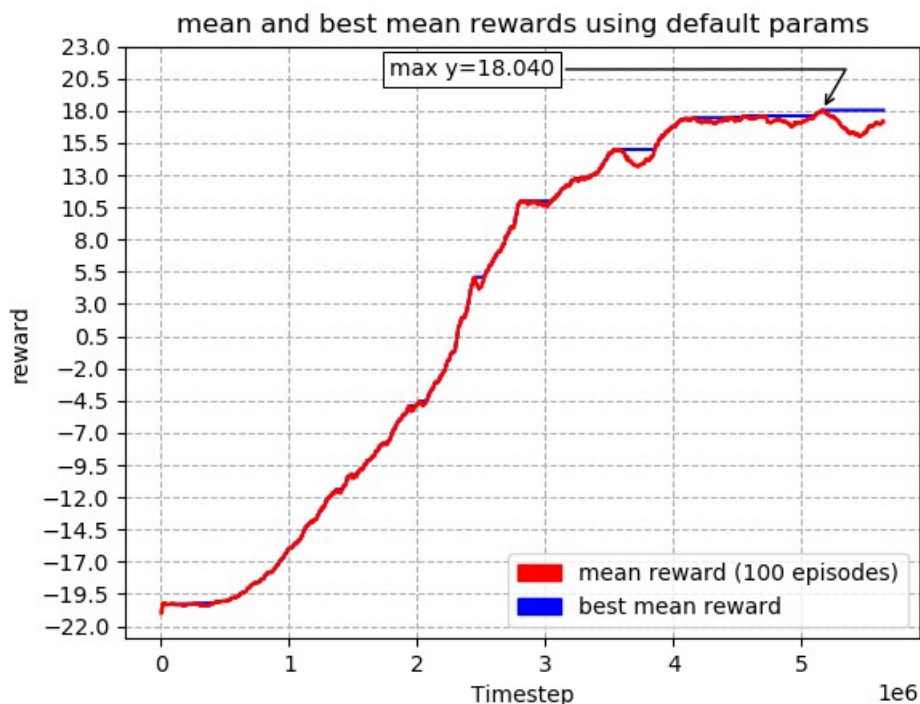


RL Project Report

057931354

Question 1:

See the plot below. Reached score (best mean reward) of 18.04 with default hyperparameters. Number of steps > 4 million.



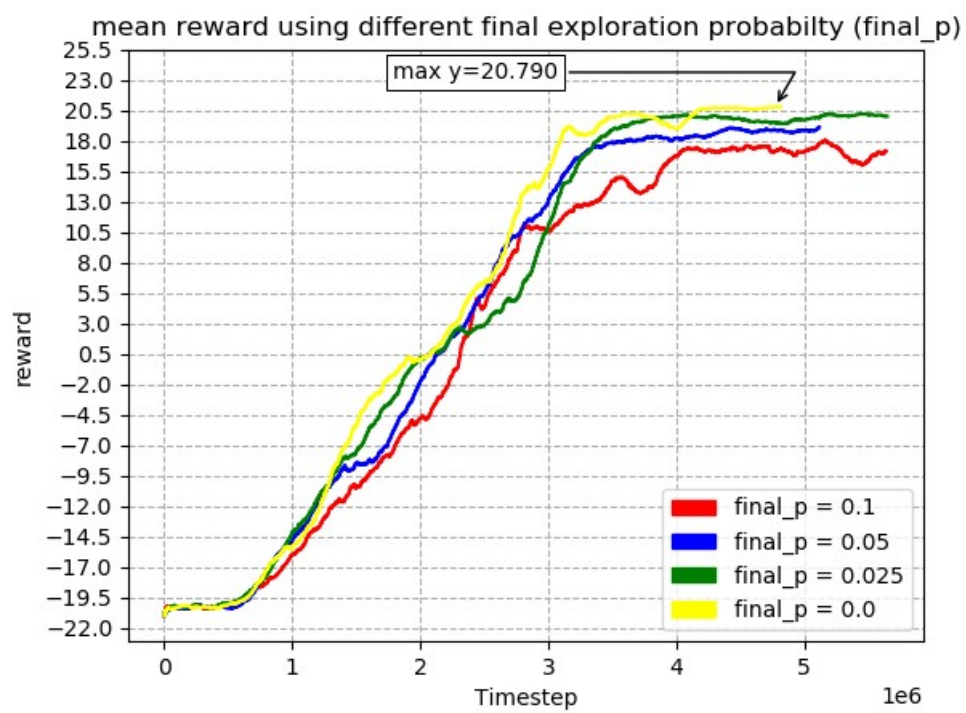
Question 2:

See the two plots below where the first one shows the mean reward of four experiments (three new ones plus the one used in Question 1) and the second shows the best mean reward of the same four experiments.

I have experimented with the exploration schedule scheme, in fact, with exactly a single parameter – the final exploration probability (*final_p*) in the linear interpolation scheme. Just to recall epsilon (*eps*), the probability

to select a random action a , is equal to $\text{fraction} * \text{final_p} + (1 - \text{fraction}) * \text{initilal_p}$ where $\text{fraction} = t / \text{schedule_timepoints}$. In the default setting $\text{schedule_timepoints} = 1,000,000$, $\text{final_p} = 0.1$ and $\text{initial_p} = 1.0$. This means that we start with $\text{eps} = 1.0$, namely, fully randomized exploration, and as we are making more and more steps we are decreasing gradually the epsilon such that at million steps it reaches 0.1, and from this point it stays fixed to 0.1. In terms of expolarion and exploitation – as we play more and more games, the time we spend on exploration decreases and the time we spend on exploitation increases. It is a good strategy since as the agent is experiencing with more and more games, as it is approaching to the goal, the estimated Q function is becoming more and more accurate, and therefore we can feel more and more confident on it. Now, if we go back to graph of Question 1, one can easily observe that at around 4 million steps the implementation reached its highest score (which is less than the maximum possible score) and from this point it stays fixed to that score. Since it is evident from the graph (of Question 1) that gradually decreasing epsilon improves the reward, one can hope that by decreasing the epsilon, even to a lower value than 0.1, the extra improvement to the reward would be gained, and, indeed, it this is what happened.

As the two graphs below show, results are improved when final_p is decreased. The graphs plot three additional experiments with different values of final_p : 0.05, 0.025 and 0. Highest score of 20.79 achieved when $\text{final_p} = 0$. To summarize, best setting is gained when we start with fully exploration, gradually changing relative proportion between exploration and exploitation until exploitation becomes *fully dominant* after million steps.



best mean reward using different final exploration probability (final_p)

