

EECS 498: Reinforcement Learning

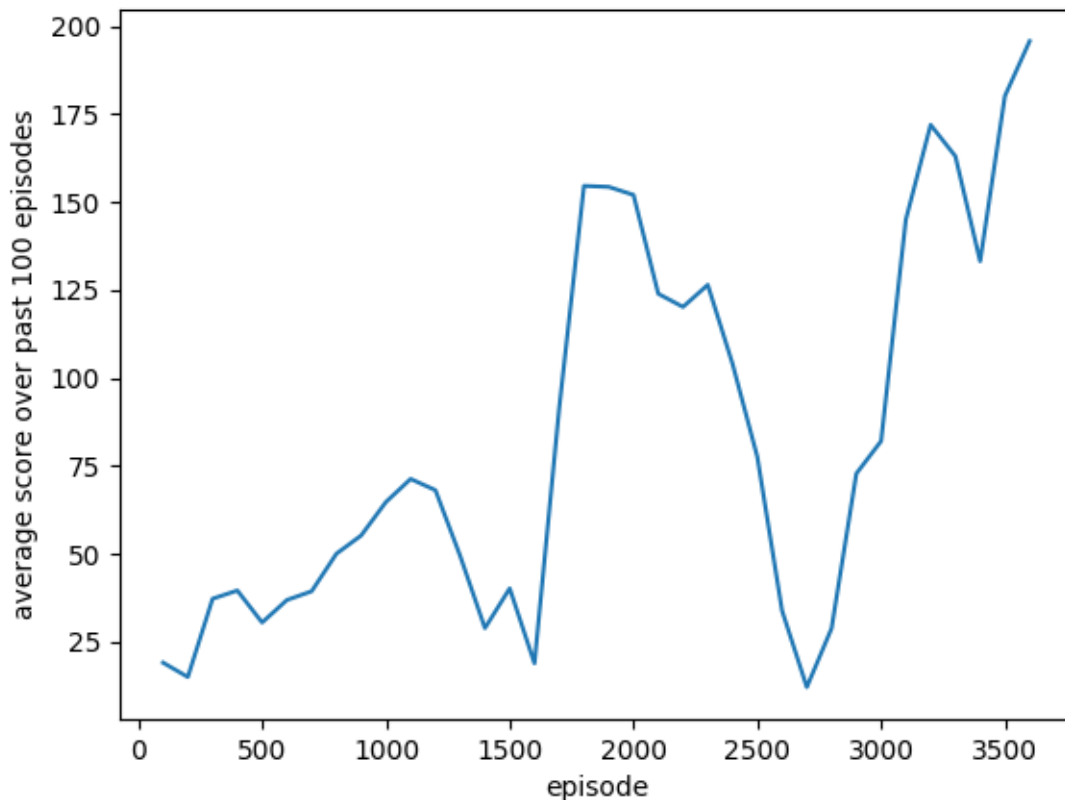
Homework 5 Responses

Tejas Jha
tjha

December 12, 2018

This document includes my responses to Homework 5 questions. Responses that involved the use of coding will provide references to specific lines of code to provide a better overview of how the problem was approached. The code can either be referenced in the Appendix or in the accompanied python script submitted with this assignment.

Question 1



(a)

(b) To try and get the model working on MountainCar-v0,

Suppose the reward function for an MDP is a linear function of d features for a state s

$$R(s) = \alpha_1\phi_1(s) + \alpha_2\phi_2(s) + \dots + \alpha_d\phi_d(s)$$

where the $\phi_1 \dots \phi_d$ are fixed, known and bounded basis functions mapping from state space S to the reals.

As presented in the Algorithms for Inverse Reinforcement Learning paper, we can define a value function for a policy π that maps $S \mapsto A$ for any state s_1 as the following:

$$V^\pi(s_1) = E[R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots | \pi]$$

where the expectation is over the distribution of the state sequence (s_1, s_2, \dots) .

Additionally from the same paper, we can use the notation V_i^π to denote the value function of the policy π in the MDP when the reward function is $R = \phi_i$. We can prove that for any policy π , the value function can be defined as $V^\pi(s_1) = \alpha_1 V_1^\pi + \alpha_2 V_2^\pi + \dots + \alpha_d V_d^\pi$

We can substitute the reward function into the expectation in the value function equation and use the linearity of expectation to expand it to a sum of expectations of fixed and known values. These expectations would simplify to just the terms. This is shown below:

$$\begin{aligned} V^\pi(s_1) &= E[R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots | \pi] \\ V^\pi(s_1) &= E[\alpha_1\phi_1(s_1) + \alpha_2\phi_2(s_1) + \dots + \alpha_d\phi_d(s_1) \\ &\quad + \gamma * (\alpha_2\phi_1(s_2) + \alpha_2\phi_2(s_2) + \dots + \alpha_d\phi_d(s_2)) + \\ &\quad \gamma^2 * (\alpha_1\phi_1(s_3) + \alpha_2\phi_2(s_3) + \dots + \alpha_d\phi_d(s_3)) + \dots | \pi] \\ \\ V^\pi(s_1) &= \alpha_1\phi_1(s_1) + \alpha_2\phi_2(s_1) + \dots + \alpha_d\phi_d(s_1) \\ &\quad + \gamma * (\alpha_2\phi_1(s_2) + \alpha_2\phi_2(s_2) + \dots + \alpha_d\phi_d(s_2)) + \\ &\quad \gamma^2 * (\alpha_1\phi_1(s_3) + \alpha_2\phi_2(s_3) + \dots + \alpha_d\phi_d(s_3)) + \dots \end{aligned}$$

The terms can now be regrouped using the alpha terms.

$$\begin{aligned} V^\pi(s_1) &= \alpha_1 * (\phi_1(s_1) + \gamma\phi_1(s_2) + \gamma^2\phi_1(s_3) + \dots) \\ &\quad + \alpha_2 * (\phi_2(s_1) + \gamma\phi_2(s_2) + \gamma^2\phi_2(s_3) + \dots) \\ &\quad + \dots \\ &\quad + \alpha_d * (\phi_d(s_1) + \gamma\phi_d(s_2) + \gamma\phi_d(s_3) + \dots) \\ \\ V^\pi(s_1) &= \alpha_1 * E[\phi_1(s_1) + \gamma\phi_1(s_2) + \gamma^2\phi_1(s_3) + \dots | \pi] \\ &\quad + \alpha_2 * E[\phi_2(s_1) + \gamma\phi_2(s_2) + \gamma^2\phi_2(s_3) + \dots | \pi] \\ &\quad + \dots \end{aligned}$$

$$+\alpha_d * E[\phi_d(s_1) + \gamma\phi_d(s_2) + \gamma^2\phi_d(s_3) + \dots|\pi]$$

Using the definition of V_i^π and the equation for the value function, we can simplify the above expression to:

$$V^\pi(s_1) = \alpha_1 V_1^\pi + \alpha_2 V_2^\pi + \dots + \alpha_d V_d^\pi$$

Since s_1 is a variable to represent any arbitrary state, we can also just write the above as:

$$V^\pi(s) = \alpha_1 V_1^\pi + \alpha_2 V_2^\pi + \dots + \alpha_d V_d^\pi$$

This concludes the proof.

Question 2

Appendix: Relevant Code - tjha_hw5.py

```

1  # Tejas Jha
2  # EECS 498 Reinforcement Learning HW 5
3
4  import random
5  import gym
6  import math
7  import numpy as np
8  from collections import deque
9  from keras.models import Sequential
10 from keras.layers import Dense
11 from keras.optimizers import Adam
12
13 class Memory():
14     def __init__(self, max_size=1000):
15         self.buffer = deque(maxlen=max_size)
16
17     def add(self, experience):
18         self.buffer.append(experience)
19
20     def sample(self, batch_size):
21         idx = np.random.choice(np.arange(len(self.buffer)),
22                                size=batch_size,
23                                replace=True)
24         return [self.buffer[ii] for ii in idx]
25
26 class DQNCartPoleSolver():
27     def __init__(self, n_episodes=2000, C=1, n_win_ticks=195,
28                  max_env_steps=None, gamma=1.0, epsilon=0.5, epsilon_max=0.5,
29                  epsilon_min=0.01, epsilon_decay=0.0001, alpha=0.0001, batch_size
30                  =32, monitor=False, quiet=False):

```

```

28     self.memory = Memory(10000)
29     self.env = gym.make('CartPole-v0')
30     if monitor: self.env = gym.wrappers.Monitor(self.env, '../data/
        cartpole-1', force=True)
31     self.gamma = gamma
32     self.epsilon = epsilon
33     self.epsilon_max = epsilon_max
34     self.epsilon_min = epsilon_min
35     self.epsilon_decay = epsilon_decay
36     self.alpha = alpha
37     self.n_episodes = n_episodes
38     self.n_win_ticks = n_win_ticks
39     self.batch_size = batch_size
40     self.quiet = quiet
41     self.C = C
42     if max_env_steps is not None: self.env._max_episode_steps =
        max_env_steps
43
44     # Init model
45     self.model = Sequential()
46     self.model.add(Dense(64, input_dim=4, activation='relu',
        use_bias=True))
47     self.model.add(Dense(2, activation='linear'))
48     self.model.compile(loss='mse', optimizer=Adam(lr=self.alpha))
49
50     self.target_model = Sequential()
51     self.target_model.add(Dense(64, input_dim=4, activation='relu',
        use_bias=True))
52     self.target_model.add(Dense(2, activation='relu'))
53     self.target_model.compile(loss='mse', optimizer=Adam(lr=self.
        alpha))
54     self.target_model.set_weights(self.model.get_weights())
55
56     def remember(self, state, action, reward, next_state, done):
57         self.memory.add((state, action, reward, next_state, done))
58
59     def choose_action(self, state, epsilon):
60         return self.env.action_space.sample() if (np.random.random() <=
            epsilon) else np.argmax(self.model.predict(state))
61
62     def get_epsilon(self, t):
63         return self.epsilon_max - min(self.epsilon_decay * t, self.
            epsilon_max - 0.01)
64
65     def preprocess_state(self, state):
66         return np.reshape(state, [1, 4])
67
68     def replay(self, batch_size):

```

```

69     x_batch, y_batch = [], []
70     minibatch = self.memory.sample(self.batch_size)
71     for state, action, reward, next_state, done in minibatch:
72         y_target = self.model.predict(state)
73         y_target[0][action] = reward if done else reward + self.
            gamma * np.max(self.target_model.predict(next_state)[0])
74         x_batch.append(state[0])
75         y_batch.append(y_target[0])
76
77     self.model.fit(np.array(x_batch), np.array(y_batch), batch_size
            =len(x_batch), verbose=0)
78
79     def run(self):
80
81         #preTrain(self.batch_size)
82
83         scores = deque(maxlen=100)
84
85         store_avg = []
86
87         for e in range(self.n_episodes):
88             state = self.preprocess_state(self.env.reset())
89             done = False
90             i = 0
91             total_reward = 0
92             while not done:
93                 action = self.choose_action(state, self.get_epsilon(e))
94                 next_state, reward, done, _ = self.env.step(action)
95                 next_state = self.preprocess_state(next_state)
96                 self.remember(state, action, reward, next_state, done)
97                 state = next_state
98                 total_reward += reward
99                 i += 1
100
101             scores.append(total_reward)
102             mean_score = np.mean(scores)
103
104             if mean_score >= self.n_win_ticks and e >= 100:
105                 if not self.quiet: print('Ran_{}_episodes_{}_Solved_after
                    _{}_trials_'.format(e, e - 100))
106                 return e - 100
107             if e % 100 == 0 and not self.quiet:
108                 print('[Episode_{}]_{}_Mean_survival_time_over_last_100_
                    episodes_was_{}_ticks.'.format(e, mean_score))
109                 store_avg.append(mean_score)
110
111             self.replay(self.batch_size)
112

```

```

113         if e % self.C == 0:
114             self.target_model.set_weights(self.model.
                get_weights())
115
116         if not self.quiet: print('Did not solve after {} episodes '.
                .format(e))
117         return e
118
119
120
121 if __name__ == '__main__':
122     agent = DQNCartPoleSolver()
123     agent.run()

```