# Report

Lecture „Fundamental 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 19, 2019

## Abstract

# 1 Motivation and overview

# 2 State

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 de avaliable cells to move, is constructed by mean of an array of four boleans.

## 2.1 States construction:

### 2.1.1 Avaliable cells array (ACA)

The first set, a boolean one, provides the information of the available cell surrounding the agent. We describe the first one, by mean of some examples of the available moves and the related set. Showed in the next table 1:

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

| Avaliable moves | Abbreviation | First set examples |
|---|---|---|
| Up | $U$ | $(1,0,0,0)$ |
| Down | $D$ | $(0,1,0,0)$ |
| Left | $L$ | $(0,0,1,0)$ |
| Right | $R$ | $(0,0,0,1)$ |
| Up/Left | $UL$ | $(1,1,0,0)$ |
| Down/Left | $DL$ | $(0,1,1,0)$ |
| Down/Right | $DR$ | $(0,1,0,1)$ |
| Up/Right | $UR$ | $(1,0,0,1)$ |
| Up/Left/Down | $URD$ | $(1,1,0,1)$ |
| .... etc | .... etc | .... etc |

## 2.2 States construction: Dysfunctional version

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

The unsuccessful results, are related with the selection of the parameters that makes up the state. Because the second array set were calculated from the regions surrounding the agent. And those parameter are correlated. And for this reason we find solutions that never converged. In order to show that first idea, we define in the next subsection, the regions definition:

### 2.2.1 Region definition

In this subsection, we show the definition of the regions. The regions where we describe all possible surrounding regions in connection with the actual situation of our agent. In figure 1, are depicted this regions.

### 2.2.2 Normalized potential rewards

And in table 2 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: |
|---|
| $(\omega_{R_U}, \omega_{R_D}, \omega_{R_L}, \omega_{R_R}, \omega_{R_{UL}}, \omega_{R_{DL}}, \omega_{R_{DR}}, \omega_{R_{UR}})$ |

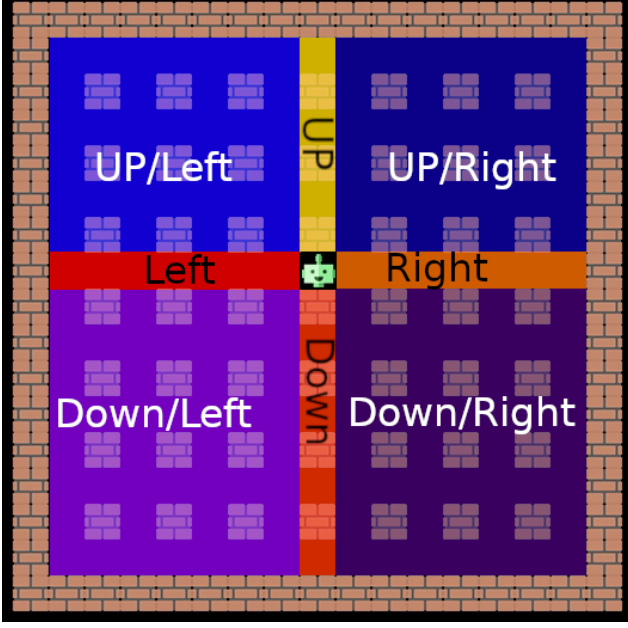Where:

$$\omega_{R_*} \in [0,1] \tag{1}$$

Figure 1: Definition of the regions in the maze.

### 2.2.3 Normalized potential danger

The third set, is similar to the second one. Each element, takes into acount the *Normalized Potential Danger* (NPD) in each region. And each weights $\omega_{D_*}$, is a measure of the potential danger in all those eigth 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: |
| --- |
| $(\omega_{D_U}, \omega_{D_D}, \omega_{D_L}, \omega_{D_R}, \omega_{D_{UL}}, \omega_{D_{DL}}, \omega_{D_{DR}}, \omega_{D_{UR}})$ |

Where:

$$\omega_{R_*} \in [0, 1] \tag{2}$$

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

And 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
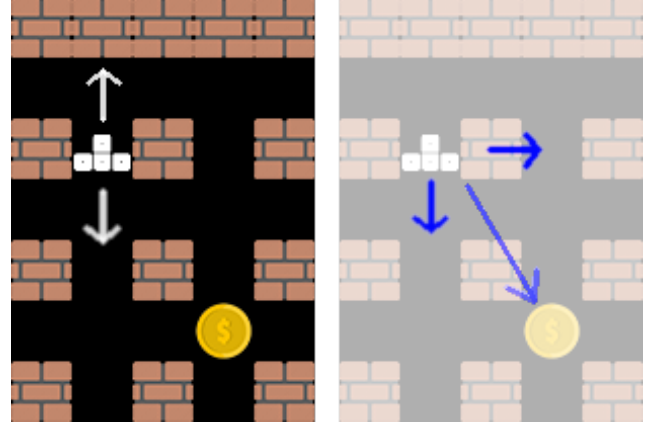


Figure 2: Example of state formation, of ACA and NPR simplified: In right side the avalible moves is showed, then ACA is: (1100), and in the left side, the direction of probable reward (DPR), and is given by: (0101). from operate with the $\wedge$ operator, we get the $NPR = (1100) \wedge (0101) = (0100)$

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

Where:

$$\omega_{B_*} \in [0, 1] \tag{3}$$

### 2.2.5 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:

## 2.3 Second attempt: Simplifying the state

With the aim to avoid correlated features, and achieve a functional version to complete the task 1 (2.4). 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

## 2.4 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 effciently.*

Table 5: State summarized: Unsuccessfully attempt

| Abbreviation | Description | Elements | Info | Type |
|---|---|---|---|---|
| **ACA** | Avalible 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 3: Total reward accumulated in episodes



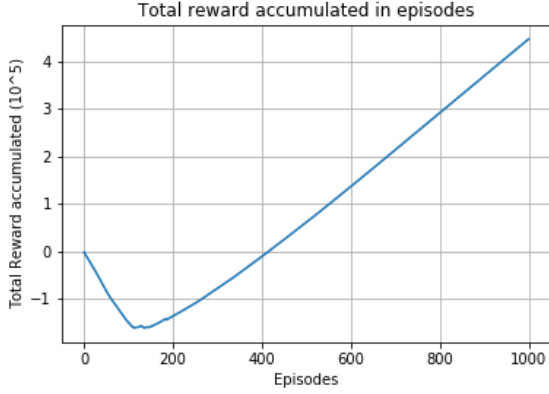Figure 4: Total reward in episodes. Plot inside: detailed behavior (in order to show the total reward tendence))

For this task we use the state made up by the **ACA** array (see section 2.1.1) and the simplified **NRP** (see second 2.3). Summarised by the table 6.

For this task, the agent must to learn two diferent skills. The fisrt one, is the skill to avoid the invalid actions, and the second one is to chose effciently the path to the coins. In order to measure those skills, we show [1]

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

    listScore = []
    listPasos = []
    listEpsil = []
    listRewaA = []
    listQMean = []
```

# References

[1] Mnih, Volodymyr ; Kavukcuoglu, Koray ; Silver, David ; Graves, Alex ; Antonoglou, 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
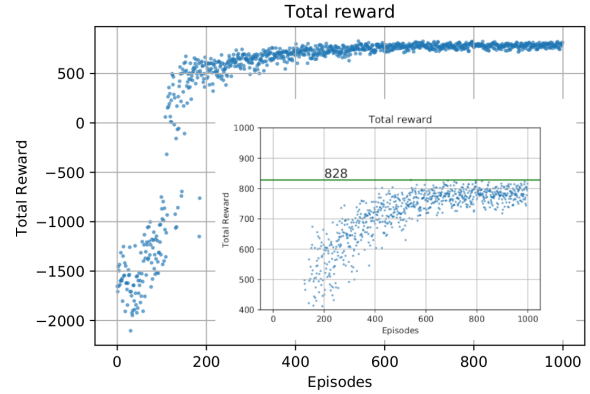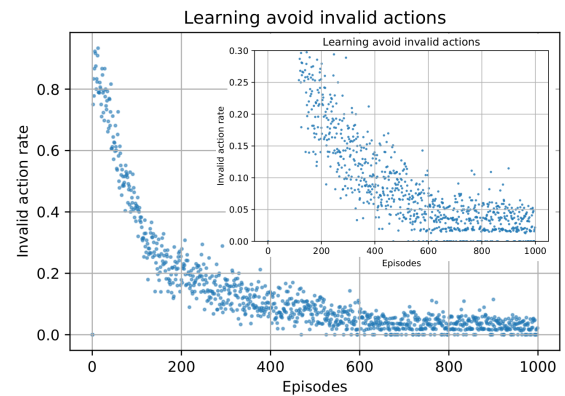
Figure 5: Meassuring the learning to avoid invalid actions. Plot inside: detailed behavior (in order to show the xtendence))

Table 6: State summarized

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