

Hi

Deep Reinforcement Learning

Deep Q-Network



Azzeddine CHENINE
@azzeddineCH_



AI EVERYTIME
EVERWHERE



 OpenAI

15 October 2019

Quick Reminder Q-Learning



Quick Reminder

Learning by experience

Quick Reminder

Learning by experience

Episodic tasks and continuous task

Quick Reminder

Learning by experience

Episodic tasks and continuous task

An action has a reward and results a new state

Quick Reminder

The value based methods try to reach an optimal policy by using a value function

Quick Reminder

The value based methods try to reach an optimal policy by using a value function

The Q function of an optima policy maximise the expected return G

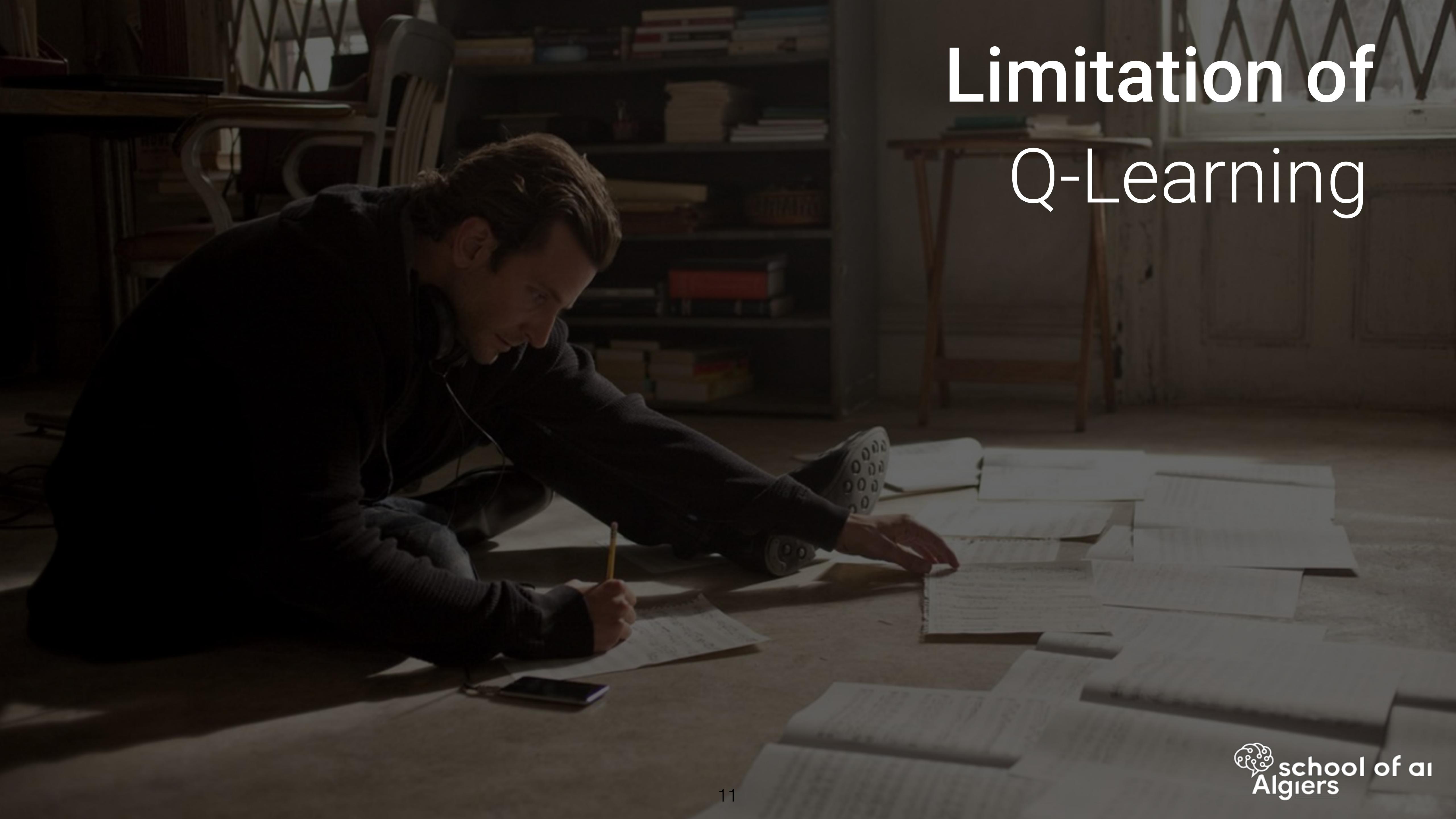
Quick Reminder

The value based methods try to reach an optimal policy by using a value function

The Q function of an optima policy maximise the expected return G

Epsilon greedy policies solve the exploration and exploitation problem

Limitation of Q-Learning



Limitation

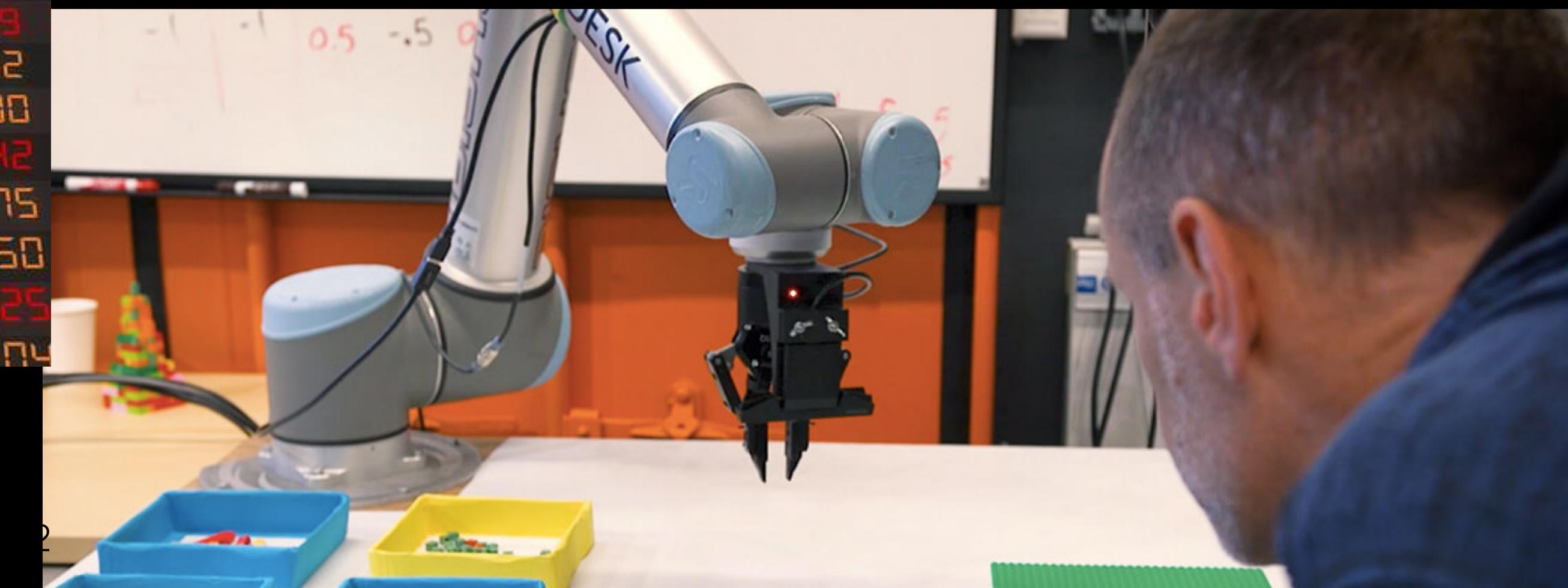
Non practical with environments with large or infinite state or action space

Games



100.00	114.00	112.00	110.00	108.00
110.00	112.00	115.00	109.50	115.00
129	▼ -2.64	▲ 1.71	▼ -5.25	▲ 5.26
110	3.06	3.52	3.04	3.85
110	3.60	3.90	3.52	4.25
110	▲ 0.54	▲ 0.37	▲ 0.48	▲ 0.39
110	57.65	70.07	59.50	65.81
110	56.50	65.84	59.93	62.72
110	▼ -6.14	▼ -4.22	▲ 0.42	▼ -3.09
102.30	103.27	103.59	102.42	
98.70	102.00	106.00	100.00	
110	▼ -3.59	▼ -1.26	▲ 2.40	▼ -2.42
116.00	114.64	113.28	114.75	
106.70	112.00	115.00	109.50	
110	▼ -9.29	▼ -2.64	▲ 1.71	▼ -5.25
129	3.06	3.52	3.04	3.85

Robots



Stock market

**Modern problems
require
Modern solutions**



Deep Q-Networks

Deepminds

- 2015 -

No more given definitions

No more given definitions

**figuring it out by
ourselves**

Neural networks can be trained to fit on a linear or a non linear function.

1 - How to replace the Q-table ?

Neural networks can be trained to fit on a linear or a non linear function.

Replace the Q-value function with a deep neural network

1 - How to replace the Q-table ?

Neural networks can be trained to fit on a linear or a non linear function.

1 - How to replace the Q-table ?

Replace the Q-value function with a deep neural network

The outputs layer hold the Q-value of each possible action

⇒ make the actions space discrete

Neural networks can be trained to fit on a linear or a non linear function.

1 - How to replace the Q-table ?

Replace the Q-value function with a deep neural network

The outputs layer hold the Q-value of each possible action

⇒ make the actions space discrete

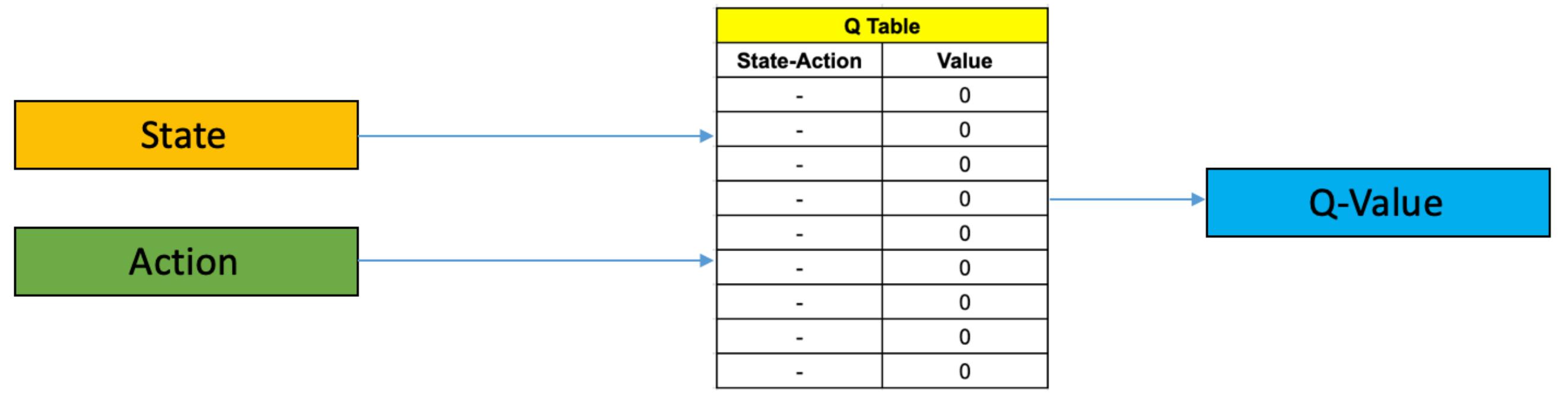
The input layer takes the desired representation of the state

⇒ a pipeline of convolutions

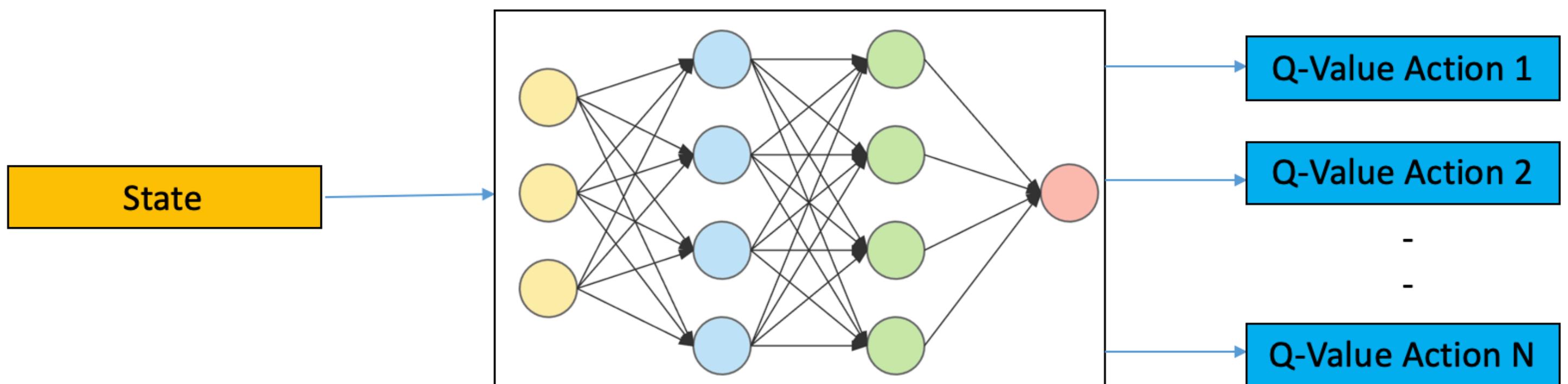
⇒ a Dense layer

⇒ a sequence of LSTM/RNN cells

1 - How to replace the Q-table ✓



Q Learning



Deep Q Learning

How did we used update the Q-table

1 - How to replace the
Q-table 

The agent took an action (a) in a state (S)

2 - How to optimize the
model

How did we used update the Q-table

1 - How to replace the
Q-table 

The agent took an action (a) in a state (S)

2 - How to optimize the
model

The agent moved to the next state (S') and
received a reward (r)

How did we used update the Q-table

1 - How to replace the
Q-table 

2 - How to optimize the
model

The agent took an action (a) in a state (S)

The agent moved to the next state (S') and received a reward (r)

the discount rate δ reflects how much we care about the futur reward

$\delta \approx 1 \implies$ we care too much about the future reward

$\delta \approx 0 \implies$ we care too much about the immediate reward

How did we used update the Q-table

1 - How to replace the
Q-table 

2 - How to optimize the
model

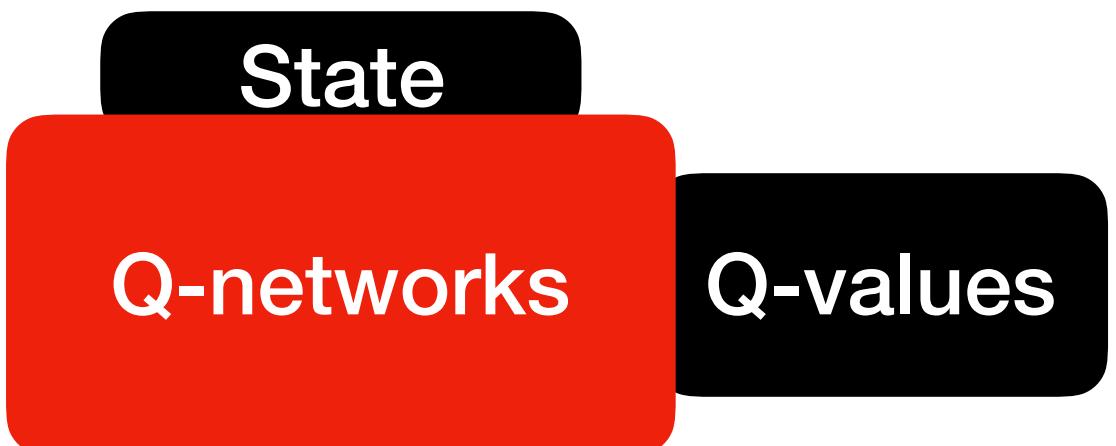
$$Q'(S, a) \leftarrow Q(S, a) + \alpha [r(S, a, S') + \delta \max Q(S', a') - Q(S, a)]$$

A good **Q** value for an action (a) given a state (S) is a good approximation to the expected discounted cumulative reward.

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

How are we going to update the Q-networks weights

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$



1 - How to replace the
Q-table ✓

2 - How to optimize the
model

How are we going to update the Q-networks weights

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

1 - How to replace the Q-table ✓

2 - How to optimize the model

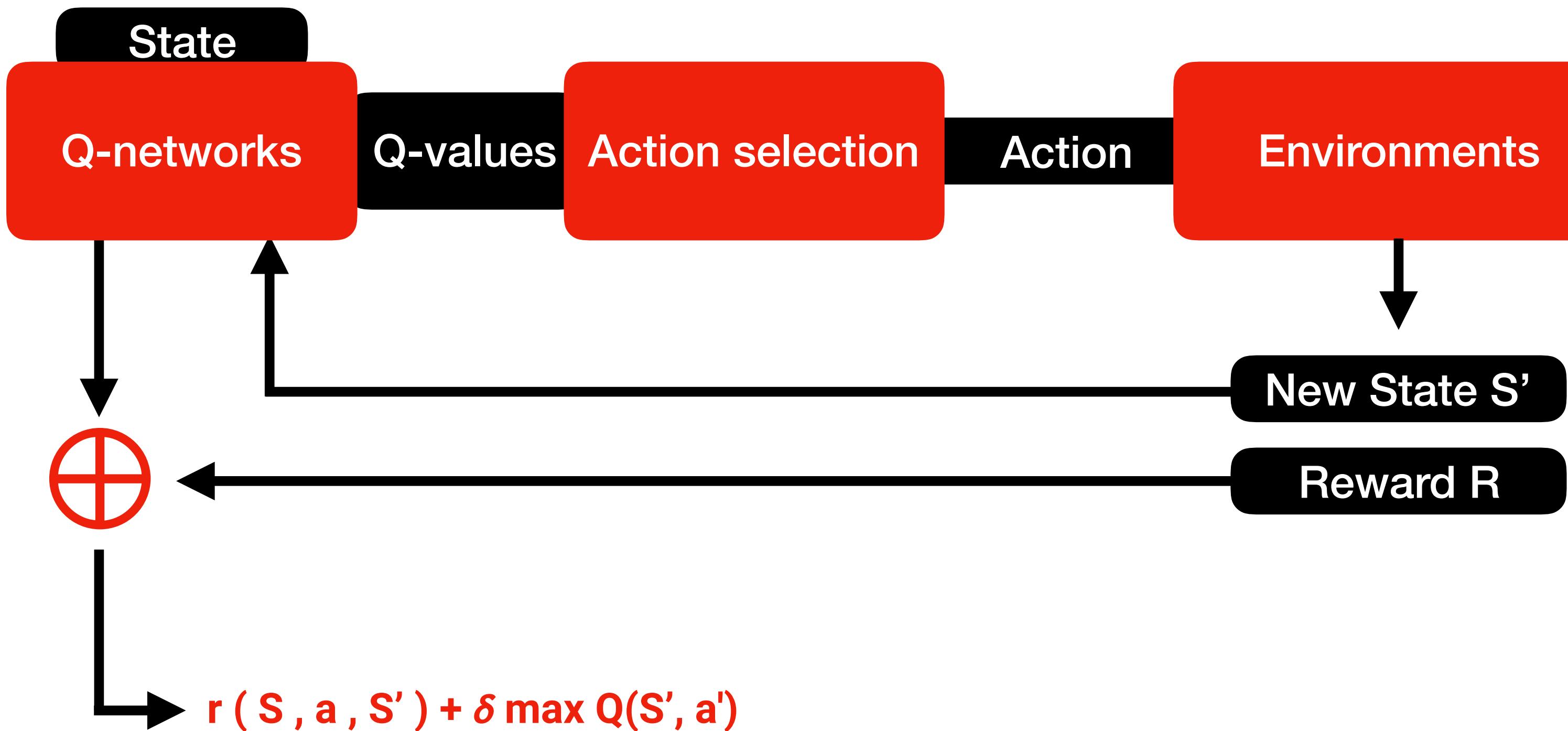


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

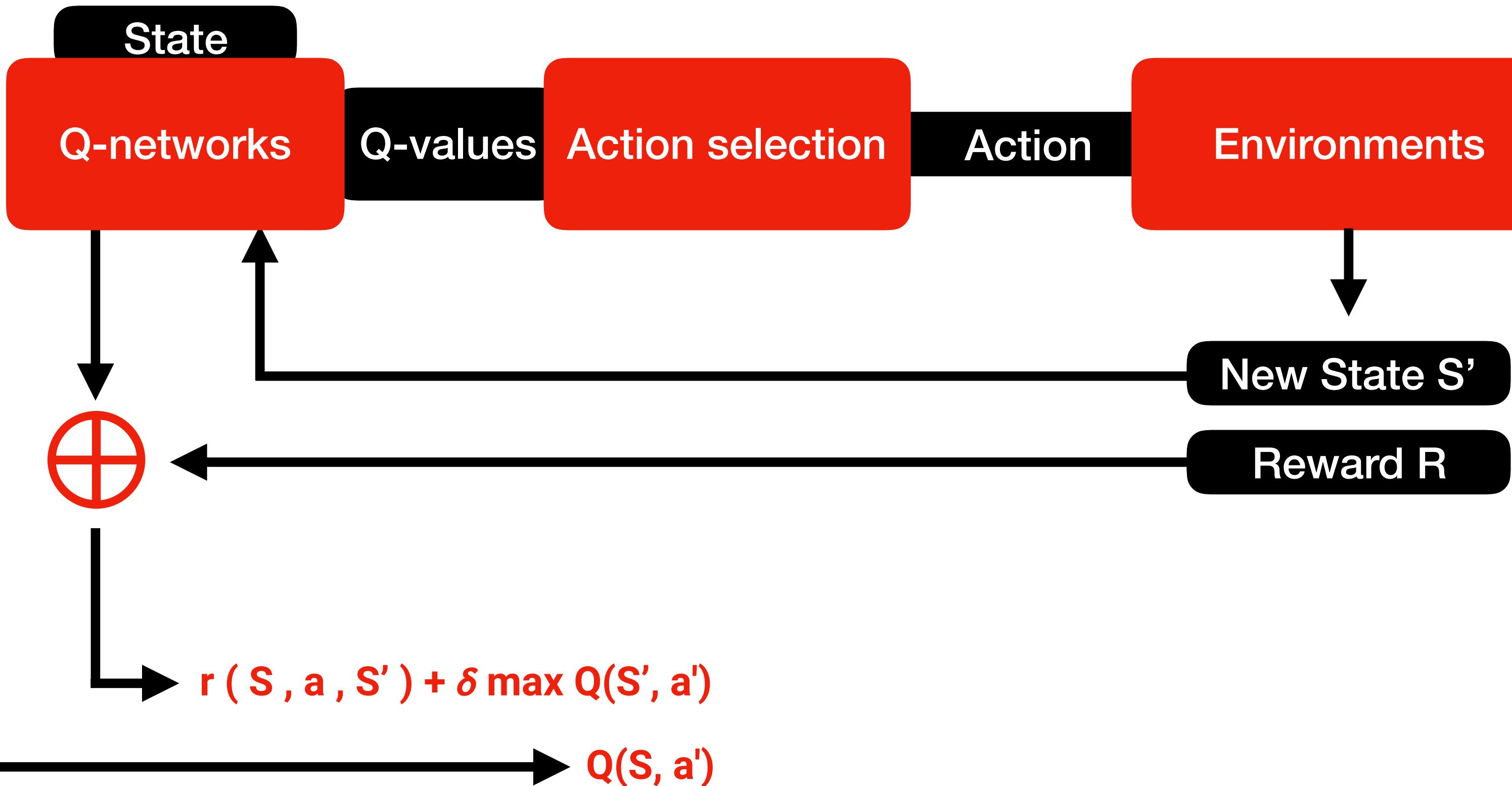


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

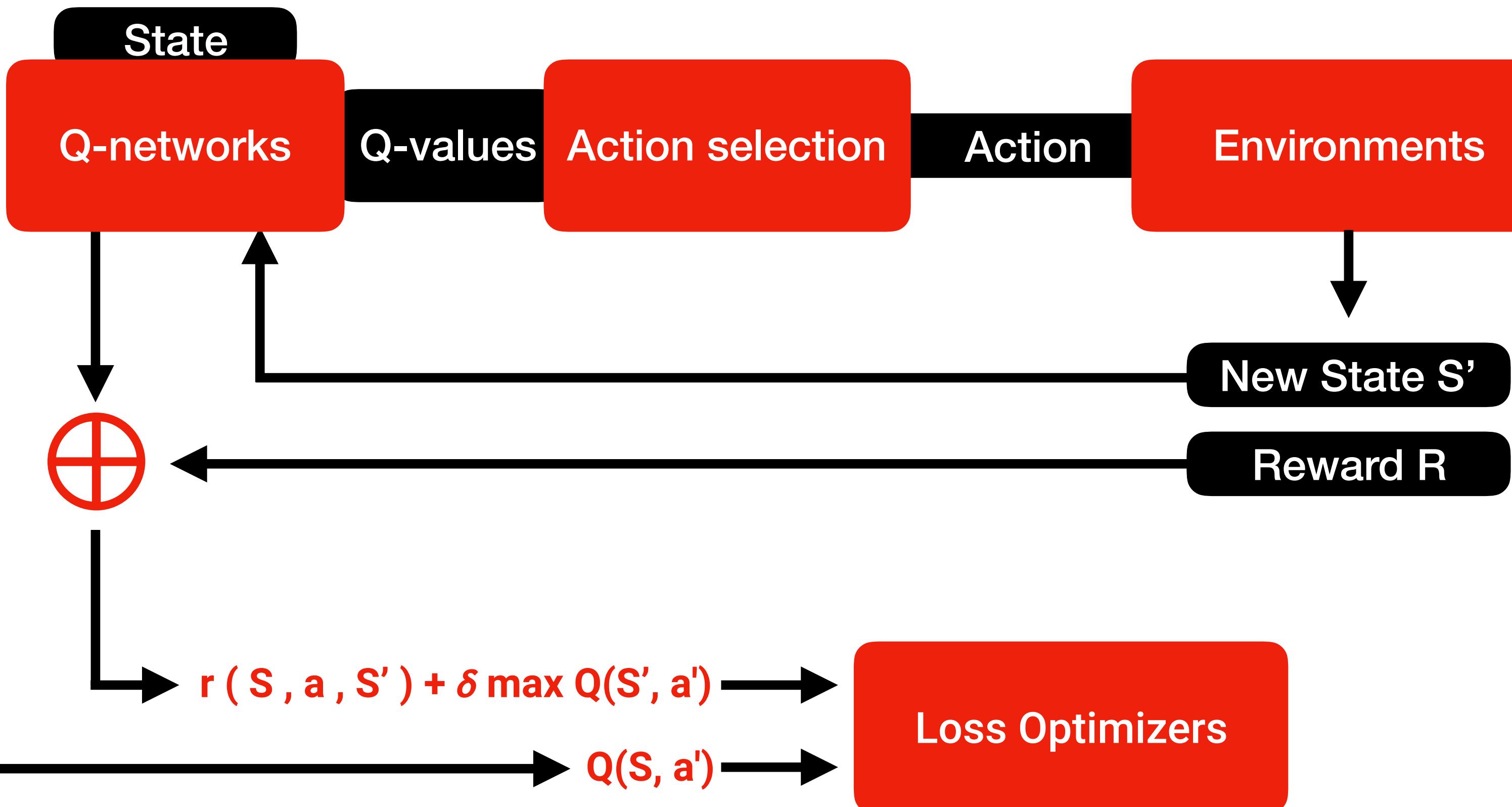


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

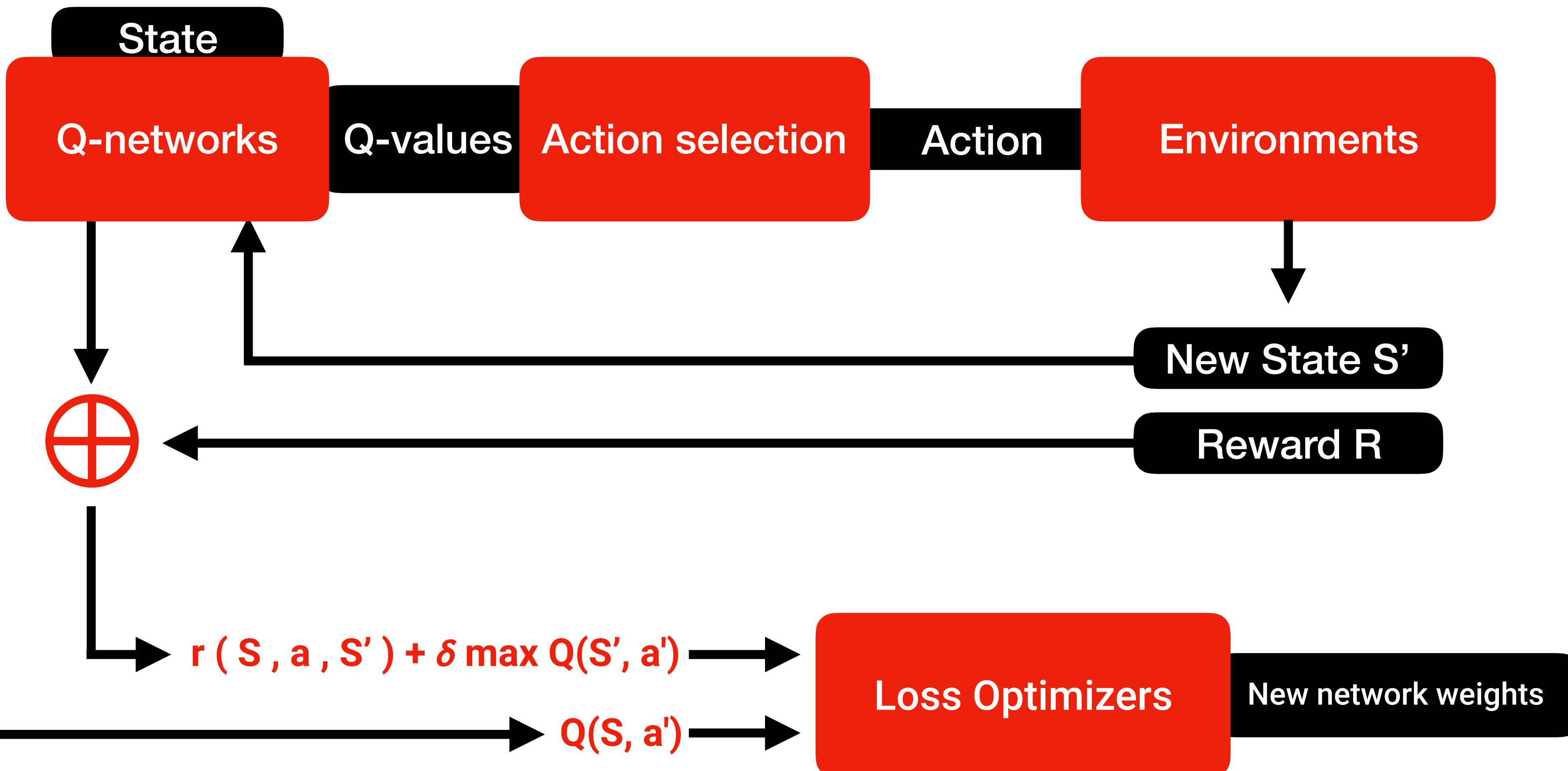


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

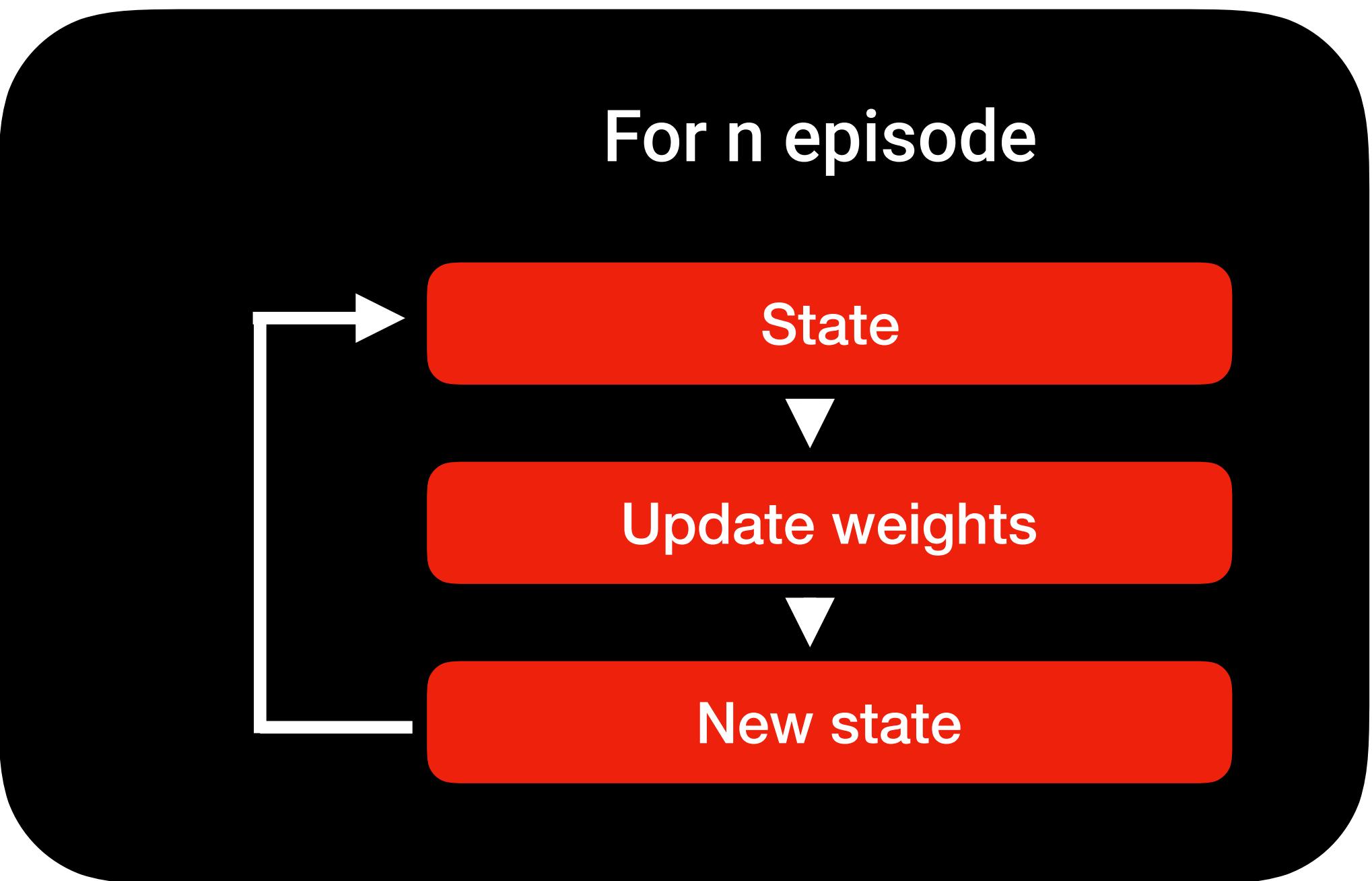


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

$$\Rightarrow Q(S, a) \approx r(S, a, S') + \delta \max Q(S', a')$$

2 - How to optimize the model





We are chasing a moving target



$r(s, a, s') + \delta \max Q(s', a')$

$Q(s, a)$

A cartoon illustration of a brown donkey's head and neck on the right side of the frame. On the left, a wooden stick is angled upwards, from which a single orange carrot hangs by a string. The donkey is looking towards the hanging carrot.

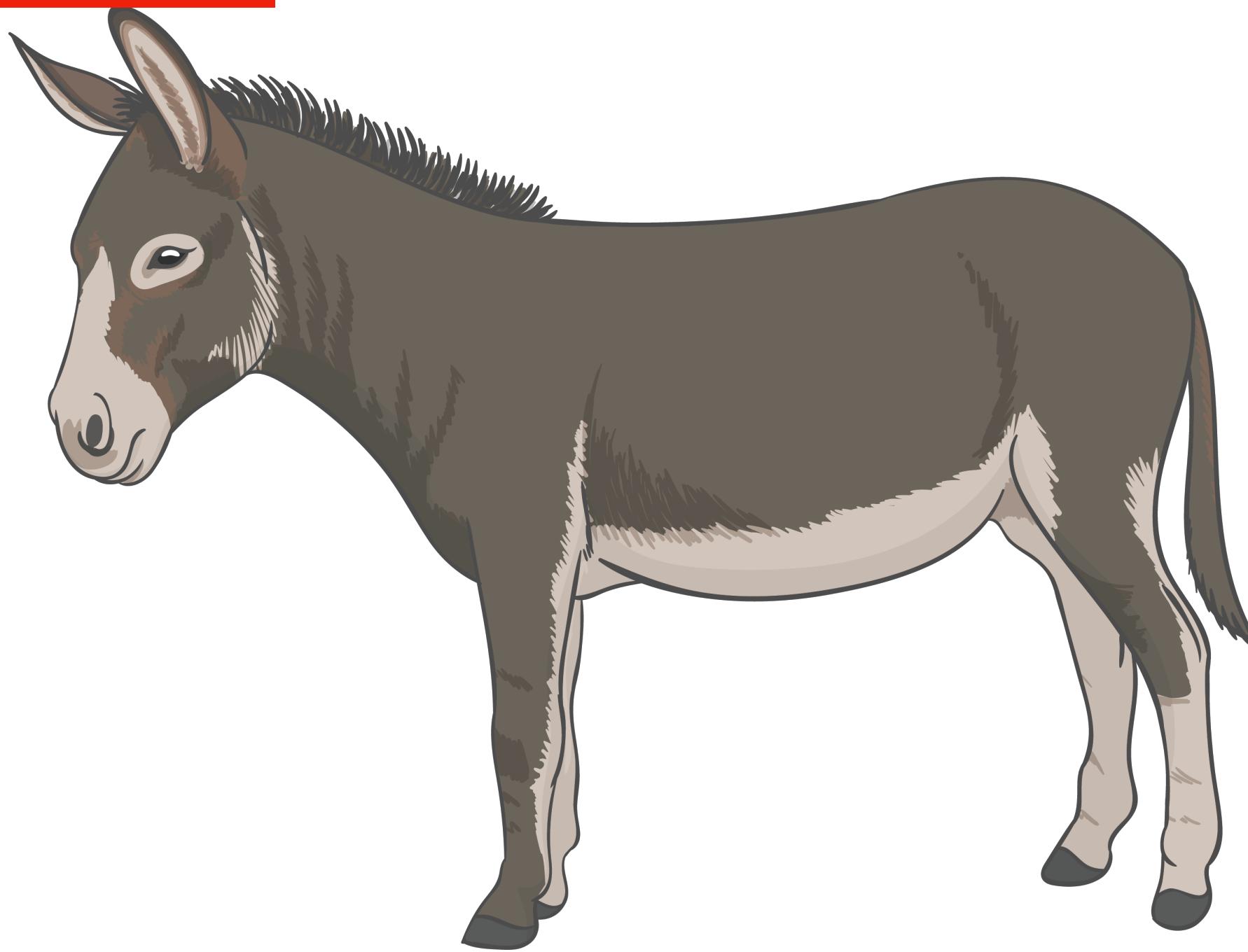
$Q(S, a)$

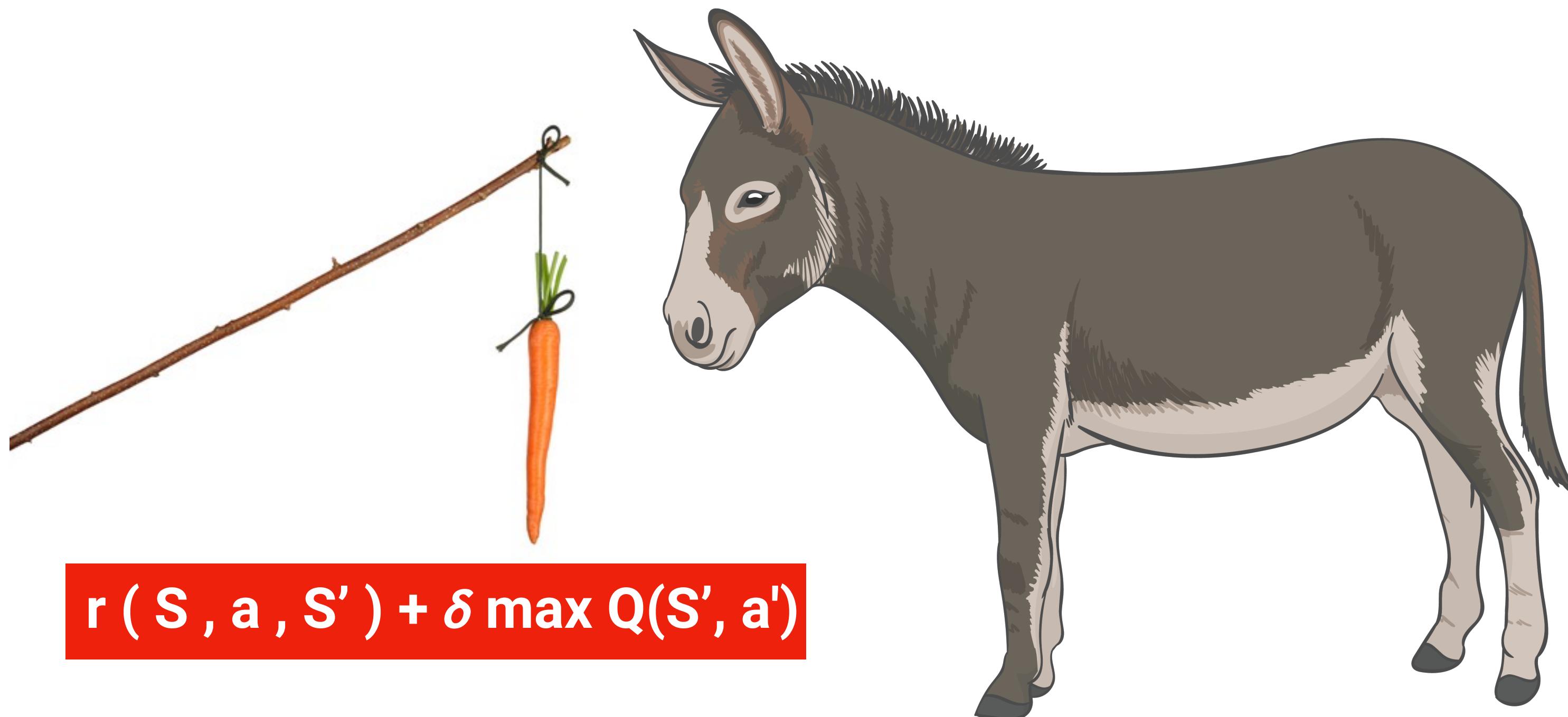
$r(S, a, S') + \delta \max Q(S', a')$



$$r(S, a, S') + \delta \max Q(S', a')$$

$$Q(S, a)$$



$Q(S, a)$  $r(S, a, S') + \delta \max Q(S', a')$

How are we going to update the Q-networks weights

We use a local network and a target network

1 - How to replace the Q-table

2 - How to optimize the model

3 - T-target

How are we going to update the Q-networks weights

We use a local network and a target network

The target neural network get slightly updated based on the local networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model ✓

3 - T-target

How are we going to update the Q-networks weights

We use a local network and a target network

The target neural network get slightly updated based on the local networks weights

The local network predict the Q-values for the current state

1 - How to replace the Q-table ✓

2 - How to optimize the model ✓

3 - T-target

How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

We use a local network and a target network

2 - How to optimize the model ✓

The target neural network get slightly updated based on the local networks weights

3 - T-target

The local network predict the Q-values for the current state

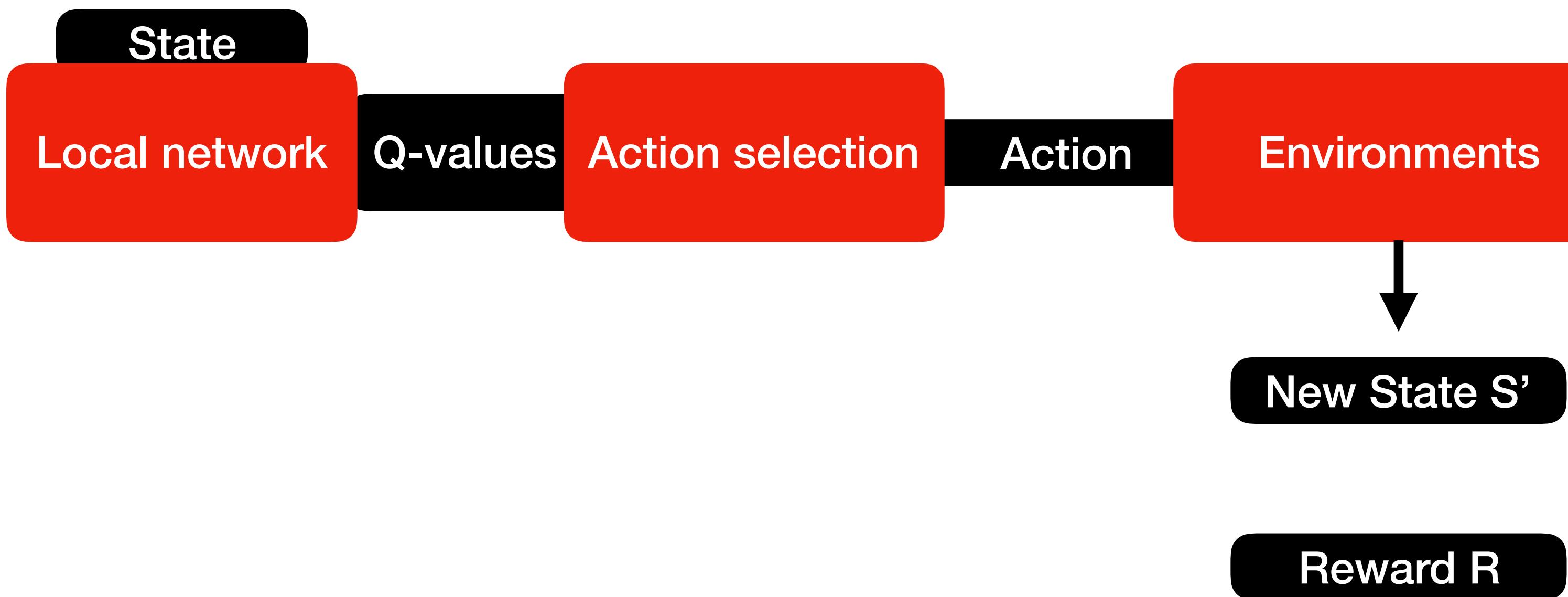
The target network predict the Q-values for the next state

How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model ✓

3 - T-target

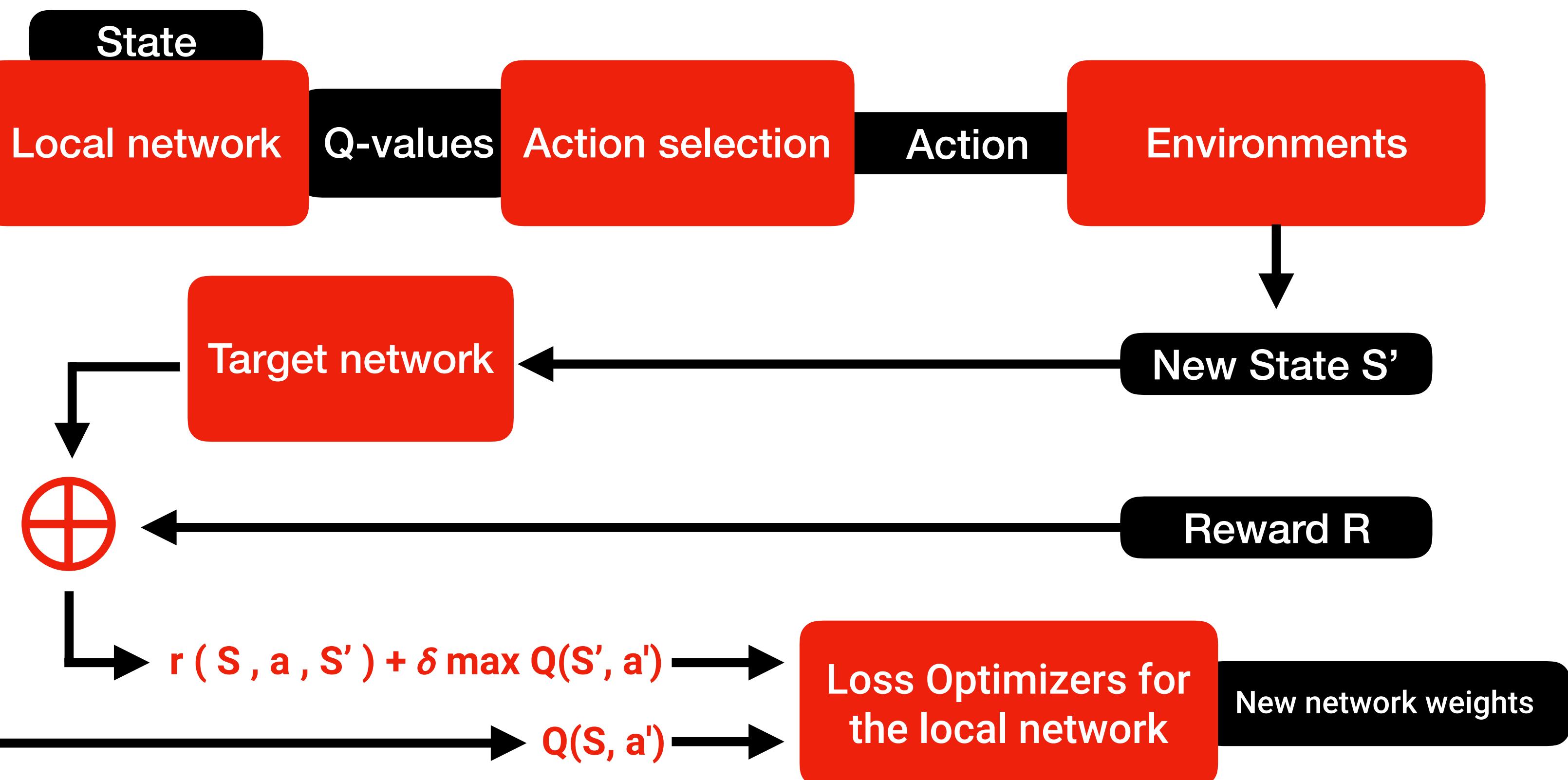


How are we going to update the Q-networks weights

1 - How to replace the Q-table ✓

2 - How to optimize the model ✓

3 - T-target



**Someone to answer the question
please**

In supervised learning, why we shuffling
the dataset gives better accuracy?



Someone to answer the question please

In supervised learning, why we shuffling
the dataset gives better accuracy?

If not shuffling data, the data can be sorted or
similar data points will lie next to each other,
which leads to slow convergence



Someone to answer the question please

In supervised learning, why we shuffling
the dataset gives better accuracy?

If not shuffling data, the data can be sorted or
similar data points will lie next to each other,
which leads to slow convergence



In our case we are having a big
correlation between consecutive states
which may lead to inefficient learning .

How to break the coloration between states

Stack the experiences in a buffer

After the specified number of steps pick an N random experiences

Proceed with optimizing the loss with the learning batch

1 - How to replace the Q-table 

2 - How to optimize the model 

3 - T-target 

4 - Experience replay

How to break the coloration between states

Stack the experiences in a buffer

After the specified number of steps pick an N
random experiences

Proceed with optimizing the loss with the
learning batch

1 - How to replace the
Q-table 

2 - How to optimize the
model 

3 - T-target 

4 - Experience
replay 

Putting it all together

Build a DQL agent



Saha ftourkoum