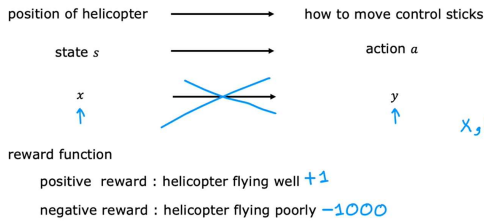


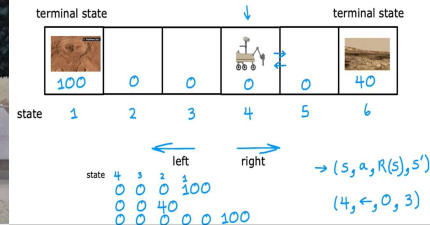
Reinforcement Learning



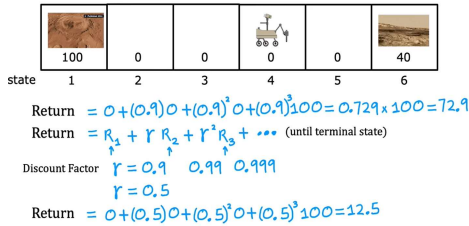
Robotic Dog Example



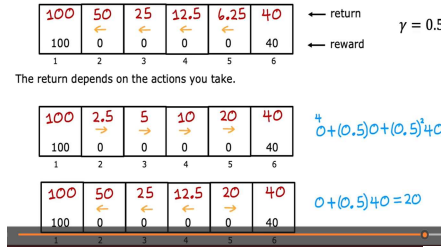
Mars Rover Example



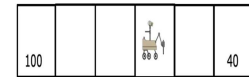
Return



Example of Return

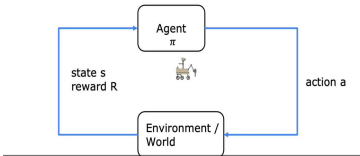


The goal of reinforcement learning

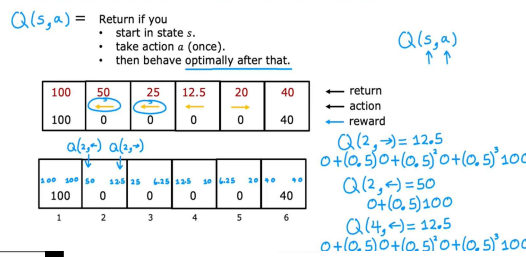


Find a policy π that tells you what action ($a = \pi(s)$) to take in every state (s) so as to maximize the return.

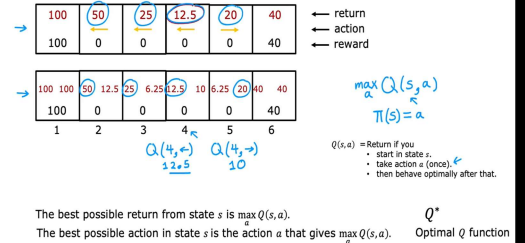
Markov Decision Process (MDP)



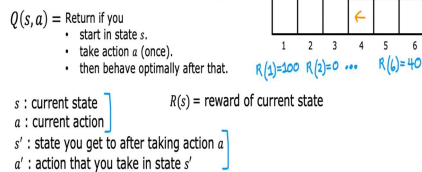
State action value function (Q-function)



Picking actions

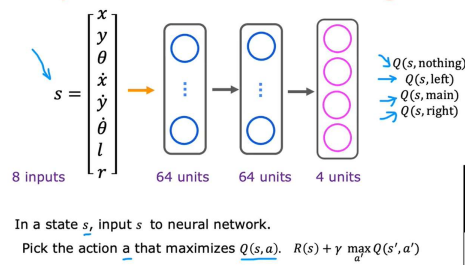


Bellman Equation

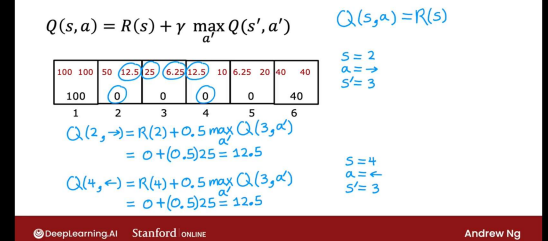


$$Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')$$

Deep Reinforcement Learning



Bellman Equation



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.

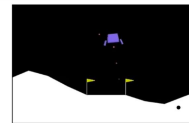
Train model:

Create training set of 10,000 examples using

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

Train Q_{new} such that $Q_{\text{new}}(s, a) = y$. $f_{W, B}(x) \approx y$

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

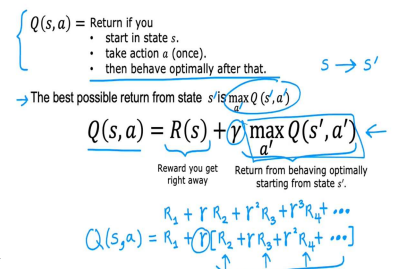


Refinement

Greedy

Policy

Explanation of Bellman Equation



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

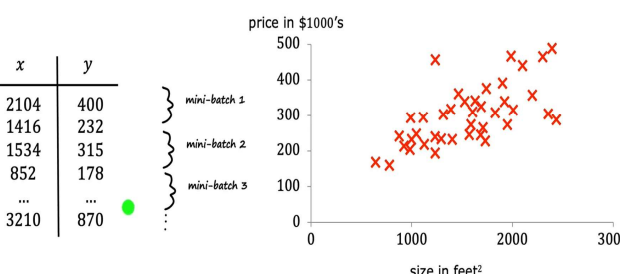
a

Start ϵ high

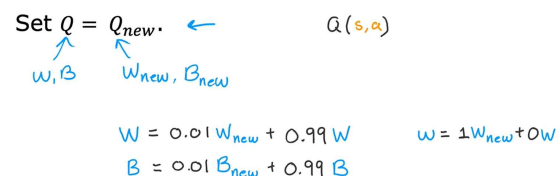
$1.0 \rightarrow 0.01$

Gradually decrease

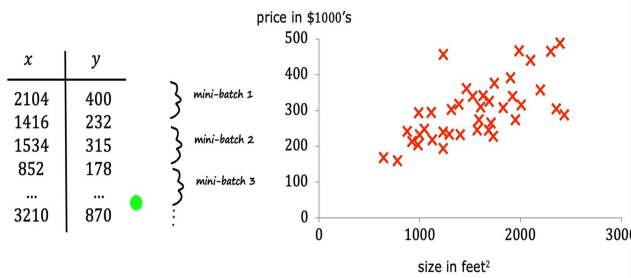
Mini-batch



Soft Update



Mini-batch



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:

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

Soft Update

Set $Q = Q_{new}$.

$Q(s, a)$

W, B

W_{new}, B_{new}

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

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

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