# MultiAgent Systems course project

Federico Vaccaro - Università degli Studi di Firenze

September 2020

**Abstract**

The purpose of this project, is to apply a reinforcement learning technique - the Deep Q-Learning - for 'teaching' to an agent how to reach a point in a given arrival point in a grid-shaped space, while trying to taking the minimum count of steps and avoiding obstacles.

## 1 Introduction

The three main paradigms of the Machine Learning are: **Unsupervised** learning, where the model try to understand the underlying distribution of the dataset without any other input aside from the data themselves; **Supervised** learning, where paired with each sample there is a label, and the learning algorithm tries to fit the distribution that bound each sample to its label; lastly, there is the **Reinforcement** learning, where given an *agent* operating within an environment $\mathcal{E}$, has to 'find' the *best policy*, i.e. the agent should determine a sequence of steps for maximizing a *reward* (and/or minimizing a punishment!) obtained at each time-step. The first difference that we observe, is the ambiguous definition of the problem: there are no dataset and labels, which makes at first glance unclear how formulating a **loss function** for guiding the training. In fact, there are various strategy for tackling the problem: the most known are **value iteration**, **policy iteration** and **Q-Learning**. In this project, we focused on the application of the Q-Learning, however the three listed strategy have in common the property of being formulated starting from the **Bellman equation of Dynamic Programmming**.

## 2 Problem formulation

We consider tasks in which an agent interacts with an environment $\mathcal{E}$, in this case the Atari emulator, in a sequence of actions, observations and rewards. At each time-step the agent selects an action $a_t$ from the set of legal game actions, $A = \{1, ..., K\}$. The action is passed to the emulator and modifies its internal state and the game score. In general $\mathcal{E}$ may be stochastic. The emulator's internal state is not observed by the agent; instead it observes an image $x_t \in \mathbf{R}^d$ from the emulator, which is a vector of raw pixel values representing the current

screen. In addition it receives a reward $r_t$ representing the change in game score. Note that in general the game score may depend on the whole prior sequence of actions and observations; feedback about an action may only be received after many time-steps have elapsed.

As previously stated, goal of the agent is to interact with the emulator by selecting actions in a way that maximises future rewards. We make the standard assumption that future rewards are discounted by a factor of $\gamma$ per time-step, and define the future discounted return at time $t$ as $R_t = \sum_{t'=t}^{T} \gamma^{t'-t} r^{t'}$ , where $T$ is the time-step at which the game terminates. We define the optimal action-value function $Q^{\star}(s, a)$ as the maximum expected return achievable by following any strategy, after seeing some state $s$ and then taking some action $a$, $Q^{\star}(s, a) = max_{\pi} \mathbf{E}[R_t | s_t = s, a_t = a, \pi]$, where $\pi$ is a policy mapping sequences to actions (or distributions over actions).

The optimal action-value function obeys an important identity known as the (previsionale mentioned) **Bellman equation** (1). This is based on the following intuition: if the optimal value $Q^{\star}(s_0, a_0)$ of the sequence $s_0$ at the next time-step was known for all possible actions $a_0$ , then the optimal strategy is to select the action $a_0$, maximising the expected value of $r + \gamma Q^{\star}(s_0, a_0)$

$$Q^{\star}(s, a) = \mathbf{E}_{s' \sim \mathcal{E}}[r + \gamma \max_{a'} Q^{\star}(s', a') | s, a] \tag{1}$$

# 3   Deep Q-Learning

We will use a function approximator (as neural networks, which we refer as *Q-networks*), with parameters $\theta$, to approximate the action-value function.

$$Q(s, a; \theta) \approx Q^{\star}(s, a). \tag{2}$$

A Q-network can be trained by minimising a sequence of loss functions $L_i(\theta_i)$ that changes at each iteration $i$.

$$L_i(\theta_i) = (y_i - Q(s, a | \theta_i))^2$$
$$y_i = \mathbf{E}_{s' \sim \mathcal{E}}[r + \gamma \max_{a'} Q^{\star}(s', a' | \theta_{i-1}) | s, a] \tag{3}$$

The formulation of the loss $L_i(\theta_i)$ in (3) is differentiable, thus we will apply a *stochastic gradient descent* for minimizing Where the parameters $\theta_{i-1}$ are held fixed during the update, othewise, with a gradient update the target would also try to match the Q-network values, instead of the opposite. Note that this method is **model-free**, i.e. it has no prior knowledge/estimation of the emulator $\mathcal{E}$, but it directly observe its "external" state. Also, the method is **off-policy**, because it does not learn by directly sampling from a policy (a sequence of steps), but rather training on different actions, even randomly sampled.

Before introducing the final Deep Q-Learning algorithm, we first need to describe the **replay memory** and the **exploration-exploitation** technique for action sampling.

## 3.1  Replay memory

At each time-step, we collect a tuple from the experience of the agent: $e_t = (s_t, a_t, s_{t+1}, r_t)$. We mantain these experiences in a pool, which we refer as *replay memory*, then as the training continues, we build a dataset of experiences. Then we are abilitated sample mini-batches of tuples from this dataset, making the reinforcement learning algorithm more similar to a "classic" supervised learning task. This is efficient for two reasons: first, it allows to build minibatches of **indipendent samples**, which are mandatory for a better estimation of the gradient via SGD; also it is **computationally-efficient**, since we can leverage the parallel-computing hardware.

## 3.2  Exploration-exploitation

During the training, we let the agent playing multiple matches, from where we collect experiences. However, especially in the first stages of the training, the agent does not know how the environment will change to a chosen action, nor the reward it will receive. Thus, it is reasonable to chose a random action instead of the action with maximum Q-score, in order to **explore** the states of the environment. Furthermore, in the very first stages the parameters are near from being random initialized, then the output Q-scores have really few sense. In order to achieving this, we select a probability $\epsilon$ with we select a random action ($1 - \epsilon$ selects the *greedy* action) during the training.



Figure 1: The Universe

# 4  Conclusion

"I always thought something was fundamentally wrong with the universe" [1]

# References

[1] D. Adams. *The Hitchhiker's Guide to the Galaxy*. San Val, 1995.