

```

1 import numpy as np    # for creating arrays and matrices
2 import gym            # for using AI labrararies.
3 import random         # generating random numbers

1 env = gym.make("Taxi-v3")
2 env.render()
3
4 # INFORMATION          POSSIBLE VALUES          Total # POSSIBLE VALUES
5
6 # TaxiRow              0, 1, 2, 3, 4              5
7 # TaxiCol              0, 1, 2, 3, 4              5
8 # Destination          R, G, B, Y                4
9 # PassngLoc            R, G, B, Y, onBoard        5
10
11
12 #####Possible Actions #####
13 #####
14 ##### North #####
15 ##### Up #####
16 ##### West <=Left <TAXI> Right => East #####
17 ##### Down #####
18 ##### South #####
19 #####
20 ##### Pickup / Dropoff #####

```

```

+-----+
|R: | : :G|
| : | : :|
|. . . .|
+-----+

```

Saving...



```

1 # Actions:( 1: go north, 2: go south, 3: go west, 4: go east,
2 #           5: pick up passenger, 6: drop off passenger )
3 action_size = env.action_space.n

```

```

4 print("Action size ", action_size)
5
6 state_size = env.observation_space.n    # state space size= 5 * 5 * 4 * 5 = 500
7 print("State size ", state_size)

```

```

    Action size  6
    State size  500

```

```

1 # For each state, action initialize the table entry Q'(s, a) to zero.
2 qtable = np.zeros((state_size, action_size))
3 print(qtable)

```

```

[[0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0.]
 ...
 [0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0.]]

```

```

1 total_episodes = 50000      # Total # of steps(episodes) for Q-Learning prob.
2 total_step_pickup = 100     # Total #of steps(episodes) since pickup passedge
3 max_steps = 99              # Max steps per episode
4
5 learning_rate = 0.7          # Learning rate
6 gamma = 0.618                # Discounting rate [0,1], it is important.
7                               # it determines the reward of immediate and future.
8
9 # Exploration parameters
10 epsilon = 1.0                # Exploration rate
                               # exploration probability at start
                               # minimum exploration probability
11 decay_rate = 0.01            # Exponential decay rate for exploration prob

```

```

1 # 2 For life or until learning is stopped
2 for episode in range(total_episodes):
3     # Reset the environment
4     state = env.reset()

```

```

5     step = 0
6     done = False
7
8     for step in range(max_steps):
9         # 3. Choose an action a in the current world state (s)
10        ## First we randomize a number
11        exp_exp_tradeoff = random.uniform(0,1)
12
13        # If this number > greater than epsilon --> exploitation
14        # (taking the biggest Q value for this state)
15        if exp_exp_tradeoff > epsilon:
16            action = np.argmax(qtable[state,:])
17
18        # Else doing a random choice --> exploration
19        else:
20            action = env.action_space.sample()
21
22        # Take the action (a) and observe the outcome state(s') and reward (r)
23        new_state, reward, done, info = env.step(action)
24
25        # Update Q(s,a):= Q(s,a) + lr [R(s,a) + gamma * max Q(s',a') - Q(s,a)]
26        qtable[state, action] = qtable[state, action] + learning_rate * (reward
27            + gamma * np.max(qtable[new_state, :]) - qtable[state, action])
28
29        # Our new state is state
30        state = new_state
31
32
33        # If done : finish episode
34        if done == True:
35            # If done : finish episode
36            # If done : finish episode
37            # If done : finish episode (need less and less exploration)
38            epsilon = min_epsilon + (max_epsilon - min_epsilon)*np.exp(-decay_rate*episode)
39
40
41        1 env.reset()
42        2 rewards = []
43        3

```

Saving...



```

4 for episode in range(total_step_pickup):
5     state = env.reset()
6     step = 0
7     done = False
8     total_rewards = 0
9     #print("*****")
10    #print("EPISODE ", episode)
11
12    for step in range(max_steps):
13        # UNCOMMENT IT IF YOU WANT TO SEE OUR AGENT PLAYING
14        env.render()
15
16        # Take the action (index) that have the maximum expected future reward
17        # given that state
18        action = np.argmax(qtable[state,:])
19
20        new_state, reward, done, info = env.step(action)
21
22        total_rewards += reward
23
24        if done:
25            rewards.append(total_rewards)
26            print ("Drop-off achieved Score", total_rewards)
27            break
28        state = new_state
29 env.close()
30 print ("Score over time: " + str(sum(rewards)/total_step_pickup))

```

```

|
|R: | : :G|
| : | : : |
| : : : : |
|

```

Saving...



(Pickup)

```

+-----+
|R: | : :G|
|_ : | : : |
| : : : : |
| | : | : |
|v| . |R. |

```

```
| . | . |  
+-----+  
(South)  
+-----+  
| R: | : :G | |
| : | : : |  
| _ | : : : |  
| | : : : |  
| Y | : | B: |  
+-----+
```

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

```
(South)  
+-----+  
| R: | : :G |  
| : | : : |  
| : : : : |  
| | : : : |  
| Y | : | B: |  
+-----+
```

(South)
Drop-off achieved Score 11

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

Saving... X

```
| R: | : :G |  
| : | : : |  
| : : : : |  
| | : : : |  
| Y | : | B: |  
+-----+  
(South)
```

```
(search)
+-----+
|R: | : :G|
| : | : : |
```

```
1 # The Program lists taxi T[s, a] and the R[s, a].
2 ##### Q-Learning Taxi Problem V3 #####
```

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.

Saving...

