# 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)

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

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

where the $\phi_1...\phi_d$ 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 $\pi$ that maps $S \mapsto A$ for any state $s_1$ as the following:

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

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

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

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. There expectations would simplify to just the terms. This is shown below:

$$V^\pi(s_1) = E[R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + ...|\pi]$$
$$V^\pi(s_1) = E[\alpha_1 \phi_1(s_1) + \alpha_2 \phi_2(s_1) + ... + \alpha_d \phi_d(s_1)$$

$$+\gamma * (\alpha_2 \phi_1(s_2) + \alpha_2 \phi_2(s_2) + ... + \alpha_d \phi_d(s_2)) +$$
$$\gamma^2 * (\alpha_1 \phi_1(s_3) + \alpha_2 \phi_2(s_3) + ... + \alpha_d \phi_d(s_3)) + ... | \pi]$$

$$V^\pi(s_1) = \alpha_1 \phi_1(s_1) + \alpha_2 \phi_2(s_1) + ... + \alpha_d \phi_d(s_1)$$
$$+\gamma * (\alpha_2 \phi_1(s_2) + \alpha_2 \phi_2(s_2) + ... + \alpha_d \phi_d(s_2)) +$$
$$\gamma^2 * (\alpha_1 \phi_1(s_3) + \alpha_2 \phi_2(s_3) + ... + \alpha_d \phi_d(s_3)) + ...$$

The terms can now be regrouped using the alpha terms.

$$V^\pi(s_1) = \alpha_1 * (\phi_1(s_1) + \gamma \phi_1(s_2) + \gamma^2 \phi_1(s_3) + ...)$$
$$+\alpha_2 * (\phi_2(s_1) + \gamma \phi_2(s_2) + \gamma^2 \phi_2(s_3) + ...)$$
$$+...$$
$$+\alpha_d * (\phi_d(s_1) + \gamma \phi_d(s_2) + \gamma \phi_d(s_3) + ...)$$

$$V^\pi(s_1) = \alpha_1 * E[\phi_1(s_1) + \gamma \phi_1(s_2) + \gamma^2 \phi_1(s_3) + ... | \pi]$$
$$+\alpha_2 * E[\phi_2(s_1) + \gamma \phi_2(s_2) + \gamma^2 \phi_2(s_3) + ... | \pi]$$
$$+...$$
$$+\alpha_d * E[\phi_d(s_1) + \gamma \phi_d(s_2) + \gamma^2 \phi_d(s_3) + ... | \pi]$$

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

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

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

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

This concludes the proof.

## Question 2

## Appendix: Relevant Code - tjha_hw5.py

```
1  # Tejas Jha
2  # EECS 498 - Reinforcement Learning
3  # Homework 5
4  # 12 December 2018
5
6
7  import gym
```

```python
 8  from gym import wrappers
 9  import random
10  import math
11
12  from keras.models import Sequential
13  from keras.layers import Dense, Activation
14  from keras.optimizers import Adam
15  import matplotlib.pyplot as plt
16
17  import numpy as np
18
19  # hyper parameters
20  EPISODES = 100   # number of episodes
21  EPS_START = 0.5   # e-greedy threshold start value
22  EPS_END = 0.01   # e-greedy threshold end value
23  EPS_DECAY = 0.0001   # e-greedy threshold decay
24  GAMMA = 1.0   # Q-learning discount factor
25  LR = 0.0001   # NN optimizer learning rate
26  HIDDEN_LAYER = 64   # NN hidden layer size
27  BATCH_SIZE = 32   # Q-learning batch size
28
29
30  class ReplayMemory:
31      def __init__(self, capacity):
32          self.capacity = capacity
33          self.memory = []
34
35      def push(self, transition):
36          self.memory.append(transition)
37          if len(self.memory) > self.capacity:
38              del self.memory[0]
39
40      def sample(self, batch_size):
41          return random.sample(self.memory, batch_size)
42
43      def __len__(self):
44          return len(self.memory)
45
46
47  class Network(nn.Module):
48      def __init__(self):
49          nn.Module.__init__(self)
50          self.l1 = nn.Linear(4, HIDDEN_LAYER)
51          self.l2 = nn.Linear(HIDDEN_LAYER, 2)
52
53      def forward(self, x):
54          x = F.relu(self.l1(x))
55          x = self.l2(x)
```

```python
56            return x
57
58
59  env = gym.make('CartPole-v0')
60  env = wrappers.Monitor(env, './tmp/cartpole-v0-1')
61
62  model = Network()
63  if use_cuda:
64      model.cuda()
65  memory = ReplayMemory(10000)
66  optimizer = optim.Adam(model.parameters(), LR)
67  steps_done = 0
68  episode_durations = []
69
70
71  def select_action(state):
72      global steps_done
73      sample = random.random()
74      eps_threshold = EPS_END + (EPS_START - EPS_END) * math.exp(-1. *
             steps_done / EPS_DECAY)
75      steps_done += 1
76      if sample > eps_threshold:
77          return model(Variable(state, volatile=True).type(FloatTensor)).
                data.max(1)[1].view(1, 1)
78      else:
79          return LongTensor([[random.randrange(2)]])
80
81
82  def run_episode(e, environment):
83      state = environment.reset()
84      steps = 0
85      while True:
86          environment.render()
87          action = select_action(FloatTensor([state]))
88          next_state, reward, done, _ = environment.step(action[0, 0])
89
90          # negative reward when attempt ends
91          if done:
92              reward = -1
93
94          memory.push((FloatTensor([state]),
95                      action, # action is already a tensor
96                      FloatTensor([next_state]),
97                      FloatTensor([reward])))
98
99          learn()
100
101         state = next_state
```

```python
102            steps += 1
103
104            if done:
105                print("{2} Episode {0} finished after {1} steps"
106                      .format(e, steps, '\033[92m' if steps >= 195 else '
                            \033[99m'))
107                episode_durations.append(steps)
108                plot_durations()
109                break
110
111
112  def learn():
113      if len(memory) < BATCH_SIZE:
114          return
115
116      # random transition batch is taken from experience replay memory
117      transitions = memory.sample(BATCH_SIZE)
118      batch_state, batch_action, batch_next_state, batch_reward = zip(*
             transitions)
119
120      batch_state = Variable(torch.cat(batch_state))
121      batch_action = Variable(torch.cat(batch_action))
122      batch_reward = Variable(torch.cat(batch_reward))
123      batch_next_state = Variable(torch.cat(batch_next_state))
124
125      # current Q values are estimated by NN for all actions
126      current_q_values = model(batch_state).gather(1, batch_action)
127      # expected Q values are estimated from actions which gives maximum
             Q value
128      max_next_q_values = model(batch_next_state).detach().max(1)[0]
129      expected_q_values = batch_reward + (GAMMA * max_next_q_values)
130
131      # loss is measured from error between current and newly expected Q
             values
132      loss = F.smooth_l1_loss(current_q_values, expected_q_values)
133
134      # backpropagation of loss to NN
135      optimizer.zero_grad()
136      loss.backward()
137      optimizer.step()
138
139
140  def plot_durations():
141      plt.figure(2)
142      plt.clf()
143      durations_t = torch.FloatTensor(episode_durations)
144      plt.title('Training ...')
145      plt.xlabel('Episode')
```

```python
146        plt.ylabel('Duration')
147        plt.plot(durations_t.numpy())
148        # take 100 episode averages and plot them too
149        if len(durations_t) >= 100:
150            means = durations_t.unfold(0, 100, 1).mean(1).view(-1)
151            means = torch.cat((torch.zeros(99), means))
152            plt.plot(means.numpy())
153
154        plt.pause(0.001)  # pause a bit so that plots are updated
155
156
157  for e in range(EPISODES):
158      run_episode(e, env)
159
160  print('Complete')
161  env.render(close=True)
162  env.close()
163  plt.ioff()
164  plt.show()
```