

CAP 6629: Reinforcement Learning Course - Neural Network Q-Learning GridWorld

Course Project 3

- Shaun Pritchard
- April 7, 2022
- Professor Ni

Pesudo Code

This implementation uses a 5x5 gridworld example, The starting state has the value (0, 0) and the ending state has the value (4, 4). I will be working on an actor-critic approach (ADP) to approximate the Q-value function in a custom grid world problem with an implementation of a neural network to learn the Q-table.

$$e_c(t) = \alpha Q(t) - [Q(t-1) - r(t)]$$

The actor-critic (ADP) uses two networks where the critic network attempts to reduce the error function. Two perceptron networks are assimilated into the actor-critic architecture using stochastic gradient descent (SGD) as the optimization algorithm.

By reducing $e_a(t)$, the actor network leads the agent system towards U_c , the optimal Q-value function objective:

$$e_a(t) = J(t) - U_c(t)$$

While the critic network reduces the error function $e_c(t)$,

The Q-value function $Q(t)$ repeats the passes of state-action pairs through the critic network while it tunes parameters that reduce $e_c(t)$

With the Actor-critic ADP class to implement the parameters. Define 3 method classes: the policy, actor, critic to instantiate the parameters. Actor networks are multilayer neural networks that take the current state as input, and output a probability distribution of the four possible actions that an agent could take i.e. Up, Down, Left, Right. The actor, critic, and policy networks implements two fully connected hidden layers, ReLU is used in each of the fully connected layers which varies based on the number of nodes. We will implement a stochastic gradient descent(SGD) algorithm to

optimize the input through a neural network. Afterwards, we will update the weights of the network by means of backpropagation.

- **Actor** Actor networks have four nodes corresponding to each action in the action space. With its softmax activation function, it identifies a probability distribution for each action the agent will take.

$$E_a(t) = 1/2e_a^2(t)$$

- **Critic** As with the actor network, the critic network t shares the same two fully connected layers, uses the same state input as it does for the actor network, is aware of the actions taken when moving between states and updating the weights, and outputs a predicted reward-to-go or Q learning utility value.

$$E_c(t) = 1/2e_c^2(t)$$

- **Policy** As with the actor and critic networks, the policy network has the same architecture. The NN mainly exists to keep track of the actor network's current plan to reach the goal state, as it outputs a similar probability distribution to the actor network and is not used for training or updating weights.

We will run 75 episodes per test at 6 test with 10 transitions per episode. While this did take longer (48.32 minutes) in google collab using 2BG GPU. I was curious to see the convergence differential based the current parameters. The wait was worth it.

▼ Initialize Project

```
1 # Librarys and Imports
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 import tensorflow as tf
8 from tensorflow.keras import backend as K
9 from tensorflow.keras.optimizers import SGD, Adam
10 from tensorflow.keras.models import Model
11 from tensorflow.keras.layers import Dense, Input
12 tf.__version__
```

▼ NN Actor Critic ADP Agent Class

```
1 # Define Actor-Centric ADP Agent class which will build policy, actor and cri
2
3 class ADP(object):
4     def __init__(self,  $\alpha$ ,  $\beta$ ,  $\gamma=0.9$ , num_actions=4,
5                 layer1_size=2, layer2_size=1, input_dims=2):
6         self. $\alpha$ = $\alpha$  # Alpha networks actor learning rate
7         self. $\beta$ = $\beta$  # Beta networks Critic learning rate
8         self. $\gamma$ = $\gamma$  # Gamma discount factor
9         self.num_actions=num_actions #number of actions
10        self.fc1_dims=layer1_size #number of nodes in first hidden layer
11        self.fc2_dims=layer2_size #number of nodes in second hidden layer
12        self.input_dims=input_dims #dimensionality of state space
13        self.actor, self.critic,self.policy=self.actor_critic_NN() # Create a
14        self.action_space = [i for i in range(self.num_actions)] # Initialize
15
16
17 # State transition of actor
18     def step(self, state, action, reward, done):
19         if self.action_space[action]==0: # Move up
20             if state[1]<4:
21                 state[1] +=1
22             elif self.action_space[action]==1: # Move down
23                 if state[1]>0:
24                     state[1]-=1
25             elif self.action_space[action]==2: # Move left
26                 if state[0]<4:
27                     state[0]+=1
28             elif self.action_space[action]==3: # Move right
29                 if state[0]>0:
30                     state[0]-=1
31             if state[0]==4 and state[1]==4: # Verify the terminal state
32                 done=True
33             return state, reward, done
34
35 # Based on the current predicted policy, this function selects an action from
36     def actions(self, observations):
37         state=observations[np.newaxis, :] # Initilize state and position
38         prob=self.policy.predict(state)[0] # Update actions probability distr
39         action=np.random.choice(self.action_space, p=prob) # Choose an actor
40         return action
41
42 # Using Kearas and Tensorflow to implement a NN function that builds an actor
43     def actor_critic_NN(self):
44         net_input=Input(shape=(self.input_dims,))
```

```

45     dense1=Dense(units=self.fc1_dims, activation='relu')(net_input)
46     dense2=Dense(units=self.fc2_dims, activation='relu')(dense1)
47     probs=Dense(units=self.num_actions, activation='softmax')(dense2)
48     values=Dense(units=1, activation='linear')(dense2)
49     actor=Model(inputs=[net_input], outputs=[probs]) # For each action,
50     actor.compile(optimizer=Adam(learning_rate=self.α), loss='mse') # Cal
51     critic=Model(inputs=[net_input], outputs=[values]) # Critic network
52     critic.compile(optimizer=Adam(learning_rate=self.β), loss='mse') # Ca
53     policy=Model(inputs=[net_input], outputs=[probs]) # he policy netwo
54     return actor, critic, policy
55
56 # Based on an observations in the environment, this function updates the weig
57 def learn(self, state, action, reward, state_, done):
58     state=state[np.newaxis, :] # Current state in network
59     state_=state_[np.newaxis, :] # Next state in network
60     critic_value=self.critic.predict(state) # Current utility value
61     critic_value_=self.critic.predict(state_) # Next utility value
62     target=reward+self.γ*critic_value_*(1-int(done)) # Utility value afte
63     actions=np.zeros((1, self.num_actions)) # Initilize array of number a
64     actions[np.arange(1), action]=1.0 # Set positive actions
65     self.actor.fit(state, actions, verbose=0) #train actor network
66     self.critic.fit(state, target, verbose=0) #train critic network
67

```

▼ Initialize Paramters

```

1  from matplotlib import axis
2  # Implement NN Q-Learning grid world
3  rewards=np.zeros((5, 5)) # Grid World size 5 x 5
4  rewards[4][4]=1 # Reward varaiable
5  episodes=75 # Number of episode in itteration
6  steps_per_episodes_output=[] # output for calulation average steps per episod
7  for i in range(6, 12): # episodes per transition itteratins
8      print(i)
9      agent=ADP(α=0.00001, β=0.00005, layer1_size=2**i, layer2_size=2**(i-1)) :
10     steps_per_episode=[]
11     step_averages=[]
12     # df = pd.DataFrame({"Steps in episode data": [i]})
13     steps_per_episodes_output.append(steps_per_episode)
14     for j in range(episodes):
15         done=False
16         count=0
17         observations=np.array((0, 0))
18         while not done:
19             action=agent.actions(observations) # Action selected

```



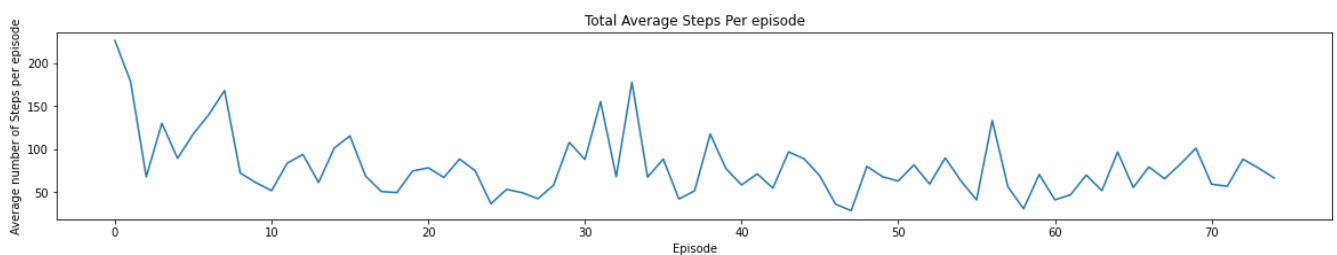
```
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
episode: 11 Number of steps: [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27,
```

▼ Analysis & Results

```
1 print('Average Steps', steps_per_episodes_output)
```

```
Average Steps [[154, 412, 169, 226, 139, 351, 543, 355, 75, 142, 27, 256, 119, 35, 197,
```

```
1 average_steps_per_episode=np.mean(steps_per_episodes_output, axis=0)
2 plt.plot(average_steps_per_episode)
3 plt.rcParams["figure.figsize"] = (20,6)
4 plt.title('Total Average Steps Per episode')
5 plt.xlabel('Episode')
6 plt.ylabel('Average number of Steps per episode')
7 plt.show()
```

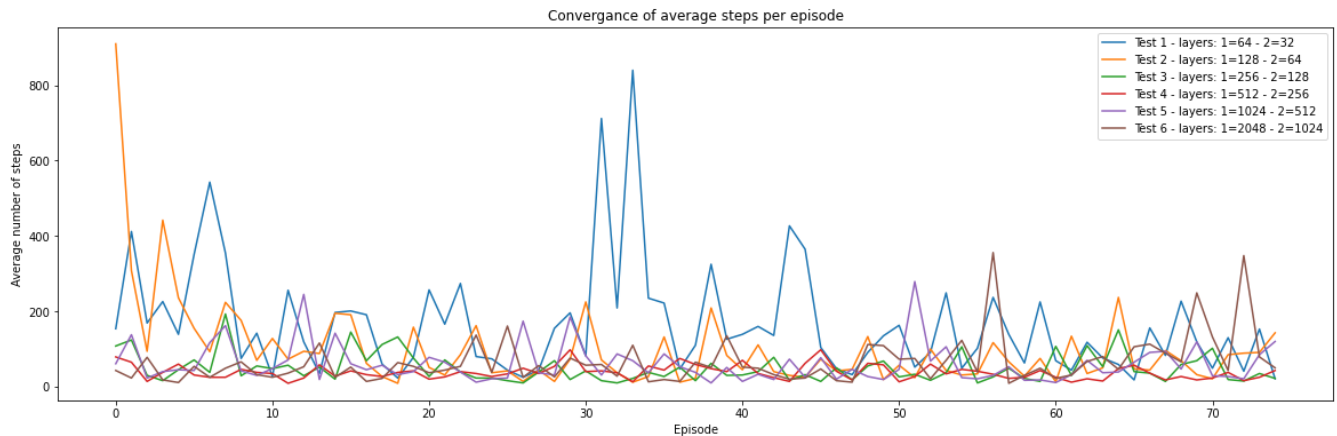


```
1 plt.plot(steps_per_episodes_output[0], label='Test 1 - layers: 1=64 - 2=32')
2 plt.plot(steps_per_episodes_output[1], label='Test 2 - layers: 1=128 - 2=64')
3 plt.plot(steps_per_episodes_output[2], label='Test 3 - layers: 1=256 - 2=128')
4 plt.plot(steps_per_episodes_output[3], label='Test 4 - layers: 1=512 - 2=256')
5 plt.plot(steps_per_episodes_output[4], label='Test 5 - layers: 1=1024 - 2=512')
```

```

5 plt.plot(steps_per_episodes_output[4],.label='Test.5---layers:.1=1024---2=512)
6 plt.plot(steps_per_episodes_output[5],.label='Test.6---layers:.1=2048---2=1024)
7 plt.title('Convergence of average steps per episode')
8 plt.rcParams["figure.figsize"] = (20,3)
9 plt.xlabel('Episode')
10 plt.ylabel('Average number of steps')
11 plt.legend()
12 plt.show()

```



Results:

The computational time complexity was very high in this initial build. Data show that networks with larger hidden layers converge to optimal policies most efficiently. Comparatively, with 1024 layer 1 nodes and 512 layer 2 nodes, the network configuration began high and then converged to the same number of optimal steps as its previous configuration. Due to the randomness in the probability distributions, it appears the configuration with 512 nodes at layer 1 and 256 nodes at layer 2 had a few outliers in terms of number of steps in the episode.

These results suggest that larger networks can create a more efficient solution than smaller ones. In addition, these results are subpar compared to regular Q-learning. If the networks are trained more, or perhaps a negative outcome occurs during implementation, then perhaps they will be comparable.

1. According to the first implementation, a fully connected layer with 64 nodes and a second layer with 32 nodes converged to an optimal policy involving 150-200 steps towards the goal.

2. According to the second implementation, where 128 nodes have been added to the first fully connected layer and 64 have been added to the second fully connected layer, the optimal policy is around 200-250 steps towards the goal.
3. According to the third implementation, which has 256 nodes for the first and 228 nodes for the second fully connected layers, converged to an optimal policy of 150-200 steps.
4. According to the fourth implementation, the first fully connected layer with 512 nodes converges to a policy of around 40-50 steps towards the goal when paired with the second fully connected layer with 256 nodes.
5. According to the fifth implementation, the first fully connected layer with 1024 nodes converges to a policy of around 50-80 steps towards the goal when paired with the second fully connected layer with 512 nodes.
6. According to the sixth implementation, the first fully connected layer with 2048 nodes converges to a policy of around 30-40 steps towards the goal when paired with the second fully connected layer with 1024 nodes.

Compare Project 2:

In project #2 faster computation and complexity with average step to action ratio when meeting the goal based on the algorithm. The algorithms implemented an agent who can enter a new situation with no prior knowledge and a bigger state space. Steps required for the agent to reach its goal are counted in each episode, then averaged over the 1000 iterations. Compared to the other results of the latter algorithms it is uncompetitive.

In comparison to the Neural network with Q-learning using ADP actor critic model. In the neural network variant, the optimal policy is not found at all, with between 20 and 50 steps per episode being an average. It can be challenging to construct a Q-table based on states and actions when the critic network only predicts one utility value in a time. This may change with more training. Nevertheless, due to limited computational resources and time constraints, each experiment was only able to run 75 episodes. In comparison to the algorithms in project 2, the neural networks' initial learning rates were quite small.

✓ 0s completed at 9:41 PM



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.