position of helicopter ⟶ how to move control sticks

state s ⟶ action a

## Option 1: →NOT good

Use supervised learning

$x$ ⟶ $y$

↳ Cannot know ideal action

## Option 2:

Reinforcement learning → What to do NOT How

## Reward function:

flying well +1

Crash −1000

## Mars Rover Example:

Terminal state | | | | | Terminal state

+100 | | | | | +40
Interesting surface to analyze | O | O | 🚙 | O | Interesting surface to analyze

State 1  2  3  4  5  6

| state | 4 | 3 | 2 | 1 | state | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 100 | | 0 | 0 | 40 |

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

↳ A new state

| state | 4 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 0 | 100 |

## The Return of RL: → How to know if a particular reward is better than another.

Return = $0 + 0(0.9) + 0(0.9)^2 + (0.9)^3 100 = 72.9$

$= R_1 + \gamma R_2 + \gamma R_3 + \cdots$ until terminal state

Discount factor $\gamma = 0.9$

- This way, getting reward sooner would result in higher return
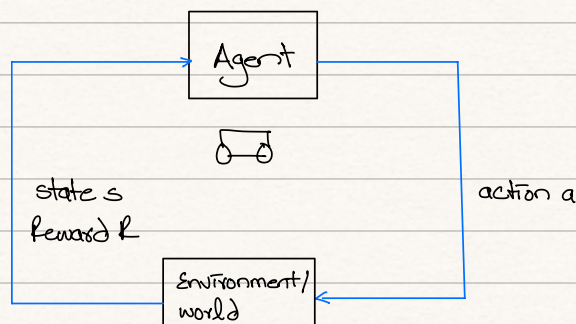
## Policy:

$$\text{state} \xrightarrow[\pi]{\text{Policy}} \text{action}$$
$$s \qquad\qquad a$$

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

## Markov Decision Process (MDP):

- states
- actions
- Rewards
- Discount factor $\gamma$
- Return
- policy $\pi$

- Future depends on where you are now NOT how you got here

Agent

state s
Reward R

action a

Environment/
world

## State action value function: (Q-function)

↳ How good is it?

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

$Q(2,←)$   $Q(2,→)$

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 100 | 100 | 50 | 12.5 | 25 | 6.25 | 12.5 | 10 | 6.25 | 20 | 40 | 40 |
| 100 | | 0 | | 0 | | 0 | | 0 | | 40 | |

State   1    2    3    4    5    6

$\gamma = 0.5$

| | | | | | |
|---|---|---|---|---|---|
| +100 Interesting surface to analyze | 50 ← | 25 ← | 12.5 ← | 20 → | +40 Interesting surface to analyze |
| | 0 | 0 | 0 | 0 | |

State   1    2    3    4    5    6

- The best possible return from state s is $\max\limits_{a} Q(s,a)$
- The best possible action in state s is the action that gives $\max Q(s,a)$

NOTE:

$Q = Q^* =$ Optimal Q function

## Bellman equation:

$s$ = current state

$R(s)$ = Reward of current state

$a$ = current action

$s'$ = state you get after action a

$a'$ = action you take in state s'

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

Reward you get right away

Return from behaving optimally starting from state $s'$

$Q(2,\leftarrow)$   $Q(2,\rightarrow)$

| 100 | 100 | 50 | 12.5 | 25 | 6.25 | 12.5 | 10 | 6.25 | 20 | 40 | 40 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 100 | | ◯ | | ◯ | | ◯ | | ◯ | | 40 | |

State    1       2        3       4       5       6

$\gamma = 0.5$

**Example 1:**

S = 2

a = →

s' = 3

**Example 2:**

S = 4

a = ←

s' = 3

$$Q(2,\rightarrow) = R(2) + 0.5 \max_{a'} Q(3,a')$$

$$= 0 + (0.5)\,25 = 12.5$$

$$Q(4,\leftarrow) = R(4) + 0.5 \max Q(3,a')$$

$$= 0 + 0.5\,(25)$$

$$= 12.5$$

For terminal state,

$$Q(s,a) = R(s)$$

• The best possible return from state $s'$ is $\max_{a'} Q(s,a)$

Random (stochastic) environment:

| ← | ← | ← | ← | → | |
|---|---|---|---|---|---|
| 100 | ◯ | ◯ | ◯ | ◯ | 40 |

State    1       2        3       4       5       6

| 4 | 3 | 2 | 1 |
|---|---|---|---|
| ◯ | ◯ | ◯ | 100 |

| 4 | 3 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|
| ◯ | ◯ | ◯ | ◯ | ◯ | 100 |

| 4 | 5 | 6 |
|---|---|---|
| ◯ | ◯ | 40 |

slips & goes to 4 instead

Expected return = Return = Average $\left( R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots \right)$

$$= \mathcal{E}\left[ R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots \right]$$

Bellman eq):

$$Q(s,a) = R(s) + \gamma \, \mathcal{E}\left[ \max_{a'} Q(s',a') \right]$$

## Continuous state spaces:

## For a truck:

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

→ How quickly is $x$ changing

→ " $y$

→ " $\theta$

## Lunar Lander:

$a =$ Do nothing

left thruster ←

Main thruster ↓

Right thruster →

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

$\left. \begin{array}{} \\ \end{array} \right\}$ left or Right

leg on ground (0 or 1)

## Reward:

Land : 100 - 140

Crash : −100

Soft landing : +100

leg grounded : +10
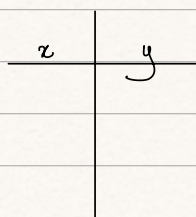
Fire main engine : −0.3

Fire side thrusters : −0.03

## Learning the state value function:

$$\vec{x} = \begin{bmatrix} s \\ a \end{bmatrix} \quad \begin{bmatrix} x \\ y \\ z \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \longrightarrow \cdots \longrightarrow \cdots \longrightarrow \bigcirc \longrightarrow \begin{array}{c} Q(s,a) \\ y \end{array}$$

$$\underset{64\ units}{\qquad} \underset{64\ units}{\qquad} \underset{\underline{1}\ unit}{\qquad}$$

In state $s$, use nueral network to compute $Q(s, nothing)$, $Q(s, left)$, $Q(s, Right)$, $Q(s, main)$. Pick the action $a$ that maximizes $Q(s,a)$.

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

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

| $x$ | $y$ |
|---|---|

{ create examples from simulation

$$\overbrace{\phantom{(s^{(1)}, a^{(1)})}}^{x} \overbrace{\phantom{R(s^{(1)}), s'^{(1)}}}^{y}$$
$$\left(s^{(1)}, a^{(1)}, R(s^{(1)}), s'^{(1)}\right)$$
$$\left(s^{(2)}, a^{(2)}, R(s^{(2)}), s'^{(2)}\right)$$
$$\vdots$$
$$\left(s^{(m)}, a^{(m)}, R(s^{(m)}), s'^{(m)}\right)$$

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

$$y^{(1)} = R(s^{(1)}) + \gamma \max_{a'} Q(s'^{(1)}, a')$$
$$y^{(2)} = R(s^{(2)}) + \gamma \max_{a'} Q(s'^{(2)}, a')$$

- Initialize nueral network randomly as guess of $\hat{Q}(s,a)$.
- Repeat {      $\uparrow$? $\Rightarrow \epsilon$ greedy policy

    Take actions in lunar lander, get $(s, a, R(s), s')$

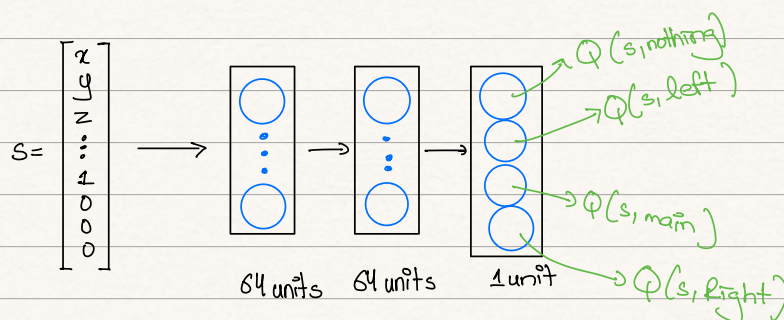    Store the $10,000$ more recent $(s, a, R(s), s')$ tuples

    $\hookrightarrow$ replay buffer

Train nueral network:

- Create training set of 10,000 examples using

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

- Train $Q_{new}$ such that $Q_{new}(s, a) \approx y$

- set $Q = Q_{new}$

}

## Improved NN:



$$S = \begin{bmatrix} x \\ y \\ z \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \rightarrow \rightarrow$$

64 units    64 units    1 unit

$Q(s, nothing)$
$Q(s, left)$
$Q(s, main)$
$Q(s, right)$

## Eplison greedy policy: How pick a while learning

## Option 1:

pick $a$, that maxs $Q(s, a)$

## Option 2:

with probability 0.95, pick a that maxs $Q(s, a)$   "Greedy, Exploitation"

with probability 0.05, pick a randomly  "Exploration"

- This is better because if $Q(s, main)$ set to low in step 1 then, it will never learn to fire main thruster.

↓

$\epsilon$ greedy policy ($\epsilon = 0.05$)
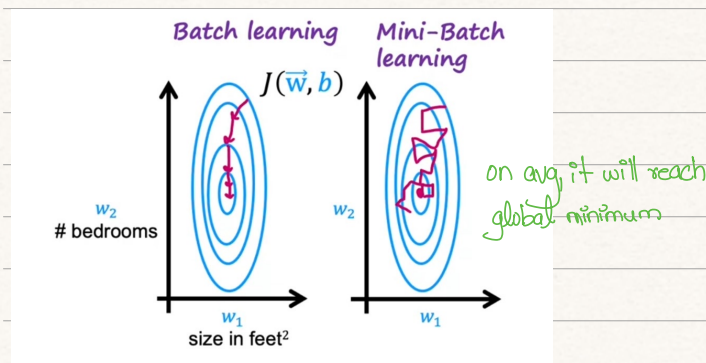
NOTE:
start $\epsilon$ high
Decrease gradually

## Mini-batches & soft updates:

If m is very large, we use a mini batch.

| $x$ | $y$ | |
|-----|-----|--|
| 2104 | 400 | } mini batch 1 |
| 1416 | 232 | |
| ⋮ | ⋮ | } mini batch 2 |
| 3210 | 870 | ⋮ |

• Every iteration look at a portion of dataset



on avg it will reach global minimum

**Soft updates.**

Set $Q = Q_{new}$

↑ ↑   ↑   ↑
W  B  Wnew  Bnew

$$W = 0.01 \, W_{new} + 0.99 \, W$$
$$B = 0.01 \, W_{new} + 0.99 \, B$$

} 99% Old W
} 1% new B

**Algorithm 1:** Deep Q-Learning with Experience Replay

1 Initialize memory buffer $D$ with capacity $N$
2 Initialize $Q$-Network with random weights $w$
3 Initialize target $\hat{Q}$-Network with weights $w^- = w$
4 **for** episode $i = 1$ **to** $M$ **do**
5     Receive initial observation state $S_1$
6     **for** $t = 1$ **to** $T$ **do**
7         Observe state $S_t$ and choose action $A_t$ using an $\epsilon$-greedy policy
8         Take action $A_t$ in the environment, receive reward $R_t$ and next state $S_{t+1}$
9         Store experience tuple $(S_t, A_t, R_t, S_{t+1})$ in memory buffer $D$
10        Every $C$ steps perform a learning update:
11        Sample random mini-batch of experience tuples $(S_j, A_j, R_j, S_{j+1})$ from $D$
12        Set $y_j = R_j$ if episode terminates at step $j + 1$, otherwise set $y_j = R_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a')$
13        Perform a gradient descent step on $(y_j - Q(s_j, a_j; w))^2$ with respect to the $Q$-Network weights $w$
14        Update the weights of the $\hat{Q}$-Network using a soft update
15     **end**
16 **end**