# Reinforcement Learning

Position is given -> How to move the control sticks?
state s -> action a

Supervised Learning?
x -> y
for x's we will find y's from our knowledge?

it is too ambigious for it, flying is unpredictable.

how to make a dog behave? Good dog! Bad dog!

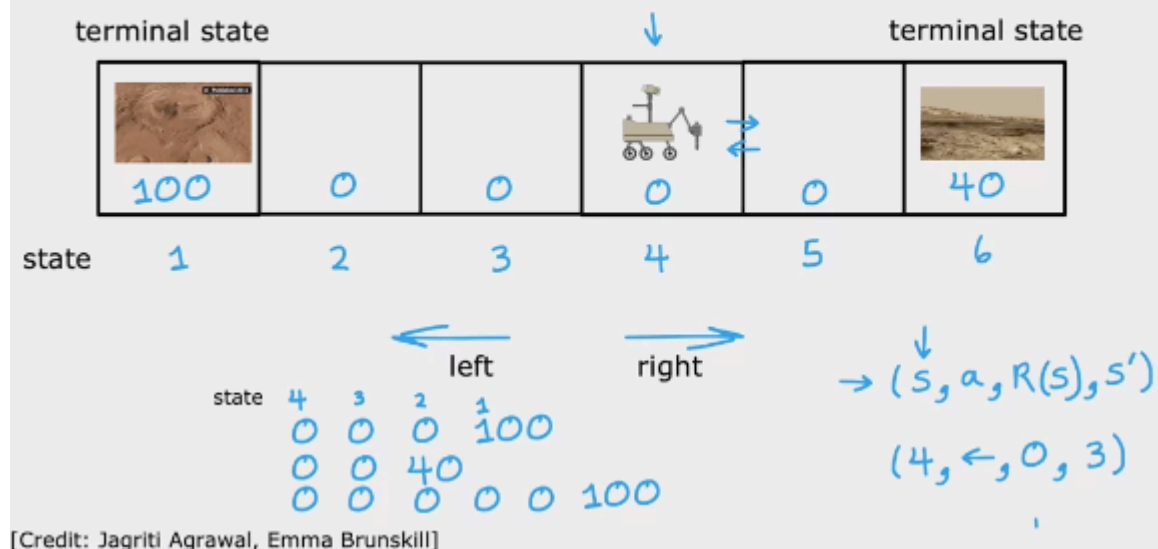you have to tell it what to do, rather than how to do it.
reward/punishment



[Thanks to Zico Kolter]

*Financial trading?*

*terminal state* when gets to the one of the *states*, day ends.

# Mars Rover Example



[Credit: Jagriti Agrawal, Emma Brunskill]

## The Return in Reinforcement Learning

*discount factor* of r = (ex) 0.9

Return = R + r.R + r^2.R + ...

Finance: Interes rate, time value of money.

# Example of Return



The return depends on the actions you take.

## Policies in reinforcement learning

can choose always go *direction*

state -> (policy) -> action

goal of RL: find a policy (pi) that tells you what action to take in every state so as to maximize the return.

states, actions, rewards, discount factor, return, policy(pi)

| | 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 $\pi(s) = a$ ↑ ↑ |

## Markov Decision Process (MDP)

*future depends on where you are now, not on how you got here.*

## State Action Value Function (Q-function)



### Picking actions

| 100 | 50 | 25 | 12.5 | 20 | 40 |
|---|---|---|---|---|---|
| 100 | 0 | 0 | 0 | 0 | 40 |

← return
← action
← reward

| 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 |

$Q(4, \leftarrow)$  $Q(4, \rightarrow)$
  12.5        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 $\max_a Q(s, a)$.

## Bellman Equation:

# 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 |

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

$s$ : 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

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

---

# 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 $s'$.

$$R_1 + r R_2 + r^2 R_3 + r^3 R_4 + \cdots$$

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

**Discrete vs Continuous State**

how would you represent the position of the helicopter/truck/...

pos, pitch, angle. with a vector:

# 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*
actions:
do nothing
left thruster
main thruster
right thruster
s=
[
x
y
velocity
angle
angular velocity
*l* sitting on ground?
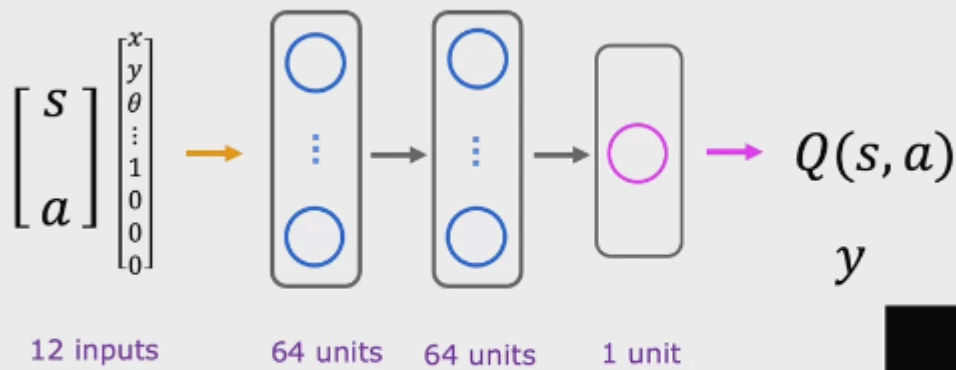*r* sitting on ...?
]

custom reward function:
crush -100
soft land +100
leg grounded +, fire main engine ---, side thrust -, get to landing pad ++

# Deep Reinforcement Learning



$$\begin{bmatrix} s \\ a \end{bmatrix} \begin{bmatrix} x \\ y \\ \theta \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \cdots \rightarrow Q(s,a)$$

12 inputs      64 units    64 units    1 unit

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 $Q(s, a)$

Create a supervised learning data for the NN output ( Learning the Q(s, a) function )

# Bellman Equation

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

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

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

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

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

$(s^{(3)}, a^{(3)}, R(s^{(3)}), s'^{(3)}) \leftarrow$

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

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

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

# 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

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

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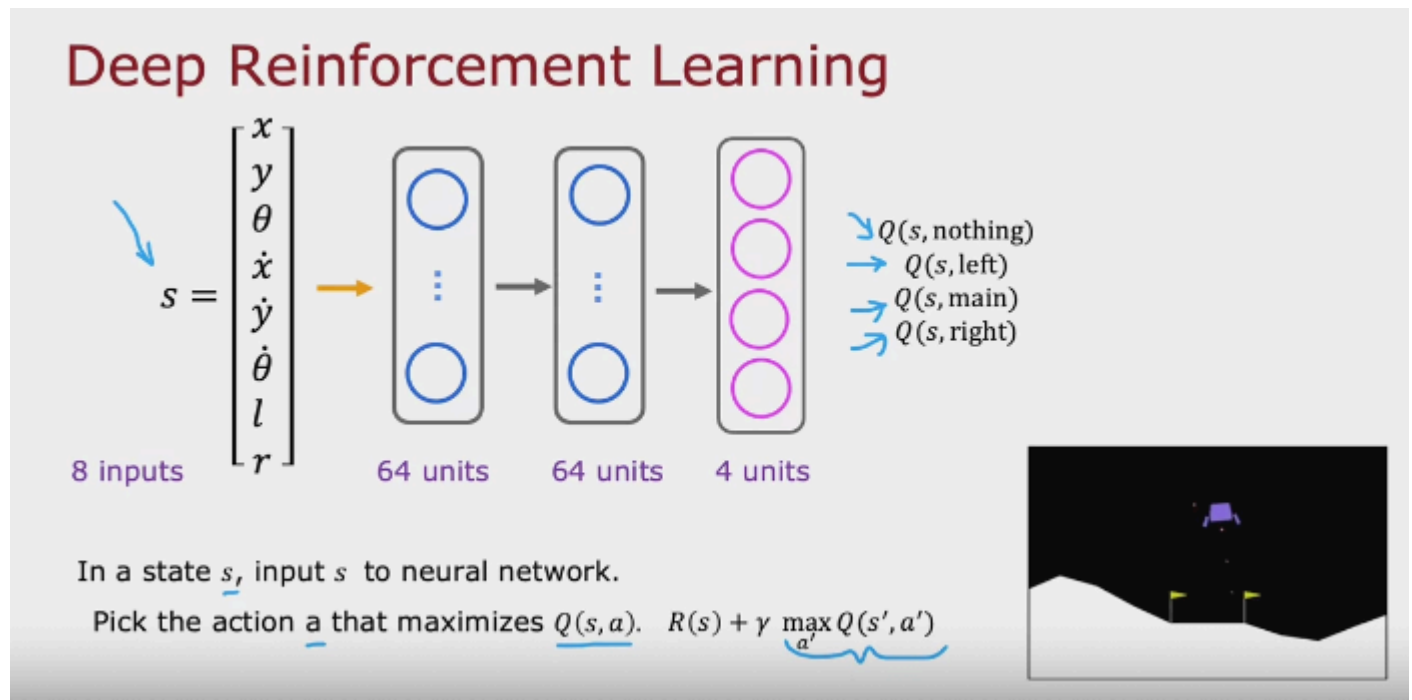
     Set $Q = Q_{new}$.

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

$x, y$

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

$\vdots$

$x^{(10,000)}, y^{(10,000)}$

Continuous State Spaces:
Improved NN Architecture

## Deep Reinforcement Learning



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

8 inputs     64 units    64 units    4 units

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

In a state $s$, input $s$ to neural network.
Pick the action $a$ that maximizes $Q(s, a)$.   $R(s) + \gamma \max_{a'} Q(s', a')$

**Algorithm refinement: ε-greedy policy**
How to choose actions while we are still learning?
at some probability, pick an action *a* randomly ( %5 )
epsilon greedy, epsilon = 0.05
Some strategies may never be tried by the NN in some cases, (stuck by chance)

start e high and gradually decrease (we know more, possibility of being stuck is lower.)

**Algorithm refinement: Mini-batch and soft updates**

## 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 $Q = Q_{new}$.

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

# Limitations

much easier in simulations rather than real world scenarios.
fewer applications compared to (un)supervised learning.

potential for future.