

EECS 498: Reinforcement Learning

Homework 4 Responses

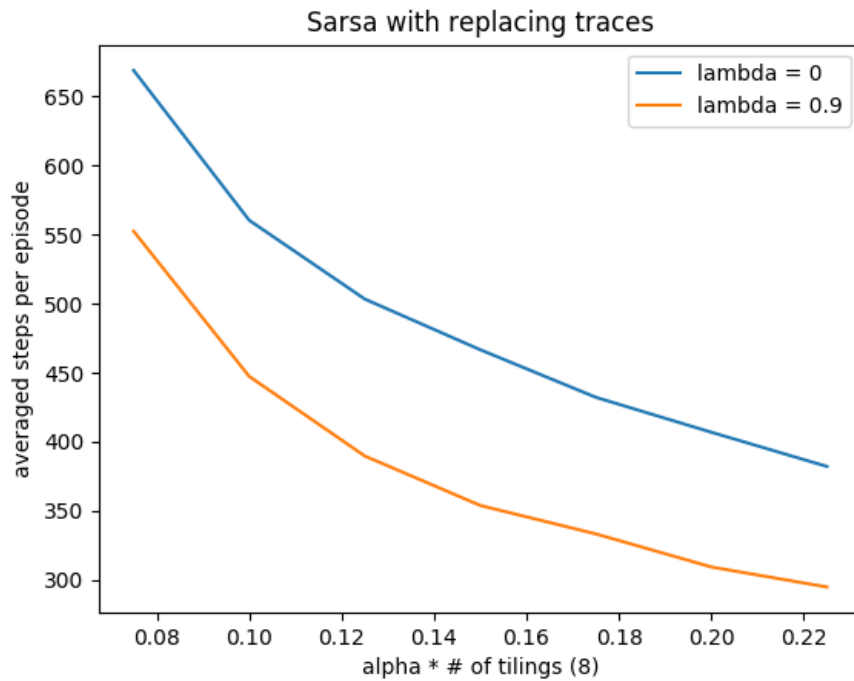
Tejas Jha
tjha

November 27, 2018

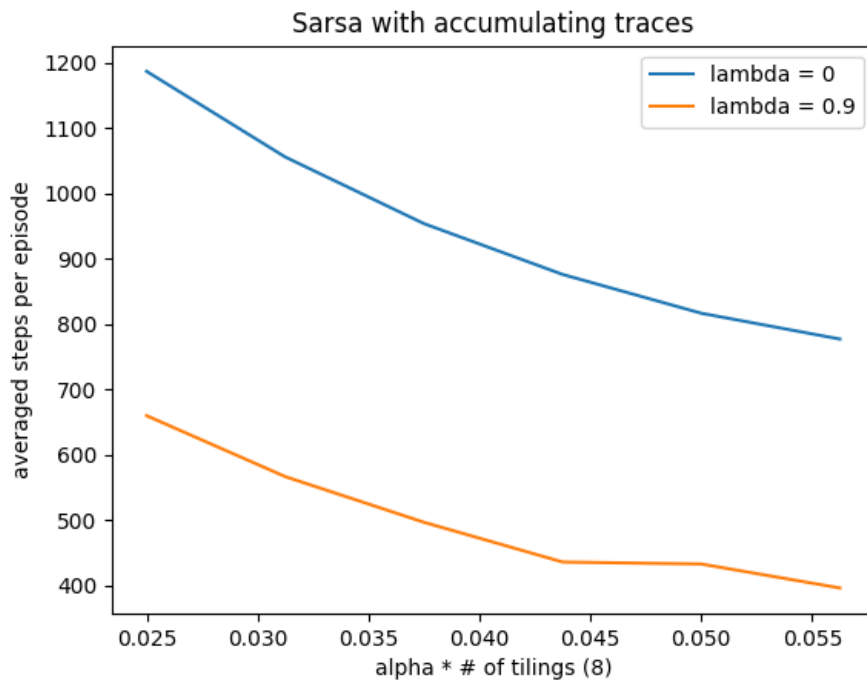
This document includes my responses to Homework 4 questions. Responses that involved the use of coding will provide references to specific lines of code to provide a better overview of how the problem was approached. The code can either be referenced in the Appendix or in the accompanied python script submitted with this assignment.

Question 1

- (a) In this specific problem, we are able to set epsilon to 0 because we begin with optimistic initial values. We covered optimistic initial values in the bandit setting, but similarly in this specific problem, these values will encourage initial exploration until values are brought down to a point where the estimated values approach the expected values for the greedy actions. The optimistic initial value we use for the values is 0, which is greater than the true values since a reward of -1 is received for every time step where the car is not at the terminal state.
- (b) Below is a generated plot similar to the left part of Figure 12.10 from the textbook for the requested trace and alpha values using replacing traces.



- (c) Below is a generated plot similar to the left part of Figure 12.10 from the textbook for the requested trace and alpha values using accumulating traces.



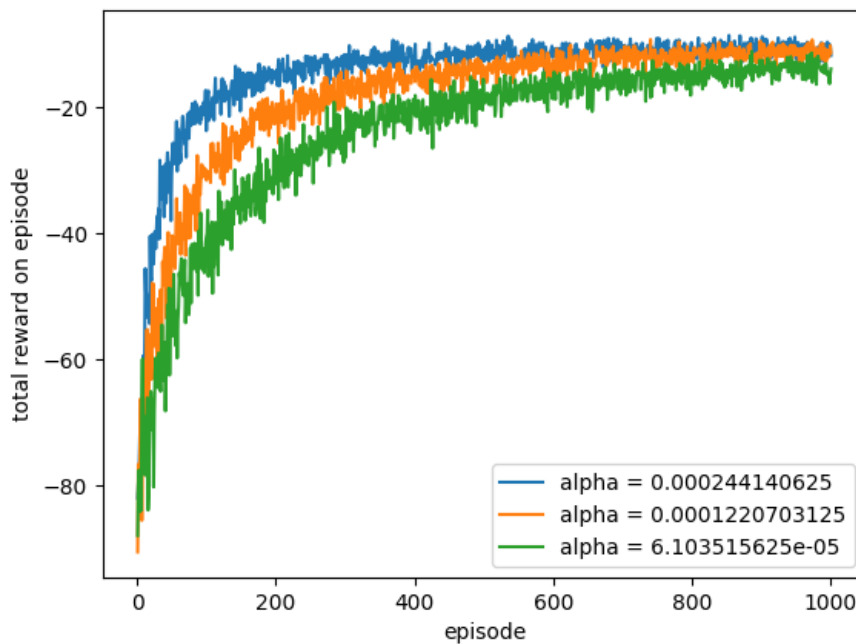
- (d) If we change the range of the alphas to be same as that which we used to generate the plot in part (b) for accumulating traces, we can expect a plot with lower average steps than the plot in part (c) where accumulating traces was also used. This will be due to the larger alpha values being looked

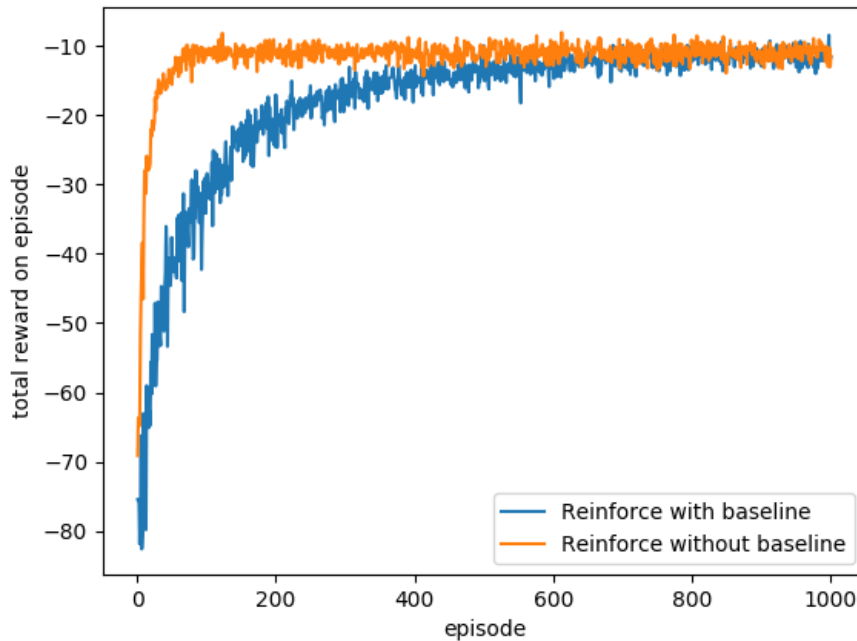
at which will now focus at the points on the curve that are lower. However, if we were to compare the plot in part (b) that used replacing traces, we will see that the curve for 0 lambda would be the same but the curve for 0.9 lambda would be lower for the accumulating traces. This result would suggest that using accumulating traces for this range of alphas seems to perform better as the averaged step per episode is lower, indicating a more optimal policy. While some searching online seems to indicate that replacing traces generally perform better than accumulating traces (<http://www.incompleteideas.net/book/ebook/node80.html>) as accumulating traces could end up selecting the "wrong" action more times resulting in a larger trace corresponding to the "wrong" action. As a result of this, accumulating traces generally take longer to learn to account for the larger trace for the "wrong" action. However, from the given alpha values, accumulating traces seems to work better in this situation, indicating that the optimal action could have been chosen earlier and then greedily selected more often with accumulating traces.

Question 2

Below are the reproduced Figures 13.1 and 13.2 from the textbook. The corresponding code used to generate these figures can be found in the appendix.

For fig1: Note that the blue curve is for 2^{-12} , the orange curve is for 2^{-13} and the green curve is for 2^{-14} . Also the data has been averaged over 100 trials.





Question 3

Below are the attempts made to recreate Figures 4 and 5 from the Options paper. The corresponding code used to generate the plots for these figures can be found in the appendix.

Appendix: Relevant Code - tjha_hw4.py

```

1 # Tejas Jha
2 # EECS 498 Reinforcement Learning HW 4
3 #
4 #####
5
6 # This python script contains all of the relevant code needed to
7 # recreate the desired
8 # plots on the assignment. Comments will separate the code to better
9 # indicate which
10 # portions corresponding to which parts of the assignment.
11
12 import gym
13 from gym.envs.registration import register
14 import numpy as np
15 import matplotlib
16 matplotlib.use('Agg')
17 import matplotlib.pyplot as plt

```

```

14 from tqdm import tqdm
15 from math import floor, log
16 from itertools import zip_longest
17 from tiles3 import IHT
18
19 # Environment Generation
20 register(
21     id = 'MountainCar-v1',
22     entry_point = 'gym.envs.classic_control:MountainCarEnv',
23     max_episode_steps = 5000
24 )
25
26 env = gym.make('MountainCar-v1')
27 #
28 #####
28 #
29 # cp code
30 # all possible actions
31 ACTION_REVERSE = -1
32 ACTION_ZERO = 0
33 ACTION_FORWARD = 1
34 # order is important
35 ACTIONS = [ACTION_REVERSE, ACTION_ZERO, ACTION_FORWARD]
36
37 # bound for position and velocity
38 POSITION_MIN = -1.2
39 POSITION_MAX = 0.5
40 VELOCITY_MIN = -0.07
41 VELOCITY_MAX = 0.07
42
43 # discount is always 1.0 in these experiments
44 DISCOUNT = 1.0
45
46 # use optimistic initial value, so it's ok to set epsilon to 0
47 EPSILON = 0
48
49 # maximum steps per episode
50 STEP_LIMIT = 5000
51
52 # take an @action at @position and @velocity
53 # @return: new position, new velocity, reward (always -1)
54 def step(position, velocity, action):
55     new_velocity = velocity + 0.001 * action - 0.0025 * np.cos(3 *
56         position)
57     new_velocity = min(max(VELOCITY_MIN, new_velocity), VELOCITY_MAX)
58     new_position = position + new_velocity
59     new_position = min(max(POSITION_MIN, new_position), POSITION_MAX)

```

```

59     reward = -1.0
60     if new_position == POSITION_MIN:
61         new_velocity = 0.0
62     return new_position, new_velocity, reward
63
64     # accumulating trace update rule
65     # @trace: old trace (will be modified)
66     # @activeTiles: current active tile indices
67     # @lam: lambda
68     # @return: new trace for convenience
69     def accumulating_trace(trace, active_tiles, lam):
70         trace *= lam * DISCOUNT
71         trace[active_tiles] += 1
72     return trace
73
74     # replacing trace update rule
75     # @trace: old trace (will be modified)
76     # @activeTiles: current active tile indices
77     # @lam: lambda
78     # @return: new trace for convenience
79     def replacing_trace(trace, activeTiles, lam):
80         active = np.in1d(np.arange(len(trace)), activeTiles)
81         trace[active] = 1
82         trace[~active] *= lam * DISCOUNT
83     return trace
84
85     # wrapper class for Sarsa(lambda)
86     class Sarsa:
87         # @maxSize: the maximum # of indices
88         def __init__(self, step_size, lam, trace_update=accumulating_trace,
89             num_of_tilings=8, max_size=2048):
89             self.max_size = max_size
90             self.num_of_tilings = num_of_tilings
91             self.trace_update = trace_update
92             self.lam = lam
93
94             # divide step size equally to each tiling
95             self.step_size = step_size / num_of_tilings
96
97             self.hash_table = IHT(max_size)
98
99             # weight for each tile
100            self.weights = np.zeros(max_size)
101
102            # trace for each tile
103            self.trace = np.zeros(max_size)
104

```

```

105     # position and velocity needs scaling to satisfy the tile
106     software
107     self.position_scale = self.num_of_tilings / (POSITION_MAX -
108         POSITION_MIN)
109     self.velocity_scale = self.num_of_tilings / (VELOCITY_MAX -
110         VELOCITY_MIN)
111
112     # get indices of active tiles for given state and action
113     def get_active_tiles(self, position, velocity, action):
114         # I think positionScale * (position - position_min) would be a
115         good normalization.
116         # However positionScale * position_min is a constant, so it's
117         ok to ignore it.
118         active_tiles = tiles(self.hash_table, self.num_of_tilings,
119             [self.position_scale * position, self.
120                 velocity_scale * velocity],
121             [action])
122         return active_tiles
123
124     # estimate the value of given state and action
125     def value(self, position, velocity, action):
126         if position == POSITION_MAX:
127             return 0.0
128         active_tiles = self.get_active_tiles(position, velocity, action
129             )
130         return np.sum(self.weights[active_tiles])
131
132     # learn with given state, action and target
133     def learn(self, position, velocity, action, target):
134         active_tiles = self.get_active_tiles(position, velocity, action
135             )
136         estimation = np.sum(self.weights[active_tiles])
137         delta = target - estimation
138         if self.trace_update == accumulating_trace or self.trace_update
139             == replacing_trace:
140             self.trace_update(self.trace, active_tiles, self.lam)
141         else:
142             raise Exception('Unexpected_Trace_Type')
143         self.weights += self.step_size * delta * self.trace
144
145     # get # of steps to reach the goal under current state value
146     function
147     def cost_to_go(self, position, velocity):
148         costs = []
149         for action in ACTIONS:
150             costs.append(self.value(position, velocity, action))
151         return -np.max(costs)
152

```

```

143 # get action at @position and @velocity based on epsilon greedy policy
    and @valueFunction
144 def get_action(position, velocity, valueFunction):
145     if np.random.binomial(1, EPSILON) == 1:
146         return np.random.choice(ACTIONS)
147     values = []
148     for action in ACTIONS:
149         values.append(valueFunction.value(position, velocity, action))
150     return np.argmax(values) - 1
151
152 # play Mountain Car for one episode based on given method @evaluator
153 # @return: total steps in this episode
154 def play(evaluator):
155     position = np.random.uniform(-0.6, -0.4)
156     velocity = 0.0
157     action = get_action(position, velocity, evaluator)
158     steps = 0
159     while True:
160         next_position, next_velocity, reward = step(position, velocity,
            action)
161         next_action = get_action(next_position, next_velocity,
            evaluator)
162         steps += 1
163         target = reward + DISCOUNT * evaluator.value(next_position,
            next_velocity, next_action)
164         evaluator.learn(position, velocity, action, target)
165         position = next_position
166         velocity = next_velocity
167         action = next_action
168         if next_position == POSITION_MAX:
169             break
170         if steps >= STEP_LIMIT:
171             print('Step_Limit_Exceeded!')
172             break
173     return steps
174
175 # end cp code
176 #
    #####

177
178 # Question 1 plots
179 def q1plots():
180     runs = 5
181     episodes = 50
182     lambdas = [0, 0.9]
183
184     # Part (b) – Generation of plot using replacing traces

```



```

185     alphas = np.arange(0.6,2.0,0.2) / 8.0
186     steps = np.zeros((len(lambdas), len(alphas), runs, episodes))
187     for lambdaIdx, lam in enumerate(lambdas):
188         for alphaIdx, alpha in enumerate(alphas):
189             for run in tqdm(range(runs)):
190                 evaluator = Sarsa(alpha, lam, replacing_trace, max_size
                                   =4096)
191                 for ep in range(episodes):
192                     step = play(evaluator)
193                     steps[lambdaIdx, alphaIdx, run, ep] = step
194
195     # average over episodes
196     steps = np.mean(steps, axis=3)
197     # average over runs
198     steps = np.mean(steps, axis=2)
199
200     for lamdaIdx, lam in enumerate(lambdas):
201         plt.plot(alphas, steps[lamdaIdx, :], label='lambda_=%s' % (str
                                   (lam)))
202     plt.xlabel('alpha_*#_of_tilings_(8)')
203     plt.ylabel('averaged_steps_per_episode')
204     plt.title('Sarsa_with_replacing_traces')
205     #plt.ylim([180, 300])
206     plt.legend()
207
208     plt.savefig('replacing_traces.png')
209     plt.close()
210
211     print("Completed_Problem_1_Part_(b)")
212
213     # Part (c) — Generation of plot using accumulating traces
214     alphas = np.arange(0.2,0.5,0.05) / 8.0
215     steps = np.zeros((len(lambdas), len(alphas), runs, episodes))
216     for lambdaIdx, lam in enumerate(lambdas):
217         for alphaIdx, alpha in enumerate(alphas):
218             for run in tqdm(range(runs)):
219                 evaluator = Sarsa(alpha, lam, accumulating_trace,
                                   max_size=4096)
220                 for ep in range(episodes):
221                     step = play(evaluator)
222                     steps[lambdaIdx, alphaIdx, run, ep] = step
223
224     # average over episodes
225     steps = np.mean(steps, axis=3)
226     # average over runs
227     steps = np.mean(steps, axis=2)
228
229     for lamdaIdx, lam in enumerate(lambdas):

```

```

230         plt.plot(alphas, steps[lamdaIdx, :], label='lambda_=%s' % (str
                (lam)))
231     plt.xlabel('alpha_*#_of_tilings_(8)')
232     plt.ylabel('averaged_steps_per_episode')
233     plt.title('Sarsa_with_accumulating_traces')
234     #plt.ylim([180, 300])
235     plt.legend()
236
237     plt.savefig('accumulating_traces.png')
238     plt.close()
239
240     print("Completed_Problem_1_Part_(c)")
241
242     #
    #####

243 #
244 # Work for Q2
245
246 class ShortCorridor:
247     """
248     Short corridor environment, see Example 13.1
249     """
250     def __init__(self):
251         self.reset()
252
253     def reset(self):
254         self.state = 0
255
256     def step(self, go_right):
257         """
258         Args:
259             go_right (bool): chosen action
260         Returns:
261             tuple of (reward, episode terminated?)
262         """
263         if self.state == 0 or self.state == 2:
264             if go_right:
265                 self.state += 1
266             else:
267                 self.state = max(0, self.state - 1)
268         else:
269             if go_right:
270                 self.state -= 1
271             else:
272                 self.state += 1
273
274         if self.state == 3:

```

```

275         # terminal state
276         return 0, True
277     else:
278         return -1, False
279
280 def softmax(x):
281     t = np.exp(x - np.max(x))
282     return t / np.sum(t)
283
284 class ReinforceAgent:
285     """
286     ReinforceAgent that follows algorithm
287     'REINFORNCE Monte-Carlo Policy-Gradient Control (episodic)'
288     """
289     def __init__(self, alpha, gamma):
290         # set values such that initial conditions correspond to left-
291         epsilon greedy
292         self.theta = np.array([-1.47, 1.47])
293         self.alpha = alpha
294         self.gamma = gamma
295         # first column - left, second - right
296         self.x = np.array([[0, 1],
297                             [1, 0]])
298         self.rewards = []
299         self.actions = []
300
301     def get_pi(self):
302         h = np.dot(self.theta, self.x)
303         t = np.exp(h - np.max(h))
304         pmf = t / np.sum(t)
305         # never become deterministic,
306         # guarantees episode finish
307         imin = np.argmin(pmf)
308         epsilon = 0.05
309
310         if pmf[imin] < epsilon:
311             pmf[:] = 1 - epsilon
312             pmf[imin] = epsilon
313
314         return pmf
315
316     def get_p_right(self):
317         return self.get_pi()[1]
318
319     def choose_action(self, reward):
320         if reward is not None:
321             self.rewards.append(reward)

```

```

322         pmf = self.get_pi()
323         go_right = np.random.uniform() <= pmf[1]
324         self.actions.append(go_right)
325
326         return go_right
327
328     def episode_end(self, last_reward):
329         self.rewards.append(last_reward)
330
331         # learn theta
332         G = np.zeros(len(self.rewards))
333         G[-1] = self.rewards[-1]
334
335         for i in range(2, len(G) + 1):
336             G[-i] = self.gamma * G[-i + 1] + self.rewards[-i]
337
338         gamma_pow = 1
339
340         for i in range(len(G)):
341             j = 1 if self.actions[i] else 0
342             pmf = self.get_pi()
343             grad_ln_pi = self.x[:, j] - np.dot(self.x, pmf)
344             update = self.alpha * gamma_pow * G[i] * grad_ln_pi
345
346             self.theta += update
347             gamma_pow *= self.gamma
348
349         self.rewards = []
350         self.actions = []
351
352     class ReinforceBaselineAgent(ReinforceAgent):
353         def __init__(self, alpha, gamma, alpha_w):
354             super(ReinforceBaselineAgent, self).__init__(alpha, gamma)
355             self.alpha_w = alpha_w
356             self.w = 0
357
358         def episode_end(self, last_reward):
359             self.rewards.append(last_reward)
360
361             # learn theta
362             G = np.zeros(len(self.rewards))
363             G[-1] = self.rewards[-1]
364
365             for i in range(2, len(G) + 1):
366                 G[-i] = self.gamma * G[-i + 1] + self.rewards[-i]
367
368             gamma_pow = 1
369

```

```

370         for i in range(len(G)):
371             self.w += self.alpha_w * gamma_pow * (G[i] - self.w)
372
373             j = 1 if self.actions[i] else 0
374             pmf = self.get_pi()
375             grad_ln_pi = self.x[:, j] - np.dot(self.x, pmf)
376             update = self.alpha * gamma_pow * (G[i] - self.w) *
                 grad_ln_pi
377
378             self.theta += update
379             gamma_pow *= self.gamma
380
381             self.rewards = []
382             self.actions = []
383
384     def trial(num_episodes, agent_generator):
385         env = ShortCorridor()
386         agent = agent_generator()
387
388         rewards = np.zeros(num_episodes)
389         for episode_idx in range(num_episodes):
390             rewards_sum = 0
391             reward = None
392             env.reset()
393
394             while True:
395                 go_right = agent.choose_action(reward)
396                 reward, episode_end = env.step(go_right)
397                 rewards_sum += reward
398
399                 if episode_end:
400                     agent.episode_end(reward)
401                     break
402
403             rewards[episode_idx] = rewards_sum
404
405         return rewards
406
407     def q2plot1():
408         num_trials = 100
409         num_episodes = 1000
410         alphas = [2**(-12), 2**(-13), 2**(-14)]
411         gamma = 1
412
413         for alpha in alphas:
414
415             print(alpha)
416

```

```

417     rewards = np.zeros((num_trials, num_episodes))
418     agent_generator = lambda : ReinforceAgent(alpha=alpha, gamma=
        gamma)
419
420     for i in tqdm(range(num_trials)):
421         reward = trial(num_episodes, agent_generator)
422         rewards[i, :] = reward
423
424     plt.plot(np.arange(num_episodes) + 1, rewards.mean(axis=0),
        label='alpha_=%s' % (str(alpha)))
425
426     #plt.plot(np.arange(num_episodes) + 1, -11.6 * np.ones(num_episodes
        ), ls='dashed', color='red', label='-11.6')
427     #plt.plot(np.arange(num_episodes) + 1, rewards.mean(axis=0), color
        ='blue')
428     plt.ylabel('total_reward_on_episode')
429     plt.xlabel('episode')
430     plt.legend(loc='lower_right')
431
432     plt.savefig('fig1.png')
433     plt.close()
434
435 def q2plot2():
436     num_trials = 100
437     num_episodes = 1000
438     alpha = 2**(-13)
439     gamma = 1
440     agent_generators = [lambda : ReinforceAgent(alpha=alpha, gamma=
        gamma),
441
        lambda : ReinforceBaselineAgent(alpha=2**(-9),
            gamma=gamma, alpha_w=2**(-6))]
442     labels = ['Reinforce_with_baseline',
443         'Reinforce_without_baseline']
444
445     rewards = np.zeros((len(agent_generators), num_trials, num_episodes
        ))
446
447     for agent_index, agent_generator in enumerate(agent_generators):
448         for i in tqdm(range(num_trials)):
449             reward = trial(num_episodes, agent_generator)
450             rewards[agent_index, i, :] = reward
451
452     #plt.plot(np.arange(num_episodes) + 1, -11.6 * np.ones(num_episodes
        ), ls='dashed', color='red', label='-11.6')
453     for i, label in enumerate(labels):
454         plt.plot(np.arange(num_episodes) + 1, rewards[i].mean(axis=0),
            label=label)
455     plt.ylabel('total_reward_on_episode')

```

```

456     plt.xlabel('episode')
457     plt.legend(loc='lower_right')
458
459     plt.savefig('fig2.png')
460     plt.close()
461
462 if __name__ == '__main__':
463     # Ensuring environment is reset to begin
464     env.reset()
465     #


---


466     # Question 1 – Implementation of replacing traces and accumulating
467     #                 traces
468     #                 to produce plots similar to that of Figure 12.10 in
469     #                 the textbook
470     #


---


471     #q1plots()
472


---


473     # Quesiton 2 – Reproduction of Figures 13.1 and 13.2 from the
474     #                 textbook.
475     #                 Plots are generated using example 13.1
476     #


---


477     # Figure 13.1 Reproduciton
478     #q2plot1()
479     # Figure 13.2 Reproduction
480     q2plot2()
481     #


---


482     # Quesiton 3 – Attempts to reproduce the work behind Figures 4 and
483     #                 5 in the
484     #                 Options paper. Heatmap figure will be created to
485     #                 represent the results
486     #


---


487     #

```

```
487     # Ensuring environment is closed at the end to avoid compilation  
         issues  
488     env.close()
```