

Reinforcement Learning - Project report

RL agent for Bomberman

Robin Charleuf
Hugo Thevenet

April 5, 2024

1 Introduction

This project aims to develop a reinforcement learning (RL) agent capable of playing the classic game Bomberman at a competitive level. Bomberman is a known game where the agents are in a maze and they can drop bombs. The goal is to kill the other agents in the explosion and to be the last survivor. In the real Bomberman game, some blocks in the maze are destroyable. To simplify, we did not include this feature.

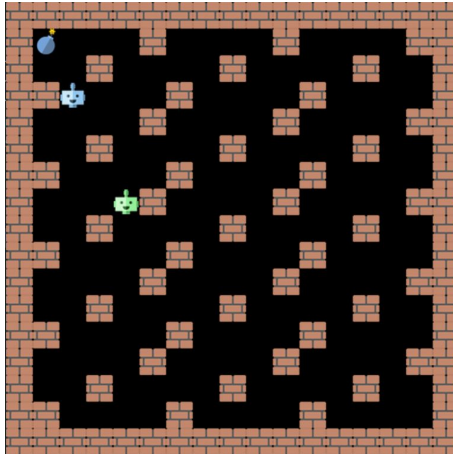


Figure 1: Our implementation of the game with two agents

2 Description of the environment

2.1 Introduction

We created the entire environment ourselves, without copying any repository. We only utilized the .png icons from an existing environment. Our game features two modes: a display mode and a classical mode. Check our Github repository [here](#).

The display mode employs Pygame for rendering. To optimize code execution speed, we added a version without rendering.

Essentially, all state information is stored in a STATE dictionary. Each agent retrieves an observation from this dictionary to make decisions about actions. Subsequently, the state is updated based on these actions. If the display is set to True, we utilize the dictionary to render everything on the screen using Pygame.

Below is the pseudo code for the main loop of our program:

Algorithm 2 Pseudo code of the main script

```
1: while running do
2:   if display then
3:     Use render functions to display the state
4:   end if
5:   for each agent do
6:     Get its observation, from this observation choose an action
7:   end for
8:   Perform the actions for all agents. Get the new state and the rewards
9:   for each agent do
10:    Update its policy with the state, the new state, and its reward
11:   end for
12:   if done then
13:     running  $\leftarrow$  False
14:   end if
15: end while
```

2.2 State Space

The state space is a dictionary containing all of these elements. It represents exhaustively everything useful in the game.

The state space contains:

- **Data Agents:** This contains information about the agents present in the environment. Each agent has attributes such as its position (**agent_x**, **agent_y**), whether it's alive (**alive**), the number of bombs available (**bombs_available**), the type of agent (**type_agent**), and a unique identifier (**agent_id**). For example:
- **Data Bombs:** This represents the bombs present in the environment. Each bomb is associated with the agent that placed it (**agent_id**), the time until explosion (**bomb_time**), and its position (**bomb_x**, **bomb_y**). For example:
- **Data Maze:** This represents the layout of the maze environment. It consists of a grid where '#' represents a wall and ' ' represents an empty space.
- **Data Explosions:** This represents explosions caused by bombs. Each explosion is associated with the agent that caused it (**agent_id**), the time until the explosion ends (**explosion_time**), and the zone affected by the explosion (**explosion_zone**). For example:

2.3 Observation space

To make a decision, the agent has to observe the state. We encode the STATE dictionary into a matrix.

So we created a function used by the agent. This function takes the observation dictionary as input and encodes the maze environment into a numerical representation, in the personal point of view of the agent.

Tile Type	Encoded Value
Wall ('#')	1
Empty Space	0
Bomb	4
Explosion	5
Enemy	3
Agent Position	2
Agent on Bomb	6

With this section, we understand that the observed state can be very high dimensional for the agent. So, we'll have to create complex agent which could handle this high dimensional state as an input.

2.4 Action space

Our action space has a length of 6:

Action	Description
0	Move Left
1	Move Right
2	Move Up
3	Move Down
4	Drop Bomb
5	No Operation

2.5 Rewards

To give the rewards, we need to get some information about the agent and its self awareness. In order to make an informed decision, the agent constantly monitors its own actions (move, do nothing, or drop a bomb), calculates the Manhattan distances to the closest living enemy, the nearest bomb, and the closest point of an explosion, and keeps track of the elapsed game time to guide its strategy within the environment. This is achieved through the following reward system

Action wanted	Reward	Points awarded
Drop bomb close to the enemy	Positive reward when approaching the enemy , based on the difference of distance between two steps	$20 * (\text{new closest distance to enemy} - \text{previous closest distance to enemy})$
	Positive reward when dropping the bomb close to the enemy	$\text{Max}(0, 5 - \text{bomb distance to enemy}) * 100$
Dodge bombs & explosions	<ul style="list-style-type: none"> Positive reward when distance between the agent and the bomb increases and the agent is too close to the bomb, negative if the distance decreases. Same idea for the explosions 	If distance to bomb < 5 / explosion < 3 : $(\text{new closest distance} - \text{previous closest distance}) * 100$
Perform actions that he's allowed to do	Positive reward if the agent can perform the action, negative otherwise	± 100

Figure 2: Reward system

3 Agents implemented and results

For the experiments, we trained each of our agents like this :

- One agent is the learning agent. It's the implemented agent.
- The competitor agent isn't a learning agent but an randomly moving agent. It is the target.

3.1 1st step : the tracking agent

In this section, we present our first implemented version. We struggled to tune the rewards policy, so as a first try, we implemented a agent that follows the other agent. The closer he is to the other ones, the highest rewards he gets. We also made the maze smaller so even a Q-learning agent could learn faster.

3.1.1 Q-Learning

We implemented a simple Q-Learning agent with an exploration factor of 10%. Even with a smaller maze, the Q-Learning agent takes a large number of environments to learn given the complexity of the game, but exhibits a clear learning pattern, with total rewards increasing throughout the learning.

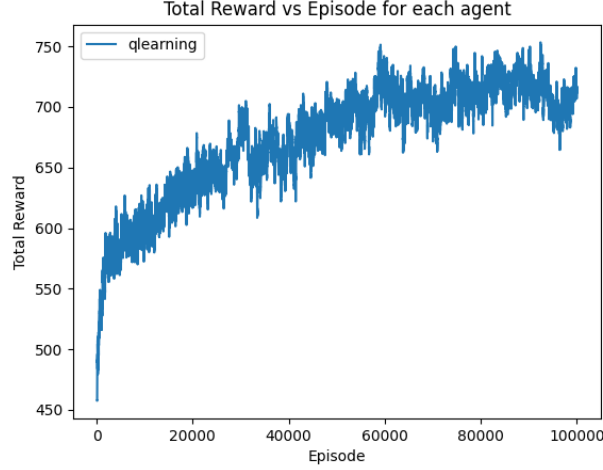


Figure 3: Follower Q-Learning agent rewards by episodes

3.1.2 DQN

We implemented a DQN agent with the following parameters:

- 4 layers of sizes (102,96),(96,64),(64,32),(32,6) and ReLU activation layers in between.
- 5% exploration
- Adam Optimizer with $\gamma = 0.1$ and a 0.005 learning rate.



Figure 4: Follower DQN agent rewards by episodes

This agent learns much faster than the Q-Learning agent, which was expected as the DQN agent learns to approximate the Q

3.2 2nd step : the bomber agent

3.2.1 Q-Learning

Given the complexity added by adding the possibility to drop a bomb, our Q-Learning agent wasn't successful, which could be expected given the very high number of episodes needed only for the follower Q-Learning Agent. This is one more proof that Q-Learning reaches limits when complexity increases.

3.2.2 DQN

We implemented a DQN agent with the following parameters:

- 4 layers of sizes (102,96),(96,64),(64,32),(32,6) and ReLU activation layers in between.
- 5% exploration
- Adam Optimizer with $\gamma = 0.1$ and a 0.005 learning rate.

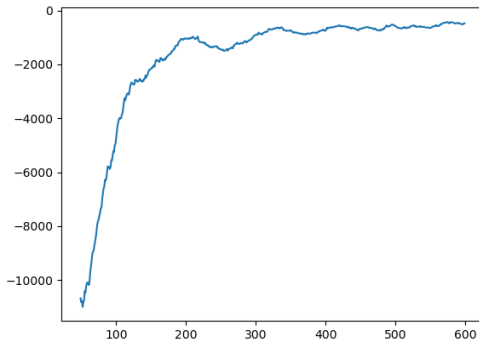


Figure 5: Bomber DQN agent rewards by episodes

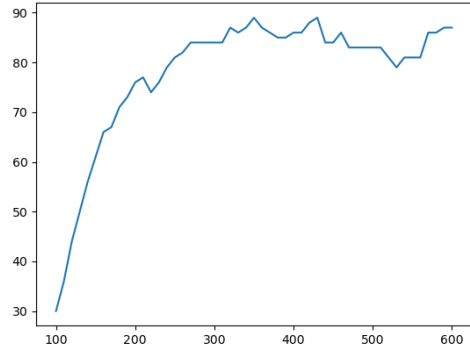


Figure 6: Win rate over the last 100 games by episode

Our agent exhibits a clear learning pattern, showing capacity to learn how to bomb a random agent in the maze, with increasing win rate over the episodes. However, it appears to reach a plateau, at around 85% win rate. We believe it to be due to the fact that our agent lacks of a true "understanding" that bombing the random agent will give him rewards, and rather has incentives to do things that will ultimately lead to this goal.

4 Conclusion

In this project, we created a scratch-built play environment. The game is based on a dictionary, which is updated every turn. This dictionary serves as a database for the agents' display and decision-making. We then created a simplified game environment with one agent simply chasing the other. This simpler system enabled us to code our agents efficiently and develop a first version of the reward system. Next, we added the ability for an agent to drop a bomb. We trained our agents to kill a random agent. To do this, we had to adapt the reward system. We obtained very good results with our agents. If we'd had to continue with the project, we would have been able to pit several types of intelligent agent against each other, and not just one agent against a random agent.

Our code and environment is completely original and was created by us from scratch. It is available [here](#)