

# Machine Learning

## Lab 9

### Simple Python Q-learning Example

Let's look at a very simple python implementation of q-learning. We create a points-list map that represents each direction our bot can take. Using this format allows us to easily create complex graphs but also easily visualize everything with networkx graphs.

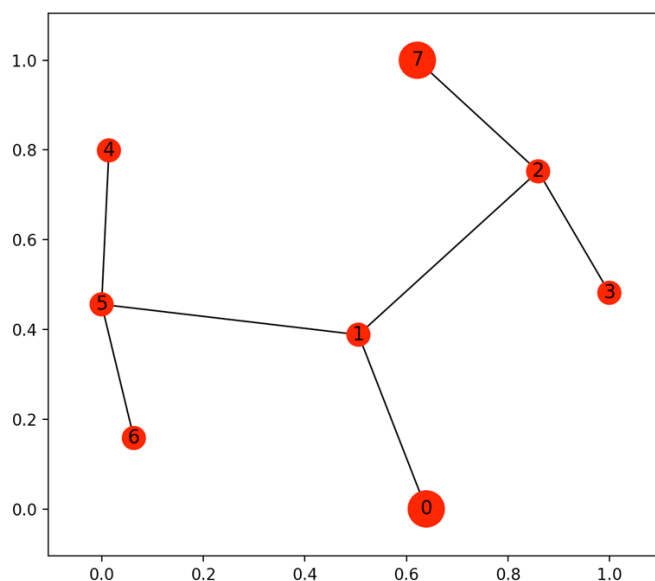
```
import numpy as np
import pylab as plt

# map cell to cell, add circular cell to goal point
points_list = [(0,1), (1,5), (5,6), (5,4), (1,2), (2,3), (2,7)]
```

Our starting point is **0**, our goal point is **7**.

```
goal = 7
```

```
import networkx as nx
G=nx.Graph()
G.add_edges_from(points_list)
pos = nx.spring_layout(G)
nx.draw_networkx_nodes(G,pos)
nx.draw_networkx_edges(G,pos)
nx.draw_networkx_labels(G,pos)
plt.show()
```



The map shows that point **0** is where our bot will start its journey and point **7** is it's final goal. The extra added points and false paths are the obstacles the bot will have to contend with.

If you look at the image, we can weave a story into this search - our bot is looking for honey, it is trying to find the hive and avoid the factory (the story-line will make sense in the second half of the article).

We then create the rewards graph - this is the matrix version of our list of points map. We initialize the matrix to be the height and width of our points list (8 in this example) and initialize all values to **-1**:

```
# how many points in graph? x points
MATRIX_SIZE = 8

# create matrix x*y
R = np.matrix(np.ones(shape=(MATRIX_SIZE, MATRIX_SIZE)))
R *= -1
```

We then change the values to be 0 if it is a viable path and 100 if it is a goal path.

```
# assign zeros to paths and 100 to goal-reaching point
for point in points_list:
    print(point)
    if point[1] == goal:
        R[point] = 100
    else:
        R[point] = 0

    if point[0] == goal:
        R[point[:-1]] = 100
    else:
        # reverse of point
        R[point[::-1]] = 0

# add goal point round trip
R[goal,goal]= 100

R

matrix([[ -1.,   0.,  -1.,  -1.,  -1.,  -1.,  -1.,  -1.],
 [   0.,  -1.,   0.,  -1.,  -1.,   0.,  -1.,  -1.],
 [  -1.,   0.,  -1.,   0.,  -1.,  -1.,  -1., 100.],
 [  -1.,  -1.,   0.,  -1.,  -1.,  -1.,  -1.,  -1.],
 [  -1.,  -1.,  -1.,  -1.,  -1.,   0.,  -1.,  -1.],
 [  -1.,   0.,  -1.,  -1.,   0.,  -1.,   0.,  -1.],
 [  -1.,  -1.,  -1.,  -1.,  -1.,   0.,  -1.,  -1.],
 [  -1.,  -1.,   0.,  -1.,  -1.,  -1.,  -1., 100.]])
```

To read the above matrix, the y-axis is the state or where your bot is currently located, and the x-axis is your possible next actions. We then build our Q-learning matrix which will hold

all the lessons learned from our bot. The Q-learning model uses a transitional rule formula and gamma is the learning parameter

```
Q = np.matrix(np.zeros([MATRIX_SIZE,MATRIX_SIZE]))

# learning parameter
gamma = 0.8

initial_state = 1

def available_actions(state):
    current_state_row = R[state,]
    av_act = np.where(current_state_row >= 0)[1]
    return av_act

available_act = available_actions(initial_state)

def sample_next_action(available_actions_range):
    next_action = int(np.random.choice(available_act,1))
    return next_action

action = sample_next_action(available_act)

def update(current_state, action, gamma):

    max_index = np.where(Q[action,] == np.max(Q[action,]))[1]

    if max_index.shape[0] > 1:
        max_index = int(np.random.choice(max_index, size = 1))
    else:
        max_index = int(max_index)
    max_value = Q[action, max_index]

    Q[current_state, action] = R[current_state, action] + gamma * max_value
    print('max_value', R[current_state, action] + gamma * max_value)

    if (np.max(Q) > 0):
        return(np.sum(Q/np.max(Q)*100))
    else:
        return (0)

update(initial_state, action, gamma)
```

Then we run the training and testing functions that will run the update function 700 times allowing the Q-learning model to figure out the most efficient path:

```

# Training
scores = []
for i in range(700):
    current_state = np.random.randint(0, int(Q.shape[0]))
    available_act = available_actions(current_state)
    action = sample_next_action(available_act)
    score = update(current_state, action, gamma)
    scores.append(score)
    print ('Score:', str(score))

print("Trained Q matrix:")
print(Q/np.max(Q)*100)

# Testing
current_state = 0
steps = [current_state]

while current_state != 7:

    next_step_index = np.where(Q[current_state,] == np.max(Q[current_state,]))[1]

    if next_step_index.shape[0] > 1:
        next_step_index = int(np.random.choice(next_step_index, size = 1))
    else:
        next_step_index = int(next_step_index)

    steps.append(next_step_index)
    current_state = next_step_index

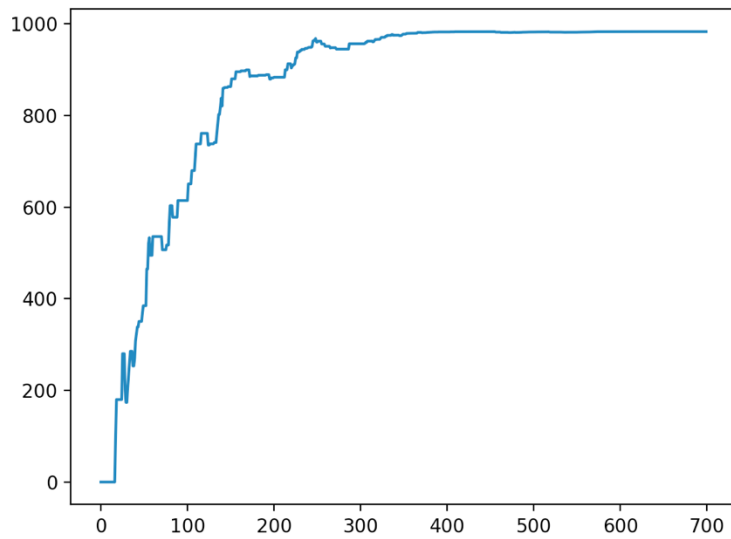
print("Most efficient path:")
print(steps)

plt.plot(scores)
plt.show()

Most efficient path:
[0, 1, 2, 7]

```

We see that the model did correctly find the most efficient path from the starting node **0** to the goal node **7** and took around 400 iterations to converge to a solution.



### Activity

Modify the codes and solve the following problem where agent want to go from 2 to 5.

