

# Reinforcement Learning

## Uses:

→ Autonomous helis → aerobatic manoeuvres  
→ drones

## Pro:

Exact: position of heli → how to move control sticks

Abstraction: states  $s$  → action  $a$

## Supervised Learning

$x \rightarrow y$

Issue: one correct action  $a$  is hard to determine

- No concept of delayed rewards
- Cannot learn from trial and error

Fig? RL is better suited to handle the complexity

RL → good dog, bad dog routine

what to do instead of how to do

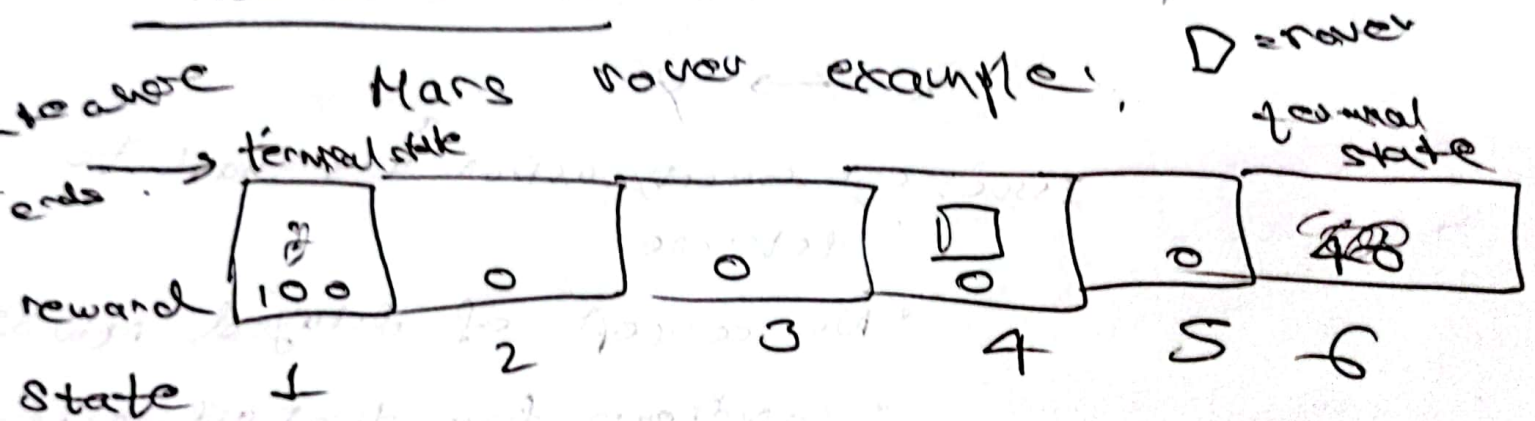
## reward function

pos reward: hell flying well +1  
neg reward: hell crashes -1000

## RL applications

- Controlling robots
- Factory optimization
- Financial (stock) trading
- Playing games (video games too).

## RL formalism:



Choice:  $\leftarrow$  left,  $\rightarrow$  right

possible paths

4 — 3 — 2 — 1

0 0 0 100

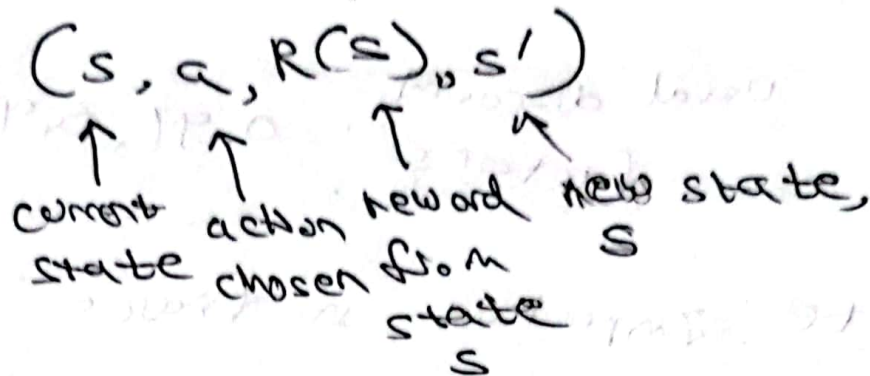
0 0 40

0 0 0 0

0 100

every  
At <sub>1</sub> time step:

robot is at state  $s$



Return in reinforcement learning:

~~Answer~~ Measure effort given for a reward.

In the example, if robot goes left for 100 reward we get:

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

discount factor

Abstract!

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

$\gamma = 0.9$



∴ Rewards obtainable earlier are more attractive

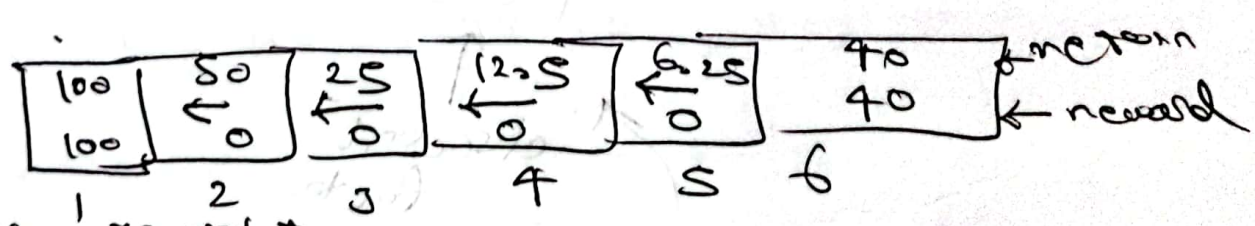
usual discount factors: 0.99, 0.999

after incorporating discount factor  
 $return = total\ reward$   
 $reward = for\ each\ state$

Note? Important in finance as a dollar today is more important than a dollar tomorrow (get e.g. car get interest if saved today)

• Reward varies when starting from diff states starting from each state and  $\gamma = 0.5$   
 Hence now r. Rewards  $r_t$

⇒ Always go left  
 $P = 0.5, P = 0.5$



⇒ Always go right



PO-8

⇒ Max <sup>return</sup> reward if we can go left or right at every state




→ There is a concept of  $-V$  rewards as well

↳ financial apps. where we want to delay payment

Policy: Function  $\pi$  which takes a state  $s$  and maps it to an action

$$\pi(s) = a$$

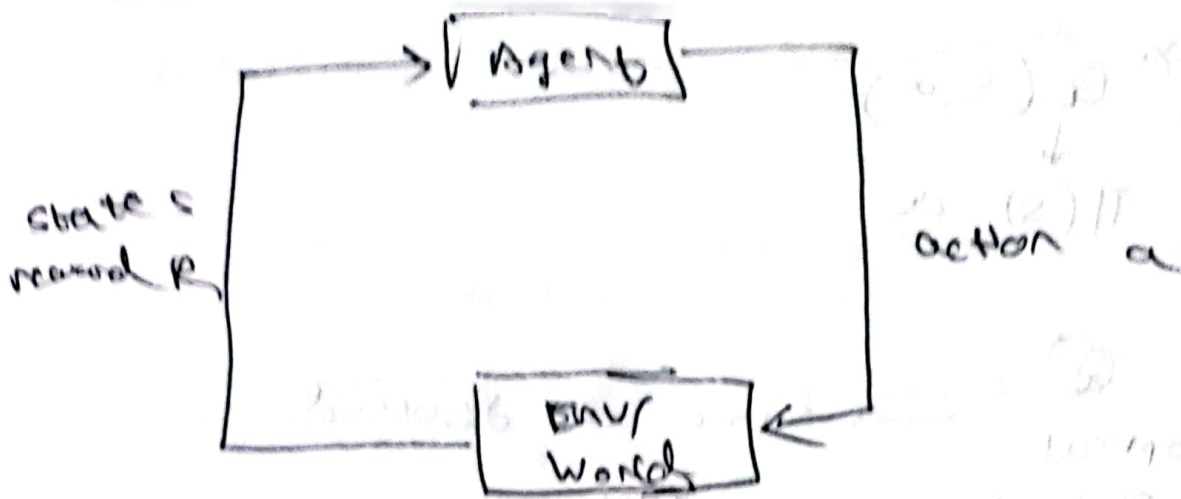
Goal of RL → Find  $\pi$  which for every state which for every state can tell us an action to maximize the return.

	Mous Rover	Helicopter	Ones
States	6	position of heli	pieces on board
actions	$\leftarrow \rightarrow$	how to move control stick	possible moves
rewards	100, 0, 40	+1, -1000	+1, -1, 0
discount factor $\gamma$	0.5	0.99	0.995
return	$R_1 + \gamma R_2 + \gamma^2 R_3 \dots$	$R_1 + \gamma R_2 + \gamma^2 R_3$	$R_1 + \gamma R_2 + \gamma^2 R_3$
policy $\pi$		Find $\pi(s) = a$	Find $\pi(s) = a$

Markov Decision Process (MDP)  $\pi(s) = a$

Future only depends on current state and not how we got to that state





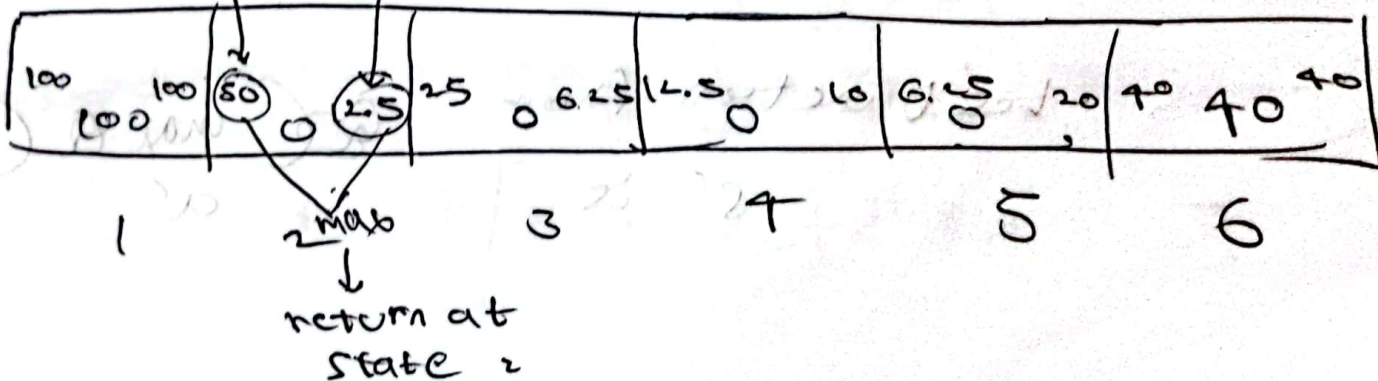
## State-action value function / Q-function

$Q(s, a)$  = return if you  
 • Start in state  $s$   
 • take action  $a$  (once)  
 • act optimally after that  
 state we are in      action we may take

We can come up with a  $Q$  before we

have even calculated  $\Pi$

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



$$\max_a Q(s, a) \downarrow \pi(s) = a$$

$Q^*$   
optimal  
Q function  $\rightarrow$  same as discussed  
here

Bellman equations:

$s$  = current state  $R(s)$  = reward of

$a$  = current state

$s'$  = state after taking  $a$

$a'$  = action that you take in  $s'$   
immediate reward  $\underbrace{\text{return for behaving optimally starting from } s'}$

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

at terminal states

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

let  $\gamma$  return from

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

$s'$  is

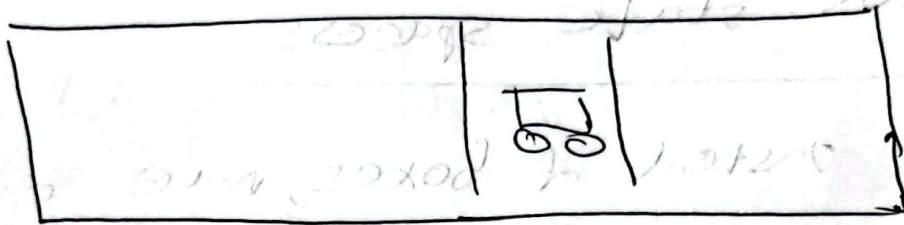


Random (stochastic) environment:

Randomness in env.  $\rightarrow$  wheels slipping, wheel sliding

$\downarrow$   
exact action or not taken

Stochastic environment:



$\xleftarrow{0.9}$   $\xrightarrow{0.1}$

tell to go left  
90% of  
going left and 10%.

Expected return.

• Return is random so we don't  
make a table

• Expected Return = Average  $(R_1 + \gamma R_2 + \gamma^2 R_3 + \dots)$

$$= E[R_1 + \gamma R_2 + \gamma^2 R_3 + \dots]$$

Modifies Bellman equations

$s'$  is random

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

# of states may be much larger

Continuous state spaces:

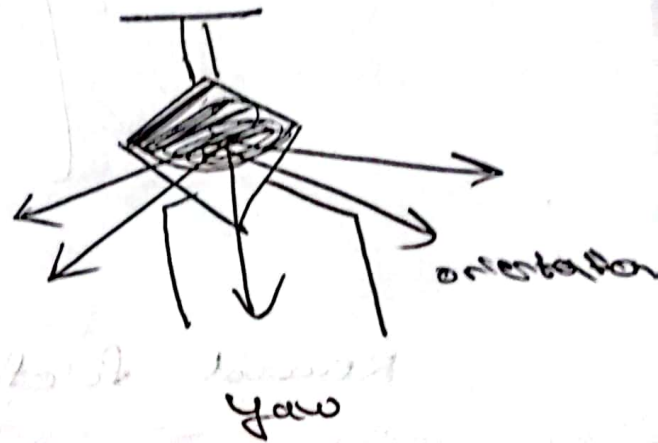
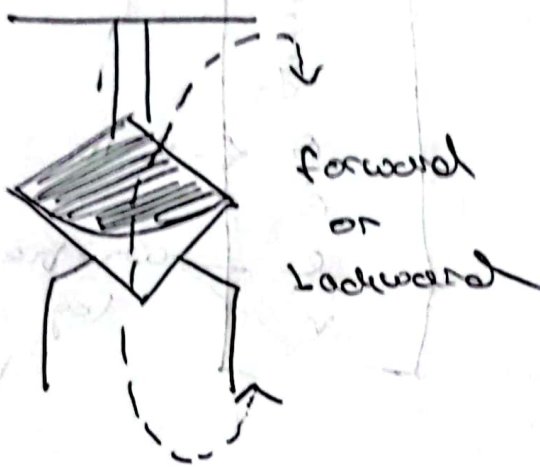
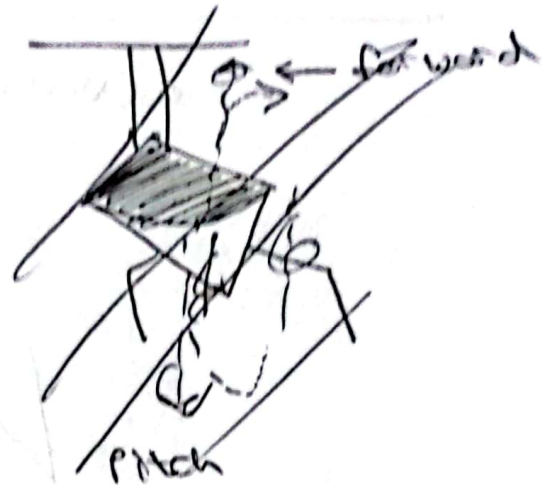
Instead of boxes, rover can be anywhere along a line

For a truck, its continuous state may be

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

velocity

For autonomous helicopters.



pitch  
pitch

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

position

orientation

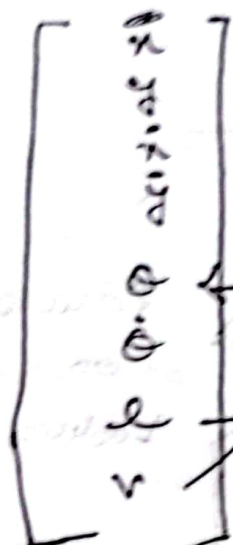
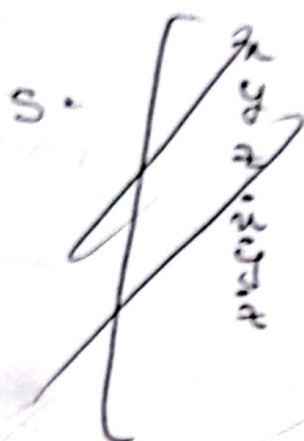
velocity

angular velocity



# Lunar Lander

actions: do nothing  
left thruster  
main thruster  
right thruster



tilt

whether left/right  
leg is sitting on  
padding or  
ground

## Reward function:

- getting to landing pad: 100 - 140
- Additional reward for moving toward/away from pad
- Crash: -100
- Soft landing: +100
- leg grounded: +10
- Fine engine management: -0.3
- Fine idle thrusters: -0.03

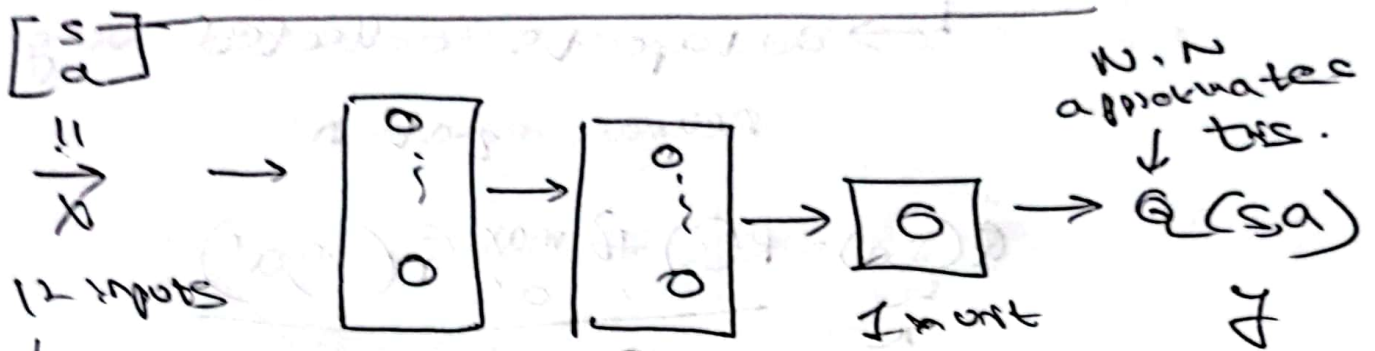
to not  
not to  
be as  
much  
as possible

Learn  $\pi$  given  $r$  &  $S$  so  $a = \pi(S)$  to maximize return

$$\gamma = 0.985$$

Learning the state-value function.

Deep reinforcement learning



$s$  from previous page  
 In a state  $s$ , use N.N to compute  $Q(s, \text{nothing})$ ,  $Q(s, \text{left})$ ,  $Q(s, \text{right})$ ,  $Q(s, \text{up})$ ,  $Q(s, \text{down})$   
 one-hot vector to rep. action taken

one-hot vector  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  (representing 'up')  
 Pick an action that maximizes  
 calculate these and choose a with largest  $Q$   
 $Q(s, a)$

The MN calculates state-action value function (~~max~~ ~~return~~ if we take action  $a$  at state  $s$  and behave optimally afterwards) for a given state-action pair  $[Q(s, a) \text{ from } s, a]$

↳ datapoints collected using Bellman equation

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

→ we try out different things in

Lunar Lander

→ If we don't have a good policy, we try random actions

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

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

$(s^2, a^2, R(s^2), s'^2)$

$(s^3, a^3, R(s^3), s'^3)$

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

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

PT = how to get this?



We can start with random  $Q$

It will get better over time

Learning algo:

Initialize  $Q$  randomly as guess of

Repeat  $[Q(s,a)]$

Take action in lower number. Get

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

Store 100 most recent  $(s, a, R(s), s')$

toples  $\leftarrow$  replay buffer

Train  $Q$ :

Create training set of 100 most recent examples

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

Train  $Q_{new}$  such that

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

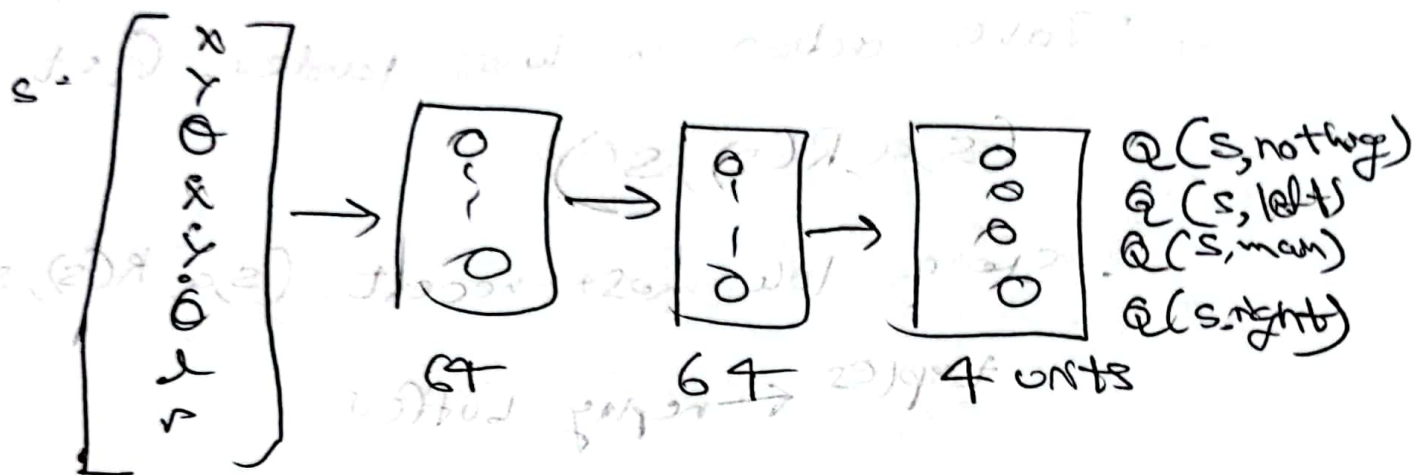
set  $Q := Q_{new}$

↑ and DQN algo.

## Refining the DQN algo!

We need to can do inferences 4 times to pick an action

Instead we change the architecture:



8 inputs

E-greedy policy

How to pick action to take with

learning  
Option 1s

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

Option 2s

→ with prob 0.95, pick a that maximizes  $Q(s, a)$

→ with prob 0.05, pick action a randomly → exploration

→ greedy/exploitation

eg. by chance N.N may break  $Q(s, a)$  is

low

→ randomness helps to break out of this.

Option-2 is  $\epsilon$ -greedy policy.

Trick: Start  $\epsilon$  high, then decrease gradually.

RL is more susceptible to hyperparameter tuning than supervised learning

MM-batches

• Same idea as <sup>seen</sup> before

Soft update

Set  $Q = Q_{new} \rightarrow Q_{new}$  may be worse

Sol<sup>n</sup>

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

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



## Limitations of RL:

• Much easier to get something working in a simulator than a real robot

• Fewer applications than supervised/unsupervised algos

$$WPP.O + w_{new} 10.0 = W$$

$$8 PPO + w_{new} 10.0 = 8$$