# ▾ Reinforcement Learning

The next and final topic in this course covers *Reinforcement Learning*. This technique is different than many of the other machine learning techniques we have seen earlier and has many applications in training agents (an AI) to interact with enviornments like games. Rather than feeding our machine learning model millions of examples we let our model come up with its own examples by exploring an enviornemt. The concept is simple. Humans learn by exploring and learning from mistakes and past experiences so let's have our computer do the same.

## Terminology

Before we dive into explaining reinforcement learning we need to define a few key peices of terminology.

**Enviornemt** In reinforcement learning tasks we have a notion of the enviornment. This is what our *agent* will explore. An example of an enviornment in the case of training an AI to play say a game of mario would be the level we are training the agent on.

**Agent** an agent is an entity that is exploring the enviornment. Our agent will interact and take different actions within the enviornment. In our mario example the mario character within the game would be our agent.

**State** always our agent will be in what we call a *state*. The state simply tells us about the status of the agent. The most common example of a state is the location of the agent within the enviornment. Moving locations would change the agents state.

**Action** any interaction between the agent and enviornment would be considered an action. For example, moving to the left or jumping would be an action. An action may or may not change the current *state* of the agent. In fact, the act of doing nothing is an action as well! The action of say not pressing a key if we are using our mario example.

**Reward** every action that our agent takes will result in a reward of some magnitude (positive or negative). The goal of our agent will be to maximize its reward in an enviornment. Sometimes the reward will be clear, for example if an agent performs an action which increases their score in the enviornment we could say they've recieved a positive reward. If the agent were to perform an action which results in them losing score or possibly dying in the enviornment then they would recieve a negative reward.

The most important part of reinforcement learning is determing how to reward the agent. After all, the goal of the agent is to maximize its rewards. This means we should reward the agent appropiatly such that it reaches the desired goal.

## Q-Learning

Now that we have a vague idea of how reinforcement learning works it's time to talk about a specific technique in reinforcement learning called *Q-Learning*.

Q-Learning is a simple yet quite powerful technique in machine learning that involves learning a matrix of action-reward values. This matrix is often reffered to as a Q-Table or Q-Matrix. The matrix is in shape (number of possible states, number of possible actions) where each value at matrix[n, m] represents the agents expected reward given they are in state n and take action m. The Q-learning algorithm defines the way we update the values in the matrix and decide what action to take at each state. The idea is that after a succesful training/learning of this Q-Table/matrix we can determine the action an agent should take in any state by looking at that states row in the matrix and taking the maximium value column as the action.

**Consider this example.**

Let's say A1-A4 are the possible actions and we have 3 states represented by each row (state 1 - state 3).

| A1 | A2 | A3 | A4 |
|----|----|----|----|
| 0  | 0  | 10 | 5  |
| 5  | 10 | 0  | 0  |
| 10 | 5  | 0  | 0  |

If that was our Q-Table/matrix then the following would be the preffered actions in each state.

> State 1: A3

> State 2: A2

> State 3: A1

We can see that this is because the values in each of those columns are the highest for those states!

## Learning the Q-Table

So that's simple, right? Now how do we create this table and find those values. Well this is where we will dicuss how the Q-Learning algorithm updates the values in our Q-Table.

I'll start by noting that our Q-Table starts of with all 0 values. This is because the agent has yet to learn anything about the enviornment.

Our agent learns by exploring the enviornment and observing the outcome/reward from each action it takes in each state. But how does it know what action to take in each state? There are two ways that our agent can decide on which action to take.

1. Randomly picking a valid action
2. Using the current Q-Table to find the best action.

Near the beginning of our agents learning it will mostly take random actions in order to explore the enviornment and enter many different states. As it starts to explore more of the enviornment it will start to gradually rely more on it's learned values (Q-Table) to take actions. This means that as our agent explores more of the enviornment it will develop a better understanding and start to take "correct" or better actions more often. It's important that the agent has a good balance of taking random actions and using learned values to ensure it does get trapped in a local maximum.

After each new action our agent wil record the new state (if any) that it has entered and the reward that it recieved from taking that action. These values will be used to update the Q-Table. The agent will stop taking new actions only once a certain time limit is reached or it has acheived the goal or reached the end of the enviornment.

## Updating Q-Values

The formula for updating the Q-Table after each action is as follows:

$$Q[state, action] = Q[state, action] + \alpha * (reward + \gamma * max(Q[newState, :]) - Q[state, action])$$

- $\alpha$ stands for the **Learning Rate**

- $\gamma$ stands for the **Discount Factor**

## Learning Rate $\alpha$

The learning rate $\alpha$ is a numeric constant that defines how much change is permitted on each QTable update. A high learning rate means that each update will introduce a large change to the current state-action value. A small learning rate means that each update has a more subtle

change. Modifying the learning rate will change how the agent explores the enviornment and how quickly it determines the final values in the QTable.

Discount Factor $\gamma$

Discount factor also know as gamma ($\gamma$) is used to balance how much focus is put on the current and future reward. A high discount factor means that future rewards will be considered more heavily.

To perform updates on this table we will let the agent explpore the enviornment for a certain period of time and use each of its actions to make an update. Slowly we should start to notice the agent learning and choosing better actions.

## ▾ Q-Learning Example

For this example we will use the Q-Learning algorithm to train an agent to navigate a popular enviornment from the [Open AI Gym](). The Open AI Gym was developed so programmers could practice machine learning using unique enviornments. Intersting fact, Elon Musk is one of the founders of OpenAI!

Let's start by looking at what Open AI Gym is.

```
import gym   # all you have to do to import and use open ai gym!
```

Once you import gym you can load an enviornment using the line `gym.make("enviornment")`.

```
env = gym.make('FrozenLake-v0')  # we are going to use the FrozenLake enviornment
```

There are a few other commands that can be used to interact and get information about the enviornment.

```
print(env.observation_space.n)   # get number of states
```

```
print(env.action_space.n)    # get number of actions

env.reset()  # reset enviornment to default state

action = env.action_space.sample()  # get a random action

new_state, reward, done, info = env.step(action)  # take action, notice it returns information about the action

env.render()    # render the GUI for the enviornment
```

## Frozen Lake Enviornment

Now that we have a basic understanding of how the gym enviornment works it's time to discuss the specific problem we will be solving.

The enviornment we loaded above `FrozenLake-v0` is one of the simplest enviornments in Open AI Gym. The goal of the agent is to navigate a frozen lake and find the Goal without falling through the ice (render the enviornment above to see an example).

There are:

- 16 states (one for each square)
- 4 possible actions (LEFT, RIGHT, DOWN, UP)
- 4 different types of blocks (F: frozen, H: hole, S: start, G: goal)

## ▾ Building the Q-Table

The first thing we need to do is build an empty Q-Table that we can use to store and update our values.

```
import gym
import numpy as np
import time

env = gym.make('FrozenLake-v0')
```

```
STATES = env.observation_space.n
ACTIONS = env.action_space.n
```

```
Q = np.zeros((STATES, ACTIONS))  # create a matrix with all 0 values
Q
```

## ▾ Constants

As we discussed we need to define some constants that will be used to update our Q-Table and tell our agent when to stop training.

```
EPISODES = 2000 # how many times to run the enviornment from the beginning
MAX_STEPS = 100  # max number of steps allowed for each run of enviornment

LEARNING_RATE = 0.81  # learning rate
GAMMA = 0.96
```

## ▾ Picking an Action

Remember that we can pick an action using one of two methods:

1. Randomly picking a valid action
2. Using the current Q-Table to find the best action.

Here we will define a new value $\epsilon$ that will tell us the probabillity of selecting a random action. This value will start off very high and slowly decrease as the agent learns more about the enviornment.

```
epsilon = 0.9  # start with a 90% chance of picking a random action

# code to pick action
if np.random.uniform(0, 1) < epsilon:  # we will check if a randomly selected value is less than epsilon.
    action = env.action_space.sample()  # take random action
```

```
else:
    action = np.argmax(Q[state, :])  # use Q table to pick best action based on current values
```

## ▾ Updating Q Values

The code below implements the formula discussed above.

```
Q[state, action] = Q[state, action] + LEARNING_RATE * (reward + GAMMA * np.max(Q[new_state, :]) - Q[state, action])
```

## ▾ Putting it Together

Now that we know how to do some basic things we can combine these together to create our Q-Learning algorithm,

```
import gym
import numpy as np
import time

env = gym.make('FrozenLake-v0')
STATES = env.observation_space.n
ACTIONS = env.action_space.n

Q = np.zeros((STATES, ACTIONS))

EPISODES = 1500 # how many times to run the enviornment from the beginning
MAX_STEPS = 100  # max number of steps allowed for each run of enviornment

LEARNING_RATE = 0.81  # learning rate
GAMMA = 0.96

RENDER = False # if you want to see training set to true
```

```python
    epsilon = 0.9

    rewards = []
    for episode in range(EPISODES):

      state = env.reset()
      for _ in range(MAX_STEPS):

        if RENDER:
          env.render()

        if np.random.uniform(0, 1) < epsilon:
          action = env.action_space.sample()
        else:
          action = np.argmax(Q[state, :])

        next_state, reward, done, _ = env.step(action)

        Q[state, action] = Q[state, action] + LEARNING_RATE * (reward + GAMMA * np.max(Q[next_state, :]) - Q[state, action])

        state = next_state

        if done:
          rewards.append(reward)
          epsilon -= 0.001
          break  # reached goal

    print(Q)
    print(f"Average reward: {sum(rewards)/len(rewards)}:")
    # and now we can see our Q values!


    # we can plot the training progress and see how the agent improved
    import matplotlib.pyplot as plt

    def get_average(values):
      return sum(values)/len(values)
```

```
avg_rewards = []
for i in range(0, len(rewards), 100):
  avg_rewards.append(get_average(rewards[i:i+100]))

plt.plot(avg_rewards)
plt.ylabel('average reward')
plt.xlabel('episodes (100\'s)')
plt.show()
```

## Sources

1. Violante, Andre. "Simple Reinforcement Learning: Q-Learning." Medium, Towards Data Science, 1 July 2019, https://towardsdatascience.com/simple-reinforcement-learning-q-learning-fcddc4b6fe56.

2. Openai. "Openai/Gym." GitHub, https://github.com/openai/gym/wiki/FrozenLake-v0.