# Reinforcement Learning

Position of helicopter $\longrightarrow$ how to move control sticks

State s $\longrightarrow$ action a

x $\longrightarrow$ y

reward function $\longrightarrow$ Positive reward : +1

negative " : -1

what to do
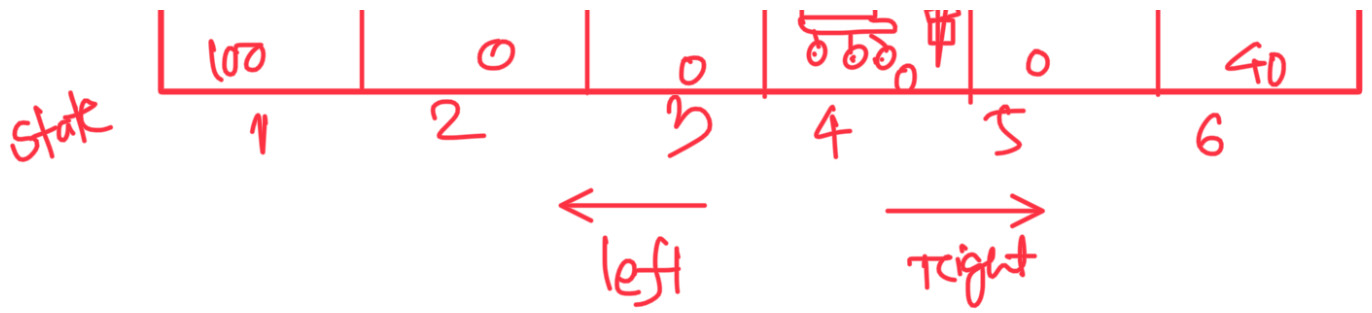
Controlling robots
factory optimization
financial trading
Playing games

terminal State

terminal State

state

| | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 100 | | 0 | | 0 | | 000,$\Psi$ | | 0 | | 40 | |

$\xleftarrow{\hspace{1cm}}$  $\xrightarrow{\hspace{1cm}}$
left        right

state  4  3  2  1
       0  0  0  100

       0  0  40

       0  0  0   0  0   100

$(S, a, R(\hat{s}), S')$

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

$\uparrow$
Reward associated with state 4

---

# Return in reinforcement learning

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

$$= R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

Discount factor $\gamma = 0.9$    0.99  0.999

$\gamma = 0.5$

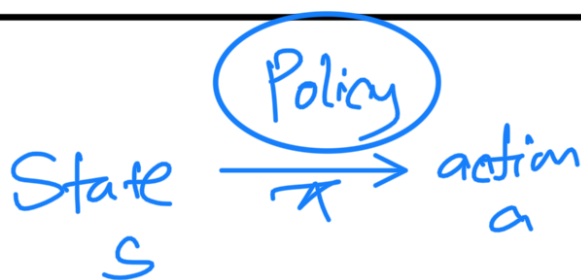$$\text{Returns} = 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 (00$$
$$= 12.5$$

| 100 | 50 | 25 | 12.5 | 625 | 40 |
|-----|-----|-----|------|-----|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |

$$\gamma = 0.5$$

| 100 | 2.5 | 5 | 10 | 20 | 40 |
|-----|-----|-----|-----|-----|-----|
| 100 | 0 | 0 | 0 | 0 | 40 |

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

Policy

State $\xrightarrow{\quad \pi \quad}$ action
$\quad S \qquad\qquad\qquad a$

Policy  Controller

$\pi(S) = a$

$\pi(2) = \Leftarrow$
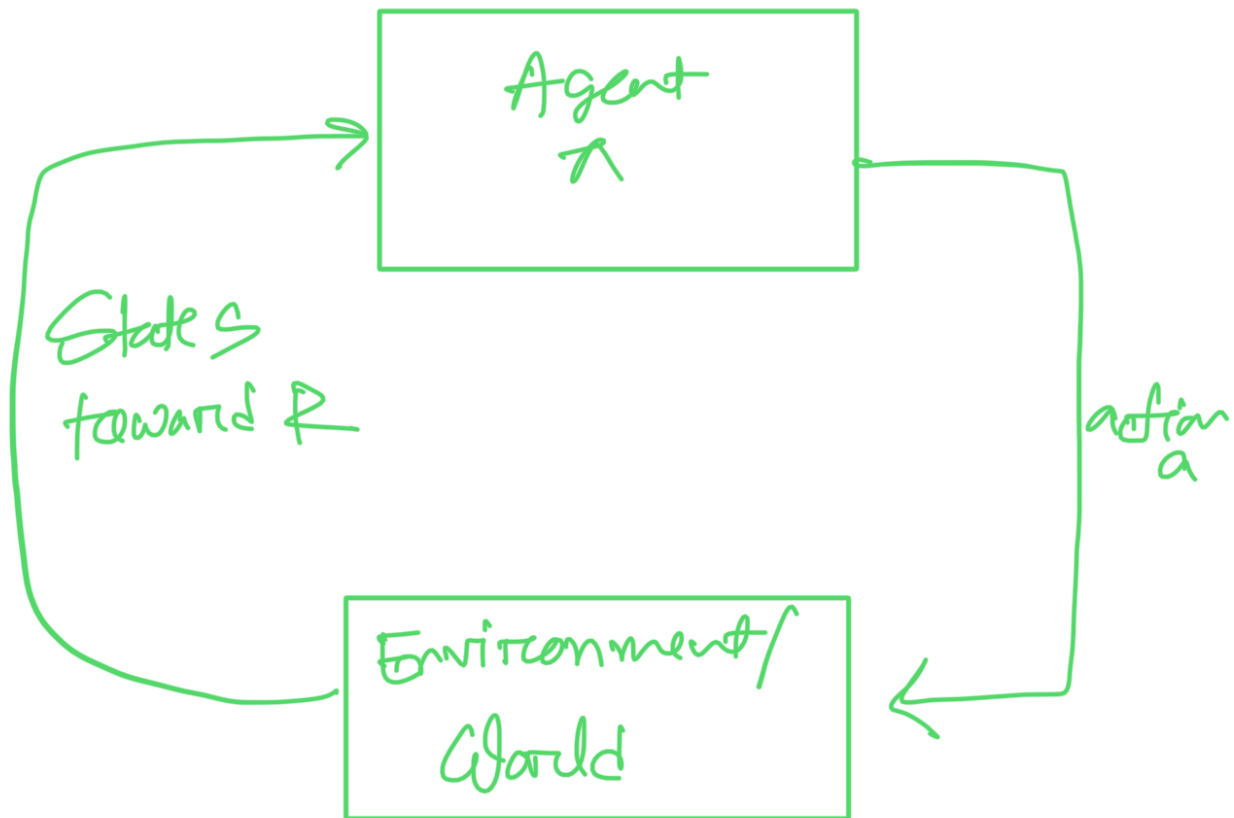
$\pi(3) = \Leftarrow$

$\pi(4) = \Leftarrow$

$\pi(5) = \Rightarrow$

# Markov Decision Process (MDP)

| | Mars rover | Helicopter | Chess |
|---|---|---|---|
| states | 6 states | position of helicopter | pieces on board |
| actions | ← → | how to move control stick | possible move |
| rewards | 100, 0, 40 | +1, −1000 | +1, 0, −1 |
| discount factor $\gamma$ | 0.5 | 0.99 | 0.995 |
| return | $R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots$ | $R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots$ | $R_1 + \gamma R_2 + \gamma^2 R_3 + \cdots$ |
| policy $\pi$ | [100] ← ← ← → [40] | Find $\pi(s) = a$ | Find $\pi(s) = a$ |

---

State action value function

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

$Q(2, \rightarrow) = 12.5$
$0 + 0(0.5) + (0.5)^2 0 + (0.5)^3 100$
$Q(2, \leftarrow) = 50$
$0 + (0.5) 100$

$$Q(4, \leftarrow) = 12.5$$

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

$$\max_{a} Q(s,a)$$

$$\pi(s) = a$$

---

# Bellman equation

$$Q(s,a) = \text{Return if you}$$
$$\Rightarrow \text{start in state } s$$
$$\Rightarrow \text{take action } a \text{ (once)}$$
$$\Rightarrow \text{Behave optimally}$$

s: current state

a: " action

s': state you get to after taking action a

a" : action you take in state s'

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

$$Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3,a')$$
$$= 0 + (0.5) \, 25 = 12.5$$

The best possible stream $s'$ is $\max_{a'} Q(s', a')$

$$R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

$$s \longrightarrow s'$$

---

Random (Stochastic) Environment

$$\text{Expected Return} = avg\left(R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots\right)$$

$$= E\left[R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots\right)$$

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

---

$$Q(s, \leftarrow) = 0 + (0.5) 0 + (0.5)^2 0 + (0.5)^4 100$$

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

$$= 0 + 0.25 0 + 0.25 0 + 0.25^3 40$$

---

Discrete vs Continuous State

$$\begin{bmatrix} x \end{bmatrix}$$   helicopter

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

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

## Lunar Lander

action:  do nothing
         left thruster
         main    "
         right   "

$$\begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix} \quad 0,1 \;\Big\{$$

pick action $a = \pi(s) \longrightarrow$ maximize return

$$\gamma = 0.985$$

$\longrightarrow \begin{bmatrix} x \\ \vdots \end{bmatrix}$

$$\vec{x} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

64    64    1 unit    $\Rightarrow Q(s,a)$    $y$

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

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

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

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

| $x$ | $y$ |
|---|---|
| $x^{(1)} = (s^{(1)}, a^{(1)})$ | $y^{(1)}$ |
| $x^{(2)} = (s^{(2)}, a^{(2)})$ | $y^{(2)}$ |
| $x^{10000}$ | $y^{10000}$ |

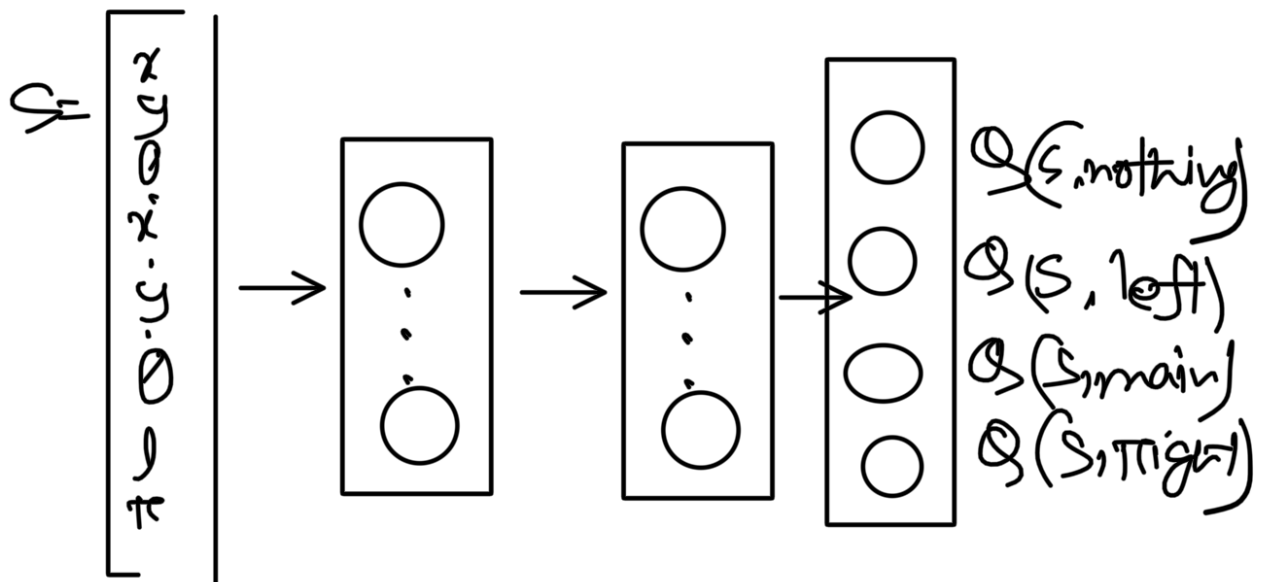$$y^{(i)} = R(s^{(i)}) + \gamma \max_{a'} Q(s'^{(i)}, a')$$

Learning algo

initialize NN randomly as guess $Q(s,a)$

Repeat {

   Take action in lunar lander, Get$(s, a, R(s), s')$

store 10k recent $(s, a, R(s), s')$ tuples

Train NN

   create training set 10k examples

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

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

Set $Q = Q_{new}$

# Deep Q Network

---

## Continous State Space



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

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

---

How to choose actions while still learning

---

In some states

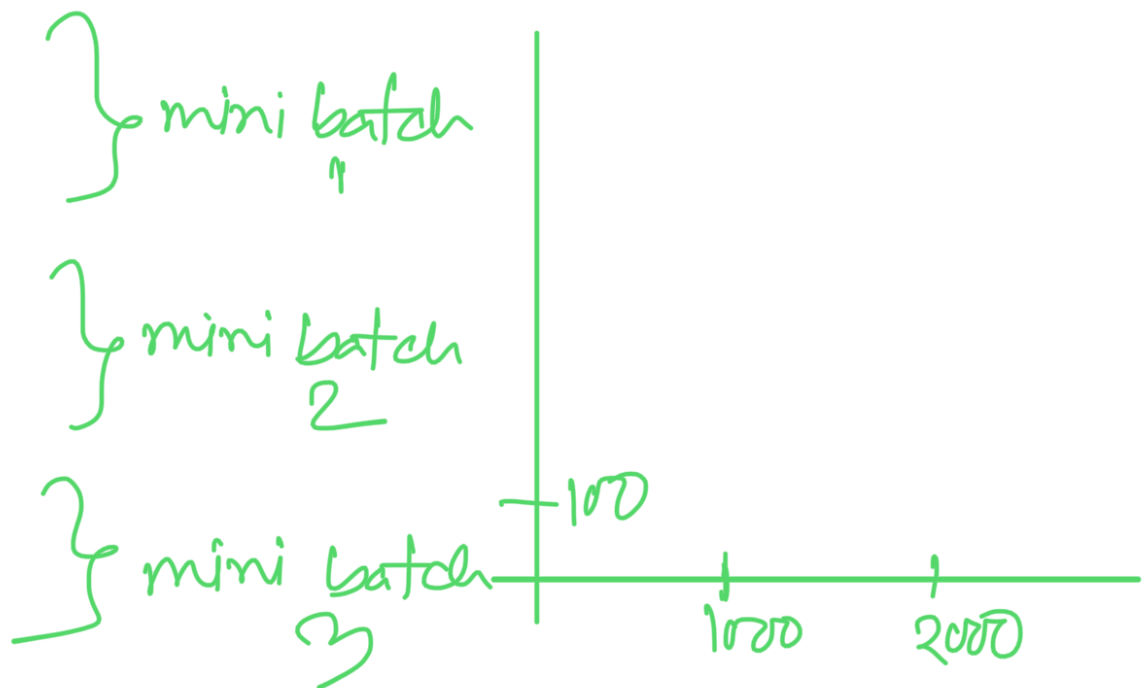Option 1

pick the action $a$ that maximizes $Q(s, a)$

Option 2

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

with prob 0.05, pick an action randomly
" exploration "

$\varepsilon$ — greedy policy

Start $\varepsilon$ high
$1.0 \rightarrow 0.01$
gradually decrease

---

## Algorithm refinement:
### Mini-batch & soft updates (optional)

$\}$ mini batch 1

$\}$ mini batch 2

$\}$ mini batch 3

100

1000    2000

Batch learning

1000 instead of 10k

Soft update

$Q = Q_{new}$

$\uparrow$

$W_{new}, B_{new}$

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

$B = 0.01 \, B_{new} + 0.99 \, B$

---

Reinforcement Learning