# Deep Reinforcement Learning AI

Maestracci V., Bardes A., Dezerces E.

January, 2019

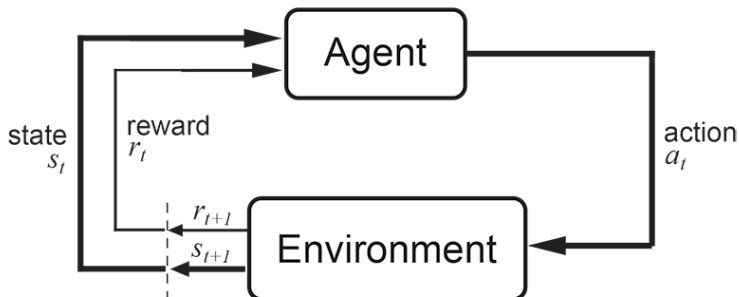# Overview

The principle of Q-learning
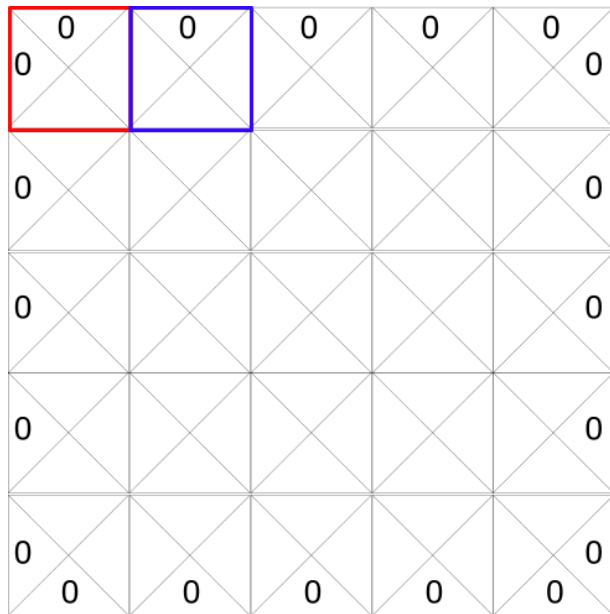
State

Action

Qtable

Q value

Q* learning

State

Deep Q Neural network

Q value action 1

Q value action 2

Q value action 3

Deep Q* learning

$$\underline{\Delta w} = \alpha[(\underline{R + \gamma \, max_a \, \hat{Q}(s', a, w)}) - \underline{\hat{Q}(s, a, w)}] \, \underline{\nabla_w \hat{Q}(s, a, w)}$$

Change in
weights

learning
rate

Maximum possible Qvalue for the
next_state (= Q_target)

Current predicted
Q-val

TD Error

Gradient of our current
predicted Q-value

## Memory

What if our AI goes in level two?

## Memory

What if our AI goes in level two?
It musn't forget about what it learned before !

## Memory

What if our AI goes in level two?
It musn't forget about what it learned before !
So let's save our previous experiences.

# Fixed Q-target

We are chasing a moving target :

# Fixed Q-target

We are chasing a moving target :



How can we make it so that the target doesn't move too much?

# Fixed Q-target

We are chasing a moving target :



How can we make it so that the target doesn't move too much?
We can have two different networks!

## Double DQN

At the beginning of the learning, our q values are noisy.
What if we favor sub-optimal choices ?

## Double DQN

At the beginning of the learning, our q values are noisy.
What if we favor sub-optimal choices ?
$\longrightarrow$ We must decouple the action selection from the q-value generation.
For that we can use our two networks!

$$Q(s,a) = r(s,a) + \gamma Q(s', \underline{argmax_a Q(s',a)})$$

TD target

DQN Network choose
action for next state

Target network calculates the Q
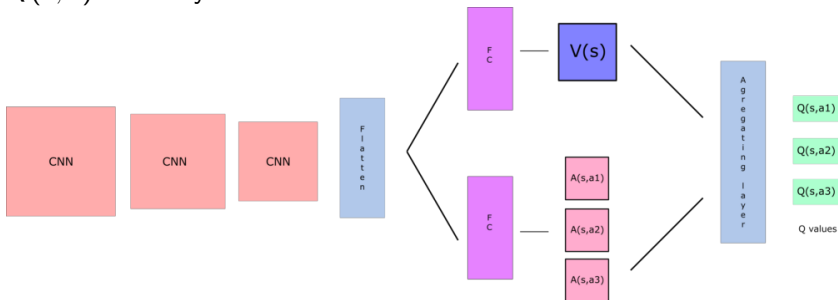value of taking that action at that
state

## Dueling DQN

A Q value represents two things :

- ▶ How good it is to be in state s.
- ▶ How good is it to take the action a in that state.

So we can split it in two : $V$ and $A$

# Dueling DQN

This is the answer to the question : what is the point of knowing $Q(s, a)$ for every $a$ when the state $s$ is bad



Also adds decoupling.

## PER

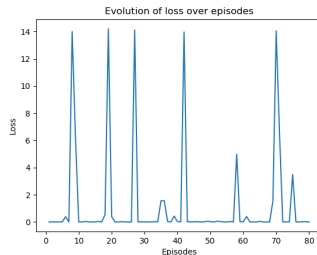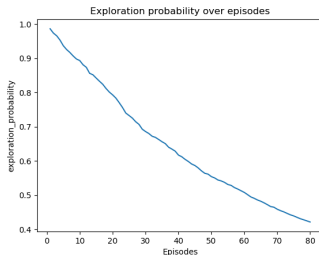Experiences with a huge loss are more important than others.

## PER

Experiences with a huge loss are more important than others. $\longrightarrow$
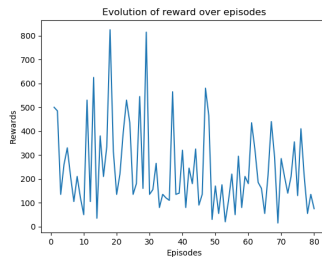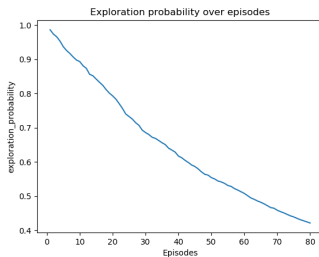So let's prioritize these.

## PER

- Add a non uniform probability to be chosen for experience replay.

- Probability of being chosen decrease if you are often chosen

# Unfortunately...

Our bot isn't really good...

Exploration probability over episodes



Evolution of reward over episodes

Our Github:
https://github.com/Adrien987k/Deep-Space-Invaders