

SolvingMazes: Q-Learning Agent

[Achi's Projects \(https://github.com/QuantumNano-AI/PROJECTS\)](https://github.com/QuantumNano-AI/PROJECTS)

The difference between SARSA and Q-Learning is in the selection of action in the next state. While SARSA follows the current policy, Q-Learning selects the next state's action in a way that maximizes utility.

```

In [1]: import numpy as np
import sklearn.preprocessing as sc
import matplotlib.pyplot as plt
import cv2
import pandas as pd
import operator

class Maze:
    def __init__(self, maze, rewards = {'goal':1000000,'wall':-15, 'other':-1}):

        """# Read file and set height and width of maze
        with open(filename) as f:
            maze = f.read()"""
        self.rewards=rewards
        # Validate start and goal
        if maze.count("A") != 1:
            raise Exception("maze must have exactly one start point")
        if maze.count("B") != 1:
            raise Exception("maze must have exactly one goal")

        self.actions = ["up", "down", "left", "right"]

        # Determine height and width of maze
        maze = maze.splitlines()
        self.height = len(maze)
        self.width = max(len(line) for line in maze)
        self.states = []
        self.gamma = .9 # This is the discount factor
        self.theta = .00001 # Small number threshold to signal convergence of the value function
        self.probs = [round(1/len(self.actions),2)] * len(self.actions)
        self.policy = list(zip(self.actions,self.probs))

        # Keep track of walls
        self.walls = []
        self.wall_cords = []
        for i in range(self.height):
            row = []
            for j in range(self.width):
                try:
                    if maze[i][j] == "A":
                        self.start = (i, j)
                        row.append(False)
                        self.states.append((i,j))
                    elif maze[i][j] == "B":
                        self.goal = (i, j)
                        row.append(False)
                        self.states.append((i,j))
                    elif maze[i][j] == " ":
                        row.append(False)
                        self.states.append((i,j))
                    else:
                        row.append(True)
                        self.wall_cords.append((i,j))
                except IndexError:
                    row.append(False)
                    self.states.append((i,j))
            self.walls.append(row)

        self.state_count = len(self.states)
        self.solution = None
        self.V = dict(zip(self.states, self.state_count*[0]))
        self.pi = dict(zip(self.states, self.state_count*[0]))
        for s in self.states:
            avail_actions = self.actions
            self.pi[s] = avail_actions[0]
        self.pi1 = dict(zip(self.states, self.state_count*[0]))

    def print(self):
        solution = self.solution[1] if self.solution is not None else None
        print()
        for i, row in enumerate(self.walls):
            for j, col in enumerate(row):
                if col:
                    print("#", end=" ")
                elif (i, j) == self.start:
                    print("A", end=" ")
                elif (i, j) == self.goal:
                    print("B", end=" ")
                elif solution is not None and (i, j) in solution:
                    print("*", end=" ")
                else:
                    print(" ", end=" ")
            print()
        print()

    def neighbors(self, state, a = None):
        """This function takes in a state and returns all available actions for that state the next state
        and reward if each action is take, with a specific transition probability"""
        row, col = state
        candidates = [
            ("up", (row - 1, col)),
            ("down", (row + 1, col)),
            ("left", (row, col - 1)),

```

```

        ("right", (row, col + 1))
    ]
    terminal = False
    result = []
    for action, (r, c) in candidates:
        if (r,c) == self.goal: terminal = True
        if 0 <= r < self.height and 0 <= c < self.width and not self.walls[r][c]:
            if (row, col) == self.goal:
                (r, c) = self.goal; terminal = True
            reward = self.rewards['goal'] if ((r,c) == self.goal) or (state == self.goal) else self.rewards['other']
            trans_prob = 1
            result.append((action, (r, c), reward, trans_prob, terminal))

    actions = [tup[0] for tup in result]

    if a:
        R = []
        if a in actions:
            inx = actions.index(a)
            R.append((result[inx]))
            return R
        else:
            R.append((a, (row,col), self.rewards['wall'], 1, terminal))
            return R
    return result

def plot_state_values(self):
    val = np.array(list(self.V.values())).reshape(-1,1)
    va = sc.MinMaxScaler(feature_range=(0, 255)).fit_transform(val).flatten()
    V = {}
    for i in range(len(va)):
        V[list(self.V.keys())[i]] = va[i]

    # create a black image
    img = np.ones((self.height,self.width,3), np.uint8)

    for item in V.items():
        (r,c),vx = item
        img[r,c] = [0,vx,0]

    for r,c in self.wall_cords:
        img[r,c] = [150,5,150]

    img[self.start[0],self.start[1]] = [255,0,0]
    img[self.goal[0],self.goal[1]] = [100,100,255]
    def showing(img):
        plt.figure(figsize = (15,15))
        plt.imshow(img, cmap='viridis')
        plt.xticks([])
        plt.yticks([])
        #plt.colorbar()
        plt.show()
    showing(img)

def policy_(s):
    row, col = s
    candidates = [
        ("up", (row - 1, col)),
        ("down", (row + 1, col)),
        ("left", (row, col - 1)),
        ("right", (row, col + 1))
    ]

    if s in self.wall_cords:
        return ('WALL!!!')
    else:
        values = {a:self.V[r,c] for a,(r,c) in candidates if 0 <= r < self.height and 0 <= c < self.width and not self.walls[r][c]}
        values = {v:k for k,v in values.items()}
        best = values[max(values)]
        return best

pi = np.zeros((self.height, self.width)).astype('str')
pi = np.where(pi=='0.0', 'wall', pi)
candidates = {
    "up": [255,0,0],
    "down": [0,0,255],
    "left": [0,255,0],
    "right": [255,255,0]
}
for item in self.pi.keys():
    action = policy_(item)
    r,c = item
    pi[r,c] = action
    img[r,c] = candidates[action]

img[self.start[0],self.start[1]] = [255,0,0]
img[self.goal[0],self.goal[1]] = [100,100,255]
def showing(img):
    plt.figure(figsize = (15,15))
    plt.imshow(img, cmap='viridis', )
    plt.yticks(list(range(self.height)))
    plt.xticks(list(range(self.width)))
    plt.title("POLICY\n\nRED => up\nBLUE => down\nGREEN => left\nYELLOW => right")
    #plt.colorbar()

```

```

plt.show()
showimg(img)
return img

class QLearning_agent(Maze):
    def __init__(self, maze, rewards = {'goal':1000000,
                                         'wall':-15,
                                         'other':-1},

                                info = {'episodes': 200,
                                         'max_steps': 1500,
                                         'alpha': 0.4,
                                         'epsilon': 0.9} ):

        Maze.__init__(self, maze, rewards)
        self.epsilon = info['epsilon']
        self.r = np.random.RandomState(seed=12345)
        self.episodes = info['episodes']
        self.max_steps = info['max_steps']
        self.alpha = info['alpha']
        self.epsilon = info['epsilon']

    def func_q(self, states,n_states,n_actions,kind = 'random'):
        if kind=='ones':
            return dict(zip(states,np.ones((n_states,n_actions)).tolist()))
        elif kind == 'zeros':
            return dict(zip(states,np.zeros((n_states,n_actions)).tolist()))
        elif kind == 'random':
            return dict(zip(states,np.round(self.r.randn(n_states,n_actions),2).tolist()))
        else : raise NameError("Wrong input: please use ['ones', 'zeros', 'random']")

    def argmax(self, test_array):
        return self.r.choice(np.flatnonzero(np.array(test_array)==np.array(test_array).max()))

    def epsilon_greedy(self, Q, epsilon, actions, state, train=False):
        current_q = Q[state]
        if self.r.rand() < epsilon:
            action = self.r.choice(actions)
            return action
        else:
            action = self.argmax(current_q)
        return actions[action]

    def run(self):
        self.Q = self.func_q(self.states,self.state_count,len(self.actions),kind = 'random')
        self.time_steps = pd.DataFrame()
        for episode in range(self.episodes):
            total_reward = 0 # This sets the total reward obtained during this episode
            s = self.states[self.r.randint(len(self.states))]
            a = self.epsilon_greedy(Q=self.Q, epsilon=self.epsilon, actions=self.actions, state=s,train=False)
            t = 0
            terminal = False
            while t < self.max_steps:
                t+=1
                _,s_, reward, p,terminal = self.neighbors(s,a)[0]
                total_reward += reward
                q_ = self.Q[s_] # Action values in the next state
                a_ = self.actions[self.argmax(q_)] # Action in the next state does not follow policy. It is rather selected to maximise utility.
                if terminal:
                    self.Q[s][self.actions.index(a)] += self.alpha * (reward - self.Q[s][self.actions.index(a)])
                else:
                    self.Q[s][self.actions.index(a)] += self.alpha * (reward + self.gamma*self.Q[s_][self.actions.index(a_)] - \
                                                                    self.Q[s][self.actions.index(a)])

                s, a = s_, a_
                if terminal:
                    self.time_steps = self.time_steps.append(pd.Series({
                        'episode':int(episode), 'steps':t, 'rewards':total_reward
                    }), ignore_index=True)

                    break
                if terminal and t%10==0:
                    print(f'.',end='')
            self.pi = {}
            self.V = {}
            for k,v in self.Q.items():
                self.pi[(k)] = self.actions[self.argmax(v)]
                self.V[(k)] = max(v)
            max_r = max([v for k,v in self.V.items()])
            self.V[self.goal] = max_r
            img = self.plot_state_values()
            #return V,pi

```

```

In [2]: maze0 = """#####
##          #
## #####  # #
##A          B#"""

```

```

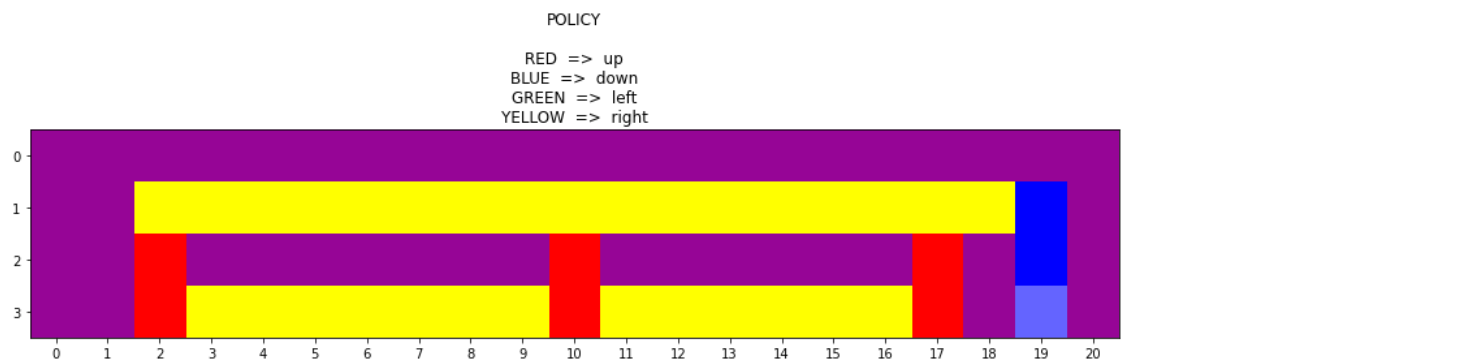
In [3]: print(maze0)

```

```

#####
##          #
## #####  # #
##A          B#

```



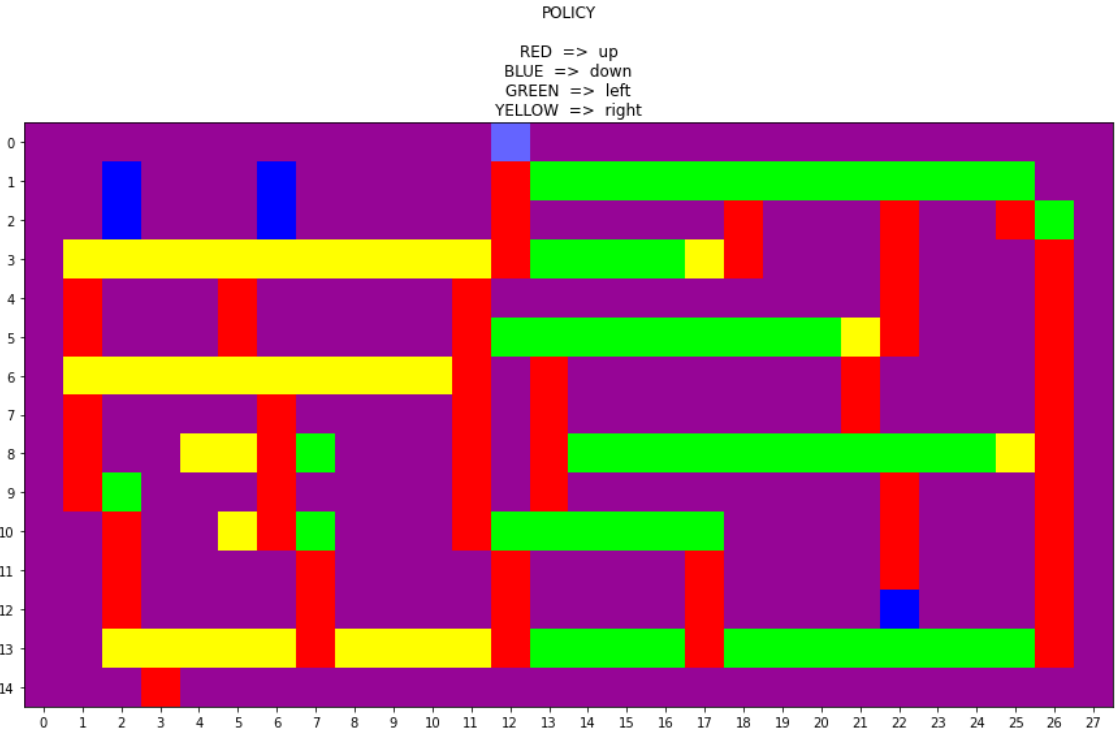
9998 rows × 3 columns

###A#####" " "

In [6]: `print(maze)`

```
#####B#####
##  ##  #####      ##
##  ##  #####  #####  ##  ##  #
#                ###  ##  #
#  ##  #####  #####  ##  #
#  ##  #####      ##  #
#          #  #####  #####  #
#  ##  #####  #  #####  #####  #
#  ##  #####  #  #####  #####  #
#  ##  #####  #  #####  #####  #
##  ##  #####  #####  #####  #
##  ##  #####  #####  #####  #
##  ##  #####  #####  #####  #
##  ##  #####  #####  #####  #
##  ##  #####  #####  #####  #
###A#####
```

```
In [7]: info = {'episodes': 10000, 'max_steps': 500, 'alpha': 0.4, 'epsilon': 0.999}
rewards = {'goal':10000, 'wall':-5, 'other':-1}
m = Qlearning_agent(maze = maze, rewards = rewards, info = info)
m.run()
t = m.time_steps; t[t.rewards!=0]
```



Out[7]:

	episode	rewards	steps
0	3.0	-926.0	498.0
1	5.0	-166.0	82.0
2	6.0	9735.0	130.0
3	7.0	9883.0	78.0
4	8.0	9881.0	84.0
...
9991	9995.0	9988.0	13.0
9992	9996.0	9981.0	20.0
9993	9997.0	9968.0	29.0
9994	9998.0	9992.0	5.0
9995	9999.0	10000.0	1.0

9996 rows × 3 columns

```
In [ ]:
```

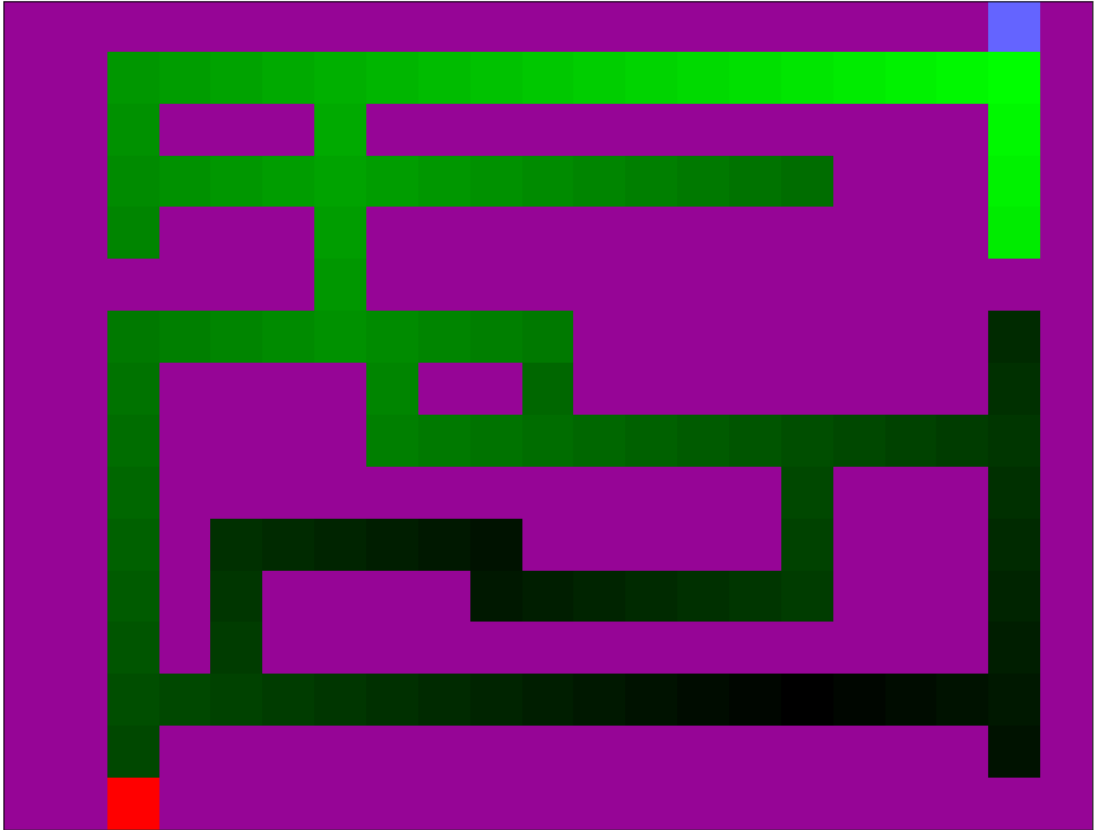
```
In [8]: maze1 = ""#####B#
##          #
## ### ##### #
##          ## #
## ### ##### #
##### #####
##          ## #
## #### # ##### #
## ####          #
## ##### ##### #
## #          ## #
## # #####      ## #
## # ##### ##### #
##          #
## ##### ##### #
##A#####
```

```
In [9]: print(maze1)
```

```
#####B#
##          #
## ### ##### #
##          ## #
## ### ##### #
##### #####
##          ## #
## #### # ##### #
## ####          #
## ##### ##### #
## #          ## #
## # #####      ## #
## # ##### ##### #
##          #
## ##### ##### #
##A#####
```

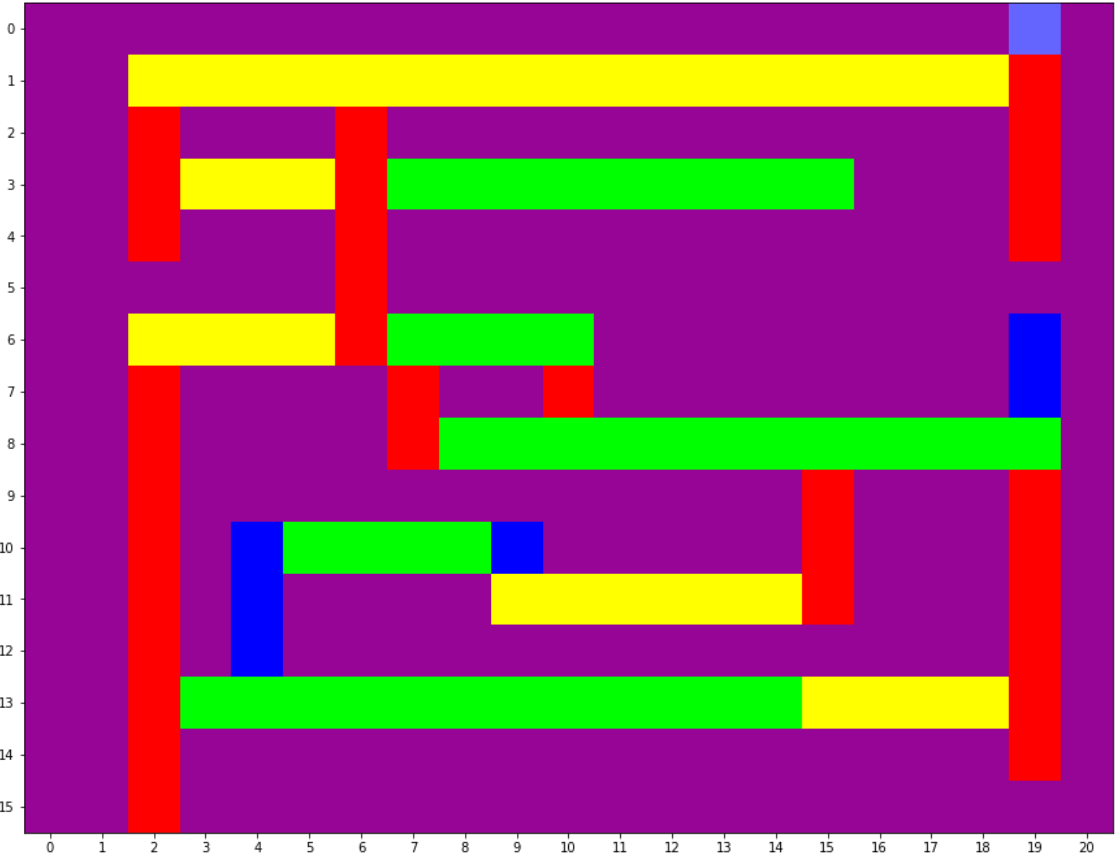


```
In [10]: info = {'episodes': 10000, 'max_steps': 500, 'alpha': 0.2, 'epsilon': 0.999}
         rewards = {'goal': 10000, 'wall': -5, 'other': -1}
         m = Qlearning_agent(maze = maze1, rewards = rewards, info = info)
         m.run()
         t = m.time_steps; t[t.rewards!=0]
```



POLICY

RED => up
BLUE => down
GREEN => left
YELLOW => right



	episode	rewards	steps
	0	7.0	-268.0
	1	10.0	9685.0
	2	12.0	9341.0
	3	15.0	9903.0
	4	18.0	9362.0

	9974	9995.0	9977.0
	9975	9996.0	9969.0
	9976	9997.0	9982.0
	9977	9998.0	9983.0
	9978	9999.0	9992.0

In []:

In []: