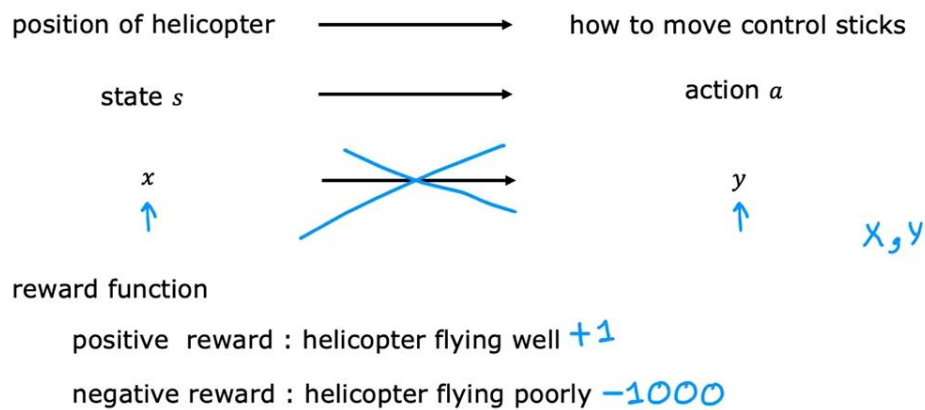


Reinforcement Learning



time flying well and
hopefully to never crash.

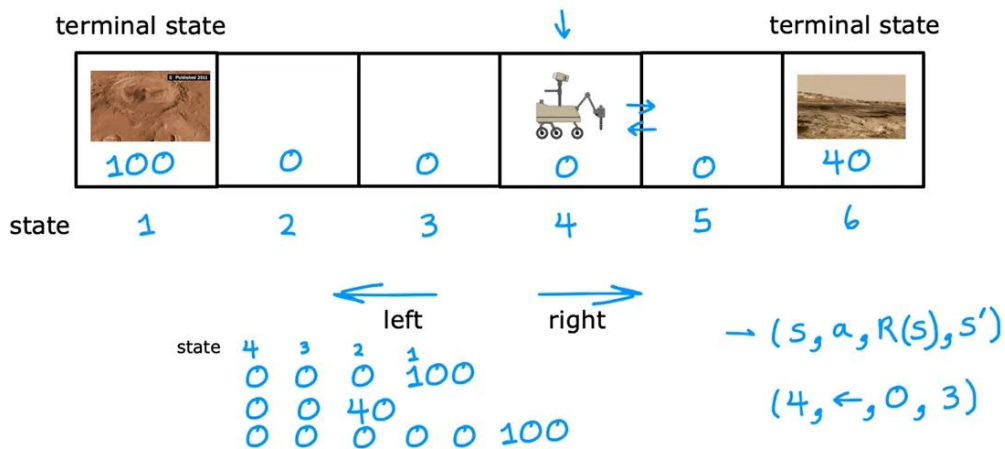
Applications

- • Controlling robots
- • Factory optimization
- • Financial (stock) trading
- Playing games (including video games)



For example, one of my friends was
working on efficient stock execution.

Mars Rover Example



[Credit: Jagriti Agrawal, Emma Brunskill]

you see that these four things,

DeepLearning.AI

Stanford ONLINE

Andrew Ng

Return



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

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

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

$$\gamma = 0.5$$

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

In financial applications,
the discount factor also has

DeepLearning.AI

Stanford ONLINE

Andrew Ng

Example of Return

100	50	25	12.5	6.25	40
100	←	←	←	←	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	→	→	→	→	40
1	2	3	4	5	6

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

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

$$0 + (0.5)40 = 20$$

Policy

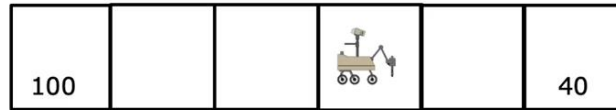
state s $\xrightarrow{\text{policy } \pi}$ action a

100	←	←	→	→	40
100	←	←	←	←	40
100	→	→	→	→	40
100	←	←	←	→	40

$$\begin{aligned}\pi(s) &= a \\ \pi(2) &= \leftarrow \\ \pi(3) &= \leftarrow \\ \pi(4) &= \leftarrow \\ \pi(5) &= \rightarrow\end{aligned}$$






A policy is a function $\pi(s) = a$ mapping from states to actions, that tells you what action a to take in a given state s .

The goal of reinforcement learning



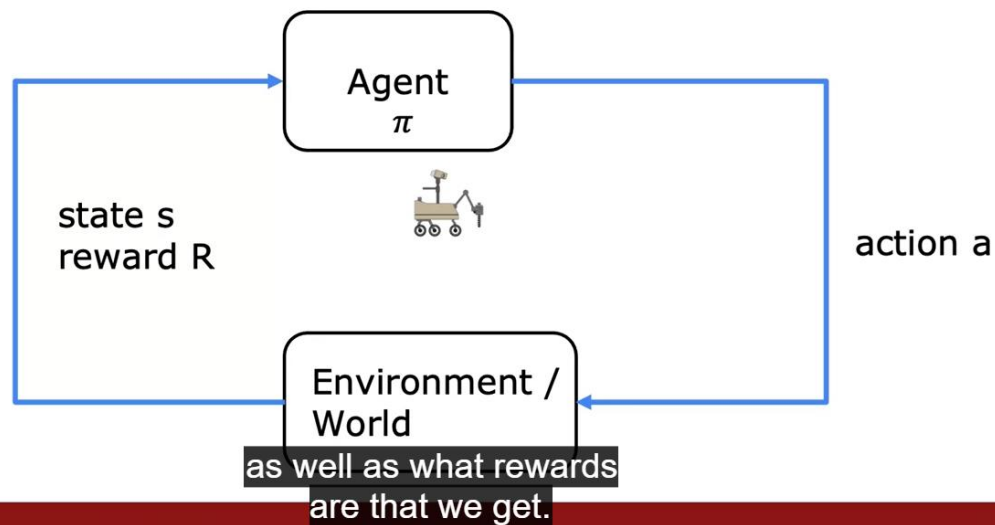
Find a policy π that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

more natural
terminology but policy

	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 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	 Find $\pi(s) = a$	Find $\pi(s) = a$	Find $\pi(s) = a$

This formalism of a reinforcement
learning application

Markov Decision Process (MDP)



State action value function

$Q(s, a)$ = Return if you

- start in state s .
- take action a (once).
- then behave optimally after that.

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

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

1 2 3 4 5 6

← return
← action
← reward

$Q(s, a)$
↑ ↑

$$Q(2, \rightarrow) = 12.5$$

$$0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

$$Q(2, \leftarrow) = 50$$

$$0 + (0.5)100$$

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

$$0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$$

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	0	0	0	0	0	40	
	1	2	3	4	5	6						

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

$$Q(4, \rightarrow) = 10$$

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

$$\pi(s) = a$$

$Q(s, a)$ = Return if you

- start in state s .
- take action a (once).
- then behave optimally after that.

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

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

Q^*
Optimal Q function

Bellman Equation

$Q(s, a)$ = Return if you

- start in state s .
- take action a (once).
- then behave optimally after that.

			←		
1	2	3	4	5	6

$$R(1)=100 \quad R(2)=0 \quad \dots \quad R(6)=40$$

s : current state
 a : current action

$R(s)$ = reward of current state

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

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

We'll come back to see why

Bellman Equation

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

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

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	0	40	40
1	2	3	4	5	6						

$$\begin{aligned} s &= 2 \\ a &= \rightarrow \\ s' &= 3 \end{aligned}$$

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

$$\begin{aligned} Q(4, \leftarrow) &= R(4) + 0.5 \max_{a'} Q(3, a') \\ &= 0 + (0.5)25 = 12.5 \end{aligned}$$

$$\begin{aligned} s &= 4 \\ a &= \leftarrow \\ s' &= 3 \end{aligned}$$

any other state action in

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

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

Reward you get right away

Return from behaving optimally starting from state s' .

$$\begin{aligned} &R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots \\ Q(s, a) &= R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots] \end{aligned}$$

the total return you get

Explanation of Bellman Equation


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

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

$Q(4, \leftarrow)$
 $= 0 + (0.5)0 + (0.5)^2 0 + (0.5)^3 100$
 $= R(4) + (0.5) [0 + (0.5) 0 + (0.5)^2 100]$
 $= R(4) + (0.5) \max_{a'} Q(3, a')$

Gamma times the returns from the next state s prime.

Expected Return

100	←	←	←	→	
100	0	0	0	0	40
1	2	3	4	5	6

$$\begin{matrix} & & 0 & & 0 & & 0 & & 100 & & 0 & & 100 \\ & & 0 & & 0 & & 0 & & 0 & & 0 & & 0 \\ & & 0 & & 0 & & 40 & & 0 & & 0 & & 0 \end{matrix}$$

Expected Return = Average($R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$)
 $= E[R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots]$

the average or the expected sum of discounted rewards.

State Action Value Function Example

In this Jupyter notebook, you can modify the mars rover example to see how the values of $Q(s,a)$ will change depending on the rewards and discount factor changing.

```

In [1]: import numpy as np
        from utils import *

In [2]: # Do not modify
        num_states = 6
        num_actions = 2

In [3]: terminal_left_reward = 100
        terminal_right_reward = 40
        each_step_reward = 0

        # Discount factor
        gamma = 0.5

        # Probability of going in the wrong direction
        misstep_prob = 0.1

In [4]: generate_visualization(terminal_left_reward, terminal_right_reward, each_step_reward, gamma, misstep_prob)

```

Optimal policy

100.0	46.06	21.25	10.49	18.52	40.0
100	←	←	←	→	40
100	0	0	0	0	40

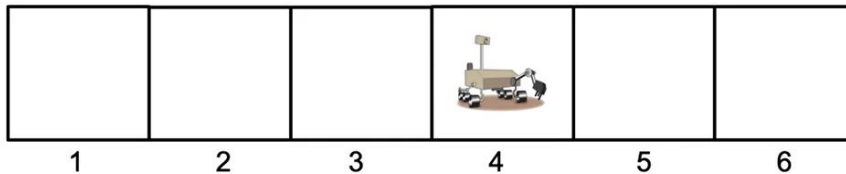
$Q(s,a)$

100.0	100.0	46.06	6.72	18.52	40.0	40.0
100					0	40

The q values, as well as the optimal returns,

Discrete vs Continuous State

Discrete State:



Continuous State:



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

between zero and 360 degrees.

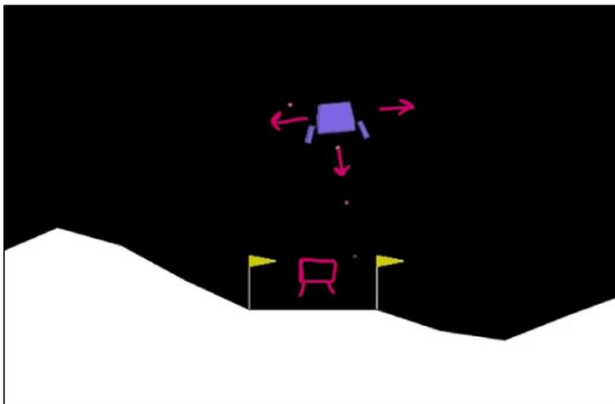
Autonomous Helicopter



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

and then the row pitch,

Lunar Lander



actions:

do nothing

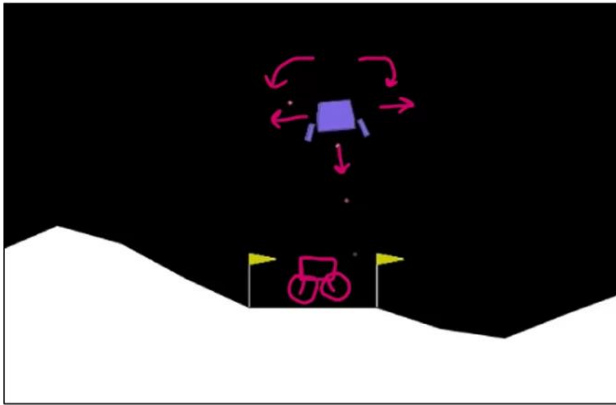
left thruster

main thruster

right thruster

So it's the lunar lander safely between these two flags here on the landing pad.

Lunar Lander



actions:

do nothing
left thruster
main thruster
right thruster

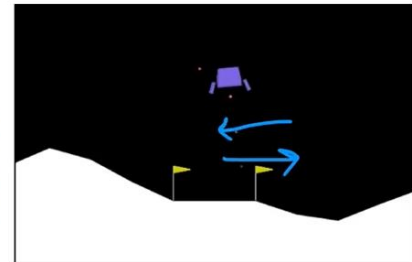
$$s = \begin{bmatrix} x \\ y \\ \dot{x} \\ \dot{y} \\ \theta \\ \dot{\theta} \\ l \\ r \end{bmatrix}$$

0 or 1

one depending on whether the left and right legs are touching the 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
- Fire main engine: -0.3
- Fire side thruster: -0.03

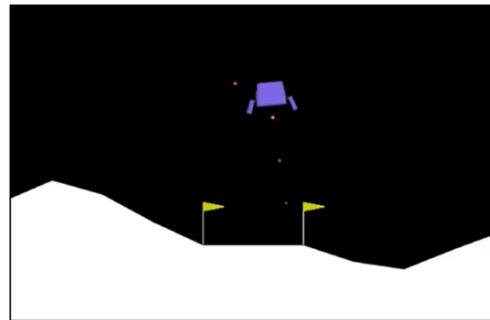


Which is much harder for this and many other reinforcement learning applications.

Lunar Lander Problem

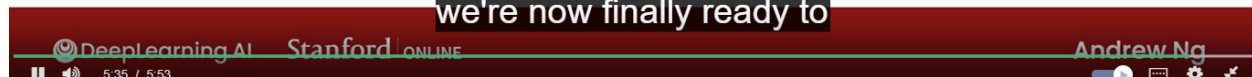
Learn a policy π that, given

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

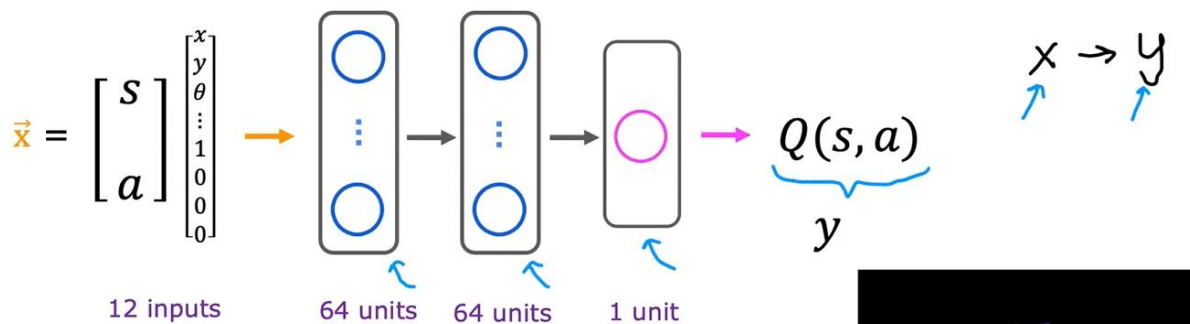


picks action $a = \pi(s)$ so as to maximize the return.

this lunar lander exciting application and
we're now finally ready to



Deep Reinforcement Learning

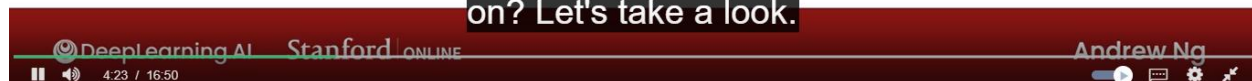


In a state s , use neural network to compute

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

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

train a neural network
on? Let's take a look.



Bellman Equation

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

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

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

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

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

$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)})$ actually take this training set where the

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

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

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

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:

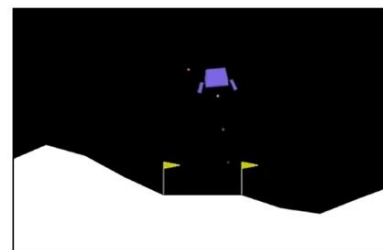
Create training set of 10,000 examples using

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

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

Set $Q = Q_{new}$.

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



x, y

x'', y''

x^{10000}, y^{10000}

How to choose actions while still learning?

In some state s

Option 1:

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

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

↑
 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 ($\epsilon = 0.05$)
0.95

Start ϵ high
 $1.0 \rightarrow 0.01$
Gradually decrease

in the Jupiter lab that
shows you how to do this.

DeepLearning.AI Stanford ONLINE

Andrew Ng

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.

machine learning systems as well.

DeepLearning.AI Stanford ONLINE

Andrew Ng

2:25 / 2:54

⏮ ⏪ ⏩ ⏭ ⚙ ⚡