

EECS 498: Reinforcement Learning

Homework 1 Responses

Tejas Jha
tjha

September 23, 2018

This document includes my responses to Homework 1 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

Question 2

Question 3

Question 4

Question 5

Appendix: Relevant Code - tjha.py

```
1 import numpy as np
2 import random
3 import math
4
5 #####
6 # Tejas Jha
7 # 24 September 2018
8 # EECS 498 Special Topics: Reinforcement Learning
9 # Homework 1
10 # Code corresponding to Q4 and Q5
11 #####
12
13
```

```

14 # For Q4:
15 #
16 # This program provides experimentations with Epsilon-Greedy,
17 # Optimistic Initial Values, and UCB algorithms to learn solutions
18 # to a 5-arm bandit problem.
19 #
20 # Consider the arm rewards to be Bernoulli Distribution with means
21 # [0.1, 0.275, 0.45, 0.625, 0.8]
22 #
23 # Each algorithm should be performed for 2000 runs for each parameter
    set
24 # Each run should choose actions over 1000 time steps
25
26 # For Q5
27 #
28 #
29 #
30 #
31
32
33 # Creation of class to store Bernoulli distribution mean and return
    reward
34 class Arm:
35     # Initialize arm given
36     def __init__(self, mean):
37         self.mean = mean
38
39     # Simulate Bernoulli distribution
40     # Get random float between 0 and 1, return 0 if greater than mean,
        1 otherwise
41     def reward(self):
42         if random.random() > self.mean:
43             return 0.0
44         else:
45             return 1.0
46
47
48
49 # Creation of a class to store details during each time step and house
50 # necessary functions needed during usage.
51 # Note: I used the code found here: https://gist.github.com/nagataka/
    c02b9acd6e8a8d7696e09f8a129d3362
52 #     for inspiration and guidance on how to construct my own
        functions and organize my code
53 class EpsilonGreedy:
54     # Initialize class values
55     # epsilon - the rate of exploration (randomly choosing an arm)
56     # counts -  $N(a)$ , number of pulls for each arm

```

```

57     # values - Q(a) average reward received from each arm (action-value
        estimates)
58     def __init__(self, epsilon, counts, values):
59         self.epsilon = epsilon
60         self.counts = counts
61         self.values = values
62
63     # Performs a single time step update through epsilon-greedy
        algorithm
64     def step(self, arms):
65         # Step 1 - Determin which arm to pick next
66         selected_arm_idx = 0
67         if random.random() > self.epsilon:
68             # EXPLOITATION
69             # Choose action with largest expected action-value (break
                ties randomly)
70             max_value = max(self.values)
71             all_max_idx = [idx for idx, val in enumerate(self.values)
                if val == max_value]
72             selected_arm_idx = all_max_idx[0]
73
74             # (Break ties randomly)
75             if len(all_max_idx) > 1:
76                 # Randomly choose an index from all_max_idx
77                 selected_arm_idx = random.choice(all_max_idx)
78         else:
79             # EXPLORATION
80             # Choose random index to explore
81             selected_arm_idx = random.choice(range(5))
82
83         # Step 2 - Get the reward for that arm
84         reward = arms[selected_arm_idx].reward()
85         # Step 3 - Update count
86         self.counts[selected_arm_idx] += 1
87         # Step 4 - Update action-value
88         self.values[selected_arm_idx] += (1 / self.counts[
            selected_arm_idx])*(reward - self.values[selected_arm_idx])
89
90     # Get the action-value estimates for EpsilonGreedy
91     def getValues(self):
92         return self.values
93
94 # Creation of another Class variable that may be fixed to be universal
    type
95 class ModelClass:
96     def __init__(self, epsilon=0, counts=np.zeros(5), values=np.zeros
        (5), c=0):
97         self.epsilon = epsilon

```

```

98         self.counts = counts
99         self.values = values
100        self.c = c
101
102    def step(self, iteration, arms):
103
104        print(iteration)
105
106        # Step 1 – Select the Action with the largest upper bound
107        # If count is 0, the index is given priority
108        all_zero_idx = [idx for idx, val in enumerate(self.counts) if
109                        val == 0]
109        selected_arm_idx = 0
110        #updated_val = 0
111        if len(all_zero_idx) > 0:
112            selected_arm_idx = random.choice(all_zero_idx)
113        else:
114            # Select the arm with the largest estimated upper bound
115            estimates = self.values
116
117            for idx in range(len(estimates)):
118                estimates[idx] += self.c*math.sqrt(math.log(iteration)/
119                    self.counts[idx])
119
120            max_value = max(estimates)
121            all_max_idx = [idx for idx, val in enumerate(estimates) if
122                          val == max_value]
122            selected_arm_idx = all_max_idx[0]
123
124            # (Break ties randomly)
125            if len(all_max_idx) > 1:
126                # Randomly choose an index from all_max_idx
127                selected_arm_idx = random.choice(all_max_idx)
128
129            #updated_val = estimates[selected_arm_idx]
130
131        # Step 2 – Get the reward for that arm
132        reward = arms[selected_arm_idx].reward()
133        # Step 3 – Update count
134        self.counts[selected_arm_idx] += 1
135        # Step 4 – Update action-value
136        self.values[selected_arm_idx] += (1 / self.counts[
137            selected_arm_idx])*(reward - self.values[selected_arm_idx])
138
139
140
141    def getValues(self):

```

```

142         return self.values
143
144
145
146 # Epsilon-Greedy Algorithm run 2000 times given epsilon value
147 def epsilon_greedy_alg(epsilon , arms):
148
149     # Print message to output describing Algorithm being run
150     print("
        *****
    ")
151     print("Performing Epsilon-Greedy Algorithm 2000 times with epsilon
        = " + str(epsilon))
152
153     # Loop through algorithm for 2000 runs
154     for current_run in range(2000):
155
156         # Each run should set a distinct random seed
157         random.seed()
158         # Create EpsilonGreedy data type initialized with epsilon and 0
            for counts and values
159         model = EpsilonGreedy(epsilon , np.zeros(5) , np.zeros(5))
160
161         # Perform a single run through algorithm using 1000 time steps
162         for i in range(1000):
163             model.step(arms)
164
165             if (current_run + 1) % 500 == 0:
166                 print("Completed Run#" + str(current_run + 1))
167                 print(model.getValues())
168
169         print("Completed Runs for epsilon = " + str(epsilon))
170         print("
        *****
    ")
171
172
173 # Optimistic Initial Value Algorithm run 2000 times given initial value
    (assume greedy arm selection)
174 def optimistic_initial_value_alg(initial_val , arms):
175
176     # Print message to output describing Algorithm being run
177     print("
        *****
    ")
178     print("Performing Optimistic Initial Value Algorithm 2000 times
        with initial value = " + str(initial_val))
179

```

```

180 # Loop through algorithm for 2000 runs
181 for current_run in range(2000):
182
183     # Each run should set a distinct random seed
184     random.seed()
185     # Create EpsilonGreedy data type initialized with epsilon = 0
186     and 0 for counts and initial values
187     initialized_values_array = np.full(5, initial_val)
188     model = EpsilonGreedy(0, np.zeros(5), initialized_values_array)
189
190     # Perform a single run through algorithm using 1000 time steps
191     for i in range(1000):
192         model.step(arms)
193
194     if (current_run + 1) % 500 == 0:
195         print("Completed_Run#" + str(current_run + 1))
196         print(model.getValues())
197
198 print("Completed_Runs_for_initial_value=" + str(initial_val))
199 print("
    *****
    ")
200
201
202 # Upper Confidence Bound (UCB) Algorithm run 2000 times with given c
203 parameter
204 def upper_confidence_bound_alg(c, arms):
205
206     # Print message to output describing Algorithm being run
207     print("
    *****
    ")
208
209     print("Performing_UCB_Algorithm_2000_times_with_c=" + str(c))
210
211     # Loop through algorithm for 2000 runs
212     for current_run in range(2000):
213
214         # Each run should set a distinct random seed
215         random.seed()
216
217         # Initialize model with c value
218         model = ModelClass(0, np.zeros(5), np.zeros(5), c)
219
220         # Perform a single run through algorithm using 1000 time steps
221         for i in range(1000):
222             model.step(i+1, arms)

```

```

222         if (current_run + 1) % 500 == 0:
223             print("Completed Run#" + str(current_run + 1))
224             print(model.getValues())
225
226     print("Completed Runs for c=" + str(c))
227     print("
    *****
    ")
228
229
230
231
232 def main():
233     # Create List of arm class variables called arms set with means for
234     Bernoulli distributions
235     means = [0.8, 0.275, 0.45, 0.625, 0.8]
236     arms = np.array(list(map(Arm, means)))
237
238     # Perform Epsilon-Greedy algorithm with QI = 0 and
239     # for each epsilon = [0.01, 0.1, 0.3]
240     epsilon_list = [0.01, 0.1, 0.3]
241     for epsilon in epsilon_list:
242         epsilon_greedy_alg(epsilon, arms)
243         print("Finished with Epsilon-Greedy Algorithm for epsilon=" +
244               str(epsilon))
245
246     # Perform Optimistic Initial Value algorithm with epsilon = 0 (
247     always greedy)
248     # for each QI = [1,5,50]
249     initial_val_list = [1.0, 5.0, 50.0]
250     for val in initial_val_list:
251         optimistic_initial_value_alg(val, arms)
252         print("Finished with Optimistic Initial Value Algorithm for
253               initial value=" + str(val))
254
255     # Perform UCB algorithm with QI = 0
256     # for each c = [0.2, 1, 2]
257     c_vals = [0.2, 1, 2]
258     for val in c_vals:
259         upper_confidence_bound_alg(val, arms)
260         print("Finished with UCB Algorithm for c=" + str(val))
261
262 if __name__ == "__main__":
263     main()

```