

# EECS 498: Reinforcement Learning

## Homework 4 Responses

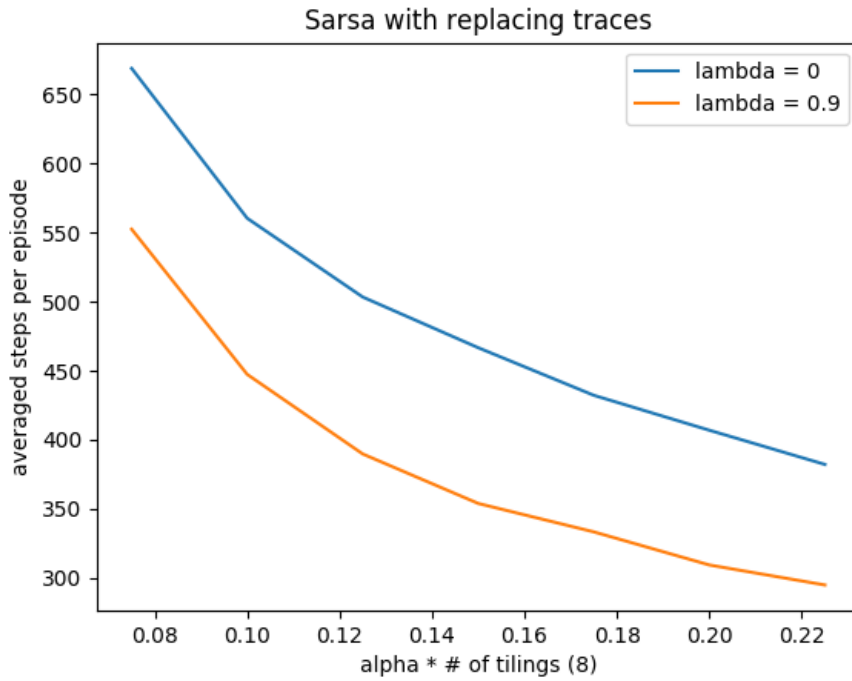
Tejas Jha  
tjha

November 27, 2018

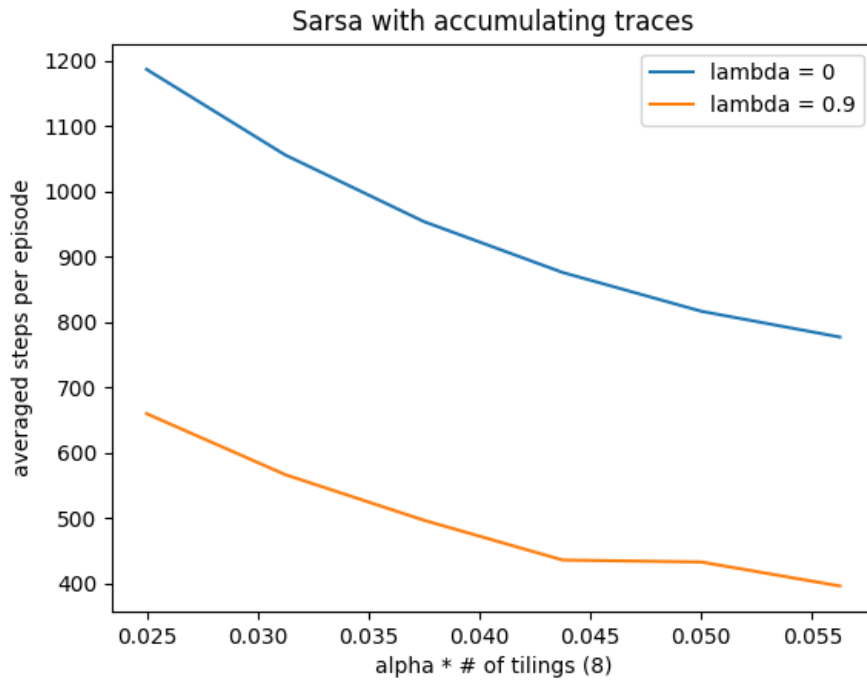
This document includes my responses to Homework 4 questions. Responses that involved the use of coding will provide references to specific lines of code to provide a better overview of how the problem was approached. The code can either be referenced in the Appendix or in the accompanied python script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



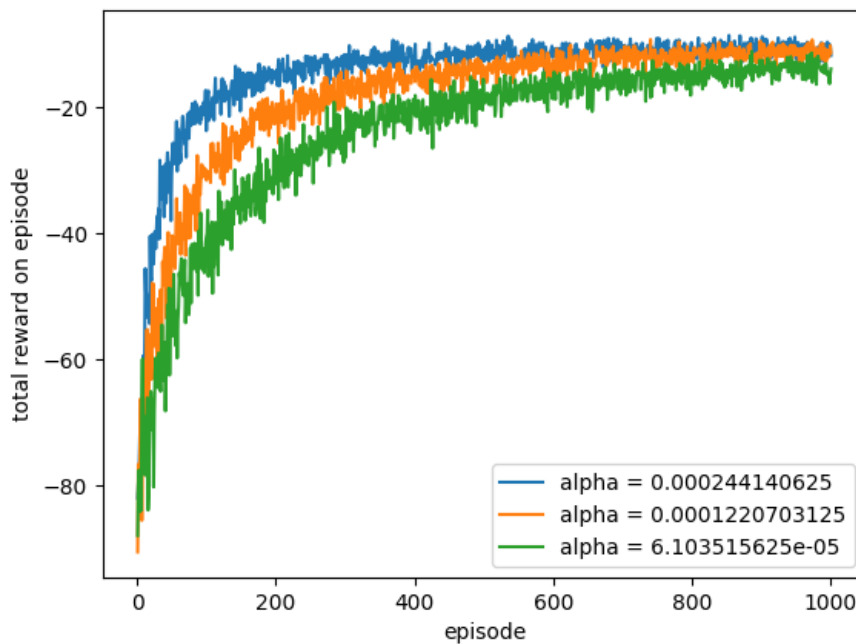
- (d) If we change the range of the alphas to be same as that which we used to generate the plot in part (b) for accumulating traces, we can expect a plot with lower average steps than the plot in part (c) where accumulating traces was also used. This will be due to the larger alpha values being looked

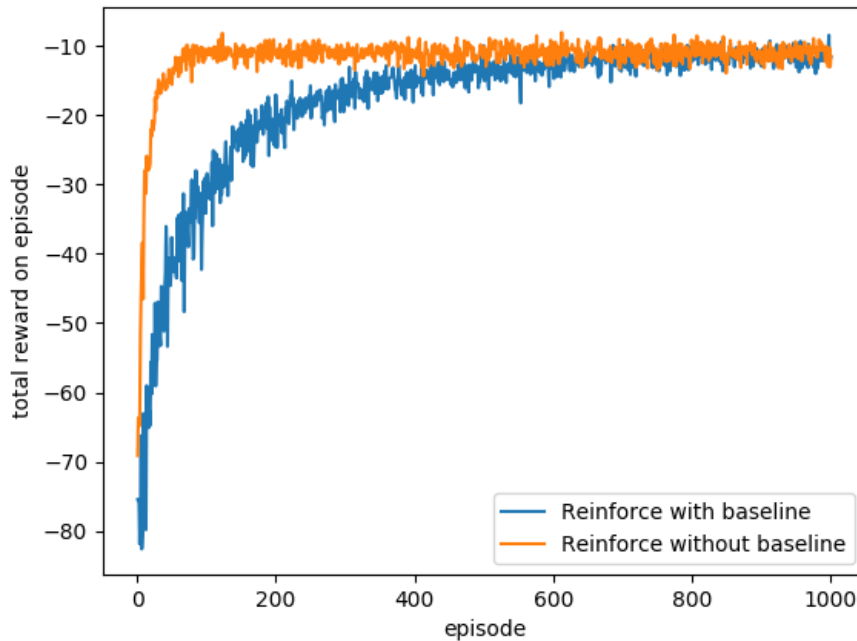
at which will now focus at the points on the curve that are lower. However, if we were to compare the plot in part (b) that used replacing traces, we will see that the curve for 0 lambda would be the same but the curve for 0.9 lambda would be lower for the accumulating traces. This result would suggest 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
18 #
19 #####
19 # The code below is from http://incompleteideas.net/tiles/tiles3.py–
20 # which is the Tile Coding Software by Rich Sutton
21 # Below is the footnote from page 246 in the textbook directing us to
22 # this code:
23 #
24 # In particular, we used the tile-coding software, available at
25 # http://incompleteideas.net/tiles/
26 # tiles3.html, with  $iht = IHT(4096)$  and  $tiles(iht, 8, [8 * x / (0.5 + 1.2), 8 * xdot / (0.07 + 0.07)], A)$  to get
27 # the indices of the ones in the feature vector for state  $(x, xdot)$  and
28 # action  $A$ .
29
30 basehash = hash
31
32 class IHT:
33     "Structure to handle collisions"
34     def __init__(self, sizeval):
35         self.size = sizeval
36         self.overfullCount = 0
37         self.dictionary = {}
38
39     def __str__(self):
40         "Prepares a string for printing whenever this object is printed"
41
42         return "Collision table:" + \
43             "\nsize:" + str(self.size) + \
44             "\noverfullCount:" + str(self.overfullCount) + \
45             "\ndictionary:" + str(len(self.dictionary)) + "\nitems"
46
47     def count(self):
48         return len(self.dictionary)
49
50     def fullp(self):
51         return len(self.dictionary) >= self.size
52
53     def getindex(self, obj, readonly=False):
54         d = self.dictionary
55         if obj in d: return d[obj]
56         elif readonly: return None

```

```

54         size = self.size
55         count = self.count()
56         if count >= size:
57             if self.overfullCount==0: print('IHT_full, _starting_to_
                    allow_collisions')
58             self.overfullCount += 1
59             return basehash(obj) % self.size
60         else:
61             d[obj] = count
62             return count
63
64 def hashcoords(coordinates, m, readonly=False):
65     if type(m)==IHT: return m.getindex(tuple(coordinates), readonly)
66     if type(m)==int: return basehash(tuple(coordinates)) % m
67     if m==None: return coordinates
68
69 def tiles (ihtORsize, numtilings, floats, ints=[], readonly=False):
70     """returns num-tilings tile indices corresponding to the floats and
        ints"""
71     qfloats = [floor(f*numtilings) for f in floats]
72     Tiles = []
73     for tiling in range(numtilings):
74         tilingX2 = tiling*2
75         coords = [tiling]
76         b = tiling
77         for q in qfloats:
78             coords.append( (q + b) // numtilings )
79             b += tilingX2
80         coords.extend(ints)
81         Tiles.append(hashcoords(coords, ihtORsize, readonly))
82     return Tiles
83
84 def tileswrap (ihtORsize, numtilings, floats, wrapwidths, ints=[],
        readonly=False):
85     """returns num-tilings tile indices corresponding to the floats and
        ints, wrapping some floats"""
86     qfloats = [floor(f*numtilings) for f in floats]
87     Tiles = []
88     for tiling in range(numtilings):
89         tilingX2 = tiling*2
90         coords = [tiling]
91         b = tiling
92         for q, width in zip_longest(qfloats, wrapwidths):
93             c = (q + b%numtilings) // numtilings
94             coords.append(c*width if width else c)
95             b += tilingX2
96         coords.extend(ints)
97         Tiles.append(hashcoords(coords, ihtORsize, readonly))

```

```

98     return Tiles
99
100 # End of Tile Coding Software
101 #
102     #####
103 #
104 # Environment Generation
105 register(
106     id = 'MountainCar-v1',
107     entry_point = 'gym.envs.classic_control:MountainCarEnv',
108     max_episode_steps = 5000
109 )
110 ENV = gym.make('MountainCar-v1')
111 #
112     #####
113 #
114 # cp code
115 # all possible actions
116 ACTION_REVERSE = -1
117 ACTION_ZERO = 0
118 ACTION_FORWARD = 1
119 # order is important
120 ACTIONS = [ACTION_REVERSE, ACTION_ZERO, ACTION_FORWARD]
121 # bound for position and velocity
122 POSITION_MIN = -1.2
123 POSITION_MAX = 0.5
124 VELOCITY_MIN = -0.07
125 VELOCITY_MAX = 0.07
126
127 # discount is always 1.0 in these experiments
128 DISCOUNT = 1.0
129
130 # use optimistic initial value, so it's ok to set epsilon to 0
131 EPSILON = 0
132
133 # maximum steps per episode
134 STEP_LIMIT = 5000
135
136 # take an @action at @position and @velocity
137 # @return: new position, new velocity, reward (always -1)
138 def step(position, velocity, action):
139     new_velocity = velocity + 0.001 * action - 0.0025 * np.cos(3 *
140         position)
141     new_velocity = min(max(VELOCITY_MIN, new_velocity), VELOCITY_MAX)

```

```

141     new_position = position + new_velocity
142     new_position = min(max(POSITION_MIN, new_position), POSITION_MAX)
143     reward = -1.0
144     if new_position == POSITION_MIN:
145         new_velocity = 0.0
146     return new_position, new_velocity, reward
147
148 # accumulating trace update rule
149 # @trace: old trace (will be modified)
150 # @activeTiles: current active tile indices
151 # @lam: lambda
152 # @return: new trace for convenience
153 def accumulating_trace(trace, active_tiles, lam):
154     trace *= lam * DISCOUNT
155     trace[active_tiles] += 1
156     return trace
157
158 # replacing trace update rule
159 # @trace: old trace (will be modified)
160 # @activeTiles: current active tile indices
161 # @lam: lambda
162 # @return: new trace for convenience
163 def replacing_trace(trace, activeTiles, lam):
164     active = np.in1d(np.arange(len(trace)), activeTiles)
165     trace[active] = 1
166     trace[~active] *= lam * DISCOUNT
167     return trace
168
169 # wrapper class for Sarsa(lambda)
170 class Sarsa:
171     # @maxSize: the maximum # of indices
172     def __init__(self, step_size, lam, trace_update=accumulating_trace,
173                 num_of_tilings=8, max_size=2048):
174         self.max_size = max_size
175         self.num_of_tilings = num_of_tilings
176         self.trace_update = trace_update
177         self.lam = lam
178
179         # divide step size equally to each tiling
180         self.step_size = step_size / num_of_tilings
181
182         self.hash_table = IHT(max_size)
183
184         # weight for each tile
185         self.weights = np.zeros(max_size)
186
187         # trace for each tile
188         self.trace = np.zeros(max_size)

```



```

188
189     # position and velocity needs scaling to satisfy the tile
       software
190     self.position_scale = self.num_of_tilings / (POSITION_MAX -
       POSITION_MIN)
191     self.velocity_scale = self.num_of_tilings / (VELOCITY_MAX -
       VELOCITY_MIN)
192
193     # get indices of active tiles for given state and action
194     def get_active_tiles(self, position, velocity, action):
195         # I think positionScale * (position - position_min) would be a
       good normalization.
196         # However positionScale * position_min is a constant, so it's
       ok to ignore it.
197         active_tiles = tiles(self.hash_table, self.num_of_tilings,
198                             [self.position_scale * position, self.
                               velocity_scale * velocity],
199                             [action])
200         return active_tiles
201
202     # estimate the value of given state and action
203     def value(self, position, velocity, action):
204         if position == POSITION_MAX:
205             return 0.0
206         active_tiles = self.get_active_tiles(position, velocity, action
       )
207         return np.sum(self.weights[active_tiles])
208
209     # learn with given state, action and target
210     def learn(self, position, velocity, action, target):
211         active_tiles = self.get_active_tiles(position, velocity, action
       )
212         estimation = np.sum(self.weights[active_tiles])
213         delta = target - estimation
214         if self.trace_update == accumulating_trace or self.trace_update
       == replacing_trace:
215             self.trace_update(self.trace, active_tiles, self.lam)
216         else:
217             raise Exception('Unexpected_Trace_Type')
218         self.weights += self.step_size * delta * self.trace
219
220     # get # of steps to reach the goal under current state value
       function
221     def cost_to_go(self, position, velocity):
222         costs = []
223         for action in ACTIONS:
224             costs.append(self.value(position, velocity, action))
225         return -np.max(costs)

```

```

226
227 # get action at @position and @velocity based on epsilon greedy policy
    and @valueFunction
228 def get_action(position , velocity , valueFunction):
229     if np.random.binomial(1, EPSILON) == 1:
230         return np.random.choice(ACTIONS)
231     values = []
232     for action in ACTIONS:
233         values.append(valueFunction.value(position , velocity , action))
234     return np.argmax(values) - 1
235
236 # play Mountain Car for one episode based on given method @evaluator
237 # @return: total steps in this episode
238 def play(evaluator):
239     position = np.random.uniform(-0.6, -0.4)
240     velocity = 0.0
241     action = get_action(position , velocity , evaluator)
242     steps = 0
243     while True:
244         next_position , next_velocity , reward = step(position , velocity ,
            action)
245         next_action = get_action(next_position , next_velocity ,
            evaluator)
246         steps += 1
247         target = reward + DISCOUNT * evaluator.value(next_position ,
            next_velocity , next_action)
248         evaluator.learn(position , velocity , action , target)
249         position = next_position
250         velocity = next_velocity
251         action = next_action
252         if next_position == POSITION_MAX:
253             break
254         if steps >= STEP_LIMIT:
255             print('Step_Limit_Exceeded!')
256             break
257     return steps
258
259 # end cp code
260 #
    #####
261
262 # Question 1 plots
263 def q1plots():
264     runs = 5
265     episodes = 50
266     lambdas = [0,0.9]
267

```

```

268 # Part (b) – Generation of plot using replacing traces
269 alphas = np.arange(0.6,2.0,0.2) / 8.0
270 steps = np.zeros((len(lambdas), len(alphas), runs, episodes))
271 for lambdaIdx, lam in enumerate(lambdas):
272     for alphaIdx, alpha in enumerate(alphas):
273         for run in tqdm(range(runs)):
274             evaluator = Sarsa(alpha, lam, replacing_trace, max_size
                               =4096)
275             for ep in range(episodes):
276                 step = play(evaluator)
277                 steps[lambdaIdx, alphaIdx, run, ep] = step
278
279 # average over episodes
280 steps = np.mean(steps, axis=3)
281 # average over runs
282 steps = np.mean(steps, axis=2)
283
284 for lamdaIdx, lam in enumerate(lambdas):
285     plt.plot(alphas, steps[lamdaIdx, :], label='lambda_=%s' % (str
        (lam)))
286 plt.xlabel('alpha_*_#_of_tilings_(8)')
287 plt.ylabel('averaged_steps_per_episode')
288 plt.title('Sarsa_with_replacing_traces')
289 #plt.ylim([180, 300])
290 plt.legend()
291
292 plt.savefig('replacing_traces.png')
293 plt.close()
294
295 print("Completed_Problem_1_Part_(b)")
296
297 # Part (c) – Generation of plot using accumulating traces
298 alphas = np.arange(0.2,0.5,0.05) / 8.0
299 steps = np.zeros((len(lambdas), len(alphas), runs, episodes))
300 for lambdaIdx, lam in enumerate(lambdas):
301     for alphaIdx, alpha in enumerate(alphas):
302         for run in tqdm(range(runs)):
303             evaluator = Sarsa(alpha, lam, accumulating_trace,
                               max_size=4096)
304             for ep in range(episodes):
305                 step = play(evaluator)
306                 steps[lambdaIdx, alphaIdx, run, ep] = step
307
308 # average over episodes
309 steps = np.mean(steps, axis=3)
310 # average over runs
311 steps = np.mean(steps, axis=2)
312

```

```

313     for lamdaIdx, lam in enumerate(lambdas):
314         plt.plot(alphas, steps[lamdaIdx, :], label='lambda_=%s' % (str
            (lam)))
315     plt.xlabel('alpha_*_#_of_tilings_(8)')
316     plt.ylabel('averaged_steps_per_episode')
317     plt.title('Sarsa_with_accumulating_traces')
318     #plt.ylim([180, 300])
319     plt.legend()
320
321     plt.savefig('accumulating_traces.png')
322     plt.close()
323
324     print("Completed_Problem_1_Part_(c)")
325
326     #
    #####

327     #
328     # Work for Q2
329
330     class ShortCorridor:
331         """
332         Short corridor environment, see Example 13.1
333         """
334         def __init__(self):
335             self.reset()
336
337         def reset(self):
338             self.state = 0
339
340         def step(self, go_right):
341             """
342             Args:
343                 go_right (bool): chosen action
344             Returns:
345                 tuple of (reward, episode terminated?)
346             """
347             if self.state == 0 or self.state == 2:
348                 if go_right:
349                     self.state += 1
350                 else:
351                     self.state = max(0, self.state - 1)
352             else:
353                 if go_right:
354                     self.state -= 1
355                 else:
356                     self.state += 1
357

```

```

358         if self.state == 3:
359             # terminal state
360             return 0, True
361         else:
362             return -1, False
363
364     def softmax(x):
365         t = np.exp(x - np.max(x))
366         return t / np.sum(t)
367
368     class ReinforceAgent:
369         """
370         ReinforceAgent that follows algorithm
371         'REINFORNCE Monte-Carlo Policy-Gradient Control (episodic)'
372         """
373         def __init__(self, alpha, gamma):
374             # set values such that initial conditions correspond to left-
375             # epsilon greedy
376             self.theta = np.array([-1.47, 1.47])
377             self.alpha = alpha
378             self.gamma = gamma
379             # first column - left, second - right
380             self.x = np.array([[0, 1],
381                                [1, 0]])
382             self.rewards = []
383             self.actions = []
384
385         def get_pi(self):
386             h = np.dot(self.theta, self.x)
387             t = np.exp(h - np.max(h))
388             pmf = t / np.sum(t)
389             # never become deterministic,
390             # guarantees episode finish
391             imin = np.argmin(pmf)
392             epsilon = 0.05
393
394             if pmf[imin] < epsilon:
395                 pmf[:] = 1 - epsilon
396                 pmf[imin] = epsilon
397
398             return pmf
399
400         def get_p_right(self):
401             return self.get_pi()[1]
402
403         def choose_action(self, reward):
404             if reward is not None:
405                 self.rewards.append(reward)

```

```

405
406     pmf = self.get_pi()
407     go_right = np.random.uniform() <= pmf[1]
408     self.actions.append(go_right)
409
410     return go_right
411
412     def episode_end(self, last_reward):
413         self.rewards.append(last_reward)
414
415         # learn theta
416         G = np.zeros(len(self.rewards))
417         G[-1] = self.rewards[-1]
418
419         for i in range(2, len(G) + 1):
420             G[-i] = self.gamma * G[-i + 1] + self.rewards[-i]
421
422         gamma_pow = 1
423
424         for i in range(len(G)):
425             j = 1 if self.actions[i] else 0
426             pmf = self.get_pi()
427             grad_ln_pi = self.x[:, j] - np.dot(self.x, pmf)
428             update = self.alpha * gamma_pow * G[i] * grad_ln_pi
429
430             self.theta += update
431             gamma_pow *= self.gamma
432
433         self.rewards = []
434         self.actions = []
435
436     class ReinforceBaselineAgent(ReinforceAgent):
437         def __init__(self, alpha, gamma, alpha_w):
438             super(ReinforceBaselineAgent, self).__init__(alpha, gamma)
439             self.alpha_w = alpha_w
440             self.w = 0
441
442         def episode_end(self, last_reward):
443             self.rewards.append(last_reward)
444
445             # learn theta
446             G = np.zeros(len(self.rewards))
447             G[-1] = self.rewards[-1]
448
449             for i in range(2, len(G) + 1):
450                 G[-i] = self.gamma * G[-i + 1] + self.rewards[-i]
451
452             gamma_pow = 1

```

```

453
454     for i in range(len(G)):
455         self.w += self.alpha_w * gamma_pow * (G[i] - self.w)
456
457         j = 1 if self.actions[i] else 0
458         pmf = self.get_pi()
459         grad_ln_pi = self.x[:, j] - np.dot(self.x, pmf)
460         update = self.alpha * gamma_pow * (G[i] - self.w) *
            grad_ln_pi
461
462         self.theta += update
463         gamma_pow *= self.gamma
464
465         self.rewards = []
466         self.actions = []
467
468 def trial(num_episodes, agent_generator):
469     env = ShortCorridor()
470     agent = agent_generator()
471
472     rewards = np.zeros(num_episodes)
473     for episode_idx in range(num_episodes):
474         rewards_sum = 0
475         reward = None
476         env.reset()
477
478         while True:
479             go_right = agent.choose_action(reward)
480             reward, episode_end = env.step(go_right)
481             rewards_sum += reward
482
483             if episode_end:
484                 agent.episode_end(reward)
485                 break
486
487             rewards[episode_idx] = rewards_sum
488
489     return rewards
490
491 def q2plot1():
492     num_trials = 100
493     num_episodes = 1000
494     alphas = [2**(-12), 2**(-13), 2**(-14)]
495     gamma = 1
496
497     for alpha in alphas:
498
499         print(alpha)

```

```

500
501     rewards = np.zeros((num_trials, num_episodes))
502     agent_generator = lambda : ReinforceAgent(alpha=alpha, gamma=
        gamma)
503
504     for i in tqdm(range(num_trials)):
505         reward = trial(num_episodes, agent_generator)
506         rewards[i, :] = reward
507
508     plt.plot(np.arange(num_episodes) + 1, rewards.mean(axis=0),
        label='alpha_=%s' % (str(alpha)))
509
510     #plt.plot(np.arange(num_episodes) + 1, -11.6 * np.ones(num_episodes
        ), ls='dashed', color='red', label='-11.6')
511     #plt.plot(np.arange(num_episodes) + 1, rewards.mean(axis=0), color
        ='blue')
512     plt.ylabel('total_reward_on_episode')
513     plt.xlabel('episode')
514     plt.legend(loc='lower_right')
515
516     plt.savefig('fig1.png')
517     plt.close()
518
519 def q2plot2():
520     num_trials = 100
521     num_episodes = 1000
522     alpha = 2**(-13)
523     gamma = 1
524     agent_generators = [lambda : ReinforceAgent(alpha=alpha, gamma=
        gamma),
525
        lambda : ReinforceBaselineAgent(alpha=2**(-9),
            gamma=gamma, alpha_w=2**(-6))]
526     labels = ['Reinforce_with_baseline',
527             'Reinforce_without_baseline']
528
529     rewards = np.zeros((len(agent_generators), num_trials, num_episodes
        ))
530
531     for agent_index, agent_generator in enumerate(agent_generators):
532         for i in tqdm(range(num_trials)):
533             reward = trial(num_episodes, agent_generator)
534             rewards[agent_index, i, :] = reward
535
536     #plt.plot(np.arange(num_episodes) + 1, -11.6 * np.ones(num_episodes
        ), ls='dashed', color='red', label='-11.6')
537     for i, label in enumerate(labels):
538         plt.plot(np.arange(num_episodes) + 1, rewards[i].mean(axis=0),
            label=label)

```



```

539     plt.ylabel('total_reward_on_episode')
540     plt.xlabel('episode')
541     plt.legend(loc='lower_right')
542
543     plt.savefig('fig2.png')
544     plt.close()
545
546 if __name__ == '__main__':
547     # Ensuring environment is reset to begin
548     ENV.reset()
549     #

```

---

```

550     # Question 1 – Implementation of replacing traces and accumulating
551     #                 traces
552     #                 to produce plots similar to that of Figure 12.10 in
553     #                 the textbook
554     #

```

---

```

553     #q1plots()
554
555     #

```

---

```

556     # Quesiton 2 – Reproduction of Figures 13.1 and 13.2 from the
557     #                 textbook.
558     #                 Plots are generated using example 13.1
559     #

```

---

```

559     # Figure 13.1 Reproduciton
560     #q2plot1()
561     # Figure 13.2 Reproduction
562     q2plot2()
563
564     #

```

---

```

565     # Quesiton 3 – Attempts to reproduce the work behind Figures 4 and
566     #                 5 in the
567     #                 Options paper. Heatmap figure will be created to
568     #                 represent the results
569     #

```

---

```

568     #
569

```

```
570
571     # Ensuring environment is closed at the end to avoid compilation
        issues
572     ENV.close()
```