

## ▼ ELEN 6885 Reinforcement Learning Coding Assignment (Part 1, 2, 3)

### Taxi Problem Overview

There are 4 locations (labeled by different letters) and your job is to pick up the passenger at one location and drop him off in another. You receive +20 points for a successful drop-off, and lose 1 point for every timestep it takes. There is also a 10 point penalty for illegal pick-up and drop-off actions.

5	R				G
4					
3					
2					
1	Y			B	
	0	2	4	6	8

- blue: passenger
- magenta: destination
- yellow: empty taxi
- green: full taxi
- other letters: locations

Please put your code into the block marked by:

#####

YOUR CODE STARTS HERE

YOUR CODE ENDS HERE

#####

You should not edit anything outside of the block.

## ▼ Playing with the environment

Run the cell below to get a feel for the environment by moving your agent(the taxi) by taking one of the actions at each step.

```
In [14]: 1 from gym.wrappers import Monitor
          2 import gym
          3 import random
          4 import numpy as np
          5 import time
          6 from helpers import utils
```

```

In [3]: 1  """
2  You can test your game now.
3  ▾ Input range from 0 to 5:
4      0 : South (Down)
5      1 : North (Up)
6      2 : East (Right)
7      3 : West (Left)
8      4: Pick up
9      5: Drop off
10     6: exit_game
11  """
12  GAME = "Taxi-v3"
13  env = gym.make(GAME)
14  env = Monitor(env, "taxi_simple", force=True)
15  s = env.reset()
16  steps = 100
17  ▾ for step in range(steps):
18      env.render()
19      action = int(input("Please type in the next action:"))
20  ▾ if action==6:
21      break
22      s, r, done, info = env.step(action)
23      print('state:',s)
24      print('reward:',r)
25      print('Is state terminal?:',done)
26      print('info:',info)
27
28  # close environment and monitor
29  env.close()

```

```

+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+

```

Please type in the next action:6

## ▾ 1.1 Incremental implementation of average

We've finished the incremental implementation of average for you. Please call the function to estimate with  $1/\text{step}$  step size and fixed step size to compare the difference between these two on a simulated Bandit problem.

```
In [4]: 1 def estimate(OldEstimate, StepSize, Target):
2         '''An incremental implementation of average.
3         OldEstimate : float
4         StepSize : float
5         Target : float
6         '''
7         NewEstimate = OldEstimate + StepSize * (Target - OldEstimate)
8         return NewEstimate
```

```
In [5]: 1 random.seed(6885)
2 numTimeStep = 10000
3 q_h = np.zeros(numTimeStep + 1) # Q Value estimate with 1/step step si
4 q_f = np.zeros(numTimeStep + 1) # Q value estimate with fixed step siz
5 FixedStepSize = 0.5 #A large number to exaggerate the difference
6 for step in range(1, numTimeStep + 1):
7     if step < numTimeStep / 2:
8         r = random.gauss(mu = 1, sigma = 0.1)
9     else:
10        r = random.gauss(mu = 3, sigma = 0.1)
11
12    #TIPS: Call function estimate defined in ./helpers/utils.py
13    #####
14    # YOUR CODE STARTS HERE
15    q_h[step] = estimate(q_h[step-1], 1/step, r)
16    q_f[step] = estimate(q_f[step-1], FixedStepSize, r)
17
18    # YOUR CODE ENDS HERE
19    #####
20
21 q_h = q_h[1:]
22 q_f = q_f[1:]
```

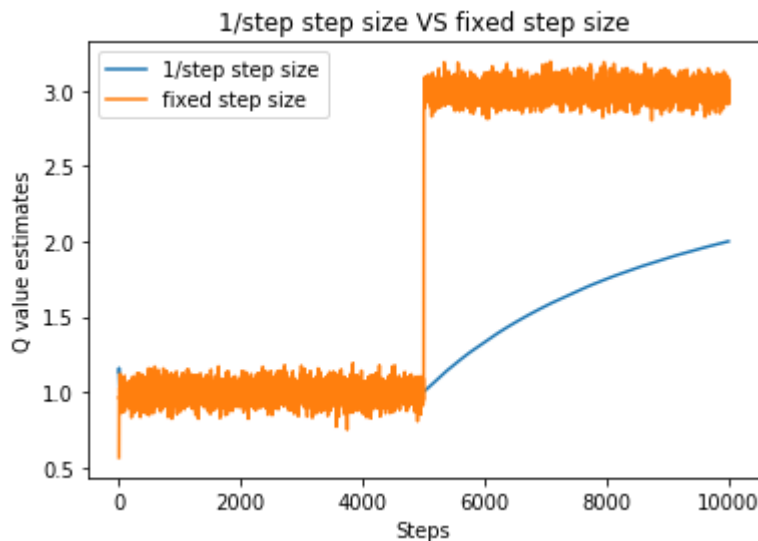
Plot the two Q value estimates. (Please include a title, labels on both axes, and legends)

```

In [5]: 1 import matplotlib.pyplot as plt
        2 #####
        3 # YOUR CODE STARTS HERE
        4 plt.title('1/step step size VS fixed step size')
        5 plt.plot(q_h, label='1/step step size')
        6 plt.plot(q_f, label='fixed step size')
        7 plt.xlabel('Steps')
        8 plt.ylabel('Q value estimates')
        9 plt.legend()
       10 # YOUR CODE ENDS HERE
       11 #####

```

Out[5]: <matplotlib.legend.Legend at 0x1dca7ce9dd8>



## 1.2 $\epsilon$ -Greedy for Exploration

In Reinforcement Learning, we are always faced with the dilemma of exploration and exploitation.  $\epsilon$ -Greedy is a trade-off between them. You are supposed to implement Greedy and  $\epsilon$ -Greedy. We combine these two policies in one function by treating Greedy as  $\epsilon$ -Greedy where  $\epsilon = 0$ . Edit the function `epsilon_greedy` the following block.

```

In [6]: 1 def epsilon_greedy(value, e, seed = None):
2         '''
3         Implement Epsilon-Greedy policy.
4
5         Inputs:
6         value: numpy ndarray
7         A vector of values of actions to choose from
8         e: float
9         Epsilon
10        seed: None or int
11        Assign an integer value to remove the randomness
12
13        Outputs:
14        action: int
15        Index of the chosen action
16        '''
17        assert len(value.shape) == 1
18        assert 0 <= e <= 1
19
20        if seed != None:
21            np.random.seed(seed)
22
23            #####
24            # YOUR CODE STARTS HERE
25            p = np.random.uniform(low=0.0, high=1.0, size=None)
26            #print(p)
27        if p < e:
28            action = random.randint(0, len(value)-1)
29        else:
30            action = random.choice(np.where(value == value.max())[0])
31
32            # YOUR CODE ENDS HERE
33            #####
34        return action

```

```

In [7]: 1 np.random.seed(6885) #Set the seed for reproducibility
2        q = np.random.normal(0, 1, size = 5)
3        #####
4        # YOUR CODE STARTS HERE
5        e_greedy_action = epsilon_greedy(q, .1, 6885)
6        greedy_action = epsilon_greedy(q, 0, 6885)
7        # YOUR CODE ENDS HERE
8        #####
9        print('Values:')
10       print(q)
11       print('Greedy Choice =', greedy_action)
12       print('Epsilon-Greedy Choice =', e_greedy_action)

```

Values:

```
[ 0.61264537  0.27923079 -0.84600857  0.05469574 -1.09990968]
```

Greedy Choice = 0

Epsilon-Greedy Choice = 0

You should get the following results:

Values:

[ 0.61264537 0.27923079 -0.84600857 0.05469574 -1.09990968]

Greedy Choice = 0

Epsilon-Greedy Choice = 0



## 1.3 Exploration VS. Exploitation

Try to reproduce Figure 2.2 (the upper one is enough) of the Sutton's book based on the experiment described in Chapter 2.3.

```

In [11]: 1  # Do the experiment and record average reward acquired in each time st
2  #####
3  # YOUR CODE STARTS HERE
4  np.random.seed(100) #Set the seed for reproducibility
5  trials = 2000
6  numTimeStep = 1000
7
8  e_final_results = np.zeros(numTimeStep + 1)
9  e2_final_results = np.zeros(numTimeStep + 1)
10 greedy_final_results = np.zeros(numTimeStep + 1)
11
12
13
14  for trial in range(1, trials+1):
15      e_gestimate = np.zeros(10)
16      e2_gestimate = np.zeros(10)
17      g_gestimate = np.zeros(10)
18      q_arms = np.random.normal(0, 1, size = 10)
19
20      e_counter = {}
21      e2_counter = {}
22      greedy_counter = {}
23
24      for step in range(1, numTimeStep + 1):
25
26          e_greedy_action = epsilon_greedy(e_gestimate, .01)
27          e_greedy_action_2 = epsilon_greedy(e2_gestimate, .1)
28          greedy_action = epsilon_greedy(g_gestimate, 0)
29
30          reward_e = q_arms[e_greedy_action] + np.random.normal(0,1)
31          reward_e2 = q_arms[e_greedy_action_2] + np.random.normal(0,1)
32          reward_greedy = q_arms[greedy_action] + np.random.normal(0,1)
33
34          if e_greedy_action in e_counter:
35              e_counter[e_greedy_action] += 1
36          else:
37              e_counter[e_greedy_action] = 1
38
39          if e_greedy_action_2 in e2_counter:
40              e2_counter[e_greedy_action_2] += 1
41          else:
42              e2_counter[e_greedy_action_2] = 1
43
44          if greedy_action in greedy_counter:
45              greedy_counter[greedy_action] += 1
46          else:
47              greedy_counter[greedy_action] = 1
48
49
50          e_gestimate[e_greedy_action] = estimate(e_gestimate[e_greedy_a
51          e2_gestimate[e_greedy_action_2] = estimate(e2_gestimate[e_gree
52          g_gestimate[greedy_action] = estimate(g_gestimate[greedy_actio
53
54          e_final_results[step] = estimate(e_final_results[step], 1/tria
55          e2_final_results[step] = estimate(e2_final_results[step], 1/t
56          greedy_final_results[step] = estimate(greedy_final_results[st

```

```

57
58 # YOUR CODE ENDS HERE
59 #####

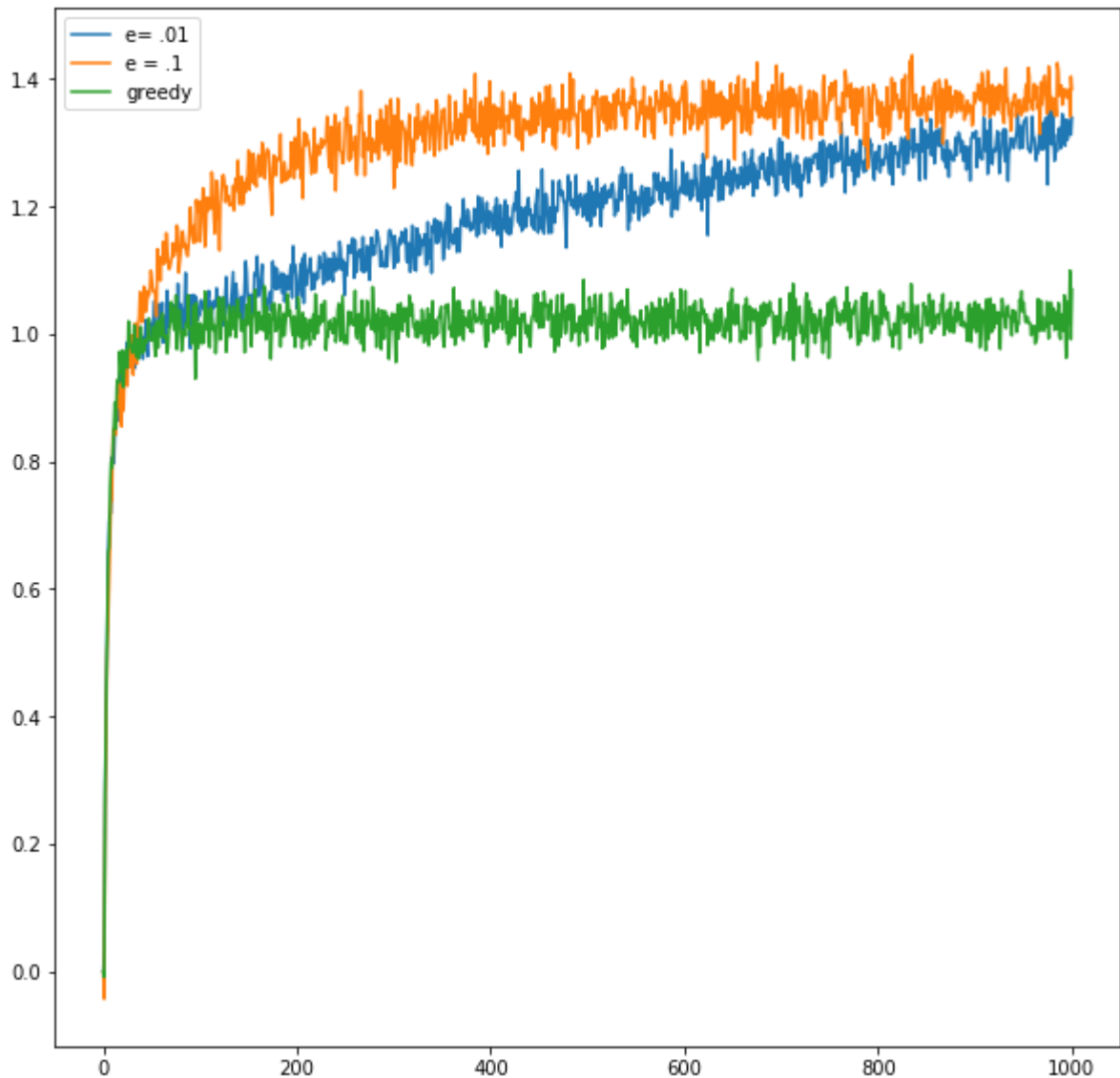
```

```

In [12]: 1 # Plot the average reward
2 #####
3 # YOUR CODE STARTS HERE
4 plt.figure(figsize=(10,10))
5 plt.plot(e_final_results, label = 'e= .01')
6 plt.plot(e2_final_results, label = 'e = .1')
7 plt.plot(greedy_final_results, label = 'greedy')
8 plt.legend()
9
10 # YOUR CODE ENDS HERE
11 #####

```

Out[12]: <matplotlib.legend.Legend at 0x1cd1765d630>



## Question 2



In this question, you will implement the value iteration and policy iteration algorithms to solve the Taxi game problem

## 2.1 Model-based RL: value iteration

For this part, you need to implement the helper functions `action_evaluation(env, gamma, v)`, and `extract_policy(env, v, gamma)` in `utils.py`. Understand `action_selection(q)` which we have implemented.

Use these helper functions to implement the `value_iteration` algorithm below.

```
In [7]: 1 import numpy as np
2 from helpers import utils
3 np.random.seed(6885) #Set the seed for reproducibility
4 def value_iteration(env, gamma, max_iteration, theta):
5     """
6     Implement value iteration algorithm. You should use extract_policy
7
8     Parameters
9     -----
10    env: OpenAI env.
11        env.P: dictionary
12            the transition probabilities of the environment
13            P[state][action] is tuples with (probability, next
14        env.nS: int
15            number of states
16        env.nA: int
17            number of actions
18    gamma: float
19        Discount factor.
20    max_iteration: int
21        The maximum number of iterations to run before stopping.
22    theta: float
23        Determines when value function has converged.
24    Returns:
25    -----
26    value function: np.ndarray
27    policy: np.ndarray
28    """
29    V = np.zeros(env.nS)
30    #####
31    # YOUR CODE STARTS HERE
32    for i in range(max_iteration):
33        q = utils.action_evaluation(env, gamma, V)
34        V_prev = V
35        V = np.max(q, axis=1)
36        if np.abs(V-V_prev).max() < theta:
37            break
38    policy = utils.extract_policy(env, V, gamma)
39    # YOUR CODE ENDS HERE
40    #####
41
42    return V, policy
43
```

After implementing the above function, read and understand the functions implemented in `evaluation_utils.py`, which we will use to evaluate our value iteration policy

```
In [9]: 1  from helpers import evaluation_utils
        2  import gym
        3  GAME = "Taxi-v3"
        4  env = gym.make(GAME)
        5  V_vi, policy_vi = value_iteration(env, gamma=0.95, max_iteration=6000,
        6  # visualize how the agent performs with the policy generated from value
        7  evaluation_utils.render_episode(env, policy_vi)
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(East)

```
+-----+
|R: | : :G|
| : | : : |
| : | : : |
+-----+
```

```
In [10]: 1  # evaluate the performance of value iteration over 100 episodes
        2  evaluation_utils.avg_performance(env, policy_vi)
```

Out[10]: 8.323232323232324

## ▼ 2.2 Model-based RL: policy iteration

In this part, you are supposed to implement policy iteration to solve the Taxi game problem.

```

In [8]: 1  from helpers import utils
2  ▾ def policy_iteration(env, gamma, max_iteration, theta):
3      """Implement Policy iteration algorithm.
4
5      You should use the policy_evaluation and policy_improvement methods to
6      implement this method.
7
8      Parameters
9      -----
10 ▾   env: OpenAI env.
11 ▾       env.P: dictionary
12           the transition probabilities of the environment
13           P[state][action] is tuples with (probability, next state, action)
14 ▾       env.nS: int
15           number of states
16 ▾       env.nA: int
17           number of actions
18 ▾   gamma: float
19           Discount factor.
20 ▾   max_iteration: int
21           The maximum number of iterations to run before stopping.
22 ▾   theta: float
23           Determines when value function has converged.
24   Returns:
25   -----
26   value function: np.ndarray
27   policy: np.ndarray
28   """
29   policy = np.zeros(env.nS, dtype=int)
30   #####
31   # YOUR CODE STARTS HERE
32 ▾   for i in range(max_iteration):
33       V = policy_evaluation(env, policy, gamma, theta)
34       policy, policy_stable = policy_improvement(env, V, policy, gamma, theta)
35 ▾       if policy_stable:
36           break
37   # YOUR CODE ENDS HERE
38   #####
39
40   return V, policy
41
42
43 ▾ def policy_evaluation(env, policy, gamma, theta):
44     """Evaluate the value function from a given policy.
45
46     Parameters
47     -----
48 ▾   env: OpenAI env.
49 ▾       env.P: dictionary
50           the transition probabilities of the environment
51           P[state][action] is tuples with (probability, next state, action)
52 ▾       env.nS: int
53           number of states
54 ▾       env.nA: int
55           number of actions
56

```

```

57 ▼     gamma: float
58         Discount factor.
59 ▼     policy: np.array
60         The policy to evaluate. Maps states to actions.
61 ▼     max_iteration: int
62         The maximum number of iterations to run before stopping.
63 ▼     theta: float
64         Determines when value function has converged.
65     Returns
66     -----
67 ▼     value function: np.ndarray
68         The value function from the given policy.
69     """
70     P = env.P
71     V = np.zeros(env.nS)
72     #####
73     # YOUR CODE STARTS HERE
74 ▼     while True:
75         delta = 0
76 ▼         for s in range(env.nS):
77             a = policy[s]
78             q_s_a = 0
79 ▼             for i in range(len(P[s][a])):
80                 next_state_tuple = P[s][a][i]
81                 v_next_state = V[next_state_tuple[1]]
82                 p_next_state = next_state_tuple[0]
83                 reward_next_state = next_state_tuple[2]
84                 q_s_a += p_next_state*(reward_next_state + gamma*v_ne
85             delta = max(delta, np.abs(q_s_a - V[s]))
86             V[s] = q_s_a
87
88 ▼         if delta < theta:
89             break
90         # YOUR CODE ENDS HERE
91         #####
92     return V
93
94
95 ▼ def policy_improvement(env, value_from_policy, policy, gamma):
96     """Given the value function from policy, improve the policy.
97
98     Parameters
99     -----
100 ▼     env: OpenAI env
101 ▼         env.P: dictionary
102                 the transition probabilities of the environment
103                 P[state][action] is tuples with (probability, nex
104 ▼         env.nS: int
105                 number of states
106 ▼         env.nA: int
107                 number of actions
108
109 ▼     value_from_policy: np.ndarray
110         The value calculated from the policy
111 ▼     policy: np.array
112         The previous policy.
113 ▼     gamma: float

```

```

114         Discount factor.
115
116     Returns
117     -----
118     new policy: np.ndarray
119         An array of integers. Each integer is the optimal action
120         in that state according to the environment dynamics and the
121         given value function.
122     stable policy: bool
123         True if the optimal policy is found, otherwise false
124     """
125     #####
126     # YOUR CODE STARTS HERE
127     nS = env.nS
128     nA = env.nA
129     q = np.zeros((nS, nA))
130     P = env.P
131     policy_stable = False
132     v = value_from_policy
133     for s in range(nS):
134         for a in range(nA):
135             #####
136             # YOUR CODE STARTS HERE
137             #[probability, nextstate, reward, terminal]=P[s][a]
138             q_s_a = 0
139             for i in range(len(P[s][a])):
140                 next_state_tuple = P[s][a][i]
141                 v_next_state = v[next_state_tuple[1]]
142                 p_next_state = next_state_tuple[0]
143                 reward_next_state = next_state_tuple[2]
144                 q_s_a += p_next_state*(reward_next_state + gamma*v_next_state)
145             q[s][a] = q_s_a
146
147     new_policy = utils.action_selection(q)
148
149     if np.array_equal(policy, new_policy):
150         policy_stable = True
151
152     # YOUR CODE ENDS HERE
153     #####
154
155     return new_policy, policy_stable
156

```

```
In [64]: 1  ## Testing out policy iteration policy for one episode
2  GAME = "Taxi-v3"
3  evaluation_utils.render_episode(env, policy_vi)
4  env = gym.make("Taxi-v3")
5  V_pi, policy_pi = policy_iteration(env, gamma=0.95, max_iteration=6000
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(North)

```
+-----+
|R: | : :G|
| : | : : |
| : | : : |
| : | : : |
+-----+
```

In [63]:

1 # visualize how the agent performs with the policy generated from poli  
2 evaluation\_utils.render\_episode(env, policy\_pi)

+-----+  
|R: | : :G|  
| : | : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+

+-----+  
|R: | : :G|  
| : | : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+

(East)

+-----+  
|R: | : :G|  
| : | : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+

(North)

+-----+  
|R: | : :G|  
| : | : : |  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+

(Pickup)

+-----+  
|R: | : :G|  
| : | : :\_|  
| : : : : |  
| | : | : |  
|Y| : |B: |  
+-----+

(South)

+-----+  
|R: | : :G|  
| : | : : |  
| : : : :\_|  
| | : | : |  
|Y| : |B: |  
+-----+

(South)

+-----+  
|R: | : :G|  
| : | : : |  
| : : : :\_|  
| | : | : |

```
| Y | : | B : |
+-----+
(West)
+-----+
|R: | : :G|
| : | : : |
| : : _ : : |
| | : | : |
| Y | : | B : |
+-----+
(West)
+-----+
|R: | : :G|
| : | : : |
| : _ : : : |
| | : | : |
| Y | : | B : |
+-----+
(West)
+-----+
|R: | : :G|
| : | : : |
| _ : : : : |
| | : | : |
| Y | : | B : |
+-----+
(West)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| _ | : | : |
| Y | : | B : |
+-----+
(South)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(South)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(Dropoff)
Episode reward: 9.000000
```



```
In [65]: 1 # evaluate the performance of policy iteration over 100 episodes  
2 print(evaluation_utils.avg_performance(env, policy_pi))
```

```
8.484848484848484
```

## ▼ Part 3: Q-learning and SARSA

### ▼ 3.1 Model-free RL: Q-learning

In this part, you will implement Q-learning.

```

In [23]: 1 def QLearning(env, num_episodes, gamma, lr, e):
2         """
3         Implement the Q-learning algorithm following the epsilon-greedy ex
4         Inputs:
5         env: OpenAI Gym environment
6         env.P: dictionary
7             P[state][action] are tuples of tuples tuples with
8             probability: float
9             nextstate: int
10            reward: float
11            terminal: boolean
12         env.nS: int
13            number of states
14         env.nA: int
15            number of actions
16         num_episodes: int
17            Number of episodes of training
18         gamma: float
19            Discount factor.
20         lr: float
21            Learning rate.
22         e: float
23            Epsilon value used in the epsilon-greedy method.
24         Outputs:
25         Q: numpy.ndarray
26         """
27
28         Q = np.zeros((env.nS, env.nA))
29
30         #####
31         # YOUR CODE STARTS HERE
32         # Set the percent you want to explore
33         for i in range(num_episodes):
34
35             s = env.reset()
36
37
38             done = False
39             while not done:
40                 if np.random.uniform(0, 1) < e:
41                     """
42                     Explore: select a random action
43                     """
44                     a = np.random.choice(range(env.nA))
45                 else:
46                     """
47                     Exploit: select the action with max value (future rewa
48                     """
49                     a = np.argmax(Q[s])
50
51                 next_s, r, done, _ = env.step(a)
52                 next_a = np.argmax(Q[next_s])
53                 Q[s][a] = Q[s][a] + lr*(r + gamma*Q[next_s][next_a] - Q[s]
54
55                 s = next_s
56

```

```

57     # YOUR CODE ENDS HERE
58     #####
59
60     return Q

```

```

In [17]: 1  def render_episode_Q(env, Q):
          2  """Renders one episode for Q function on environment.
          3
          4      Parameters
          5      -----
          6  env: gym.core.Environment
          7      Environment to play Q function on.
          8  Q: np.array of shape [env.nS x env.nA]
          9      state-action values.
         10  """
         11
         12  episode_reward = 0
         13  state = env.reset()
         14  done = False
         15  while not done:
         16      env.render()
         17      time.sleep(0.5)
         18      action = np.argmax(Q[state])
         19      state, reward, done, _ = env.step(action)
         20      episode_reward += reward
         21
         22  print ("Episode reward: %f" %episode_reward)

```

```

In [24]: 1  Q = QLearning(env = env.env, num_episodes = 1000, gamma = 1, lr = 0.1,
          2  print('Action values:'))
          3  print(Q)

```

Action values:

```

[[ 0.          0.          0.          0.          0.          0.
]
 [-3.19794405 -3.87267875 -4.05310345 -3.7965246   7.81255557 -5.1332576
]
 [-0.46050136  0.02574472 -1.7852255  -0.95707712 14.41700161 -1.9119776
7]
...
 [-1.1         -1.05841259 -1.1         -1.08586767 -2.         -2.
]
 [-2.58165982 -2.64337669 -2.6556984  -0.57262706 -2.94391081 -3.9352165
]
 [-0.2         -0.2         -0.2         1.909         -1.         -1.
]]

```

```
In [29]: 1 # Uncomment the following to evaluate your result, comment them when y
2 env = gym.make('Taxi-v3')
3 Q = QLearning(env = env.env, num_episodes = 1000, gamma = 1, lr = 0.1,
4 render_episode_Q(env.env, Q)
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(North)

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(West)

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(West)

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(West)

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
```

(South)

```
+-----+
|R: | : :G|
| : | : : |
```

```

| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(South)
+-----+
| R : | : : G |
| : | : : |
| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(Pickup)
+-----+
| R : | : : G |
| : | : : |
| : : : : |
| _ | : | : |
| Y | : | B : |
+-----+
(North)
+-----+
| R : | : : G |
| : | : : |
| _ : : : : |
| | : | : |
| Y | : | B : |
+-----+
(North)
+-----+
| R : | : : G |
| _ : | : : |
| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(North)
+-----+
| R : | : : G |
| : | : : |
| : : : : |
| | : | : |
| Y | : | B : |
+-----+
(North)
Episode reward: 9.000000

```



## 3.2 Model-free RL: SARSA

In this part, you will implement Sarsa.

```

In [19]: 1 def SARSA(env, num_episodes, gamma, lr, e):
2         """
3         Implement the SARSA algorithm following epsilon-greedy exploration
4         Inputs:
5         env: OpenAI Gym environment
6         env.P: dictionary
7             P[state][action] are tuples of tuples tuples with
8             probability: float
9             nextstate: int
10            reward: float
11            terminal: boolean
12         env.nS: int
13            number of states
14         env.nA: int
15            number of actions
16         num_episodes: int
17            Number of episodes of training
18         gamma: float
19            Discount factor.
20         lr: float
21            Learning rate.
22         e: float
23            Epsilon value used in the epsilon-greedy method.
24         Outputs:
25         Q: numpy.ndarray
26            State-action values
27         """
28         Q = np.zeros((env.nS, env.nA))
29         #####
30         # YOUR CODE STARTS HERE
31         for i in range(num_episodes):
32
33             s = env.reset()
34             if np.random.uniform(0, 1) < e:
35                 """
36                 Explore: select a random action
37                 """
38                 a = np.random.choice(range(env.nA))
39             else:
40                 """
41                 Exploit: select the action with max value (future reward)
42                 """
43                 a = np.argmax(Q[s])
44
45             done = False
46             while not done:
47
48                 next_s, r, done, _ = env.step(a)
49
50                 if np.random.uniform(0, 1) < e:
51                     """
52                     Explore: select a random action
53                     """
54                     next_a = np.random.choice(range(env.nA))
55                 else:
56                     """

```

```

57         Exploit: select the action with max value (future reward)
58         """
59         next_a = np.argmax(Q[next_s])
60         Q[s][a] = Q[s][a] + lr*(r + gamma*Q[next_s][next_a] - Q[s][a])
61
62         s = next_s
63         a = next_a
64
65     # YOUR CODE ENDS HERE
66     #####
67
68     return Q

```

```

In [20]: 1 def render_episode_Q(env, Q):
2         """Renders one episode for Q function on environment.
3
4         Parameters
5         -----
6         env: gym.core.Environment
7             Environment to play Q function on.
8         Q: np.array of shape [env.nS x env.nA]
9             state-action values.
10        """
11
12        episode_reward = 0
13        state = env.reset()
14        done = False
15        while not done:
16            env.render()
17            time.sleep(0.5)
18            action = np.argmax(Q[state])
19            state, reward, done, _ = env.step(action)
20            episode_reward += reward
21
22        print ("Episode reward: %f" % episode_reward)

```

```

In [21]: 1 Q = SARSA(env = env.env, num_episodes = 1000, gamma = 1, lr = 0.1, e = 0.01)
2         print('Action values:')
3         print(Q)

```

```

Action values:
[[ 0.          0.          0.          0.          0.          0.
]
 [-4.74823032 -4.52974963 -4.50450374 -3.6320741  -0.68727041 -7.3560936
6]
 [-2.4731041  -1.37530551 -2.33613848  0.12065845  9.84693682 -5.0100564
]
 ...
 [-0.95104891 -0.53842485 -0.98139445 -0.93855415 -4.82619345 -2.09
]
 [-3.30326141 -3.12965199 -3.32168319 -3.33724226 -5.5078639  -5.3846284
]
 [-0.281      -0.299      -0.29          8.14927588 -1.9        -1.91
]]

```

```
In [22]: 1 # Uncomment the following to evaluate your result, comment them when y
2 env = gym.make('Taxi-v3')
3 render_episode_Q(env.env, Q)
```

```
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+

+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(East)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(North)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(Pickup)
+-----+
|R: | : _:G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(West)
+-----+
|R: | : :G|
| : | : _: |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(South)
+-----+
|R: | : :G|
| : | : : |
| : : : _: |
```



```
| | : | : |
|Y| : |B: |
+-----+
(South)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : |_: |
|Y| : |B: |
+-----+
(South)
+-----+
|R: | : :G|
| : | : : |
| : : : : |
| | : | : |
|Y| : |B: |
+-----+
(South)
Episode reward: 12.000000
```