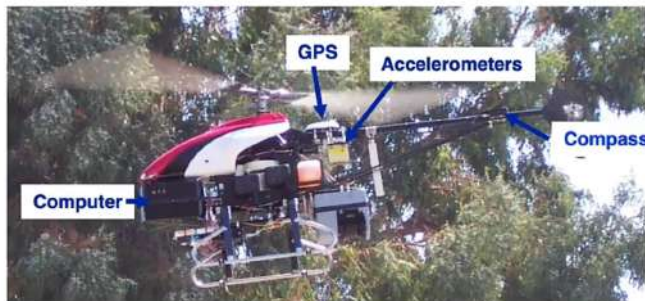


Course 3

Third Week

What is Reinforcement Learning?

Autonomous Helicopter



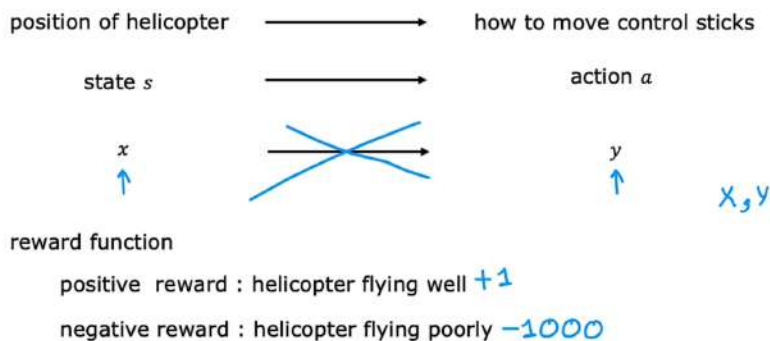
How to fly it?

学习到飞

Autonomous Helicopter



Reinforcement Learning



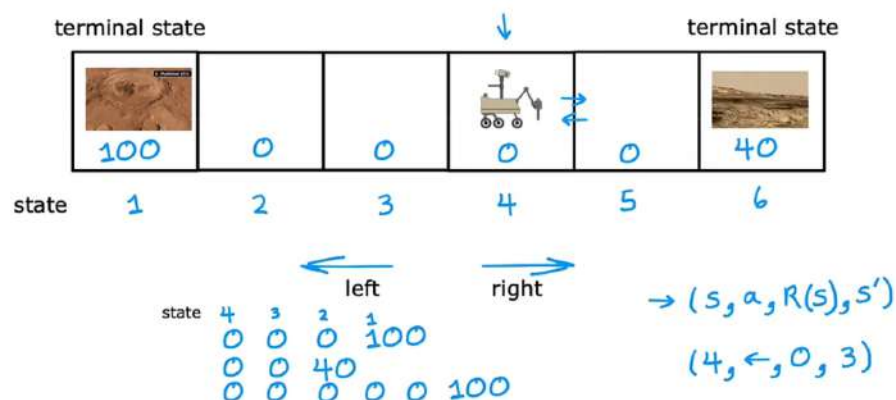
Applications

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



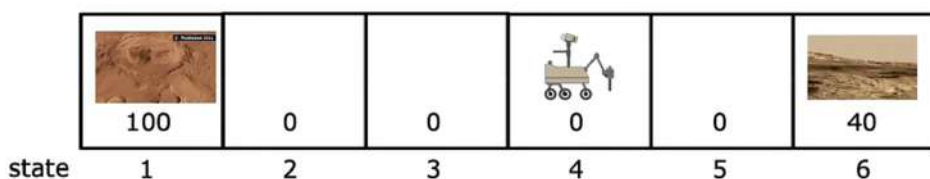
Mars rover example

Mars Rover Example



The Return in Reinforcement Learning γ

Return



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

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

Discount Factor $\gamma = 0.9$ 0.99 0.999 \rightarrow 远视

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

Example of Return

100	50	25	12.5	6.25	40
100	0	0	0	0	40
1	2	3	4	5	6

← return

$$\gamma = 0.5$$

← reward

The return depends on the actions you take.

100	2.5	5	10	20	40
100	0	0	0	0	40
1	2	3	4	5	6

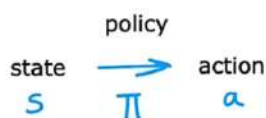
$$0 + (0.5)0 + (0.5)^2 40 = 10$$

100	50	25	12.5	20	40
100	0	0	0	0	40
1	2	3	4	5	6

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

Making decisions: Policies in reinforcement Learning

Policy




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

$\pi(1) = \rightarrow$
 $\pi(2) = \leftarrow$
 $\pi(3) = \leftarrow$
 $\pi(4) = \leftarrow$
 $\pi(5) = \rightarrow$




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

100					40
-----	--	--	---	--	----

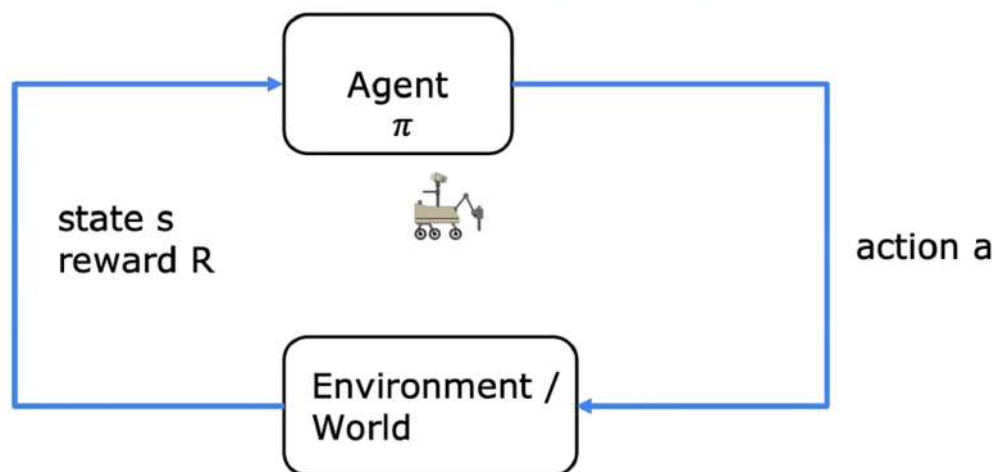
Find a policy π that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

Review of key concepts

	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	<table border="1"><tr><td>100</td><td>←</td><td>←</td><td>←</td><td>→</td><td>40</td></tr></table>	100	←	←	←	→	40	Find $\pi(s) = a$ ↑ ↑	Find $\pi(s) = a$ ↑ ↑
100	←	←	←	→	40				

总结
→

Markov Decision Process (MDP)



State-action value function definition

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

$Q(s, a)$
↑ ↑

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

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

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

$$0 + (0.5)100$$

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

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

这通常也称为 Q 函数。

Picking actions

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

← return
← action
← reward

100	50	25	12.5	20	40
100	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

State-action value function example

Optional Lab

see how γ and reward

change policy π $Q(s, a)$

Bellman Equation \rightarrow help compute

Bellman Equation

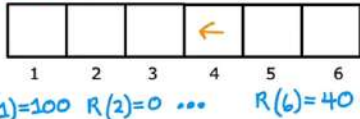
$Q(s, a)$ = Return if you

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

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

$R(s)$ = reward of current state

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



$R(1)=100$ $R(2)=0$... $R(6)=40$

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

Bellman Equation

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

$Q(s, a) = R(s)$ \rightarrow terminal T

$s=2$
 $a \rightarrow$
 $s'=3$

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						

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

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

$s=4$
 $a \leftarrow$
 $s'=3$

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

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' .

$$Q(s, a) = R_1 + \gamma [R_2 + \gamma R_3 + \gamma^2 R_4 + \dots]$$

Explanation of Bellman Equation

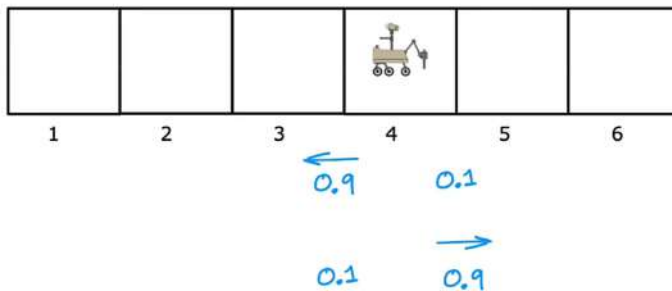
$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

100	100	50	12.5	25	6.25	12.5	10	6.25	20	40	40
100	0	0	0	0	0	0	0	0	0	40	40
1	2	3	4	5	6						

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

Random (stochastic) environment (Optional)

Stochastic Environment



Expected Return

100	←	←	←	→	Robot
1	2	3	4	5	6

$$\begin{aligned}
 &0 \quad 0 \quad 0 \quad 100 \quad 0 \quad 100 \\
 &0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \\
 &0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Expected Return} &= \text{Average}(R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots) \\
 &= E[R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots]
 \end{aligned}$$

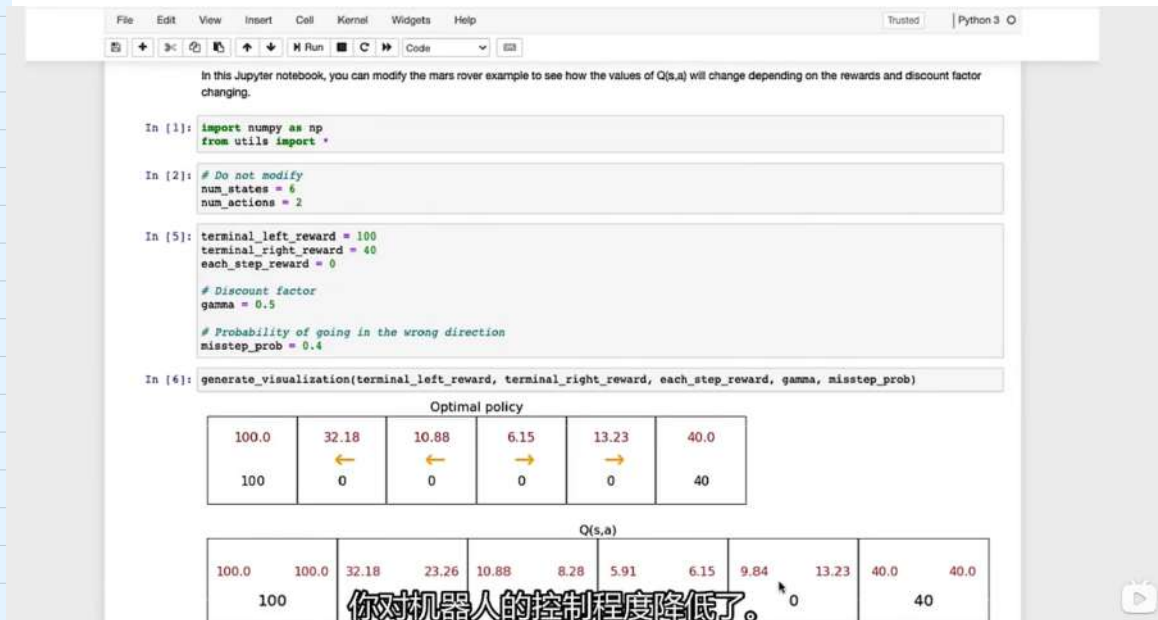
Expected Return

Goal of Reinforcement Learning:

Choose a policy $\pi(s) = a$ that will tell us what action a to take in state s so as to maximize the expected return.

Bellman Equation:

$$Q(s, a) = \underbrace{R(s)}_{\substack{\uparrow \\ 3}} + \gamma \underbrace{E[\max_{a'} Q(s', a')]}_{\substack{\uparrow \\ 2 \text{ or } 4}}$$



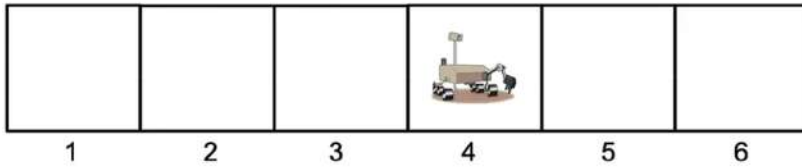
通过更改 misstep_prob 观察 Q 值

更多的 s

Example of Continuous State application

Discrete vs Continuous State

Discrete State:



Continuous State:



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

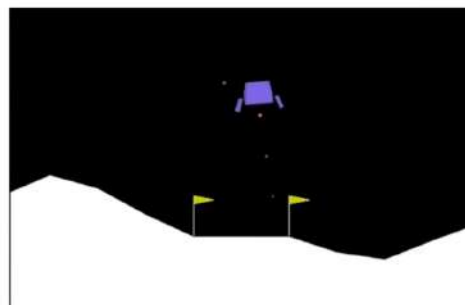
Autonomous Helicopter



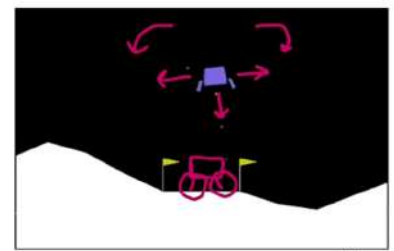
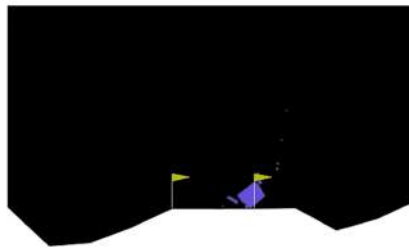
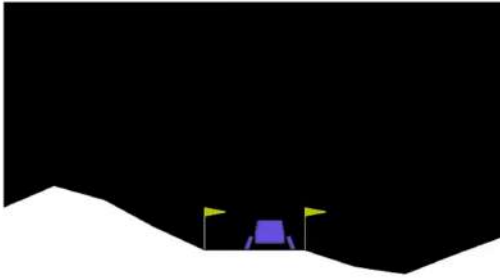
$$s = \begin{bmatrix} x \\ y \\ z \\ \phi \\ \theta \\ \omega \\ \dot{x} \\ \dot{y} \\ \dot{z} \\ \dot{\phi} \\ \dot{\theta} \\ \dot{\omega} \end{bmatrix}$$

Lunar lander

Lunar Lander



if Good Policy is bad



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

Reward Function

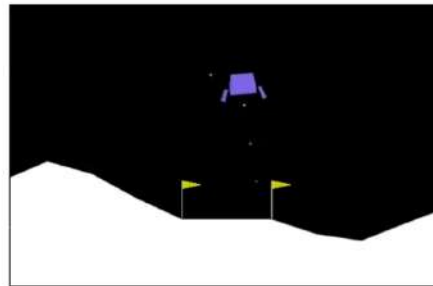
- Getting to landing pad: 100 - 140
- Additional reward for moving toward/away from pad.
- Crash: -100
- Soft landing: +100
- Leg grounded: +10
- Fire main engine: -0.3
- Fire side thruster: -0.03



Lunar Lander Problem

Learn a policy π that, given

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



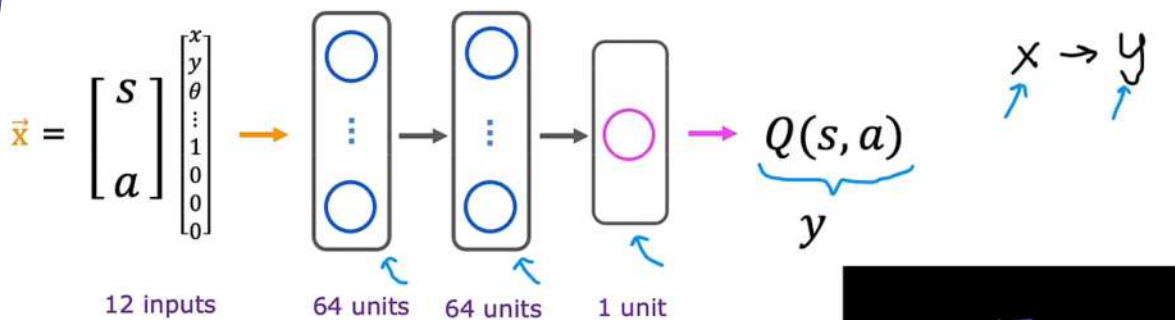
picks action $a = \pi(s)$ so as to maximize the return.

$$\gamma = 0.985$$

Learning the state-value function

DQN

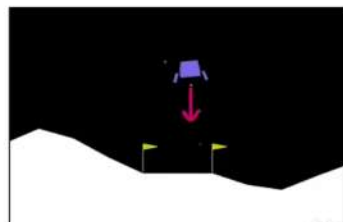
Deep Reinforcement Learning



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 $Q(s, a)$



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

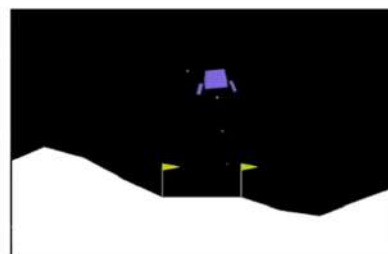
Set $Q = Q_{new}$.

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

x, y

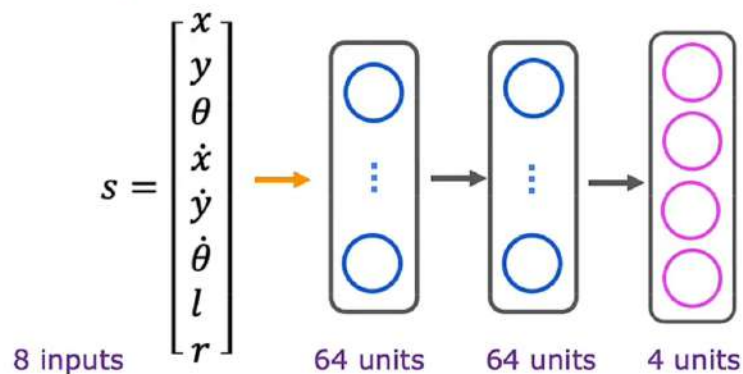
$x^{(1)}, y^{(1)}$

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



Algorithm refinement: Improved neural network architecture

Deep Reinforcement Learning



$\left. \begin{matrix} Q(s, \text{nothing}) \\ Q(s, \text{left}) \\ Q(s, \text{main}) \\ Q(s, \text{right}) \end{matrix} \right\} \text{并行计算}$

In a state s , input s to neural network.

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



Algorithm refinement: ϵ -greedy policy

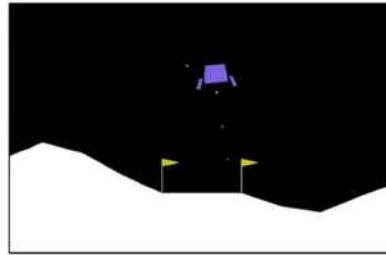
Learning Algorithm

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.



Train model:

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 $Q_{new}(s, a) \approx y$. $f_{w,b}(x) \approx y$

Set $Q = Q_{new}$.

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 ($\epsilon = 0.05$)
0.95

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

↑
 a

Start ϵ high

$1.0 \rightarrow 0.01$

Gradually decrease

Algorithm refinement: Mini-batch and soft update

How to choose actions while still learning?

x	y
2104	400
1416	232
1534	315
852	178
...	...
3210	870

100,000,000

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$m = 100,000,000$

$m' = 1,000$

repeat {

$$w = w - \alpha \frac{\partial}{\partial w} \left[\frac{1}{2m'} \sum_{i=1}^{m'} (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$$

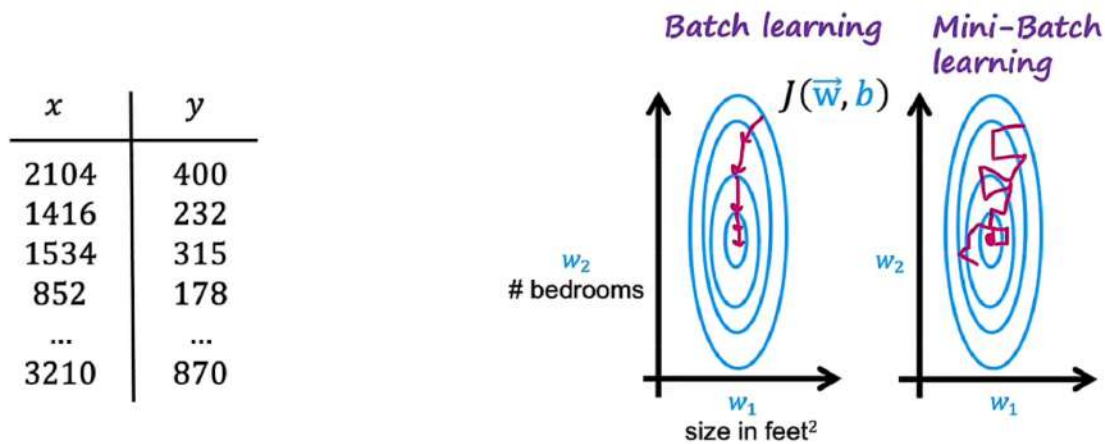
$$b = b - \alpha \frac{\partial}{\partial b} \left[\frac{1}{2m'} \sum_{i=1}^{m'} (f_{w,b}(x^{(i)}) - y^{(i)})^2 \right]$$

}

Mini-batch



Mini-batch



解决 training set 过大的问题

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:

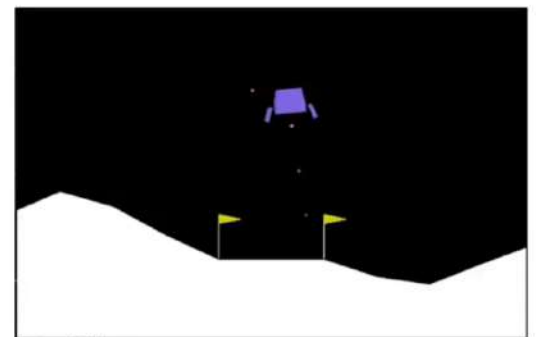
1,000

Create training set of 10,000 examples using

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

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

Set $Q = Q_{new}$.



$x^{(1)}, y^{(1)}$
 \vdots
 $x^{(1000)}, y^{(1000)}$

② minibatch 会带来 Q 突然变化

3/1 soft update

Soft Update

Set $Q = Q_{new}$ $\leftarrow Q(s, a)$
 \uparrow \uparrow
 W, B W_{new}, B_{new}

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

$$W = 1 W_{new} + 0 W$$

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

The state of reinforcement Learning

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.

Summary of the Course

Courses

- Supervised Machine Learning: Regression and Classification
Linear regression, logistic regression, gradient descent
- Advanced Learning Algorithms
Neural networks, decision trees, advice for ML
- Unsupervised Learning, Recommenders, Reinforcement Learning
Clustering, anomaly detection, collaborative filtering, content-based filtering, reinforcement learning