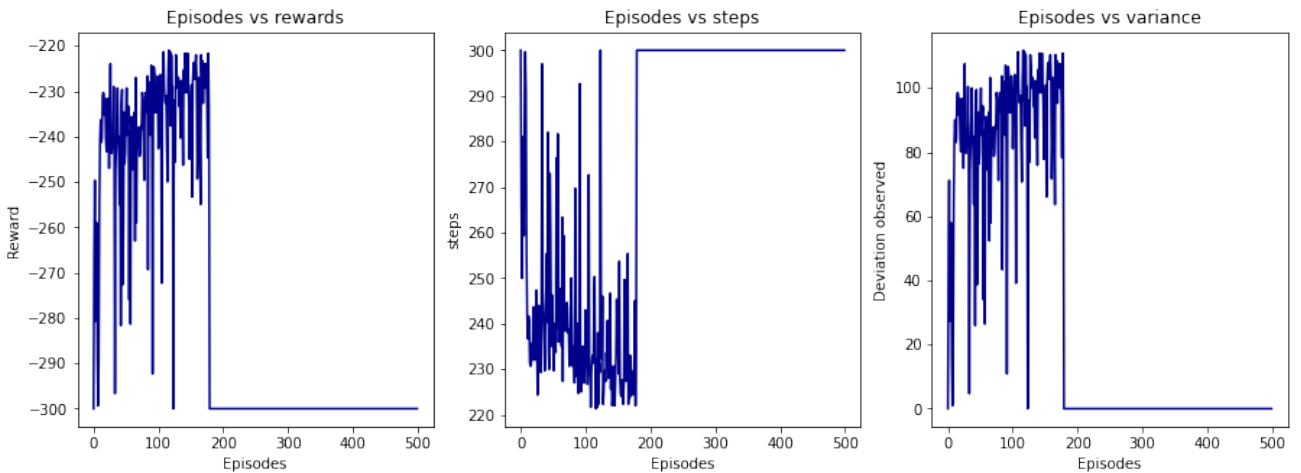
Figure 14: Actor Critic using N-step return for Cartpole-v1

4.2 Acrobot-v1

Actor-critic models have been successful on a wide range of reinforcement learning tasks, but they may not perform as well on certain tasks, such as in this case with the acrobot-v1 environment. The Acrobot-v1 environment is quite complicated and exploring to find the winning trajectory proves to be quite tough for the Actor-Critic Model with only a few runs (owing to limited computational resources and time).

Hence, we have the following trajectories, with only 1 of the 3 giving us results that result in eventual convergence:

State	Learning rate	Num_episodes	Num_layers	Num_neurons_each_layer
Tutorial agent	1e-05	1800	02	1024,512
Hyperparameter setting 1	4.8e-04	2000	08	128
Hyperparameter setting 2	1e-04	500	02	128



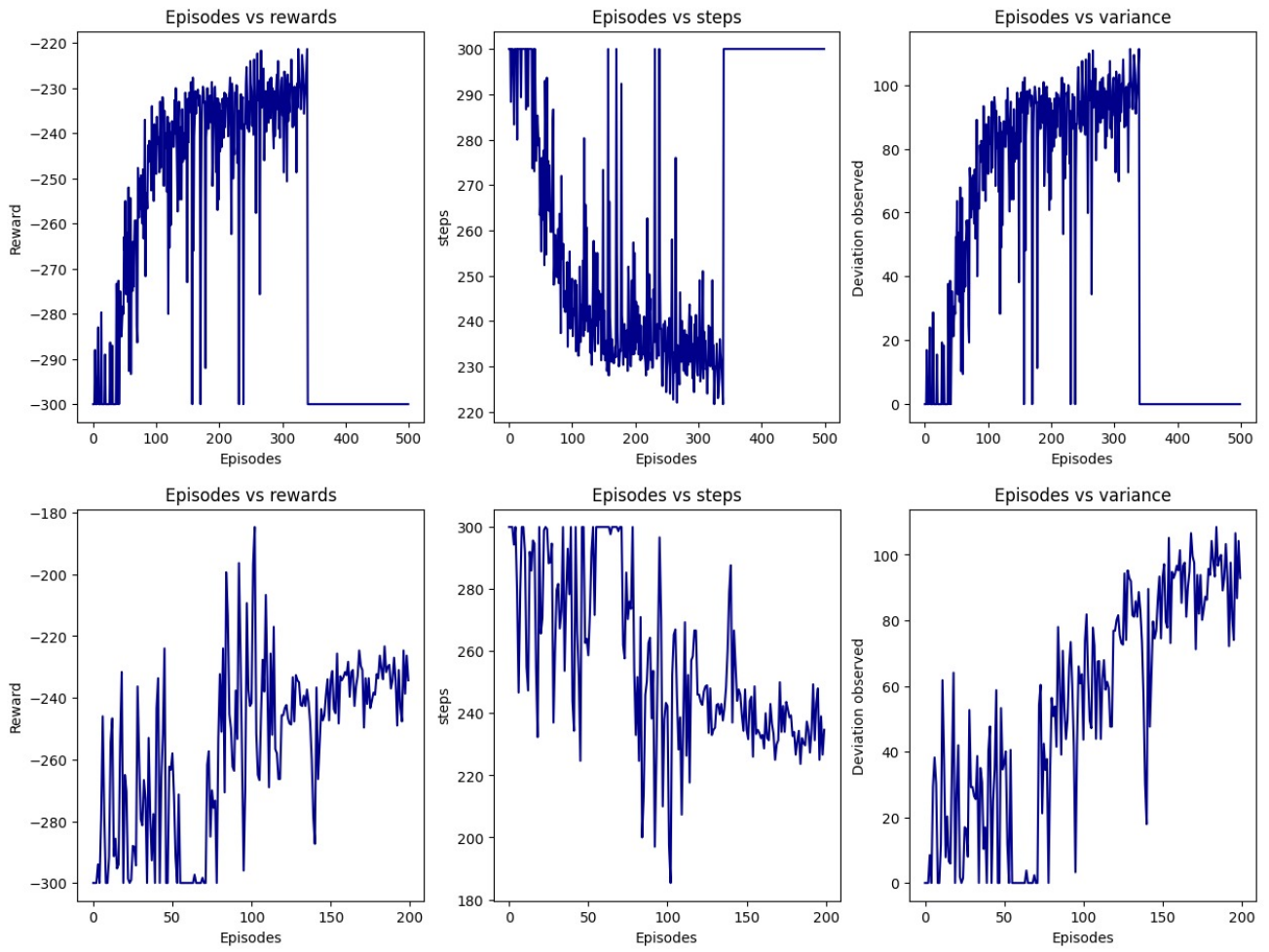


Figure 15: Actor Critic Method on the OpenAI Acrobot-v1 Environment for 3 different runs
Now we plot the best full returns for cartpole:

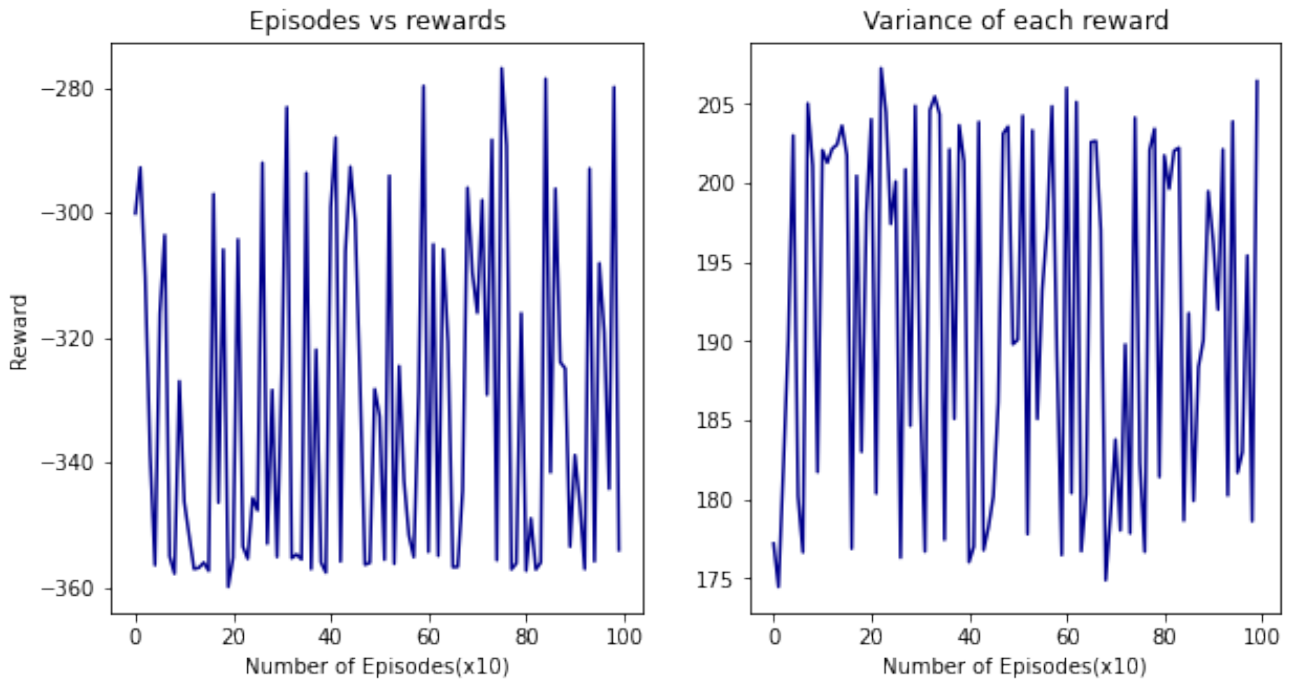


Figure 16: Actor Critic using Full-step return for Acrobot-v1

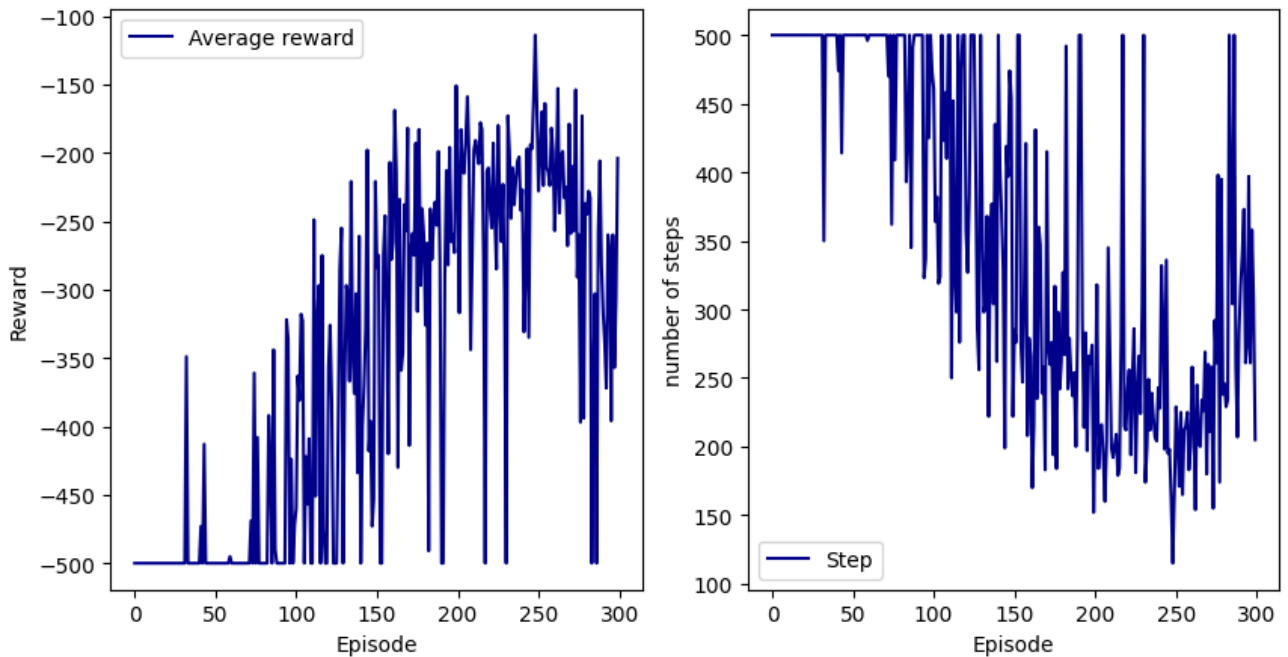


Figure 17: Actor Critic using N-step return for Acrobot-v1

4.3 MountainCar-v0

In the MountainCar-v0 environment of OpenAI Gym, the agent controls a car that must reach the top of a hill, but its engine is not powerful enough to climb the hill directly. Instead, the agent must learn to rock back and forth to gain enough momentum to get to the top. This reward function is rather tough to optimize due to the tricky shape - the actor critic model is unable to generalize to this environment as a result of this trickiness and often ends up criticizing winning algorithms.

Owing to the computational complexity, we were unable to get the Actor-Critic model to converge despite trying a large amount of hyperparameters.

Many hyperparameter settings were tried for mountain car and many of them did not even budge and had a constant reward of -200 when the system is set to truncate at 200 steps and has a constant reward of -500 when the system is set to truncate at -500.

Mountain car again did not improve its performance in one step, n-step Actor critic and full return actor critic and has been constantly maintaining a return of -200, even multiple configurations of the network and truncation criteria are used and thus has not been plotted.

5 Conclusions

We found much better results with the Deep-Q Networks than with the Actor-Critic model.

Deep-Q Learning in general is more effective in learning in discrete action environments with low-dimensional state spaces, whereas Actor-Critic would be more effective with continuous action and with more complex state spaces.

For cartpole algorithm, it has been observed that out all hyperparamter settings learning rate is the most important hyperparamter of all, followed by batch size of experiences, which is followed by the update frequency of the target network.

For Acrobot algorithm, number of neurons each layer, maximum value of gradients to which the magnitude of gradients truncated, gamma, and update frequency are some of the important parameters for acrobot environment.

The success of Moutaincar environment depends on the discretization of state space. Results

have been significantly improved in case of mountain car if discretization of state space has been employed.

The below table gives the importance of various hyperparameters in decreasing order of their correlation with the performance metric.

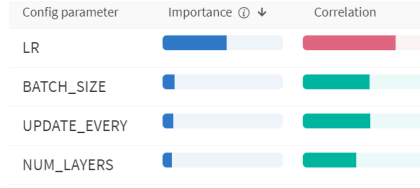


Figure 18: Relative importance of hyperparameters and its correlation with the performance metric

Bayes exploration strategy has been used to identify the best set of hyperparameters. The parameters are chosen in such a way that the distance of new parameters from the current hyperparameters is proportional to the error in current metric chosen. Mountain car has the slowest learning process and it was not able to improve its performance for any wide range of hyperparameter settings in Actor critic. From the graphs that have been plotted for actor critic methods, it is evident that the one step actor critic methods have huge variance in all the three environments namely cartpole, acrobot and mountain car. It is evident from the graphs that have been plotted above. The huge variation in one step actor critic is because of the TD(0) update and less variance in the full return actor critic methods is due to monte-carlo return update. As n-step actor critic is in between the TD(0) update and the monte carlo update, it has variance less than one step, but greater than Full return actor critic. For, Actor critic methods, the value of learning rate is a very important hyperparameter, followed by the neural network architecture. The architecture of neural network determines how well we generalized the action and state space values.

Thus DQN and Actor critic methods have been implemented on the gym environments CartPole, Acrobot and MountainCar.