# Agenda

1. Introduction To Machine Learning
2. What is Reinforcement Learning?
3. Reinforcement Learning with an Analogy
4. Reinforcement Learning Process
5. Reinforcement Learning – Counter strike example
6. Reinforcement Learning Definitions
7. Reinforcement Learning Concepts
8. Markov's Decision Process
9. Understanding Q-Learning
10. Hands-On

# What Is Machine Learning?

*Machine learning is a subset of artificial intelligence (AI) which provides machines the ability to learn automatically & improve from experience without being explicitly programmed.*



They look the same!

Data

Algorith

Cherry

Apple

Orange

# Types Of Machine Learning
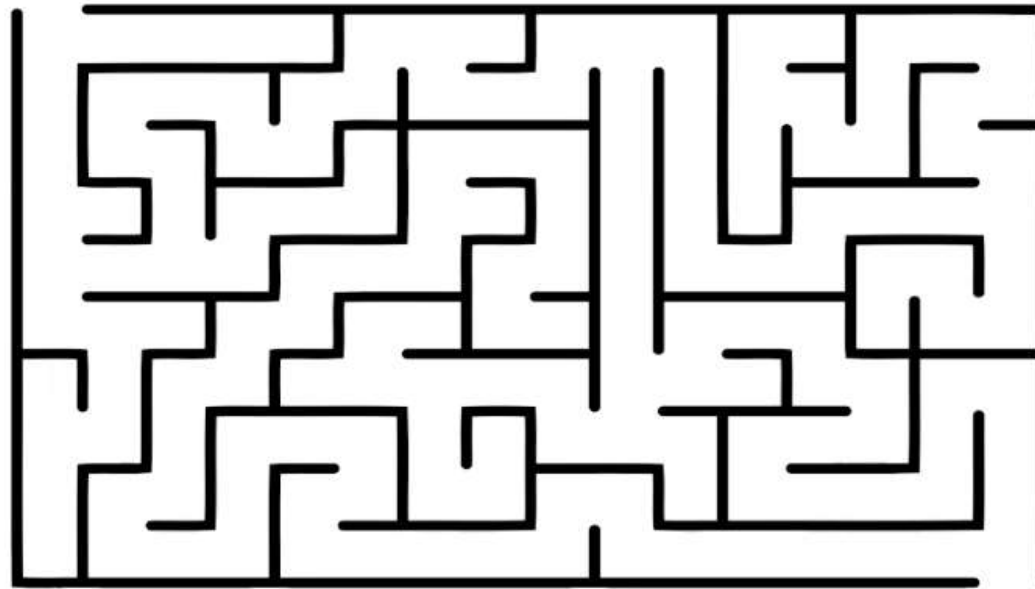


Supervised Learning



Unsupervised Learning



Reinforcement Learning

# What Is Reinforcement Learning?

*Reinforcement learning is a type of Machine Learning where an agent learns to behave in a environment by performing actions and seeing the results*

# Reinforcement Learning With An Analogy

Scenario 1: Baby starts crawling and makes it to the candy

# Reinforcement Learning With An Analogy

Scenario 1: Baby starts crawling and makes it to the candy

# Reinforcement Learning With An Analogy

Scenario 2: Baby starts crawling but falls due to some hurdle in between

# Reinforcement Learning With An Analogy

Scenario 2: Baby starts crawling but falls due to some hurdle in between

# Reinforcement Learning With An Analogy

Scenario 2: Baby starts crawling but falls due to some hurdle in between

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main
components:
- Agent
- Environment

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# Reinforcement Learning Process

Reinforcement Learning system is comprised of two main components:
- Agent
- Environment

# Counter Strike Example



1. The RL Agent (Player1) collects state $S^0$ from the environment

2. Based on the state $S^0$, the RL agent takes an action $A^0$, initially the action is random

3. The environment is now in a new state $S^1$

4. RL agent now gets a reward $R^1$ from the environment

5. The RL loop goes on until the RL agent is dead or reaches the destination

# Reinforcement Learning Definitions

Agent: The RL algorithm that learns from trial and error

Environment: The world through which the agent moves

Action (A): All the possible steps that the agent can take

State (S): Current condition returned by the environment

# Reinforcement Learning Definitions

Reward (R): An instant return from the environment to appraise the last action

Policy (π): The approach that the agent uses to determine the next action based on the current state

Value (V): The expected long-term return with discount, as opposed to the short-term reward R

Action-value (Q): This similar to Value, except, it takes an extra parameter, the current action (A)

# Reward Maximization

Reward maximization theory states that, *a RL agent must be trained in such a way that, he takes the best action so that the reward is maximum.*

# Exploration & Exploitation

*Exploitation* is about using the already known exploited information to heighten the rewards

*Exploration* is about exploring and capturing more information about an environment

# Markov Decision Process

The mathematical approach for mapping a solution in reinforcement learning is called *Markov Decision Process* (MDP)

The following parameters are used to attain a solution:

- Set of actions, A
- Set of states, S
- Reward, R
- Policy, π
- Value, V

# Markov Decision Process – Shortest Path Problem

Goal: Find the shortest path between A and D with minimum possible cost

In this problem,

- Set of states are denoted by nodes i.e. {A, B, C, D}

- Action is to traverse from one node to another {A -> B, C -> D}

- Reward is the cost represented by each edge

- Policy is the path taken to reach the destination {A -> C -> D}

# Understanding Q-Learning With An Example

*Place an agent in any one of the rooms (0,1,2,3,4) and the goal is to reach outside the building (room 5)*



- 5 rooms in a building connected by doors

- each room is numbered 0 through 4

- The outside of the building can be thought of as one big room (5)

- Doors 1 and 4 lead into the building from room 5 (outside)

# Understanding Q-Learning With An Example

Let's represent the rooms on a graph, each room as a node, and each door as a link

# Understanding Q-Learning With An Example

Next step is to associate a reward value to each door:

- doors that lead directly to the goal have a reward of 100

- Doors not directly connected to the target room have zero reward

- Because doors are two-way, two arrows are assigned to each room

- Each arrow contains an instant reward value

# Understanding Q-Learning With An Example

The terminology in Q-Learning includes the terms state and action:
- Room (including room 5) represents a state
- agent's movement from one room to another represents an action
- In the figure, a state is depicted as a node, while "action" is represented by the arrows



Example (Agent traverse from room 2 to room5):

1. Initial state = state 2

2. State 2 -> state 3

3. State 3 -> state (2, 1, 4)

4. State 4 -> state 5

# Understanding Q-Learning With An Example

We can put the state diagram and the instant reward values into a reward table, matrix R.



$$R =$$

| State | Action 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | -1 | -1 | -1 | -1 | 0 | -1 |
| 1 | -1 | -1 | -1 | 0 | -1 | 100 |
| 2 | -1 | -1 | -1 | 0 | -1 | -1 |
| 3 | -1 | 0 | 0 | -1 | 0 | -1 |
| 4 | 0 | -1 | -1 | 0 | -1 | 100 |
| 5 | -1 | 0 | -1 | -1 | 0 | 100 |

The -1's in the table represent null values

# Understanding Q-Learning With An Example

Add another matrix Q, representing the memory of what the agent has learned through experience.
- The rows of matrix Q represent the current state of the agent
- columns represent the possible actions leading to the next state
- Formula to calculate the Q matrix:

Q(state, action) = R(state, action) + Gamma * Max [Q(next state, all actions)]

**Note**

The Gamma parameter has a range of 0 to 1 (0 <= Gamma > 1).
- If Gamma is closer to zero, the agent will tend to consider only immediate rewards.
- If Gamma is closer to one, the agent will consider future rewards with greater weight

# Q – Learning Algorithm

1. Set the gamma parameter, and environment rewards in matrix R

2. Initialize matrix Q to zero

3. Select a random initial state

4. Set initial state = current state

5. Select one among all possible actions for the current state

6. Using this possible action, consider going to the next state

7. Get maximum Q value for this next state based on all possible actions

8. Compute: Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

9. Repeat above steps until current state = goal state

# Q – Learning Example

First step is to set the value of the learning parameter Gamma = 0.8, and the initial state as Room 1.

Next, initialize matrix Q as a zero matrix:
- From room 1 you can either go to room 3 or 5, let's select room 5.
- From room 5, calculate maximum Q value for this next state based on all possible actions:

*Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]*

$Q(1,5) = R(1,5) + 0.8 * Max[Q(5,1), Q(5,4), Q(5,5)] = 100 + 0.8 * 0 = 100$

$$Q = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array}\right] \end{array}$$

Action

State

$$R = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array}\right] \end{array}$$

# Q – Learning Example

For the next episode, we start with a randomly chosen initial state, i.e. state 3
- From room 3 you can either go to room 1,2 or 4, let's select room 1.
- From room 1, calculate maximum Q value for this next state based on all possible actions:

Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

Q(3,1) = R(3,1) + 0.8 * Max[Q(1,3), Q(1,5)]= 0 + 0.8 * [0, 100] = 80
The matrix Q get's updated

$$Q = \begin{array}{c} \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} 0 & 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 100 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 80 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{array}\right] \end{array}$$

$$R = \begin{array}{c} State \\ \\ 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{array}{cccccc} & & Action & & & \\ 0 & 1 & 2 & 3 & 4 & 5 \\ \left[\begin{array}{cccccc} -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & -1 & 0 & -1 & 100 \\ -1 & -1 & -1 & 0 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & -1 \\ 0 & -1 & -1 & 0 & -1 & 100 \\ -1 & 0 & -1 & -1 & 0 & 100 \end{array}\right] \end{array}$$

# Q – Learning Example

For the next episode, the next state, 1, now becomes the current state. We repeat the inner loop of the Q learning algorithm because state 1 is not the goal state.
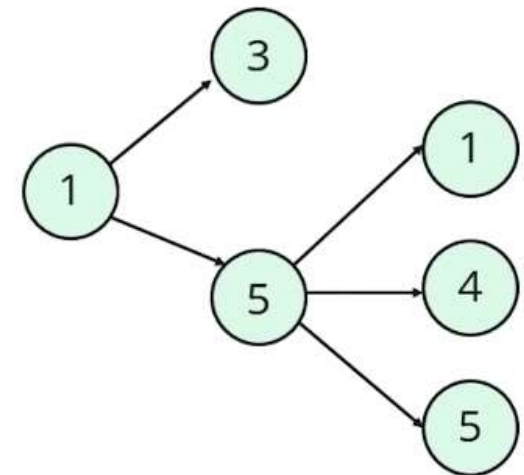- From room 1 you can either go to room 3 or 5, let's select room 5.
- From room 5, calculate maximum Q value for this next state based on all possible actions:

Q(state, action) = R(state, action) + Gamma * Max[Q(next state, all actions)]

Q(1,5) = R(1,5) + 0.8 * Max[Q(5,1), Q(5,4), Q(5,5)] = 100 + 0.8 * 0 = 100
The matrix Q remains the same since, Q(1,5) is already fed to the agent

$$Q = \begin{array}{c|cccccc} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 100 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 3 & 0 & 80 & 0 & 0 & 0 & 0 \\ 4 & 0 & 0 & 0 & 0 & 0 & 0 \\ 5 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

$$R = \begin{array}{c|cccccc} \text{State} & 0 & 1 & 2 & 3 & 4 & 5 \\ \hline 0 & -1 & -1 & -1 & -1 & 0 & -1 \\ 1 & -1 & -1 & -1 & 0 & -1 & 100 \\ 2 & -1 & -1 & -1 & 0 & -1 & -1 \\ 3 & -1 & 0 & 0 & -1 & 0 & -1 \\ 4 & 0 & -1 & -1 & 0 & -1 & 100 \\ 5 & -1 & 0 & -1 & -1 & 0 & 100 \end{array}$$

Action

```python
import numpy as np

# R matrix
R = np.matrix([[-1, -1, -1, -1, 0, -1],
               [-1, -1, -1, 0, -1, 100],
               [-1, -1, -1, 0, -1, -1],
               [-1, 0, 0, -1, 0, -1],
               [-1, 0, 0, -1, -1, 100],
               [-1, 0, -1, -1, 0, 100]])

# Q matrix
Q = np.matrix(np.zeros([6, 6]))

# Gamma (learning parameter).
gamma = 0.8

# Initial state. (Usually to be chosen at random)
initial_state = 1
```

```python
16
17     # Initial state. (Usually to be chosen at random)
18     initial_state = 1
19
20
21     # This function returns all available actions in the state given as an argument
22     def available_actions(state):
23         current_state_row = R[state,]
24         av_act = np.where(current_state_row >= 0)[1]
25         return av_act
26
27
28     # Get available actions in the current state
29     available_act = available_actions(initial_state)
30
31
32     # This function chooses at random which action to be performed within the range
33     # of all the available actions.
34     def sample_next_action(available_actions_range):
35         next_action = int(np.random.choice(available_act, 1))
```

```python
# Get available actions in the current state
available_act = available_actions(initial_state)



# This function chooses at random which action to be performed within the range
# of all the available actions.
def sample_next_action(available_actions_range):
    next_action = int(np.random.choice(available_act, 1))
    return next_action



# Sample next action to be performed
action = sample_next_action(available_act)



# This function updates the Q matrix according to the path selected and the Q
# learning algorithm
def update(current_state, action, gamma):
    max_index = np.where(Q[action,] == np.max(Q[action,]))[1]
```

Test  >  demo.py                                                                                           demo ▼   ▶   🐞   ■   Q

demo.py

```python
37
38
39      # Sample next action to be performed
40      action = sample_next_action(available_act)
41
42
43      # This function updates the Q matrix according to the path selected and the Q
44      # learning algorithm
45      def update(current_state, action, gamma):
46          max_index = np.where(Q[action,] == np.max(Q[action,]))[1]
47
48          if max_index.shape[0] > 1:
49              max_index = int(np.random.choice(max_index, size=1))
50          else:
51              max_index = int(max_index)
52          max_value = Q[action, max_index]
53
54          # Q learning formula
55          Q[current_state, action] = R[current_state, action] + gamma * max_value
56
```

update()

Run:   demo ×

▶ 4: Run     6: TODO     Terminal     Python Console                          Event Log

55:11   CRLF ÷   UTF-8 ÷   4 spaces ÷

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

Test ) demo.py

demo.py

```python
51          max_index = int(max_index)
52      max_value = Q[action, max_index]
53
54      # Q learning formula
55      Q[current_state, action] = R[current_state, action] + gamma * max_value
56
57
58  # Update Q matrix
59  update(initial_state, action, gamma)
60
61  # ------------------------------------------------------------------
62  # Training
63
64  # Train over 10 000 iterations. (Re-iterate the process above).
65  for i in range(10000):
66      current_state = np.random.randint(0, int(Q.shape[0]))
67      available_act = available_actions(current_state)
68      action = sample_next_action(available_act)
69      update(current_state, action, gamma)
70
```

for i in range(10000)

Run:  demo

4: Run    6: TODO    Terminal    Python Console                    Event Log

65:23    CRLF    UTF-8    4 spaces

```python
        update(current_state, action, gamma)


# Normalize the "trained" Q matrix
print("Trained Q matrix:")
print(Q / np.max(Q) * 100)


# ------------------------------------------------------------------
# Testing


# Goal state = 5
# Best sequence path starting from 2 -> 2, 3, 1, 5

current_state = 1
steps = [current_state]


while current_state != 5:

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

    if next_step_index.shape[0] > 1:
```

Test ) demo.py

demo.py

```python
    current_state = 1
    steps = [current_state]


while current_state != 5:

    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 selected sequence of steps
print("Selected path:")
print(steps)
```

Run:   demo

▶ 4: Run      6: TODO      Terminal      Python Console                                    Event Log

96:35   CRLF   UTF-8   4 spaces

File   Edit   View   Navigate   Code   Refactor   Run   Tools   VCS   Window   Help

Test  >  demo.py                                                    demo  ▼

demo.py

```
79     # Best sequence path starting from 2 -> 2, 3, 1, 5
80
81     current_state = 1
82     steps = [current_state]
83
84     while current_state != 5:
85
```

Run:     demo

```
C:\Users\zulaikha\PycharmProjects\Test\venv\Scripts\python.exe C:/Users/zulaikha/PycharmProjects/Test/demo
Trained Q matrix:
[[   0.    0.    0.    0.   80.    0. ]
 [   0.    0.    0.   64.    0.  100. ]
 [   0.    0.    0.   64.    0.    0. ]
 [   0.   80.   51.2   0.   80.    0. ]
 [   0.   80.   51.2   0.    0.  100. ]
 [   0.   80.    0.    0.   80.  100. ]]
Selected path:
[1, 5]

Process finished with exit code 0
```

▶ 4: Run     ⁞ 6: TODO     Terminal     Python Console                    Event Log

10:7   CRLF   UTF-8   4 spaces

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

Test  demo.py                                                          demo ▾  ▶  ⚙  ■  Q

demo.py

```python
# Best sequence path starting from 2 -> 2, 3, 1, 5

current_state = 2
steps = [current_state]

while current_state != 5:

    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 selected sequence of steps
```

while current_state != 5

Run:  demo

```
[  0.    80.    51.2    0.    80.     0.  ]
[  0.    80.    51.2    0.     0.   100.  ]
[  0.    80.     0.     0.    80.   100.  ]]
```

▶ 4: Run    ≡ 6: TODO    ⊞ Terminal    ⬥ Python Console                    ○ Event Log

86:34    CRLF ÷    UTF-8 ÷    4 spaces ÷

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

Test  demo.py

demo.py

```python
79      # Best sequence path starting from 2 -> 2, 3, 1, 5
80
81      current_state = 2
82      steps = [current_state]
83
84      while current_state != 5:
85
86          next_step_index = np.where(Q[current_state,] == np.max(Q[current_state,]))[1]
87
88          if next_step_index.shape[0] > 1:
89              next_step_index = int(np.random.choice(next_step_index, size=1))
```

while current_state != 5

Run:  demo

```
Trained Q matrix:
[[  0.     0.     0.     0.    80.     0.  ]
 [  0.     0.     0.    64.     0.   100.  ]
 [  0.     0.     0.    64.     0.     0.  ]
 [  0.    80.    51.2    0.    80.     0.  ]
 [  0.    80.    51.2    0.     0.   100.  ]
 [  0.    80.     0.     0.    80.   100.  ]]
Selected path:
[2, 3, 4, 5]
```

▶ 4: Run    6: TODO    Terminal    Python Console                              Event Log

10:13    CRLF    UTF-8    4 spaces