

# Report

Lecture „Fundamentals of Machine Learning“ im WS 2018

Juan Antonio Ruiz Leal  
Carlos Alberto Rios Rubiano  
hr229@ix.urz.uni-heidelberg.de  
sy226@ix.urz.uni-heidelberg.de

March 23, 2019

## Abstract

## 1 Reinforcement learning method and regression model

### 1.1 First attempt approach

In order to construct the states, we found three set of crucial parameters to describe the situations in the game. The first set, related to the available cells to move, is constructed by means of an array of four booleans.

### 1.2 States construction:

#### 1.2.1 Available cells array (ACA)

The first set, a boolean one, provides the information of the available cells surrounding the agent. If a bit is on, that means this direction is free to move. We provide some examples for better understanding, the first column indicates the available moves, the second the abbreviation of the direction and the third one the related set. Shown in the next table 1:

Table 1: Available moves, abbreviation and related boolean array.

Available moves	Abbreviation	First set examples
Up	<i>U</i>	(1, 0, 0, 0)
Down	<i>D</i>	(0, 1, 0, 0)
Left	<i>L</i>	(0, 0, 1, 0)
Right	<i>R</i>	(0, 0, 0, 1)
Up/Left	<i>UL</i>	(1, 1, 0, 0)
Down/Left	<i>DL</i>	(0, 1, 1, 0)
Down/Right	<i>DR</i>	(0, 1, 0, 1)
Up/Right	<i>UR</i>	(1, 0, 0, 1)
Up/Left/Down	<i>URD</i>	(1, 1, 0, 1)
.... etc	.... etc	.... etc

### 1.3 States construction: Dysfunctional version

In this section we describe the first attempt we made for the construction of the state, which we had several difficulties with, and we did not achieve the expected solution.

The unsuccessful results are due to the selection of the parameters that make up the state. The second array set were calculated from the regions surrounding the agent as well. We realized that these parameters are correlated with each other, therefore we find solutions that never converged. In order to show that, we define the regions definition in the next subsection:

#### 1.3.1 Region definition

In this subsection, we show the definition of the regions. A region is defined based on the current position of the agent and the possible directions where the agent can go to. We defined 8 possible regions, Up, Down, Left, Right, Up/Left, Down/Left, Down/Right and Up/Right. The regions we described before are depicted in figure 1.

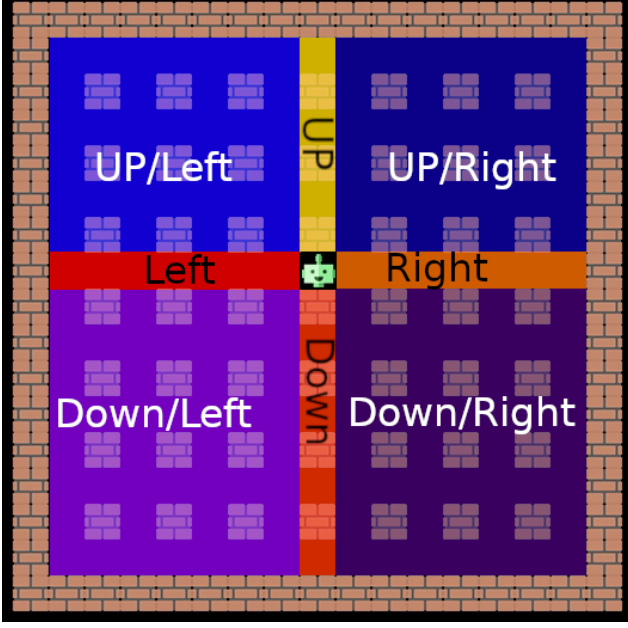


Figure 1: Definition of the regions in the maze.

### 1.3.2 Normalized potential rewards

We show the regions regarding the second set, which make up the state. Each element, have the *normalized potential rewards* (NPR) in each region, by mean of the  $\omega_{R_*}$ . And the weights  $\omega_{R_*}$ , are elements to measure the potential reward to acquire, in all those eight regions described in figure 1.

Table 2: How looks like the second array. Where the  $\omega_{R_*}$  is the weight related to the NPR in each region.

Second array of the state:
$(\omega_{RU}, \omega_{RD}, \omega_{RL}, \omega_{RR}, \omega_{RUL}, \omega_{RDL}, \omega_{RDR}, \omega_{RUR})$

Where:

$$\omega_{R_*} \in [0, 1] \quad (1)$$

### 1.3.3 Normalized potential danger

The third set, is similar to the second one. Each element takes into account the *Normalized Potential Danger* (NPD) in each region. And each weights  $\omega_{D_*}$ , is a measure of the potential danger in all those eight regions.

Similary to the table 2, we show how looks like the NPD float array.

Table 3: How looks like the second array. Where the  $\omega_{R_{D_*}}$  is the weight related to the NPD in each region.

How looks like the second array of the state:
$(\omega_{DU}, \omega_{DD}, \omega_{DL}, \omega_{DR}, \omega_{DUL}, \omega_{DDL}, \omega_{DDR}, \omega_{DUR})$

Where:

$$\omega_{R_*} \in [0, 1] \quad (2)$$

### 1.3.4 Drop a bomb (DB): Feasibility and usefulness

To complete the state, we additionally add an array of two float values, regarding the situations where is feasible and useful to drop a bomb. The first value is a measure of the surrounding situation with the crates, and the second value is a measure of the proximity of the opponents.

Table 4: How looks like the third array. Where the  $\omega_{R_{D_{Cr}}}$  is the weight that measure the usefulness to get coins by blowing up crates, and  $\omega_{R_{D_O}}$  the feasible of killing an opponent

How looks like the fourth array of the state:
$(\omega_{B_{Cr}}, \omega_{B_O})$

Where:

$$\omega_{B_*} \in [0, 1] \quad (3)$$

### 1.3.5 Summary of state components

In order to summarise the construct of the state, we describe the features selected. Which is made up of the four arrays shown in tables 1,2,3 and 4:

## 1.4 Second attempt: Simplifying the state

In order to avoid correlated features, and achieve a functional version to complete the task 1 (6.4). We simplify the NPR, in a binary version.

The new NPR array is made up by four Booleans, where we calculate the nearest coin direction to operate the ACA with.

Using the *AND* operator  $\wedge$ . We construct the *NPR* array, a picted example is showed in figure 2. And in case of  $ACA \wedge DPR = (0, 0, 0, 0)$ , we chose randomly the *NPR* form the composed *ACA* possibilities

## 1.5 Observations

### 1.5.1 Rewards

The rewards for the agent should be evolutive, that means, the agent should learn first to break the state of stillness, then it should acquire the capacity to move, then dropp bombs, survive the bombs that it has dropped, collect the coins and finally fight against their opponents. In this process the rewards play an important role. For example at the beggining the reward of "WAIT" should be negative to encourage the agent to move. The re

Table 5: State summarized: Unsuccessfully attempt

Abbreviation	Description	Elements	Info	Type
ACA	Availible cells	4	Directions	boolean
NPR	Potential rewards	8	Regions	float
NPD	Potential danger	8	Regions	float
DB	Feasibility of throwing a bomb	4	Info crates and opponents	float

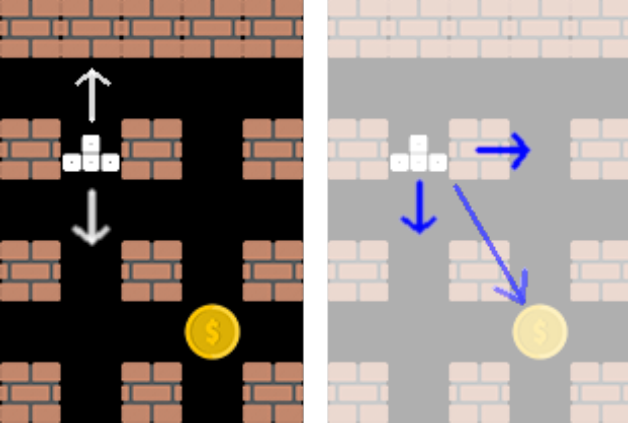


Figure 2: Example of state formation, of ACA and NPR simplified: In the left side the available moves is showed, then ACA is: (1100), and in the right side, the direction of probable reward (DPR), and is given by: (0101). from operate with the  $\wedge$  operator, we get the  $NPR = (1100) \wedge (0101) = (0100)$

## 2 Training process

## 3 Experimental Results

### 3.1 Task 1:

On a game board without any crates, collect a number of revealed coins as quickly as possible. This task does not require dropping any bombs. The agent should learn how to navigate the board efficiently.

For this task we use the state made up by the **ACA** array (see section 6.1.1) and the simplified **NRP** (see second 6.3). Summarised by the table 6.

For this task, the agent must to learn two different skills. The first one, is the skill to avoid the invalid actions, and the second one is to choose efficiently the path to the coins.

With the aim to know the learning evolution of the training, we calculate the total reward accumulated in all the training realization. At the beginning of the training the amount of random actions predominate and the total reward accumulated start to decrease, later at the 150 episode, the agent will perform better actions, and get positive rewards. Then, the total re-

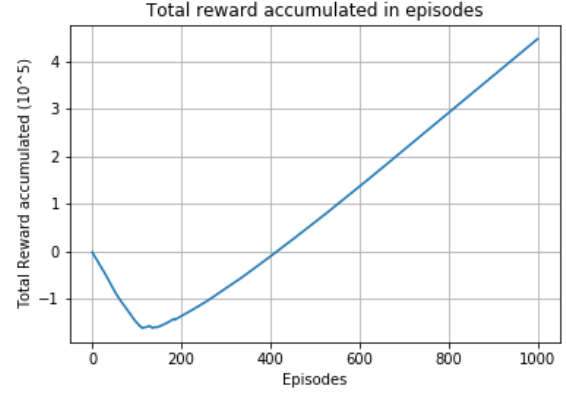


Figure 3: Total reward accumulated in episodes

ward accumulated increase. See figure 3.

This behavior becomes clarified by mean of the figure 4. We shows, a convergence tendency after episode 400. Up to get a maximal rewarded performed of 828 (green)

In order to measure those skills, we show [1]

### 3.2 Task 2:

On a game board with randomly placed crates, find all hidden coins and collect them within the step limit. The agent must drop bombs to destroy the crates. It should learn how to use bombs without killing itself, while not forgetting efficient navigation.

We tried to face the task 2 with the same strategy, but the problem was that when we performed the AND operator over the *ACA* and *DPR* the result could be (0, 0, 0, 0), because the nearest coin could be in a region where the direction was blocked, which again caused that the model did not converge to a desired result.

We solved this problem by means of getting the nearest coin with respect to an unblocked direction and re-defining *DPR* by turning on the bit in this direction.

Another problem we found when solving the task 2 was that we could not punish too much the action of killed by itself, because since the agent could die by dropping a bomb the agent avoided that action, which led to a situation where the agent survive the entire episode but without collecting any coin. So we decided

Table 6: State summarized

Abbreviation	Description	Elements	Info	Type
ACA	Avalible cells	4	Directions	Boolean
NPR simplified	Nearest coin	4	Directions	Boolean

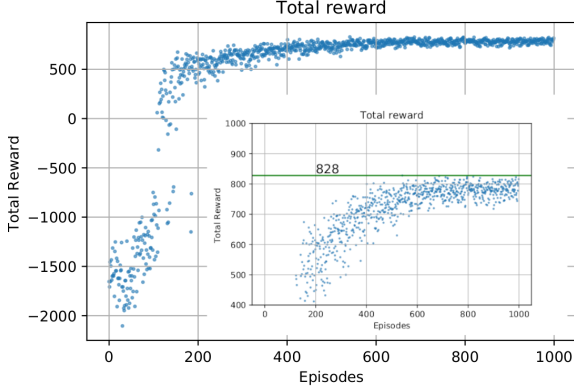


Figure 4: Total reward in episodes. Plot inside: detailed behavior (in order to show the total reward tendence))

to does not punish the action of killed by itself.

## 4 Improving the agent and feedback

### 4.1 Improving the agent

We learned many things in the development of the project, among the most importants are, try the easy solutions first and make sure that these approaches work, then continue adding more complicated situations little by little, model them, and make sure that nothing is broken, because we discarded the simple solutions without even prove that they worked and when we tested the implementation of the complicated models and these didn't work we returned to previous ideas.

### 4.2 Feedback

When we were training the model we realized that it could be handy to have the possibility to send parameters in the main to the program to speed up the the set up of the hyperparameters.

Another thing that could be improved is to have one not toy but easy example in the assignments or in the class to have a better understanding about the designing of the states, because we were pretty lost with this part at the begginig of the project.

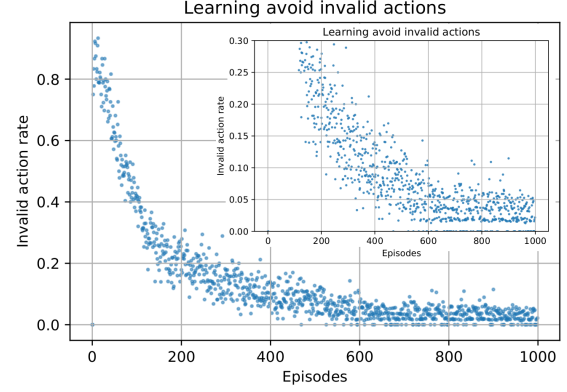


Figure 5: Measuring the learning to avoid invalid actions. Plot inside: detailed behavior (in order to show the xtendence))

### 4.3 Q-learning with regression forest prediction

### 4.4 Neural Networks

### 4.5 General definitions:

#### 4.5.1 Epsilon decay

The exploration rate is defined by mena of the follow parameters:

```
// Code
#Initial exploration rate when training
self.epsilon = 1.0
# Exploration steps
self.exploration_steps = 2000000
# Minimum value of epsilon , after this value
# epsilon does not decrease anymore
self.epsilon_min = 0.1
# This hyperparameter is to decrease the num
# of explorations as the agent gets better
self.epsilon_decay = (self.epsilon - self.ep
```

Then the *exploration\_steps* is the parameter to fit in each realization. Were *epsilon* decrease by *self.epsilon* - = *self.epsilon\_decay* with every step.

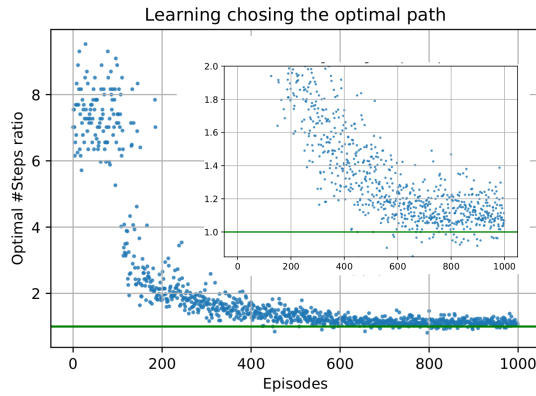


Figure 6: Measuring the learning find the optimal path to get the coins. Plot inside: detailed behavior (in order to show the tendency))

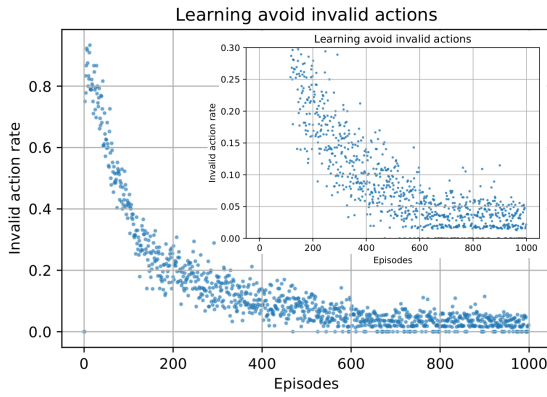


Figure 7: Measuring the learning to avoid invalid actions. Plot inside: detailed behavior (in order to show the tendency))

## 5 Motivation and overview

## 6 Model

We basically tried to use 2 models, the normal Q learning and the deep Q learning. The model used for the training was constructed based on the algorithm known as Deep Q Learning [1]

### 6.1 Algorithm Deep Q Learning

```
// Code
stringToMatch = 'Score'
matchedLine = ''

listScore = []
listPasos = []
```

```
listEpsil = []
listRewaA = []
listQMean = []
```

## References

- [1] MNIH, Volodymyr ; KAVUKCUOGLU, Koray ; SILVER, David ; GRAVES, Alex ; ANTONOGLU, Ioannis ; WIERSTRA, Daan ; RIEDMILLER, Martin A.: Playing Atari with Deep Reinforcement Learning. In: *CoRR* abs/1312.5602 (2013). <http://arxiv.org/abs/1312.5602> 3, 4