# 強化學習概念 Hw06

# Introduction
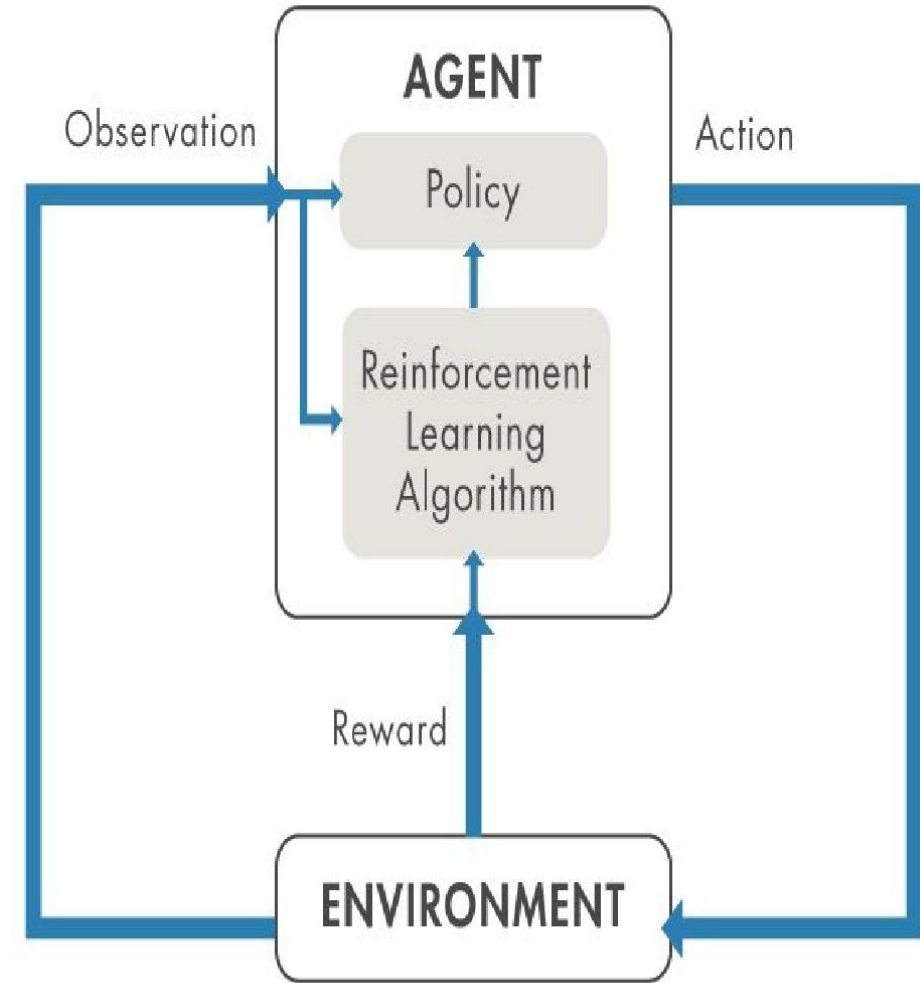
- 強化學習 (Reinforcement learning) 潛力無窮, 能解決許多開發應用上面臨的艱難決策問題, 包括產業自動化、自主駕駛、電玩競技遊戲以及機器人等, 因此備受矚目。

- 強化學習是機器學習 (Machine learning) 的一種, 指的是電腦透過與一個**動態(dynamic) 環境**不斷重複地互動, 來學習正確地執行一項任務。這種**嘗試錯誤(trial-and-error) 的學習方法**, 使電腦在沒有人類干預、沒有被寫入明確的執行任務程式下, 就能夠做出一系列的決策。

- 最著名的強化學習案例就是 **AlphaGo**, 它是第一支打敗人類圍棋比賽世界冠軍的電腦程式。

- 強化學習的運作主要是仰賴動態環境中的資料 -- 也就是會隨著外部條件變化而改變的資料。強化學習演算法的目標, 即是於找出能夠產生最佳結果的策略。強化學習之所以能達成目標, 是藉著軟體當中被稱為主體 (agent) 的部分在環境中進行探索、互動和學習的方法。

- 自助停車(self-parking) 是自動駕駛功能中極為重要的一環, 目標是要讓**車輛中的電腦(主體, agent)** 能準確地尋找位置並將車輛停入正確的停車格。

- 在以下的範例中, 環境(environment)指的是主體之外的所有事物—比如車輛本身的動態、附近的車輛、天候條件等等。訓練過程中, 主體使用從各種感測器如攝影機、GPS、光學雷達(LiDAR)以及其他感測器讀取的資料來產生**駕駛、煞車、與加速指令(動作, action)**。

- 為了學習如何從觀察去產生正確的動作(也就是策略調整, policy tuning), 主體會不斷反覆地嘗試錯誤來試著停車, 而正確的動作會得到一個**獎賞(reward)**。

1.  **Q-learning :**

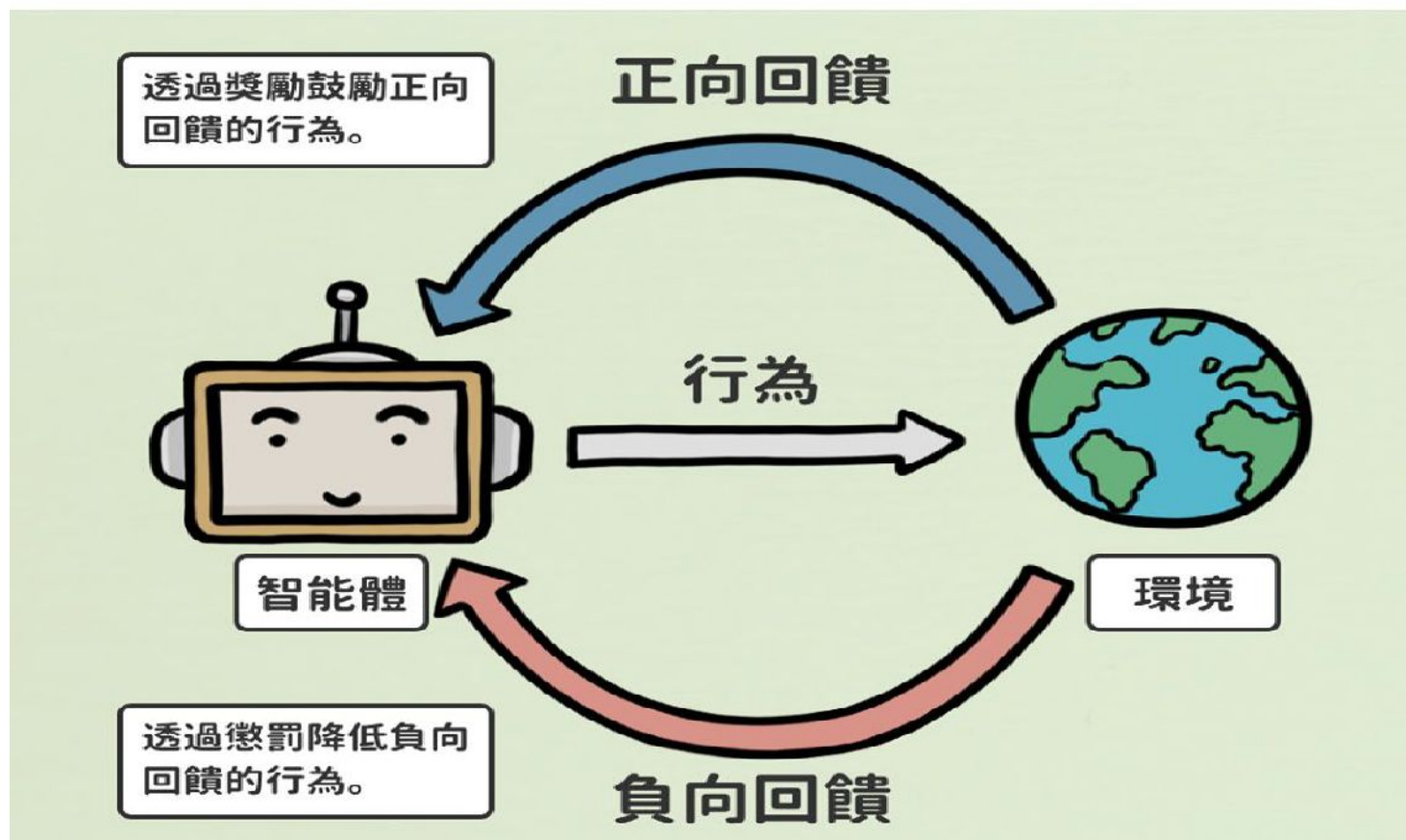    - Q-learning 是強化學習的一種方法，主要是透過記錄學習過的策略，來告訴 Agent，什麼情況下要對應採取什麼 Action 會得到最大的獎勵 Reward。
    - Q-learning 最基本的應用方式，就是將對應行動的獎勵值存在一個 Q-table 中。

2.  **Q-table :**

    - 簡單來說就是一個查詢表，用來計算某狀態下做某行為後未來可以期望得到最大的 Reward 為多少，這個表可以引導我們選出每個狀態 state 下，最好的行為 Action。

- $\alpha$ is the learning rate $(0 < \alpha < 1)$. $\gamma$ is the discount factor $(0 < \gamma < 1)$.

- When the value of $\gamma$ is larger, the long-term rewards obtained in the future will receive more attention; the smaller the value, the more current rewards will be considered.

- $R$ is the reward; it acts on each state and will receive the corresponding reward.

$$Q(s,a)^{new} = (1-\alpha) * Q(s,a)^{old} + \alpha[r + \gamma * max_{a'}Q(s',a')]$$

# Q-learning 檔案

| 檔案名稱 | 類型 |
|---|---|
| discount_factor.csv | Excel 檔案 |
| Q-learning.py | Python 檔案 |
| q_table.npy | Numpy 檔案 |
| numpy_to_txt_grid.py | Python 檔案 |
| q_table.txt | 文字檔案 |

```
 8        # set the rows and columns length
 9        BOARD_ROWS = 4
10        BOARD_COLS = 6
11
12        # initalise start, win and lose states
13        START = (0, 0)
14        WIN_STATE = (3, 5)
```

Column

|       | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|---|---|---|---|---|---|
| 0 | (0, 0) | (0, 1) | (0, 2) | (0, 3) | (0, 4) | (0, 5) |
| 1 | (1, 0) | (1, 1) | (1, 2) | (1, 3) | (1, 4) | (1, 5) |
| 2 | (2, 0) | (2, 1) | (2, 2) | (2, 3) | (2, 4) | (2, 5) |
| 3 | (3, 0) | (3, 1) | (3, 2) | (3, 3) | (3, 4) | (3, 5) |

START

Row

WIN_STATE

```python
19  ∨  class State:
20         def __init__(self, state=START):
21             self.state = state
22             self.isEnd = False
23
24  ∨      def getReward(self):
25             if self.state == WIN_STATE:
26                 return 100
27             else:
28                 return 0
29
30         def isEndFunc(self):
31             if self.state == WIN_STATE:
32                 self.isEnd = True
33
34  ∨      def nxtPosition(self, action):
35             if action == 0:
36                 nxtState = (self.state[0] - 1, self.state[1])        # up
37             elif action == 1:
38                 nxtState = (self.state[0] + 1, self.state[1])        # down
39             elif action == 2:
40                 nxtState = (self.state[0], self.state[1] - 1)        # left
41             elif action == 3:
42                 nxtState = (self.state[0], self.state[1] + 1)        # right
43             elif action == 4:
44                 nxtState = (self.state[0] - 1, self.state[1] - 1)   # up-left
45             elif action == 5:
46                 nxtState = (self.state[0] - 1, self.state[1] + 1)   # up-right
47             elif action == 6:
48                 nxtState = (self.state[0] + 1, self.state[1] - 1)   # down-left
49             elif action == 7:
50                 nxtState = (self.state[0] + 1, self.state[1] + 1)   # down-right
51
52             if (nxtState[0] >= 0) and (nxtState[0] < BOARD_ROWS) and (nxtState[1] >= 0) and (nxtState[1] < BOARD_COLS):
53                 return nxtState
54
55             return self.state
```

```python
58    # class agent to implement reinforcement learning through grid
59    class Agent:
60        def __init__(self):
61            # inialise states and actions
62            self.states = []
63            self.actions = [0, 1, 2, 3, 4, 5, 6, 7]  # up, down, left, right, up-left, up-right, down-left, down-right
64            self.State = State()
65            self.alpha = 0.8
66            self.epsilon = 0.5
67            self.isEnd = self.State.isEnd
68
69            # array to retain reward values for plot
70            self.plot_reward = []
71
72            # initalise Q values as a dictionary for current and new
73            self.Q = {}
74            self.new_Q = {}
75            self.rewards = 0
76
77            # initalise all Q values across the board to 0, print these values
78            for i in range(BOARD_ROWS):
79                for j in range(BOARD_COLS):
80                    for k in range(len(self.actions)):
81                        self.Q[(i, j, k)] = 0
82                        self.new_Q[(i, j, k)] = 0
83            # print(self.Q)
```

```
+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+
|            (0,0)           |            (0,1)           |            (0,2)           |            (0,3)           |            (0,4)           |            (0,5)           |
|↑ Q(0,0,0) = 0              |↑ Q(0,1,0) = 0              |↑ Q(0,2,0) = 0              |↑ Q(0,3,0) = 0              |↑ Q(0,4,0) = 0              |↑ Q(0,5,0) = 0              |
|↓ Q(0,0,1) = 0              |↓ Q(0,1,1) = 0              |↓ Q(0,2,1) = 0              |↓ Q(0,3,1) = 0              |↓ Q(0,4,1) = 0              |↓ Q(0,5,1) = 0              |
|← Q(0,0,2) = 0              |← Q(0,1,2) = 0              |← Q(0,2,2) = 0              |← Q(0,3,2) = 0              |← Q(0,4,2) = 0              |← Q(0,5,2) = 0              |
|→ Q(0,0,3) = 0              |→ Q(0,1,3) = 0              |→ Q(0,2,3) = 0              |→ Q(0,3,3) = 0              |→ Q(0,4,3) = 0              |→ Q(0,5,3) = 0              |
|↖ Q(0,0,4) = 0              |↖ Q(0,1,4) = 0              |↖ Q(0,2,4) = 0              |↖ Q(0,3,4) = 0              |↖ Q(0,4,4) = 0              |↖ Q(0,5,4) = 0              |
|↗ Q(0,0,5) = 0              |↗ Q(0,1,5) = 0              |↗ Q(0,2,5) = 0              |↗ Q(0,3,5) = 0              |↗ Q(0,4,5) = 0              |↗ Q(0,5,5) = 0              |
|↙ Q(0,0,6) = 0              |↙ Q(0,1,6) = 0              |↙ Q(0,2,6) = 0              |↙ Q(0,3,6) = 0              |↙ Q(0,4,6) = 0              |↙ Q(0,5,6) = 0              |
|↘ Q(0,0,7) = 0              |↘ Q(0,1,7) = 0              |↘ Q(0,2,7) = 0              |↘ Q(0,3,7) = 0              |↘ Q(0,4,7) = 0              |↘ Q(0,5,7) = 0              |
+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+
|            (1,0)           |            (1,1)           |            (1,2)           |            (1,3)           |            (1,4)           |            (1,5)           |
|↑ Q(1,0,0) = 0              |↑ Q(1,1,0) = 0              |↑ Q(1,2,0) = 0              |↑ Q(1,3,0) = 0              |↑ Q(1,4,0) = 0              |↑ Q(1,5,0) = 0              |
|↓ Q(1,0,1) = 0              |↓ Q(1,1,1) = 0              |↓ Q(1,2,1) = 0              |↓ Q(1,3,1) = 0              |↓ Q(1,4,1) = 0              |↓ Q(1,5,1) = 0              |
|← Q(1,0,2) = 0              |← Q(1,1,2) = 0              |← Q(1,2,2) = 0              |← Q(1,3,2) = 0              |← Q(1,4,2) = 0              |← Q(1,5,2) = 0              |
|→ Q(1,0,3) = 0              |→ Q(1,1,3) = 0              |→ Q(1,2,3) = 0              |→ Q(1,3,3) = 0              |→ Q(1,4,3) = 0              |→ Q(1,5,3) = 0              |
|↖ Q(1,0,4) = 0              |↖ Q(1,1,4) = 0              |↖ Q(1,2,4) = 0              |↖ Q(1,3,4) = 0              |↖ Q(1,4,4) = 0              |↖ Q(1,5,4) = 0              |
|↗ Q(1,0,5) = 0              |↗ Q(1,1,5) = 0              |↗ Q(1,2,5) = 0              |↗ Q(1,3,5) = 0              |↗ Q(1,4,5) = 0              |↗ Q(1,5,5) = 0              |
|↙ Q(1,0,6) = 0              |↙ Q(1,1,6) = 0              |↙ Q(1,2,6) = 0              |↙ Q(1,3,6) = 0              |↙ Q(1,4,6) = 0              |↙ Q(1,5,6) = 0              |
|↘ Q(1,0,7) = 0              |↘ Q(1,1,7) = 0              |↘ Q(1,2,7) = 0              |↘ Q(1,3,7) = 0              |↘ Q(1,4,7) = 0              |↘ Q(1,5,7) = 0              |
+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+
|            (2,0)           |            (2,1)           |            (2,2)           |            (2,3)           |            (2,4)           |            (2,5)           |
|↑ Q(2,0,0) = 0              |↑ Q(2,1,0) = 0              |↑ Q(2,2,0) = 0              |↑ Q(2,3,0) = 0              |↑ Q(2,4,0) = 0              |↑ Q(2,5,0) = 0              |
|↓ Q(2,0,1) = 0              |↓ Q(2,1,1) = 0              |↓ Q(2,2,1) = 0              |↓ Q(2,3,1) = 0              |↓ Q(2,4,1) = 0              |↓ Q(2,5,1) = 0              |
|← Q(2,0,2) = 0              |← Q(2,1,2) = 0              |← Q(2,2,2) = 0              |← Q(2,3,2) = 0              |← Q(2,4,2) = 0              |← Q(2,5,2) = 0              |
|→ Q(2,0,3) = 0              |→ Q(2,1,3) = 0              |→ Q(2,2,3) = 0              |→ Q(2,3,3) = 0              |→ Q(2,4,3) = 0              |→ Q(2,5,3) = 0              |
|↖ Q(2,0,4) = 0              |↖ Q(2,1,4) = 0              |↖ Q(2,2,4) = 0              |↖ Q(2,3,4) = 0              |↖ Q(2,4,4) = 0              |↖ Q(2,5,4) = 0              |
|↗ Q(2,0,5) = 0              |↗ Q(2,1,5) = 0              |↗ Q(2,2,5) = 0              |↗ Q(2,3,5) = 0              |↗ Q(2,4,5) = 0              |↗ Q(2,5,5) = 0              |
|↙ Q(2,0,6) = 0              |↙ Q(2,1,6) = 0              |↙ Q(2,2,6) = 0              |↙ Q(2,3,6) = 0              |↙ Q(2,4,6) = 0              |↙ Q(2,5,6) = 0              |
|↘ Q(2,0,7) = 0              |↘ Q(2,1,7) = 0              |↘ Q(2,2,7) = 0              |↘ Q(2,3,7) = 0              |↘ Q(2,4,7) = 0              |↘ Q(2,5,7) = 0              |
+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+
|            (3,0)           |            (3,1)           |            (3,2)           |            (3,3)           |            (3,4)           |            (3,5)           |
|↑ Q(3,0,0) = 0              |↑ Q(3,1,0) = 0              |↑ Q(3,2,0) = 0              |↑ Q(3,3,0) = 0              |↑ Q(3,4,0) = 0              |↑ Q(3,5,0) = 0              |
|↓ Q(3,0,1) = 0              |↓ Q(3,1,1) = 0              |↓ Q(3,2,1) = 0              |↓ Q(3,3,1) = 0              |↓ Q(3,4,1) = 0              |↓ Q(3,5,1) = 0              |
|← Q(3,0,2) = 0              |← Q(3,1,2) = 0              |← Q(3,2,2) = 0              |← Q(3,3,2) = 0              |← Q(3,4,2) = 0              |← Q(3,5,2) = 0              |
|→ Q(3,0,3) = 0              |→ Q(3,1,3) = 0              |→ Q(3,2,3) = 0              |→ Q(3,3,3) = 0              |→ Q(3,4,3) = 0              |→ Q(3,5,3) = 0              |
|↖ Q(3,0,4) = 0              |↖ Q(3,1,4) = 0              |↖ Q(3,2,4) = 0              |↖ Q(3,3,4) = 0              |↖ Q(3,4,4) = 0              |↖ Q(3,5,4) = 0              |
|↗ Q(3,0,5) = 0              |↗ Q(3,1,5) = 0              |↗ Q(3,2,5) = 0              |↗ Q(3,3,5) = 0              |↗ Q(3,4,5) = 0              |↗ Q(3,5,5) = 0              |
|↙ Q(3,0,6) = 0              |↙ Q(3,1,6) = 0              |↙ Q(3,2,6) = 0              |↙ Q(3,3,6) = 0              |↙ Q(3,4,6) = 0              |↙ Q(3,5,6) = 0              |
|↘ Q(3,0,7) = 0              |↘ Q(3,1,7) = 0              |↘ Q(3,2,7) = 0              |↘ Q(3,3,7) = 0              |↘ Q(3,4,7) = 0              |↘ Q(3,5,7) = 0              |
+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+----------------------------+
```

```python
79 ∨      def Action(self):
80            # random value vs epsilon
81            rnd = random.random()
82            # set arbitraty low value to compare with Q values to find max
83            mx_nxt_reward = -10
84            action = None
85
86            # find max Q value over actions
87            if rnd > self.epsilon:
88                # iterate through actions, find Q-value and choose best
89                for k in self.actions:
90                    i, j = self.State.state
91                    nxt_reward = self.Q[(i, j, k)]
92
93                    if nxt_reward >= mx_nxt_reward:
94                        action = k
95                        mx_nxt_reward = nxt_reward
96            # else choose random action
97            else:
98                action = np.random.choice(self.actions)
99
100           # select the next state based on action chosen
101           position = self.State.nxtPosition(action)
102           return position, action
```

```
104            # Q-learning Algorithm
105    ∨       def Q_Learning(self, episodes):
106                df_factor_path = "discount_factor.csv"
107                df_factor = pd.read_csv(df_factor_path, index_col="Grid 座標")
108                x = 0
109                # iterate through best path for each episode
110                while x < episodes:
111                    # check if state is end
112                    if self.isEnd:
113                        # get current rewrard and add to array for plot
114                        reward = self.State.getReward()
115                        self.rewards += reward
116                        self.plot_reward.append(self.rewards)
117
118                        # get state, assign reward to each Q_value in state
119                        i, j = self.State.state
120                        for a in self.actions:
121                            self.new_Q[(i, j, a)] = round(reward, 3)
122
123                        # reset state
124                        self.State = State()
125                        self.isEnd = self.State.isEnd
126
127                        # set rewards to zero and iterate to next episode
128                        self.rewards = 0
129                        x += 1
```

# discount_factor.csv 檔案說明

| Grid 座標 | discount factor |
|---|---|
| 0_0 | 0.3/0.9/0.3/0.69449/0.3/0.3/0.3/0.3 |
| 0_1 | 0.3/0.43638/0.68073/0.55197/0.3/0.55602/0.3/0.32028 |
| 0_2 | 0.3/0.53606/0.55858/0.5093/0.3/0.3/0.3/0.3 |
| 0_3 | 0.3/0.3/0.72106/0.60534/0.3/0.32053/0.3/0.36304 |
| 0_4 | 0.3/0.54984/0.50086/0.61222/0.3/0.3/0.3/0.3 |
| 0_5 | 0.3/0.34667/0.9/0.3/0.3/0.36411/0.3/0.3 |
| 1_0 | 0.41314/0.9/0.3/0.46518/0.3/0.3/0.3/0.39871 |
| 1_1 | 0.40229/0.59559/0.72101/0.58091/0.3/0.53315/0.30577/0.32393 |
| 1_2 | 0.44158/0.51077/0.50963/0.44838/0.3/0.37714/0.3/0.3 |

$$Q(s,a)^{new} = (1 - \alpha) * Q(s,a)^{old} + \alpha[r + \gamma * max_{a'}Q(s',a')]$$

```python
130                else:
131                    mx_nxt_value = -10
132                    next_state, action = self.Action()
133                    i, j = self.State.state
134                    reward = self.State.getReward()
135                    # add reward to rewards for plot
136                    self.rewards += reward
137
138                    # iterate through actions to find max Q value for action based on next state action
139                    for a in self.actions:
140                        now_index = str(i) + "_" + str(j)
141                        df = df_factor.loc[now_index, "discount factor"]
142                        df = df.split("/")
143
144                        nxtStateAction = (next_state[0], next_state[1], a)
145                        q_value = (1 - self.alpha) * self.Q[(i, j, action)] + self.alpha * (
146                            reward + float(df[a]) * self.Q[nxtStateAction]
147                        )
148
149                        if q_value >= mx_nxt_value:
150                            mx_nxt_value = q_value
151
152                    # next state is now current state, check if end state
153                    self.State = State(state=next_state)
154                    self.State.isEndFunc()
155                    self.isEnd = self.State.isEnd
156                    self.new_Q[(i, j, action)] = round(mx_nxt_value, 3)
157                self.Q = self.new_Q.copy()
158            # print(self.Q)
```

16

```
                (0,0)          |              (0,1)            |              (0,2)            |              (0,3)            |              (0,4)            |              (0,5)
↑ Q(0, 0, 0) = 2.286           |↑ Q(0, 1, 0) = 1.785           |↑ Q(0, 2, 0) = 3.030           |↑ Q(0, 3, 0) = 6.154           |↑ Q(0, 4, 0) = 7.931           |↑ Q(0, 5, 0) = 6.195
↓ Q(0, 0, 1) = 2.540           |↓ Q(0, 1, 1) = 3.086           |↓ Q(0, 2, 1) = 4.902           |↓ Q(0, 3, 1) = 10.685          |↓ Q(0, 4, 1) = 14.424          |↓ Q(0, 5, 1) = 14.620
← Q(0, 0, 2) = 2.286           |← Q(0, 1, 2) = 1.556           |← Q(0, 2, 2) = 1.654           |← Q(0, 3, 2) = 3.667           |← Q(0, 4, 2) = 5.875           |← Q(0, 5, 2) = 6.957
→ Q(0, 0, 3) = 2.777           |→ Q(0, 1, 3) = 3.235           |→ Q(0, 2, 3) = 5.728           |→ Q(0, 3, 3) = 8.418           |→ Q(0, 4, 3) = 8.039           |→ Q(0, 5, 3) = 6.261
↖ Q(0, 0, 4) = 2.286           |↖ Q(0, 1, 4) = 1.785           |↖ Q(0, 2, 4) = 3.030           |↖ Q(0, 3, 4) = 6.154           |↖ Q(0, 4, 4) = 7.931           |↖ Q(0, 5, 4) = 6.261
↗ Q(0, 0, 5) = 2.286           |↗ Q(0, 1, 5) = 1.785           |↗ Q(0, 2, 5) = 3.030           |↗ Q(0, 3, 5) = 6.154           |↗ Q(0, 4, 5) = 7.931           |↗ Q(0, 5, 5) = 6.261
↙ Q(0, 0, 6) = 2.286           |↙ Q(0, 1, 6) = 2.018           |↙ Q(0, 2, 6) = 2.847           |↙ Q(0, 3, 6) = 5.826           |↙ Q(0, 4, 6) = 10.807          |↙ Q(0, 5, 6) = 14.008
↘ Q(0, 0, 7) = 3.883           |↘ Q(0, 1, 7) = 5.313           |↘ Q(0, 2, 7) = 10.102          |↘ Q(0, 3, 7) = 16.952          |↘ Q(0, 4, 7) = 23.189          |↘ Q(0, 5, 7) = 6.261

                (1,0)          |              (1,1)            |              (1,2)            |              (1,3)            |              (1,4)            |              (1,5)
↑ Q(1, 0, 0) = 2.286           |↑ Q(1, 1, 0) = 1.879           |↑ Q(1, 2, 0) = 3.030           |↑ Q(1, 3, 0) = 6.408           |↑ Q(1, 4, 0) = 8.117           |↑ Q(1, 5, 0) = 7.377
↓ Q(1, 0, 1) = 2.330           |↓ Q(1, 1, 1) = 4.064           |↓ Q(1, 2, 1) = 6.705           |↓ Q(1, 3, 1) = 18.845          |↓ Q(1, 4, 1) = 23.319          |↓ Q(1, 5, 1) = 42.174
← Q(1, 0, 2) = 2.097           |← Q(1, 1, 2) = 2.124           |← Q(1, 2, 2) = 2.606           |← Q(1, 3, 2) = 7.096           |← Q(1, 4, 2) = 10.528          |← Q(1, 5, 2) = 14.008
→ Q(1, 0, 3) = 3.657           |→ Q(1, 1, 3) = 5.591           |→ Q(1, 2, 3) = 9.625           |→ Q(1, 3, 3) = 17.652          |→ Q(1, 4, 3) = 23.561          |→ Q(1, 5, 3) = 21.280
↖ Q(1, 0, 4) = 2.097           |↖ Q(1, 1, 4) = 1.648           |↖ Q(1, 2, 4) = 1.594           |↖ Q(1, 3, 4) = 4.223           |↖ Q(1, 4, 4) = 5.969           |↖ Q(1, 5, 4) = 7.278
↗ Q(1, 0, 5) = 2.777           |↗ Q(1, 1, 5) = 3.327           |↗ Q(1, 2, 5) = 5.457           |↗ Q(1, 3, 5) = 8.766           |↗ Q(1, 4, 5) = 8.167           |↗ Q(1, 5, 5) = 21.280
↙ Q(1, 0, 6) = 2.097           |↙ Q(1, 1, 6) = 2.910           |↙ Q(1, 2, 6) = 3.137           |↙ Q(1, 3, 6) = 11.025          |↙ Q(1, 4, 6) = 17.450          |↙ Q(1, 5, 6) = 21.062
↘ Q(1, 0, 7) = 3.351           |↘ Q(1, 1, 7) = 8.687           |↘ Q(1, 2, 7) = 14.955          |↘ Q(1, 3, 7) = 27.513          |↘ Q(1, 4, 7) = 46.695          |↘ Q(1, 5, 7) = 21.280

                (2,0)          |              (2,1)            |              (2,2)            |              (2,3)            |              (2,4)            |              (2,5)
↑ Q(2, 0, 0) = 2.618           |↑ Q(2, 1, 0) = 2.615           |↑ Q(2, 2, 0) = 4.616           |↑ Q(2, 3, 0) = 11.820          |↑ Q(2, 4, 0) = 17.009          |↑ Q(2, 5, 0) = 20.916
↓ Q(2, 0, 1) = 2.539           |↓ Q(2, 1, 1) = 3.262           |↓ Q(2, 2, 1) = 8.618           |↓ Q(2, 3, 1) = 26.452          |↓ Q(2, 4, 1) = 41.742          |↓ Q(2, 5, 1) = 83.584
← Q(2, 0, 2) = 3.586           |← Q(2, 1, 2) = 3.227           |← Q(2, 2, 2) = 3.355           |← Q(2, 3, 2) = 10.014          |← Q(2, 4, 2) = 18.159          |← Q(2, 5, 2) = 20.701
→ Q(2, 0, 3) = 5.009           |→ Q(2, 1, 3) = 6.996           |→ Q(2, 2, 3) = 14.955          |→ Q(2, 3, 3) = 24.989          |→ Q(2, 4, 3) = 37.320          |→ Q(2, 5, 3) = 41.453
↖ Q(2, 0, 4) = 3.586           |↖ Q(2, 1, 4) = 1.887           |↖ Q(2, 2, 4) = 2.681           |↖ Q(2, 3, 4) = 6.445           |↖ Q(2, 4, 4) = 10.022          |↖ Q(2, 5, 4) = 14.008
↗ Q(2, 0, 5) = 4.003           |↗ Q(2, 1, 5) = 4.502           |↗ Q(2, 2, 5) = 8.466           |↗ Q(2, 3, 5) = 15.776          |↗ Q(2, 4, 5) = 18.831          |↗ Q(2, 5, 5) = 41.453
↙ Q(2, 0, 6) = 3.586           |↙ Q(2, 1, 6) = 2.063           |↙ Q(2, 2, 6) = 3.345           |↙ Q(2, 3, 6) = 12.031          |↙ Q(2, 4, 6) = 22.149          |↙ Q(2, 5, 6) = 27.489
↘ Q(2, 0, 7) = 4.993           |↘ Q(2, 1, 7) = 8.405           |↘ Q(2, 2, 7) = 18.946          |↘ Q(2, 3, 7) = 49.852          |↘ Q(2, 4, 7) = 56.315          |↘ Q(2, 5, 7) = 41.453

                (3,0)          |              (3,1)            |              (3,2)            |              (3,3)            |              (3,4)            |              (3,5)
↑ Q(3, 0, 0) = 2.547           |↑ Q(3, 1, 0) = 2.715           |↑ Q(3, 2, 0) = 6.802           |↑ Q(3, 3, 0) = 14.955          |↑ Q(3, 4, 0) = 27.785          |↑ Q(3, 5, 0) = 100.000
↓ Q(3, 0, 1) = 2.292           |↓ Q(3, 1, 1) = 2.707           |↓ Q(3, 2, 1) = 8.173           |↓ Q(3, 3, 1) = 20.961          |↓ Q(3, 4, 1) = 55.428          |↓ Q(3, 5, 1) = 100.000
← Q(3, 0, 2) = 2.292           |← Q(3, 1, 2) = 1.585           |← Q(3, 2, 2) = 3.172           |← Q(3, 3, 2) = 9.534           |← Q(3, 4, 2) = 29.411          |← Q(3, 5, 2) = 100.000
→ Q(3, 0, 3) = 3.547           |→ Q(3, 1, 3) = 6.974           |→ Q(3, 2, 3) = 17.968          |→ Q(3, 3, 3) = 39.504          |→ Q(3, 4, 3) = 74.450          |→ Q(3, 5, 3) = 100.000
↖ Q(3, 0, 4) = 2.292           |↖ Q(3, 1, 4) = 1.959           |↖ Q(3, 2, 4) = 3.182           |↖ Q(3, 3, 4) = 7.935           |↖ Q(3, 4, 4) = 18.604          |↖ Q(3, 5, 4) = 100.000
↗ Q(3, 0, 5) = 3.558           |↗ Q(3, 1, 5) = 5.805           |↗ Q(3, 2, 5) = 14.955          |↗ Q(3, 3, 5) = 19.802          |↗ Q(3, 4, 5) = 30.862          |↗ Q(3, 5, 5) = 100.000
↙ Q(3, 0, 6) = 2.292           |↙ Q(3, 1, 6) = 2.707           |↙ Q(3, 2, 6) = 8.173           |↙ Q(3, 3, 6) = 20.961          |↙ Q(3, 4, 6) = 55.428          |↙ Q(3, 5, 6) = 100.000
↘ Q(3, 0, 7) = 2.292           |↘ Q(3, 1, 7) = 2.707           |↘ Q(3, 2, 7) = 8.173           |↘ Q(3, 3, 7) = 20.961          |↘ Q(3, 4, 7) = 55.428          |↘ Q(3, 5, 7) = 100.000
```

```python
164 ∨        def showValues(self, arr):
165            outArr = [[0 for _ in range(0, BOARD_COLS)] for _ in range(0, BOARD_ROWS)]
166            for i in range(0, BOARD_ROWS):
167                print("--------------------------------------------------------------------------------")
168                out = "| "
169                for j in range(0, BOARD_COLS):
170                    mx_nxt_value = -10
171                    for a in self.actions:
172                        nxt_value = self.Q[(i, j, a)]
173                        if nxt_value >= mx_nxt_value:
174                            mx_nxt_value = nxt_value
175                    out += str(mx_nxt_value).ljust(6) + " | "
176                    outArr[i][j] = str(mx_nxt_value)
177                print(out)
178            print("--------------------------------------------------------------------------------")
```

| 3.883 | 5.313 | 10.102 | 16.952 | 23.189 | 14.62 |
| 3.657 | 8.687 | 14.955 | 27.513 | 46.695 | 42.174 |
| 5.009 | 8.405 | 18.946 | 49.852 | 56.315 | 83.584 |
| 3.558 | 6.974 | 17.968 | 39.504 | 74.45 | 100 |

`[-1.        3.657 -1.        5.313 -1.        -1.        -1.        8.687]`

```python
190         # up
191         if START[0] - 1 < 0:
192             arr[0] = -1
193         else:
194             arr[0] = outArr[START[0] - 1][START[1]]
195
196         # down
197         if START[0] + 1 >= BOARD_ROWS:
198             arr[1] = -1
199         else:
200             arr[1] = outArr[START[0] + 1][START[1]]
201
202         # left
203         if START[1] - 1 < 0:
204             arr[2] = -1
205         else:
206             arr[2] = outArr[START[0]][START[1] - 1]
207
208         # right
209         if START[1] + 1 >= BOARD_COLS:
210             arr[3] = -1
211         else:
212             arr[3] = outArr[START[0]][START[1] + 1]
```

```python
214         # up-left
215         if START[0] - 1 < 0 or START[1] - 1 < 0:
216             arr[4] = -1
217         else:
218             arr[4] = outArr[START[0] - 1][START[1] - 1]
219
220         # up-right
221         if START[0] - 1 < 0 or START[1] + 1 >= BOARD_COLS:
222             arr[5] = -1
223         else:
224             arr[5] = outArr[START[0] - 1][START[1] + 1]
225
226         # down-left
227         if START[0] + 1 >= BOARD_ROWS or START[1] - 1 < 0:
228             arr[6] = -1
229         else:
230             arr[6] = outArr[START[0] + 1][START[1] - 1]
231
232         # down-right
233         if START[0] + 1 >= BOARD_ROWS or START[1] + 1 >= BOARD_COLS:
234             arr[7] = -1
235         else:
236             arr[7] = outArr[START[0] + 1][START[1] + 1]
237
238         print(arr)
```

```python
202    if __name__ == "__main__":
203        # create agent for 15,000 episodes implementing a Q-learning algorithm plot and show values.
204        episodes = 10000
205        ag = Agent()
206
207        filename = "q_table"
208        q_value = np.zeros((8), dtype=np.float64)
209
210        print(START, WIN_STATE)
211
212        ag.Q_Learning(episodes)
213        ag.showValues(q_value)
214        np.save(os.path.join(filename), q_value)
```

```python
1   import numpy as np
2   import os
3
4   fileName = 'q_table.txt'
5
6   try:
7       os.remove(fileName)
8   except:
9       pass
10
11  test = np.load('q_table.npy')
12  fo = open("q_table.txt", "a+")
13
14  string = " ".join(map(str, test))
15  fo.write(string + "\n")
16
17  fo.close()
```



q_table.npy    q_table.txt

# 作業繳交 說明

# 作業說明

✔ **作業要求：**

1.  將 row 改為 8、column 改為 10

2.  改變起點、終點

3.  嘗試產生自己的 discount factor 資料 (隨機生成八個方位的 discount factor, 範圍介於 0 到 1 區間, 數值取至小數點後第三位。)

4.  紀錄最佳路徑策略圖

5.  紀錄完整的 Q-table (不單單紀錄起點的 Q-value)

6.  紀錄 ε (epsilon) 變化 (x 軸：ε 參數變化, y 軸：執行時間)

7.  檔名內容：學號_HW6.zip  (包含  Q-learning.py、discount_factor.csv、numpy_to_txt_grid.py、q_table.npy、q_table.txt、學號_HW6.pdf) □ 補交檔名：學號_HW6-2.zip

繳交檔案：1. 紙本檔案、2. 電子檔案上傳 FTP、3. 檔名為：學號_HW6.zip

上傳：120.107.172.19  使用者名稱：1132VANET  密碼：1132student

```
1   0_0 -1.0 15.041 -1.0 20.294 -1.0 -1.0 -1.0 20.857
2   0_1 -1.0 20.857 15.629 17.87 -1.0 -1.0 15.041 21.384
3   0_2 -1.0 21.384 20.294 16.429 -1.0 -1.0 20.857 28.823
4   0_3 -1.0 28.823 17.87 15.479 -1.0 -1.0 21.384 24.929
5   0_4 -1.0 24.929 16.429 26.493 -1.0 -1.0 28.823 24.116
6   0_5 -1.0 24.116 15.479 18.48 -1.0 -1.0 24.929 30.278
7   0_6 -1.0 30.278 26.493 18.366 -1.0 -1.0 24.116 38.868
8   0_7 -1.0 38.868 18.48 25.769 -1.0 -1.0 30.278 34.67
9   0_8 -1.0 34.67 18.366 27.285 -1.0 -1.0 38.868 25.25
10  0_9 -1.0 25.25 25.769 -1.0 -1.0 -1.0 34.67 -1.0
11  1_0 15.629 20.946 -1.0 20.857 -1.0 20.294 -1.0 25.712
12  1_1 20.294 25.712 15.041 21.384 15.629 17.87 20.946 27.151
13  1_2 17.87 27.151 20.857 28.823 20.294 16.429 25.712 32.608
14  1_3 16.429 32.608 21.384 24.929 17.87 15.479 27.151 36.393
15  1_4 15.479 36.393 28.823 24.116 16.429 26.493 32.608 36.546
16  1_5 26.493 36.546 24.929 30.278 15.479 18.48 36.393 39.661
17  1_6 18.48 39.661 24.116 38.868 26.493 18.366 36.546 30.996
18  1_7 18.366 30.996 30.278 34.67 18.48 25.769 39.661 31.238
19  1_8 25.769 31.238 38.868 25.25 18.366 27.285 30.996 29.889
20  1_9 27.285 29.889 34.67 -1.0 25.769 -1.0 31.238 -1.0
21  2_0 15.041 23.164 -1.0 25.712 -1.0 20.857 -1.0 23.331
22  2_1 20.857 23.331 20.946 27.151 15.041 21.384 23.164 16.95
23  2_2 21.384 16.95 25.712 32.608 20.857 28.823 23.331 22.081
24  2_3 28.823 22.081 27.151 36.393 21.384 24.929 16.95 26.916
25  2_4 24.929 26.916 32.608 36.546 28.823 24.116 22.081 43.171
26  2_5 24.116 43.171 36.393 39.661 24.929 30.278 26.916 45.172
27  2_6 30.278 45.172 36.546 30.996 24.116 38.868 43.171 37.925
28  2_7 38.868 37.925 39.661 31.238 30.278 34.67 45.172 39.79
29  2_8 34.67 39.79 30.996 29.889 38.868 25.25 37.925 29.869
30  2_9 25.25 29.869 31.238 -1.0 34.67 -1.0 39.79 -1.0
31  3_0 20.946 30.545 -1.0 23.331 -1.0 25.712 -1.0 31.95
32  3_1 25.712 31.95 23.164 16.95 20.946 27.151 30.545 34.775
33  3_2 27.151 34.775 23.331 22.081 25.712 32.608 31.95 39.862
34  3_3 32.608 39.862 16.95 26.916 27.151 36.393 34.775 34.346
35  3_4 36.393 34.346 22.081 43.171 32.608 36.546 39.862 45.491
36  3_5 36.546 45.491 26.916 45.172 36.393 39.661 34.346 32.237
37  3_6 39.661 32.237 43.171 37.925 36.546 30.996 45.491 37.678
38  3_7 30.996 37.678 45.172 39.79 39.661 31.238 32.237 53.734
39  3_8 31.238 53.734 37.925 29.869 30.996 29.889 37.678 55.279
40  3_9 29.889 55.279 39.79 -1.0 31.238 -1.0 53.734 -1.0
```

```
41  4_0 23.164 30.817 -1.0 31.95 -1.0 23.331 -1.0 35.053       61  6_0 30.817 33.594 -1.0 41.777 -1.0 35.053 -1.0 31.049
42  4_1 23.331 35.053 30.545 34.775 23.164 16.95 30.817 31.502  62  6_1 35.053 31.049 38.184 44.825 30.817 31.502 33.594 32.842
43  4_2 16.95 31.502 31.95 39.862 23.331 22.081 35.053 38.812   63  6_2 31.502 32.842 41.777 44.87 35.053 38.812 31.049 35.684
44  4_3 22.081 38.812 34.775 34.346 16.95 26.916 31.502 40.428  64  6_3 38.812 35.684 44.825 41.492 31.502 40.428 32.842 55.102
45  4_4 26.916 40.428 39.862 45.491 22.081 43.171 38.812 57.657 65  6_4 40.428 55.102 44.87 56.42 38.812 57.657 35.684 63.918
46  4_5 43.171 57.657 34.346 32.237 26.916 45.172 40.428 39.246 66  6_5 57.657 63.918 41.492 72.8 40.428 39.246 55.102 62.531
47  4_6 45.172 39.246 45.491 37.678 43.171 37.925 57.657 49.288 67  6_6 39.246 62.531 56.42 79.429 57.657 49.288 63.918 78.873
48  4_7 37.925 49.288 32.237 53.734 45.172 39.79 39.246 55.866  68  6_7 49.288 78.873 72.8 70.8 39.246 55.866 62.531 92.9
49  4_8 39.79 55.866 37.678 55.279 37.925 29.869 49.288 59.44   69  6_8 55.866 92.9 79.429 77.7 49.288 59.44 78.873 100.0
50  4_9 29.869 59.44 53.734 -1.0 39.79 -1.0 55.866 -1.0         70  6_9 59.44 100.0 70.8 -1.0 55.866 -1.0 92.9 -1.0
51  5_0 30.545 38.184 -1.0 35.053 -1.0 31.95 -1.0 41.777        71  7_0 38.184 -1.0 -1.0 31.049 -1.0 41.777 -1.0 -1.0
52  5_1 31.95 41.777 30.817 31.502 30.545 34.775 38.184 44.825  72  7_1 41.777 -1.0 33.594 32.842 38.184 44.825 -1.0 -1.0
53  5_2 34.775 44.825 35.053 38.812 31.95 39.862 41.777 44.87   73  7_2 44.825 -1.0 31.049 35.684 41.777 44.87 -1.0 -1.0
54  5_3 39.862 44.87 31.502 40.428 34.775 34.346 44.825 41.492  74  7_3 44.87 -1.0 32.842 55.102 44.825 41.492 -1.0 -1.0
55  5_4 34.346 41.492 38.812 57.657 39.862 45.491 44.87 56.42   75  7_4 41.492 -1.0 35.684 63.918 44.87 56.42 -1.0 -1.0
56  5_5 45.491 56.42 40.428 39.246 34.346 32.237 41.492 72.8    76  7_5 56.42 -1.0 55.102 62.531 41.492 72.8 -1.0 -1.0
57  5_6 32.237 72.8 57.657 49.288 45.491 37.678 56.42 79.429    77  7_6 72.8 -1.0 63.918 78.873 56.42 79.429 -1.0 -1.0
58  5_7 37.678 79.429 39.246 55.866 32.237 53.734 72.8 70.8     78  7_7 79.429 -1.0 62.531 92.9 72.8 70.8 -1.0 -1.0
59  5_8 53.734 70.8 49.288 59.44 37.678 55.279 79.429 77.7      79  7_8 70.8 -1.0 78.873 100.0 79.429 77.7 -1.0 -1.0
60  5_9 55.279 77.7 55.866 -1.0 53.734 -1.0 70.8 -1.0           80  7_9 77.7 -1.0 92.9 -1.0 70.8 -1.0 -1.0 -1.0
```