



## COSC 2753 | Machine Learning

### Week 8+1 Lab Exercises: Reinforcement learning

```
1 !pip install gym[box2d]==0.17.* pyvirtualdisplay==0.2.* PyOpenGL==3.1.* PyOpenGL-accelerate==3.1.*
```

#### Introduction

In this lab you will be:

1. Learning how to use OpenAI gym.
2. Implement Q-learning to solve a well-known toy reinforcement learning problem called MountainCar problem or Cartpole problem (<https://gym.openai.com/envs/CartPole-v1/>).

\*\* Please run this lab on the anaconda environment on your PC.\*\*

#### Mountain Car with RL

Mountain Car is a classic control Reinforcement Learning problem that was first introduced by A. Moore in 1991 [1]. A car is on a one-dimensional track, positioned between two “mountains”. The goal is to drive up the mountain on the right; however, the car’s engine is not strong enough to scale the mountain in a single pass. Therefore, the only way to succeed is to drive back and forth to build up momentum. It can be tricky to find this optimal solution due to the sparsity of the reward.

##### Mountain Car Problem definition:

- Objective: Get the car to the top of the right hand side mountain.
- State: Car’s horizontal position and velocity (can be negative).
- Action: Direction of push (left, nothing or right).
- Reward: -1 for every time step until success, which incentivises quick solutions.

More information about Mountain Car can be found: [https://en.wikipedia.org/wiki/Mountain\\_car\\_problem](https://en.wikipedia.org/wiki/Mountain_car_problem) or <https://www.toptal.com/machine-learning/deep-dive-into-reinforcement-learning>

#### OpenAI Gym

OpenAI Gym is a Python package comprising a selection of RL environments, ranging from simple “toy” environments to more challenging environments, including simulated robotics environments and Atari video game environments. It was developed with the aim of becoming a standardized environment and benchmark for RL research. In this Lab, we will use the OpenAI Gym Cartpole environment to demonstrate how to get started in using this exciting tool and show how Q-learning can be used to solve this problem.

▼

#### Setting up the environment

Lets first import the libraries required for the implementation.

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3 from IPython import display
4 !pip install gym
5 import numpy as np
6 import gym

Requirement already satisfied: gym in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (0.17.3)
Requirement already satisfied: pygame<=1.5.0,>=1.4.0 in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (from gym) (1.5.0)
Requirement already satisfied: cloudpickle<1.7.0,>=1.2.0 in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (from gym) (1.6.0)
Requirement already satisfied: scipy in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (from gym) (1.7.1)
Requirement already satisfied: numpy>=1.10.4 in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (from gym) (1.20.3)
Requirement already satisfied: future in /Users/thienbao/opt/anaconda3/lib/python3.9/site-packages (from pygame<=1.5.0,>=1.4.0->gym) (0.18.2)
```

To start using the mountain car environment initialize it as follows:

```
1 import gym
2 print(gym.__version__)# for me: 0.15.4
3 env = gym.make("MountainCar-v0")
4 obs = env.reset()
5 for i in range(1000):# it's changable
6     env.step(env.action_space.sample())
7     env.render()# won't work in Google Colab
8 env.close()

0.17.3
```

The `env.reset()` command resets the environemnt and return the initial state

Lets explore the state space and the action space og the Cartpole environment

```
1 print(obs)

[-0.4160749  0.          ]

1 print('State space: ', env.observation_space)
2 print('Action space: ', env.action_space)

State space:  Box(-1.2000000476837158, 0.6000000238418579, (2,), float32)
Action space:  Discrete(3)
```

1. `print('State space: ', env.observation_space):` This line is printing the state space of the environment. The state space, represented by `env.observation_space`, is the set of all possible states that the agent can be in. In a game, for example, this could include the positions of all the objects on the screen. The structure of the state space depends on the specific environment. It could be a multi-dimensional array for complex environments or a simple numeric value for simpler ones.
2. `print('Action space: ', env.action_space):` This line is printing the action space of the environment. The action space, represented by `env.action_space`, is the set of all possible actions that the agent can take. In a game, this could include moving left, moving right, jumping, etc. Similar to the state space, the structure of the action space depends on the specific environment. It could be a simple discrete value for environments with a fixed number of actions or a continuous value for environments with a range of possible actions.

The `print` function is a built-in Python function that writes the specified message to the screen, or other standard output device. The message can be a string, or any other object, the object will be converted into a string before written to the screen. In this case, it's being used to display the state and action spaces of the environment.

This tells us that the state space is a 2-dimensional space, so each state observation is a vector of 2 (float) values, and that the action space comprises three discrete actions (left, nothing or right). By default, the three actions are represented by the integers 0, 1 and 2. How about the state space? What are the limits of the state space?

```
1 print('State space Low: ', env.observation_space.low)
2 print('State space High: ', env.observation_space.high)
```

```
State space Low:  [-1.2  -0.07]
State space High:  [0.6   0.07]
```

State: Two-dimensional continuous state space.

```
Velocity=(-0.07,0.07)

Position=(-1.2,0.6)
```

Actions: One-dimensional discrete action space.

```
action=(left,neutral,right)
```

Reward: For every time step:

```
reward=-1
```

Update function: For every time step:

```
Action=[-1,0,1]

Velocity=Velocity+(Action)0.001+cos((3Position)*(-0.0025))

Position=Position+Velocity
```

Starting condition: Optionally, many implementations include randomness in both parameters to show better generalized learning.

```
Position=-0.5

Velocity=0.0
```

Termination condition: End the simulation when:

```
Position >= 0.6
```

This shows that the first state variable (horizontal position) has a range [-1.2, 0.6] and the second state variable(speed) has a range [-0.07, 0.07]. The state space of the environment is a continuous state space, which means that there are infinitely many state-action pairs, making it impossible to build a Q table. As a solution to this problem we can discretize the state space. One simple discretization is to cover the state space to a grid with spacing of 0.1 along first element and 0.01 along second element in the state space. The states can be given integer indexes multiply the first element by 10 and the second by 100. lets see the size of discretized state space.

```
1 num_states = (env.observation_space.high - env.observation_space.low)*np.array([10, 100])
2 num_states = np.round(num_states, 0).astype(int) + 1
3 print(num_states)

[19 15]
```

We can also write a function that will convert a continuous state vector to a discrete one.

```
1 # Discretize state
2 def discretize_state(state, env_low):
3     state_adj = (state - env_low)*np.array([10, 100])
4     state_adj = np.round(state_adj, 0).astype(int)
5     return state_adj
```

Let's now make some random actions in the environment and see what the output will be. For this we need a function to plot the output of the environment.

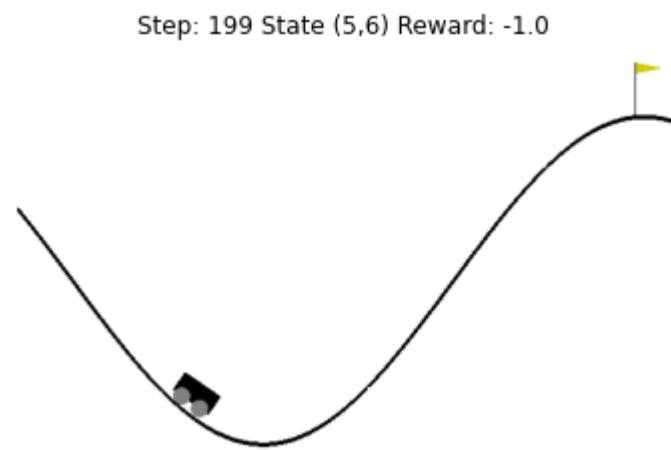
```
1 def show_state(env, step=0, info=""):
2     plt.figure(1)
3     plt.clf()
4     plt.imshow(env.render(mode='rgb_array'))
5     plt.title("Step: %d %s" % (step, info))
6     plt.axis('off')
7     display.clear_output(wait=True)
8     display.display(plt.gcf())
```

Now we take some random actions:

```

1 env.reset()
2 done = False
3 step_index = 0
4 while done != True:
5     action = env.action_space.sample() # get a random action from the set of actions
6     state, reward, done, info = env.step(action) # perform the action and receive new state and reward
7     d_state = discretize_state(state, env.observation_space.low)
8     show_state(env, step=step_index, info='State ({}{}) Reward: {}'.format(d_state[0], d_state[1], reward))
9     step_index = step_index + 1

```



The provided Python code is a simple loop that runs a reinforcement learning episode in an environment represented by the `env` object. Here's a step-by-step explanation of what the code does:

1. `env.reset()`: This line resets the environment to its initial state. This is typically done at the beginning of each episode to ensure that the environment is in a known state.
2. `done = False`: This line initializes the `done` variable to `False`. This variable is used to indicate whether the episode has ended. The loop continues to run as long as `done` is `False`.
3. `while done != True`: This line starts a while loop that continues to run until `done` is `True`, which indicates that the episode has ended.
4. `action = env.action_space.sample()`: This line selects a random action from the action space of the environment. The `env.action_space.sample()` function returns a random action.
5. `state, reward, done, info = env.step(action)`: This line performs the selected action in the environment, which results in a new state and a reward. The `env.step(action)` function takes an action as input and returns the new state, the reward for performing the action, a boolean indicating whether the episode has ended (`done`), and an info dictionary containing extra information.
6. `d_state = discretize_state(state, env.observation_space.low)`: This line discretizes the continuous state space into a discrete state space. This is typically done when you want to apply a discrete-space algorithm, like Q-learning, to a continuous-space problem.
7. `show_state(env, step=step_index, info='State ({}{}) Reward: {}'.format(d_state[0], d_state[1], reward))`: This line visualizes the current state of the environment. The `show_state` function is not defined in the provided code, but it likely takes the environment, the current step index, and some info text as input and displays the state of the environment in some way.
8. `step_index = step_index + 1`: This line increments the step index by one. The step index is used to keep track of the number of steps that have been taken in the current episode.

In summary, this code runs a single episode of a reinforcement learning experiment in a given environment. It selects actions randomly, performs them in the environment, and visualizes the resulting states.

Run the following block if the visualization of the environment gives an error. On mac you need to install pyglet version 1.5.11 to get the gym environment to render. The installation will give an error, but it will work.

```

1 #!pip install pyglet==1.5.11

```

Now let's develop a function for Q learning. The function prototype is given below and the algorithm for Q learning is given in Algorithm 1. Assume that `Q` is a numpy matrix with dimensions (number of elements for state 1, number of elements for state 2, number of actions).

```

1 # Define Q-learning function
2 def QLearning(env, Q, learning, discount, epsilon, episodes):
3     # Env: The OpenAI gym environment
4     # Q: Initial Q table
5     # learning: Learning Rate of Q learning
6     # discount: discount factor (gamma)
7     # epsilon: epsilon for exploration vs exploitation
8     # episodes: number of episodes to run when learning the Q table
9
10    # Initialize variables to hold rewards
11    reward_list = []
12
13    # Calculate reduction in epsilon per episode
14    epsilon_d = (epsilon)/episodes
15
16    for i in range(episodes):
17        done = False
18        tot_reward, reward = 0,0
19        state = env.reset()
20
21        state_adj = discretize_state(state, env.observation_space.low)
22
23        while done != True:
24
25            # Determine next action - epsilon greedy strategy for explore vs exploitation
26            if np.random.random() < 1 - epsilon:

```

```

27         # select the best action according to Qtable (exploitation)
28         # TODO
29     else:
30         # select a random action (exploration)
31         # TODO
32
33     # Step and Get the next state and reward
34     # TODO
35
36     # Allow for terminal states
37     if done and state2[0] >= 0.5:
38         Q[state_adj[0], state_adj[1], action] = reward
39
40     # Update the Q table
41     else:
42         # TODO
43
44     # Update variables
45     tot_reward += reward
46     state_adj = state2_adj
47
48     # Update epsilon
49     if epsilon > 0:
50         epsilon -= epsilon_d
51
52     # Track rewards
53     reward_list.append(tot_reward)
54
55     if (i+1) % 100 == 0:
56         ave_reward = np.mean(reward_list)
57         reward_list = []
58
59     # Average reward is the average number of steps that the agent spent to win
60     if (i+1) % 100 == 0:
61         print('Episode {} Average Reward: {} Epsilon {}'.format(i+1, ave_reward, np.round(ep
62
63 env.close()
64
65 return Q

```

```
1 # Define Q-learning function
2 def QLearning(env, Q, learning, discount, epsilon, episodes):
3     # Env: The OpenAI gym environment
4     # Q: Initial Q table
5     # learning: Learning Rate of Q learing
6     # discount: discount factor (gamma)
7     # epsilon: epsilon for exploration vs exploitation
8     # episodes: number of episodes to run when learing the Q table
9
10    # Initialize variables to hold rewards
11    reward_list = []
12    ave_reward_list = []
13
14    # Calculate reduction in epsilon per episode
15    epsilon_d = (epsilon)/episodes
16
17    for i in range(episodes):
```

The provided Python code defines a function `QLearning` that implements the Q-learning algorithm, a type of reinforcement learning algorithm. The function takes as input an OpenAI Gym environment (`env`), an initial Q-table (`Q`), a learning rate (`learning`), a discount factor (`discount`), an epsilon value for the epsilon-greedy policy (`epsilon`), and the number of episodes to run (`episodes`).

The function begins by initializing two lists to hold the rewards for each episode (`reward_list`) and the average rewards over every 100 episodes (`ave_reward_list`). It also calculates the amount by which epsilon should be decreased after each episode (`epsilon_d`) to gradually shift from exploration to exploitation.

The function then enters a loop that runs for the specified number of episodes. For each episode, it initializes some variables, resets the environment to its initial state, and discretizes this state. It then enters another loop that continues until the episode ends.

In this inner loop, the function first decides whether to take the best action according to the current Q-table (exploitation) or a random action (exploration), based on a random number and the current epsilon value. It then performs the chosen action in the environment using the `env.step(action)` function, which returns the new state, the reward for the action, a boolean indicating whether the episode has ended, and an info dictionary.

The function then discretizes the new state and updates the Q-table. If the episode has ended and the agent has reached the goal state, it sets the Q-value for the current state-action pair to the reward. Otherwise, it calculates the difference between the current Q-value and the expected Q-value based on the reward and the maximum Q-value for the new state, scales this difference by the learning rate, and adds the result to the current Q-value.

After updating the Q-table, the function updates the total reward for the episode and the current state, and exits the inner loop. It then decreases epsilon by `epsilon_d`, adds the total reward for the episode to `reward_list`, and calculates and prints the average reward every 100 episodes.

Finally, after all episodes have been run, the function closes the environment and returns the final Q-table and the list of average rewards.

```
43         if done and state[0] < 0.5:
```

Sample Solutions

If you are struggling with the above function, a sample solution has been provided. Only use this if you have made your absolute best attempts at implementing the function yourself. The purpose of this lab is to understand common aspects of RL algorithm, though the Q-learning algorithm. You will gain significantly less out of this lab if you don't try to solve the problems yourself.

```
55         state2 = env.reset()
56         state2 = discretize(state2)
57         state2_action_list = []
```

> Learning & testing the model

Now we have all the elements required. Let's learn the model.

```
[ ] ↳ 6 cells hidden
```

> Cartpole with RL

`Cartpole` is a classic control Reinforcement Learning problem that was first introduced by by Barto, Sutton, and Anderson [Barto83]. A pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The system is controlled by applying a force of +1 or -1 to the cart. The pendulum starts upright, and the goal is to prevent it from falling over. A reward of +1 is provided for every timestep that the pole remains upright. The episode ends when the pole is more than 15 degrees from vertical, or the cart moves more than 2.4 units from the center.

```
[ ] ↳ 16 cells hidden
```

```
71         reward_list = []
```

> Learning

We are going to use Q-learning for this task. Lets first define some hyper parameters. You may change them to get better performance later.

```
[ ] ↳ 7 cells hidden
```

```
70
79         return Q, ave_reward_list
```