

g) Autonomous Helicopter

state $s \rightarrow$ action a

if supervised learning

x - if this state $\rightarrow y$

~~difficult to get states~~

Applications;

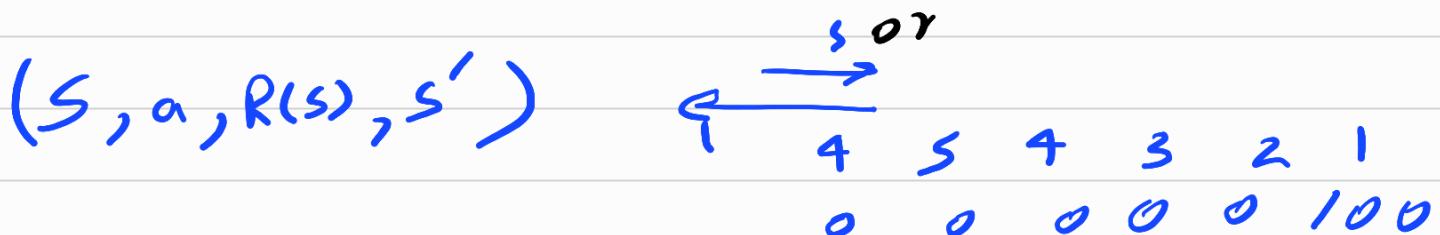
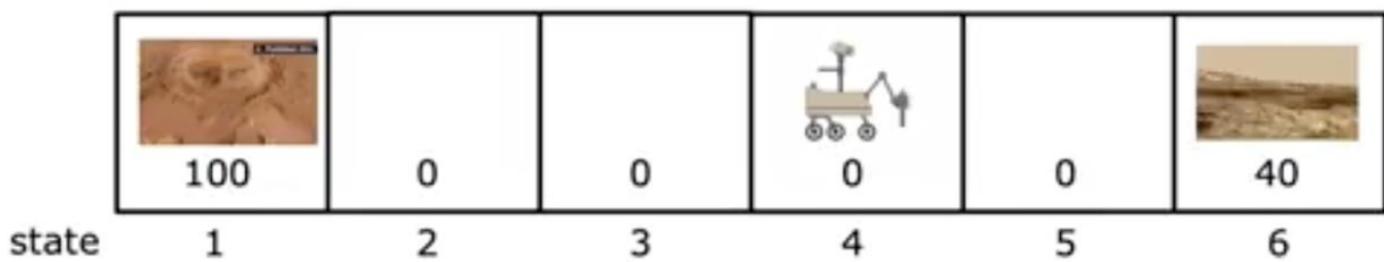
Controlling robots

Factory optimization

Financial (stock) trading

Games

e.g. - Mars Rover Example;



Return in Reinforcement Learning?

	0	0		0	
state 1	2	3	4	5	6

no dis. first state

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

discount factor

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

(mostly γ close to 1)

if always left

dis. factor = 0.5

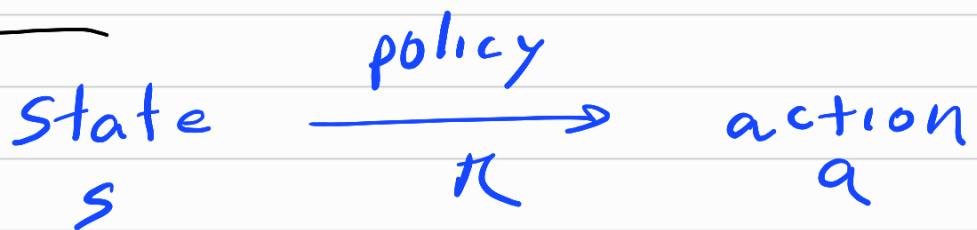
100	50	25	12.5	6.25	40
100	0	0	0	0	40

if always right $\frac{100}{2.5} \rightarrow \frac{5}{10} \rightarrow \frac{20}{40}$

Both sides considered

100	50	25	12.5	20	40
100	0	0	0	0	40

Policies



$$\begin{aligned}\kappa(2) &= \leftarrow \\ \kappa(3) &= \leftarrow \\ \kappa(4) &= \rightarrow\end{aligned}$$

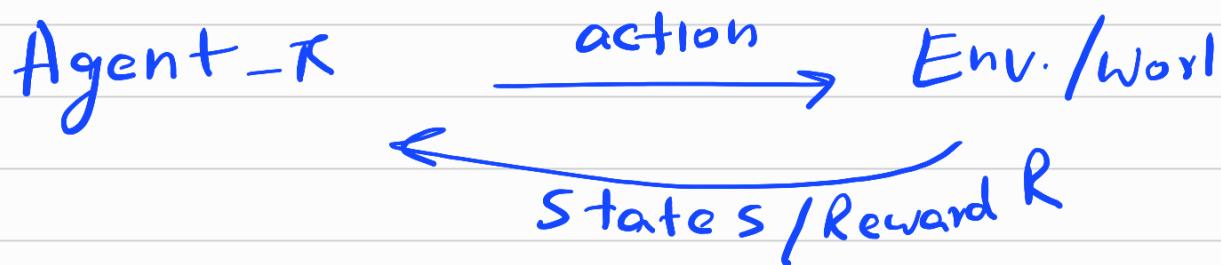
$$\left\{ \begin{array}{l} \kappa(s) \rightarrow a \end{array} \right.$$

find policy to maximize return,

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

Markov Decision Process (MDP)

Future depends on current state now
(not how you got here)



State Action Value function;
(Q-Function)

$Q(S, a) = \text{Return if Start state } s \text{ take action}(a) \text{ once optimal after that}$

100	50	25	12.5	20	40
100	0	0	0	0	40

1 2 3 4 5 6

$$\begin{aligned}
 Q(4, \leftarrow) &= 0 + \alpha^{0.5} + 0 + 100(0.5)^2 \\
 &= 12.5
 \end{aligned}$$

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

1 2 3 4 5 6

$$Q(4, \leftarrow) > Q(4, \rightarrow)$$

Bellman Equation

s = current state

a - current action

$R(s)$ = reward of current state

a' = action taken in s'

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

eg:- $\underbrace{\text{reward right away}}_{\text{a'}}$ $\underbrace{\text{best possible return from state}}_{\gamma = 0.5}$

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

$$Q(2, \rightarrow)$$

$$s = 2$$

$$a = \rightarrow$$

$$s' = 3$$

$$\textcircled{1} \quad Q(2, \rightarrow) = R(2) + 0.5 \max_{a'} Q(3, a')$$

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

$$= 12.5$$

$$\textcircled{2} Q(4, \leftarrow) = R(4) + 0.5 \max_{a'} Q(3, a')$$

$$= 0 + 0.5 \times (2s)$$

$$= 12.5$$

$$\text{eg!- } = R(4) + 0.5 \left[0 + (0.5)0 + \underbrace{0.5^2}_{10s} \right]$$

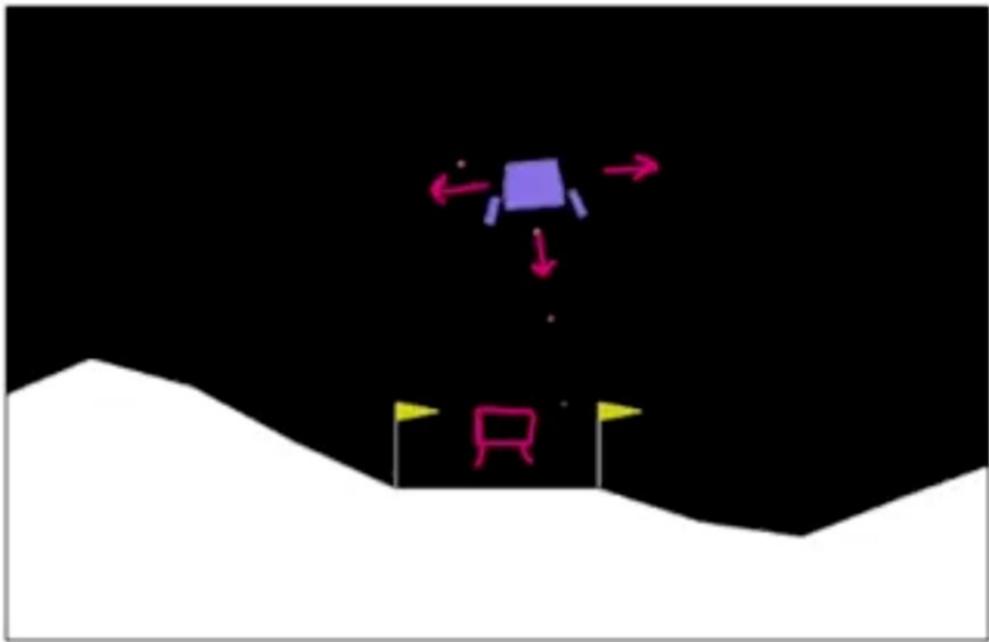
If Stochastic
 Expected \rightarrow sequence of rewards
 Return = $\text{Avg } (R_1 + \gamma R_2 + \gamma^2 R_3 + \dots)$
 $= E [R_1 + \gamma R_2 + \gamma^2 R_3 + \dots]$

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

2 & 4

for Continuous state space;

$$\text{eg!- } s = \begin{vmatrix} x \\ y \\ z \\ b \\ \theta \\ \omega \\ x_- \end{vmatrix} \downarrow$$



actions:

- do nothing
- left thruster
- main thruster
- right thruster

$$S = \begin{bmatrix} x \\ y \\ z \\ \dot{x} \\ \dot{y} \\ \dot{z} \\ \theta_x \\ \theta_y \\ \theta_z \end{bmatrix}$$

velocity

angle

Reward Function

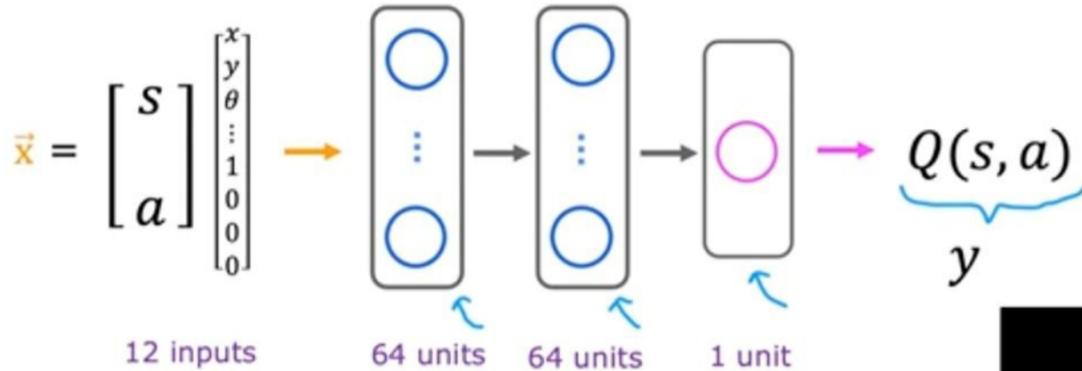
- Getting to landing pad: 100 - 140
- Additional reward for moving toward/away from pad.
- Crash: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03

Learn Policy π given s_j

Pick action $a \leftarrow \pi(s)$ to maximize return

$$\gamma = 0.985$$

Deep reinforcement Learning



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 $Q(s, a)$



Bellman; $Q(s, a) = R(s) + \gamma \max Q(s', a')$

$$f_{w, b}(x) \approx y \quad \text{using } \underbrace{\gamma}_{\text{supervised learning}}$$

try different things &

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

$$(s^{(2)}, \dots, \dots) \rightarrow y^{(2)}$$

Learning Algorithm

DQN
Deep Queue Network

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

Set $Q = Q_{\text{new}}$.

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

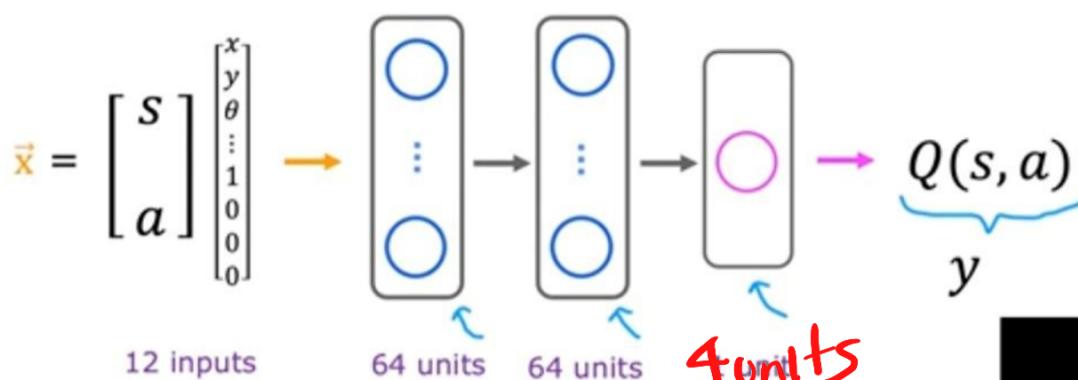
x, y

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

\vdots
 $x^{(10000)}, y^{(10000)}$

If separate 4 actions
 → have to run NN 4 separate times

Output all Q values at once
 only inference one time



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 $Q(s, a)$

to get all
 Q values one time



$Q(s, \text{nothing})$

$Q(s, \text{left})$

$Q(s, \text{right})$

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

ε -greedy policy ($\varepsilon = 0.05$)
 0.95

$Q(s, \text{main})$ is low



Start ε high

Gradually decrease

Mini batch and soft-updates;

making the examples

x	y		<u>learning</u>
2104	400	y mini batch 1	y much faster
1416	232	y mini batch 2	
--	--		
--	--		

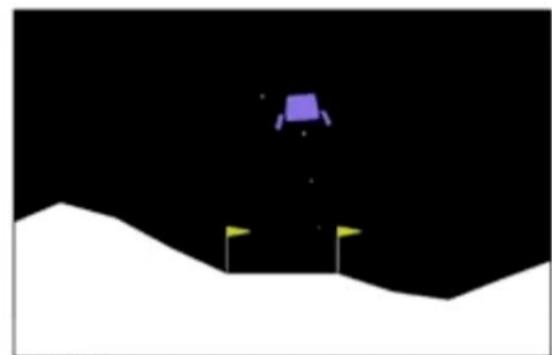
Learning Algorithm

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 model:

1,000

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

Set $Q = Q_{\text{new}}$.

$$\begin{aligned} &x^{(1)}, y^{(1)} \\ &\vdots \\ &x^{(1000)}, y^{(1000)} \end{aligned}$$

Soft Update \rightarrow Causes P-L. algorithm to converge more likely

Set $Q = Q_{\text{new}}$

$w_{\text{new}}, b_{\text{new}}$

$$w = 0.01 w_{\text{new}} + 0.99 w$$

$$b = 0.01 b_{\text{new}} + 0.99 b$$