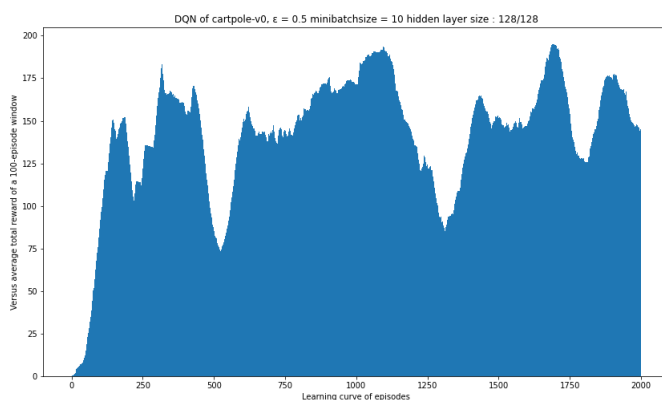


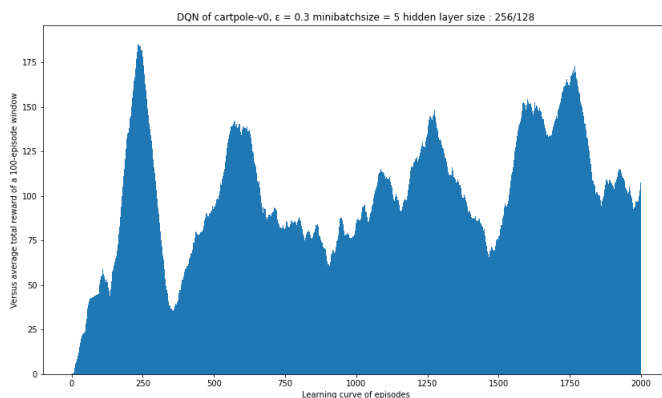
2 Deep RL

2.1 Plot of the DQN algorithm

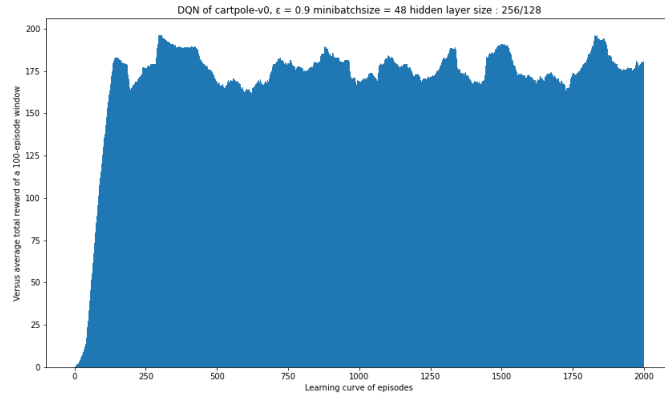
I did a lot of tests (more than 30) to look for the best set of hyper parameters. You can find the results in the boardin section 2.3. I decided to show only the more interesting plots here. I first started with the given set of hyper parameters :



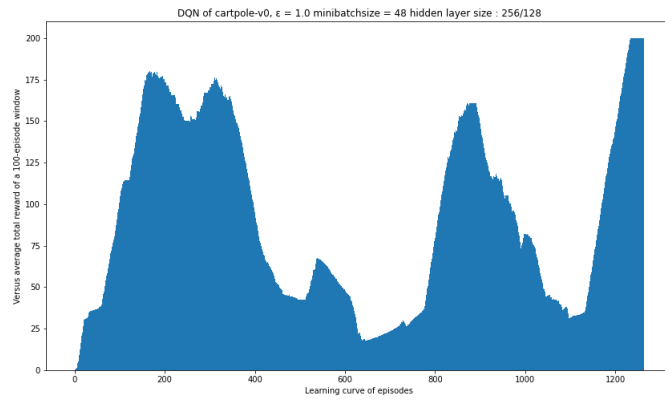
Then, I tried to change the epsilon/hidden layer/ minibatch size which helped me finding better parameters.

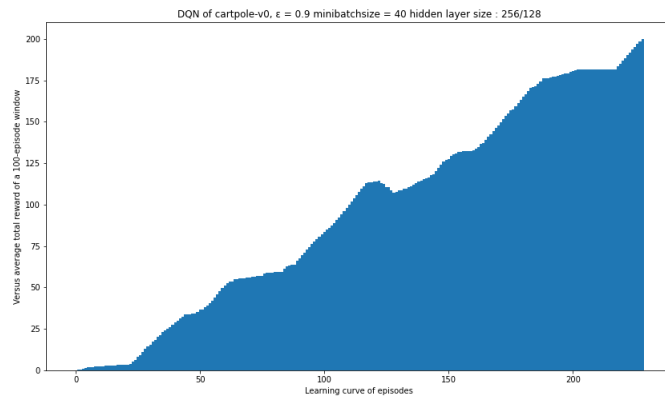
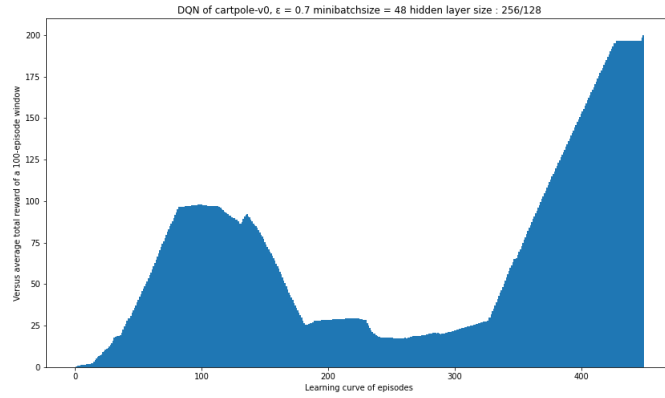


I finally found a good set of hyper parameters :



And then, I changed the discount factor from 0.9 to 1.0, the plot's form really changed. Before, it was more or less homogeneous but it never solved the problem and after, it was more up and downs and it solved the problem :





2.2 Best set of hyper parameters

The set of hyper parameters which is working the best for me is :

REPLAY MEMORY SIZE = 10000

EPSILON = 0.9

EPSILON DECAY = 0.99

EPISODES NUM = 2000

MINIBATCH SIZE = 40

DISCOUNT FACTOR = 1.0

TARGET UPDATE FREQ = 100

HIDDEN1 SIZE = 256

HIDDEN2 SIZE = 128

LEARNING RATE = 0.0001

indeed, with this set of hyper parameters, the problem is always solve and often it doesn't use more than 500 episode to solve it.

2.3 Observations on variation of certain hyper parameters

| γ | Hidden layer(s) | ϵ | Minibatch | Solved | Max consecutive eps | Last Episode nb | Final Mean |
|----------|-----------------|------------|-----------|--------|---------------------|-----------------|------------|
| 0.9 | 128 / 128 | 0.5 | 10 | False | 18 | 1100 | 200 |
| 0.9 | 128 / 128 | 0.8 | 10 | False | 12 | 1018 | 83 |
| 0.9 | 128 / 128 | 0.3 | 10 | False | 29 | 1517 | 200 |
| 0.9 | 256 / 256 | 0.3 | 10 | False | 26 | 1037 | 199.72 |
| 0.9 | 512 / 256 | 0.3 | 10 | False | 12 | 502 | 123.66 |
| 0.9 | 256 / 128 | 0.3 | 10 | False | 28 | 429 | 153 |
| 0.9 | 256 / 64 | 0.3 | 10 | False | 14 | 1770 | 125.16 |
| 0.9 | 256 / 128 | 0.3 | 5 | False | 44 | 206 | 199.86 |
| 0.9 | 256 / 128 | 0.3 | 48 | False | 63 | 933 | 200 |
| 0.9 | 256 / 128 | 0.3 | 64 | False | 60 | 81 | 177.4 |
| 0.9 | 256 / 128 | 0.3 | 128 | False | 58 | 773 | 200 |
| 0.9 | 256 / 128 | 0.1 | 48 | False | 45 | 1714 | 108.32 |
| 0.9 | 256 / 128 | 0.9 | 48 | False | 47 | 88 | 199.18 |
| 0.9 | 256 / 128 | 1.0 | 48 | False | 19 | 1529 | 142.72 |
| 1.0 | 256 / 128 | 1.0 | 48 | True | 100 | 1263 | 200 |
| 1.0 | 256 / 128 | 0.9 | 48 | True | 100 | 371 | 200 |
| 1.0 | 256 / 128 | 0.5 | 48 | True | 100 | 875 | 196.04 |
| 1.0 | 256 / 128 | 0.8 | 48 | True | 100 | 175 | 200 |
| 1.0 | 256 / 128 | 0.7 | 48 | True | 100 | 449 | 200 |
| 1.0 | 256 / 128 | 0.8 | 25 | True | 100 | 356 | 200 |
| 1.0 | 256 / 128 | 0.8 | 64 | True | 100 | 1838 | 200 |
| 1.0 | 256 / 128 | 0.8 | 55 | True | 100 | 1461 | 200 |
| 1.0 | 256 / 128 | 0.8 | 40 | True | 100 | 187 | 200 |
| 1.0 | 256 / 128 | 0.8 | 45 | True | 100 | 856 | 200 |
| 1.0 | 256 / 128 | 0.8 | 48 | True | 100 | 623 | 200 |
| 1.0 | 256 / 128 | 0.9 | 40 | True | 100 | 218 | 200 |
| 1.0 | 256 / 128 | 0.9 | 40 | True | 100 | 228 | 200 |

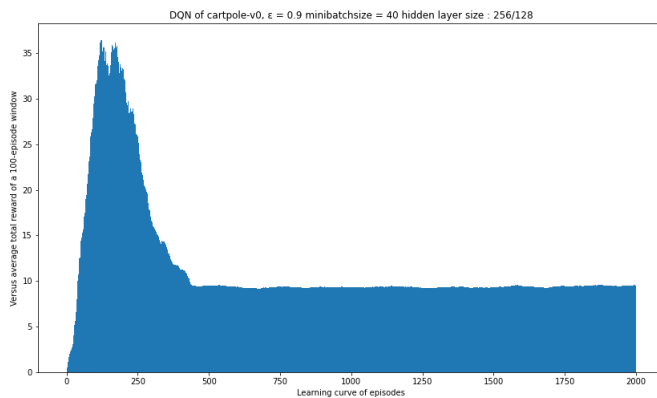
Actually, in this set of hyper parameters, one of the most important thing is the discount factor. Indeed, if it is less than 1 it is almost never going to solve the problem, if it is equal to 1, it will usually solve the problem even if the other parameters are a bit different from optimal ones. It will just be slower. From what I noticed, on the other hyper parameters, it seems like choosing small values won't give good results, but if values are too big then the results are not good either. It takes a really long time to find a value which is not too big but not too small. For example, for the minibatch size, I started with 10 and then I tried 5 which was really too small and then I also tried with 128 which was too big this time. I finally found that the best values were 48 and 64 and I chose

to keep going with 48 because it had less computing. At the end, when I found a set of hyper parameters which was solving the problem, I decide to readjust the minibatch size and I found that 40 was even working better. I did the same method to find all my parameters. That is how I finally starting from epsilon = 0.5, I chose epsilon = 0.3 for most of my test but readjusting at the end, it happens that the best value was 0.9. And for the hidden layers, I found that with both layers with big values, it was not working well and moreover that too big value (512) was not helping so I finally chose 256/128.

2.4 Removing experience replay and/or target network

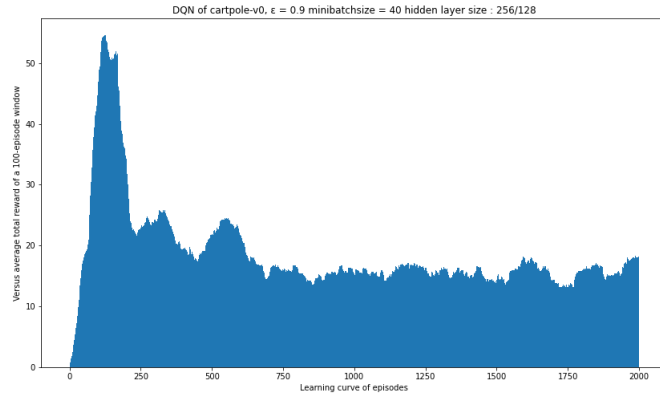
2.4.1 Removing experience replay

If I remove the experience replay and I make my target only on the state that I just compute, my model is going to learn a little bit a the beginning but then it will stop learning and even decreasing. It will never reach more than 100 steps and on the 100 consecutive episode more than 35.



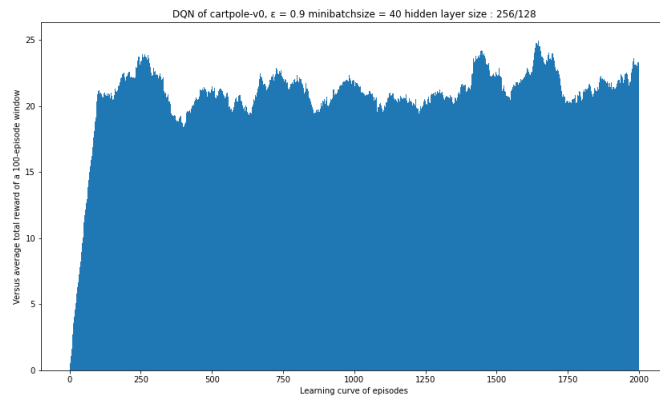
2.4.2 Removing target network

If I remove the target network and instead, my target is the reward if the episode is done and 0 otherwise. It will be able to go up to 200 steps in one episode but it will not learn on the long term, and of course it won't solve the problem.



2.4.3 Removing experience replay and target network

If we remove the experience replay and target network, our model is not learning anymore. It is computing the 2000 episodes very quickly but in the end, it never reaches 200 steps, his final mean is 9.22.



References

<https://pylessons.com/CartPole-reinforcement-learning/>