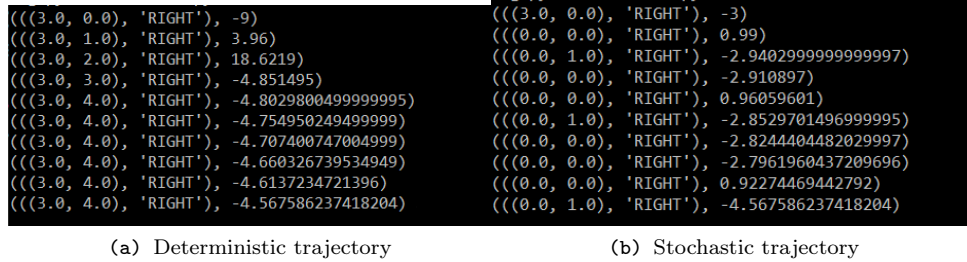


## 1 General structure of the code

In the following report, the basic chosen policy is to move always right. All the other basic movement (always move left, up,down) and a uniform random policy can be called for section 2 and section 3. Three different classes are used. The agent class contains all the data needed to go from one state to another and information about the current state. The grid class contains the environment in which the agent is. The game class is the manager code for both classes

## 2 Section 2

The trajectory of the always go right policy is given by :

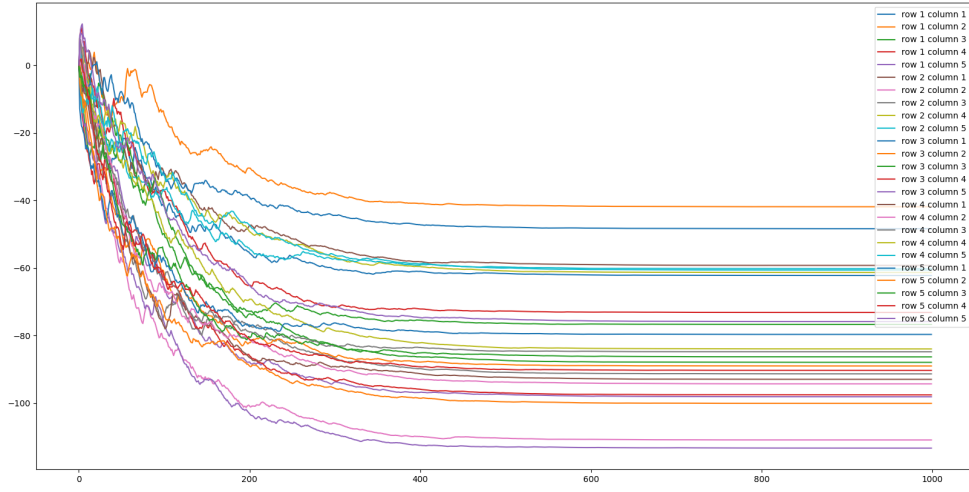


**Figure 1** – Evolution of the trajectory for 10 steps starting at position (3,0) in deterministic and stochastic case ( $\beta = 0.5$ ). Trajectory contains the position before the move, the chosen move and the reward it got after the move

It can be seen that as expected in the deterministic case, the agent moves to the right and stays in the position (3,4) as he can not go further right. Whereas in the stochastic case, a probability of 0.5 sends the agent to position (0,0) when it wanted to go right.

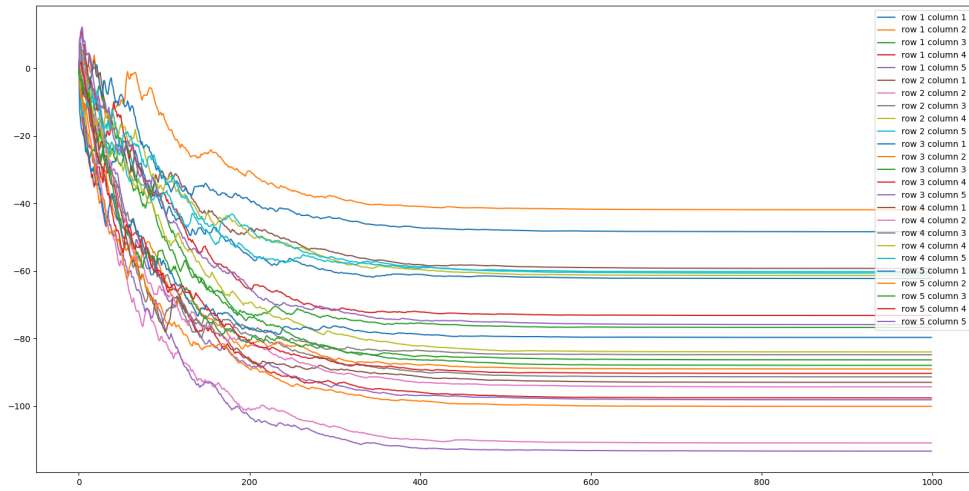
## 3 Section 3

Similarly to section 2, the always right. In this section, the agents reward is simulated 10 times for 1000 movement and the mean of this result is plotted for each starting position. The value 1000 for the number of timestep to have a good estimation was chosen arbitrarily in the beginning but , as can be seen in Figure 2 and 3,  $N = 500$  could also have been a good choice. Before 500, there are a lot of oscillations still and after that the rewards do not evolve anymore.



**Figure 2** – Deterministic evolution of the rewards over 1000 timesteps for each starting position

In Figure 2, it can be seen that for a deterministic always right policy the column at which the agent starts has a little influence on the expected final reward received as the agent will stay in column 4 once it is there and the rewards to get from the initial state to the column 4 are also different. Whereas, the row has a huge influence on the expected final reward as we can clearly see. If another policy such as always up was used the column would have a huge influence on the final result and the influence of the row would be lesser.



**Figure 3** – Stochastic evolution of the rewards over 1000 timesteps for each starting position

In the stochastic case, the starting position does not influence much the expected final reward. For each action that is taken, we have a probability of  $\beta = 0.5$  to end up in state  $(0,0)$ . As we have,  $\sum_{n=1}^{\infty} (\frac{1}{2})^n = 1$  the agent will always end up in the  $(0,0)$  position which explain why in this case the initial starting position is not that important.

## 4 Section 4