# Week 3

# Reinforcement Learning

## Why use Reinforcement Learning?
- Supervised learning is not a good approach → Very difficult to get use a dataset to get an ideal action → Many ambiguities/nuances on what is the right action to take
- You don't need to know the exact actions to take → Robot will automatically learn based on the reward function



## Overview of how Reinforcement Learning Works
- Robot will be aware of:
  **(state, action, reward, new state)**

# Mars Rover Example



terminal state      ↓      terminal state

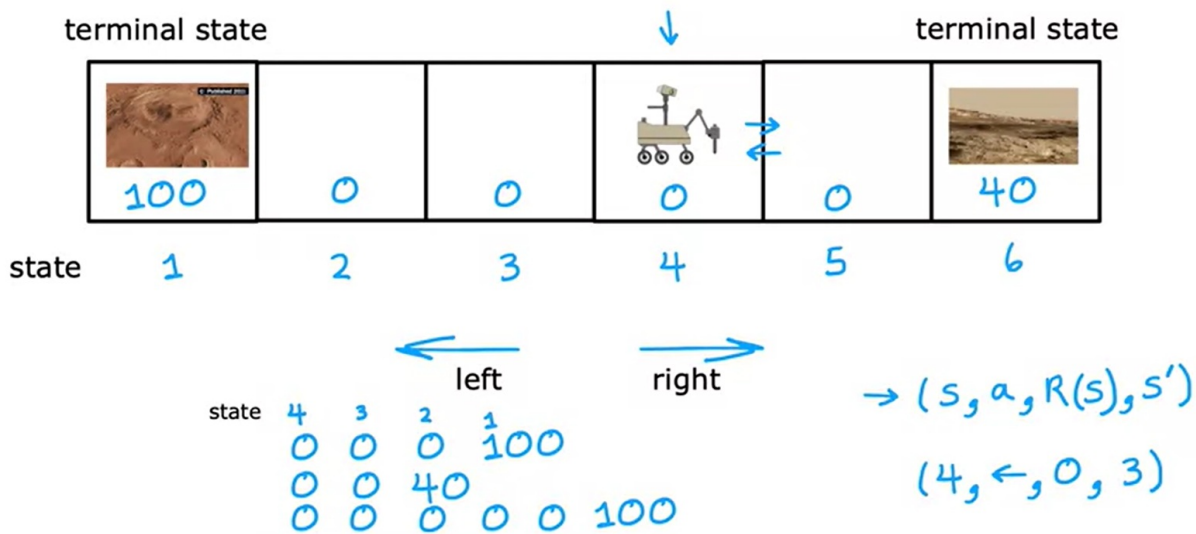| 100 | 0 | 0 | 0 | 0 | 40 |
|------|------|------|------|------|------|

state   1    2    3    4    5    6

← left    right →

```
state   4    3    2    1
        0    0    0   100
        0    0   40
        0    0    0    0    0   100
```

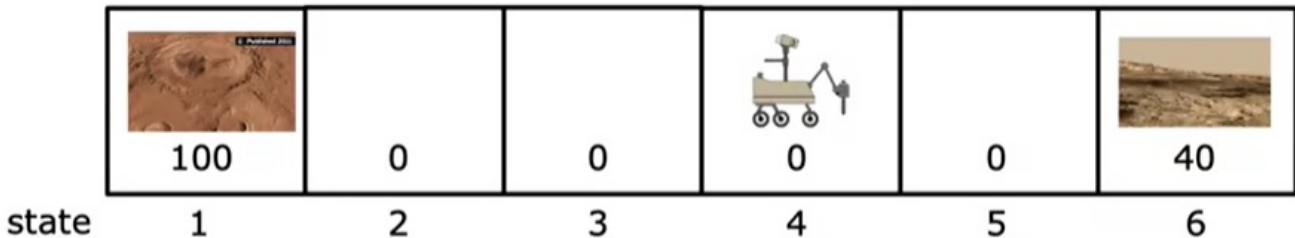$$\rightarrow (s, a, R(s), s')$$

$$(4, \leftarrow, 0, 3)$$

[Credit: Jagriti Agrawal, Emma Brunskill]

## Return

- Idea: Penalize the additional states taken to reach the terminal state by adding a discount factor $(\gamma)(\gamma)(\gamma)$
- high $\gamma$ = agent seeks long-term rewards

# Return

| 100 | 0 | 0 | 0 | 0 | 40 |
|------|------|------|------|------|------|

state   1    2    3    4    5    6

$$\text{Return} = 0 + (0.9)0 + (0.9)^2 0 + (0.9)^3 100 = 0.729 \times 100 = 72.9$$

$$\text{Return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots \quad \text{(until terminal state)}$$

$$\text{Discount Factor} \quad \gamma = 0.9 \quad 0.99 \quad 0.999$$

$$\gamma = 0.5$$

$$\text{Return} = 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100 = 12.5$$

## Example

- Reward at each state (depending on the action) if discount factor γγγ **is 0.5**
- Keep in mind that a state can have **negative** reward

# Example of Return

| 100 | 50 ← | 25 ← | 12.5 ← | 6.25 ← | 40 |
|-----|------|------|--------|--------|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |
| 1 | 2 | 3 | 4 | 5 | 6 |

← return

← reward

$\gamma = 0.5$

The return depends on the actions you take.

| 100 | 2.5 → | 5 → | 10 → | 20 → | 40 |
|-----|-------|-----|------|------|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |
| 1 | 2 | 3 | 4 | 5 | 6 |

$0 + (0.5)0 + (0.5)^2 40 = 10$

| 100 | 50 ← | 25 ← | 12.5 ← | 20 → | 40 |
|-----|------|------|--------|------|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |
| 1 | 2 | 3 | 4 | 5 | 6 |

$0 + (0.5)40 = 20$

## Goal of Reinforcement Learning - Policy
- Create a policy that maps a state to an action

# The goal of reinforcement learning

| | | | | | |
|---|---|---|---|---|---|
| 100 | | | 🤖 | | 40 |

Find a policy $\pi$ that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

## State-action Value Function

### Q-learning function
- $Q(s,a)$ represents the **maximum potential reward** possible if a particular action $a$ is taken at state $s$, assuming that steps taken after are **optimal** (which depends on the policy).

# Picking actions

| 100 | 50 | 25 | 12.5 | 20 | 40 | ← return |
|-----|-----|-----|------|-----|-----|----------|
| 100 | 0 | 0 | 0 | 0 | 40 | ← action |
|     |    |   |      |    |    | ← reward |

| 100 100 | 50 12.5 | 25 6.25 | 12.5 10 | 6.25 20 40 | 40 |
|---------|---------|---------|---------|------------|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |
| 1 | 2 | 3 | 4 | 5 | 6 |

$$\max_a Q(s,a)$$
$$\Pi(s) = a$$

$Q(4, \leftarrow)$   $Q(4, \rightarrow)$
12.5            10

$Q(s,a)$ = Return if you
- start in state $s$.
- take action $a$ (once).
- then behave optimally after that.

The best possible return from state $s$ is $\max\limits_a Q(s,a)$.

The best possible action in state $s$ is the action $a$ that gives $\max\limits_a Q(s,a)$.

## Bellman Equation

Notations:

$$Q(s,a) = \text{Return if you}$$
- start in state $s$.
- take action $a$ (once).
- then behave optimally after that.

$R(1) = 100 \quad R(2) = 0$

$s$ : current state
$a$ : current action
$s'$ : state you get to after taking action $a$
$a'$ : action that you take in state $s'$

$R(s)$ = reward of current state

1   2   3

- Bellman equation is a recursive function
- maxaQ(s',a')max_aQ(s',a') maxaQ(s',a') represents the return from optimal behavior from s's's'

# Explanation of Bellman Equation

$Q(s, a) =$ Return if you
- start in state $s$.
- take action $a$ (once).
- then behave optimally after that.

$s \rightarrow s'$

→ The best possible return from state $s'$ is $\max\limits_{a'} Q(s', a')$

$$Q(s, a) = R(s) + \gamma \max\limits_{a'} Q(s', a')$$

Reward you get right away

Return from behaving optimally starting from state $s'$.

$$R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \cdots$$
$$Q(s, a) = R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \cdots]$$

## Random Stochastic Environment
- Misstep probability → The next state can be different from the optimal based on probability
- Maximize average expected return

# Expected Return

Goal of Reinforcement Learning:

Choose a policy $\pi(s) = a$ that will tell us what action $a$ to take in state $s$ so as to maximize the expected return.

Bellman Equation:

$$Q(s, a) = R(s) + \gamma E[\max\limits_{a'} Q(s', a')]$$

3 ←

2 or 4

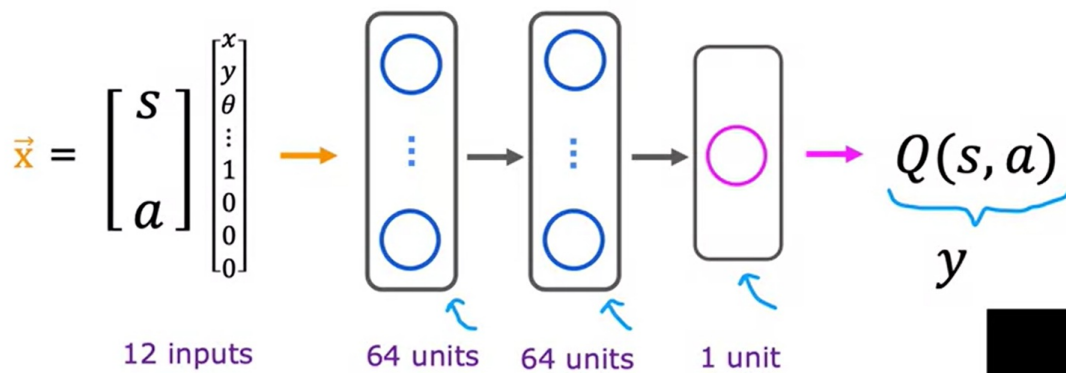# Continuous State Spaces

## Discrete vs Continuous states

- Instead of simply moving from state 1 - 6, think about helicopters moving x, y z, rotation, which is continuous
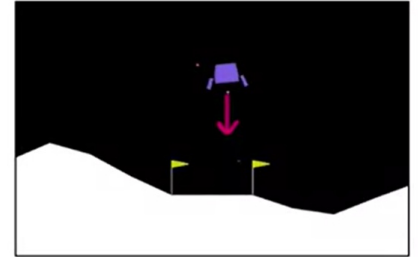
## Deep Reinforcement Learning
- In the example below, the 12 inputs come from:
  - 8 states
  - 4 q-learning functions (1 if the function is used, 0 otherwise) → actions

- How do we train a neural network to output $Q(s,a) Q(s,a) Q(s,a)$? Use Bellman Equation to create a training set of x and y.

# Deep Reinforcement Learning

$$\vec{x} = \begin{bmatrix} s \\ a \end{bmatrix} \begin{bmatrix} x \\ y \\ \theta \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{[64 units]} \rightarrow \text{[64 units]} \rightarrow \text{[1 unit]} \rightarrow Q(s,a)$$

12 inputs     64 units   64 units    1 unit

$Q(s,a)$ is $y$

In a state $s$, use neural network to compute
$Q(s, \text{nothing})$, $Q(s, \text{left})$, $Q(s, \text{main})$ $Q(s, \text{right})$
Pick the action $a$ that maximizes $Q(s,a)$

## Creating training set for Reinforcement Learning
- Note that s'1s'^1s'1 is not fixed → Q is from a random guess

# Bellman Equation

$$Q(s,a) = \underbrace{R(s)}_{} + \underbrace{\gamma \max_{a'} Q(s',a')}_{y}$$

$\underbrace{Q(s,a)}_{x}$

$f_{w,B}(x) \approx y$

$(s, a, R(s), s')$

$(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)}) \leftarrow$

$(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)}) \leftarrow$

$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)}) \leftarrow$

$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$

$y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$

| $x$ | $y$ |
|---|---|
| $x^{(1)} = (s^{(1)}, a^{(1)})$ | $y^{(1)}$ |
| $x^{(2)} = (s^{(2)}, a^{(2)})$ | $y^{(2)}$ |
| $x^{(10,000)}$ | $y^{(10,000)}$ |

## Learning Algorithm
- As you run this algorithm for many iterations, the accuracy will improve

# Learning Algorithm

Initialize neural network randomly as guess of $Q(s,a)$.

Repeat {

Take actions in the lunar lander. Get $(s, a, R(s), s')$.

Store 10,000 most recent $(s, a, R(s), s')$ tuples.

Replay Buffer

Train neural network:

Create training set of 10,000 examples using

$$x = (s,a) \text{ and } y = R(s) + \gamma \max_{a'} Q(s', a')$$

Train $Q_{new}$ such that $Q_{new}(s,a) \approx y$.
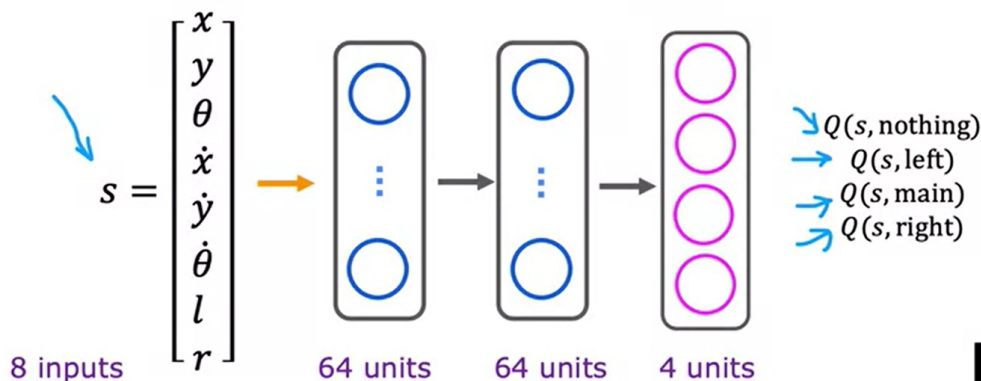
Set $Q = Q_{new}$.

$$f_{W,B}(x) \approx y$$

$x, y$

$x^{(1)}, y^{(1)}$
:
$x^{10000}, y^{10000}$

## Algorithm Refinement

- Inefficient → The original requires you to carry out inference 4 times in the neural network for some state
- Train neural network to output all 4 values simultaneously
  - 8 inputs only (for each state)
  - Given a state s, we only need to run the inference once!
  - Bellman Equation is also more efficient due to the 4 Q functions outputted in each inference

# Deep Reinforcement Learning

$$s = \begin{bmatrix} x \\ y \\ \theta \\ \dot{x} \\ \dot{y} \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

8 inputs    64 units    64 units    4 units

$Q(s, \text{nothing})$
$Q(s, \text{left})$
$Q(s, \text{main})$
$Q(s, \text{right})$

In a state $s$, input $s$ to neural network.

Pick the action a that maximizes $Q(s,a)$.    $R(s) + \gamma \max_{a'} Q(s', a')$

## How to choose actions while still learning?

- Option 1 can be bad as it can be influenced easily by the random Q()
  - E.g Q(s,main) generated could be low, and using option 1, main thrusters will never be used
- Hence an exploration step is needed to find an alternate Q() which can sometimes be a good option to take
- **ε-greedy policy (ε = 0.05)**

# How to choose actions while still learning?

In some state s

Option 1:

Pick the action $a$ that maximizes $Q(s, a)$.

Option 2:

→ With probability 0.95, pick the action a that maximizes $Q(s, a)$. Greedy, "Exploitation"

→ With probability 0.05, pick an action $a$ randomly. "Exploration"

$\varepsilon$-greedy policy $(\varepsilon = 0.05)$

0.95

*Q(s,main) is low*

↑

a

*Start $\varepsilon$ high*

1.0 → 0.01

*Gradually decrease*

---

**Limitations of Reinforcement Learning**

# Limitations of Reinforcement Learning

- Much easier to get to work in a simulation than a real robot!
- Far fewer applications than supervised and unsupervised learning.
- But … exciting research direction with potential for future applications.