

Prior Experiments

Experiment 1

```
A. L. F. Arcade Learning Environment (version 0.11.2+ecce1138)
(Powered by Stella)
Loading trained model
100%|#####| 100/100 [01:16:00.00, 1.311k/s]
Run 100 episodes for 5 lives each
Mean: 239.06
rewards [310.0, 383.0, 268.0, 258.0, 315.0, 302.0, 196.0, 277.0, 225.0, 298.0, 225.0, 298.0, 234.0, 325.0, 351.0, 336.0, 363.0, 392.0, 225.0, 215.0, 303.0, 153.0, 335.0, 279.0, 251.0, 137.0, 225.0, 225.0, 294.0, 194.0, 246.0, 151.0, 225.0, 262.0, 296.0, 315.0, 225.0, 22
0.0, 165.0, 273.0, 211.0, 176.0, 301.0, 315.0, 251.0, 294.0, 363.0, 362.0, 379.0, 263.0, 224.0, 258.0, 261.0, 193.0, 189.0, 257.0, 256.0, 248.0, 268.0, 327.0, 225.0, 245.0, 258.0, 258.0, 347.0, 268.0, 317.0, 268.0, 294.0, 305.0, 225.0, 336.0, 307.0, 169.0, 287.0, 231.0,
303.0, 227.0, 296.0, 238.0, 253.0, 258.0, 223.0, 288.0, 286.0, 295.0, 183.0, 225.0, 258.0, 225.0, 225.0, 225.0, 294.0, 258.0, 225.0, 307.0, 276.0, 386.0, 262.0, 165.0]
running time 76.37722683864575
running time 77.19755285123451
```

Experiment Details

Parameters

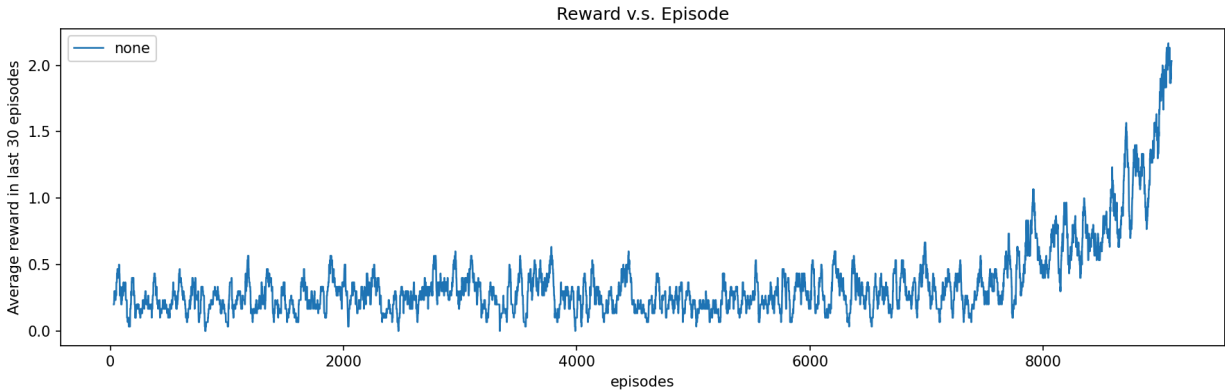
- Episodes = 100000
- epsilon
 - start = 0.01
 - end = 0.001
- decay episodes = 100000
- 200–500 million frames

These were the parameters I configured to tune my model. When first running the experiment I noticed that the step count was far too quick to decay the epsilon value. Therefore, I decided to make it dependent on the decay steps so that I could alter how long it takes for the epsilon decay. I did this because as we know, the higher the epsilon, the higher exploration is, so I wanted to make it that 20% in the beginning was the exploration phase and the rest was more of an exploitation approach. I set decay episodes to be 100000, and with that I also set `eps_decay_steps` to increase to at least 1M. I did this since we need more experience in order for the agent to learn how to score highly on breakout. So if we have more step for the decay when training, we are able to have more available steps to gain more experience. Also, for my double dqn approach I first decided to take a very conservative approach of making the epsilon be 0.01~0.001. This is because when I did more of an exploration approach, the learning progression was not as good as with a more conservative approach. The only drawback of this was that since the decay rate was extremely slow, it led to a very slow training process. Therefore I decided to take a faster learning approach leading to my final results for this project.

Experiment 2

mean: 73.94

```
Run 100 episodes for 5 Lives each
Mean: 73.94
rewards [100.0, 85.0, 63.0, 65.0, 67.0, 67.0, 70.0, 98.0, 181.0, 63.0, 69.0, 101.0, 53.0, 58.0, 62.0, 80.0, 46.0, 100.0, 70.0, 80.0, 55.0, 55.0, 64.0, 60.0, 91.0, 87.0, 71.0, 93.0, 81.0, 113.0, 74.0, 64.0, 61.0, 71.0, 65.0, 55.0, 76.0, 75.0, 65.0, 63.0, 77.0, 77.0, 71.0,
73.0, 96.0, 107.0, 47.0, 66.0, 67.0, 57.0, 73.0, 82.0, 62.0, 76.0, 113.0, 64.0, 111.0, 77.0, 99.0, 93.0, 73.0, 99.0, 101.0, 85.0, 53.0, 70.0, 66.0, 109.0, 39.0, 54.0, 68.0, 82.0, 82.0, 62.0, 50.0, 93.0, 49.0, 85.0, 86.0, 51.0, 86.0, 60.0, 74.0, 67.0, 71.0, 96.0, 61.0,
66.0, 72.0, 76.0, 86.0, 73.0, 69.0, 70.0, 64.0, 70.0, 70.0, 96.0, 34.0, 80.0]
running time 76.2333955378055
running time 77.65039142635277
```



Experiment Details

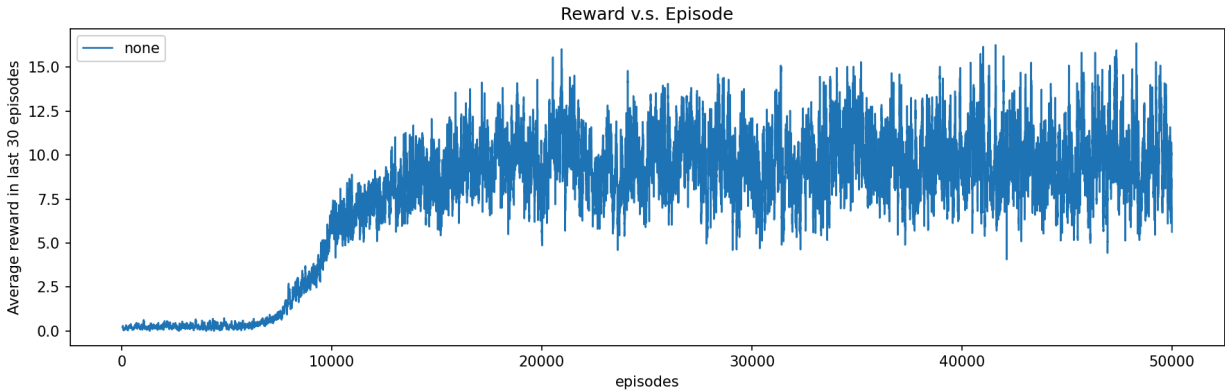
Parameters

- Episodes = 10000
- epsilon
 - start = 1.0
 - end = 0.01
- decay episodes = 10000
- 40 million frames in total

Unfortunately due to limited time, I could not train this model as long as I would've idealistically would like to. However, from the results, the average was above for most of them and less training time. Therefore I decided to further investigate this by training it for longer with these settings since from the first experiment having a longer training time leads to better results.

Final Results

```
(venv) hq@diehaiz:/mnt/storage/elyssa/Project $ !1977
python main.py --test_dqn
A.L.E. Arcade Learning Environment (version 0.11.2+eccl138)
(Powered by Stella)
Loading trained model
100%
Run 100 episodes for 5 Lives each
Mean: 355.35
rewards: 1599.0, 290.0, 344.0, 379.0, 398.0, 379.0, 389.0, 311.0, 428.0, 364.0, 360.0, 370.0, 399.0, 364.0, 377.0, 372.0, 368.0, 335.0, 328.0, 331.0, 293.0, 370.0, 336.0, 383.0, 369.0, 325.0, 346.0, 399.0, 360.0, 341.0, 385.0, 348.0, 346.0, 400.0, 391.0, 401.0, 412.0, 361.0, 360.0, 341.0, 367.0, 243.0, 356.0, 373.0, 399.0, 360.0, 292.0, 350.0, 352.0, 355.0, 307.0, 370.0, 366.0, 348.0, 301.0, 330.0, 377.0, 341.0, 337.0, 407.0, 279.0, 389.0, 323.0, 341.0, 362.0, 367.0, 302.0, 404.0, 343.0, 416.0, 334.0, 387.0, 366.0, 338.0, 390.0, 372.0, 384.0, 378.0, 360.0, 342.0, 313.0, 346.0, 360.0, 357.0, 338.0, 338.0, 345.0, 377.0, 348.0, 370.0, 404.0, 69.0, 336.0, 376.0, 341.0, 358.0, 383.0, 390.0, 338.0]
running time: 106.35912489891852
running time: 107.65813912558091
```



Experiment Details

Parameters

- Episodes = 50000
- epsilon
 - start = 1.0
 - end = 0.01
- decay episodes = 10000
- 200,000,000 frames
- Exploration for the 20% beginning of training
- Number of layers:
- Optimizer: RMSprop
- Loss function: Huber loss

From our first two experiments, for the final model I wanted to have a target of 20% exploration for the beginning and the rest to take an exploitation approach. Since I also wanted to train longer I changed the episodes to 50000 which took approximately 5 hours to train. I kept my epsilon to have a very explorative to a good exploitation rate. Like stated before, in my code I depended the amount of decay episodes to spend a certain amount of episodes with a certain decay rate. I set it to 10000 in order to have a slower decay rate of change. Therefore I was able to hit my target of 20% exploration and the rest have more of an exploitation approach. Overall, with these parameters, I was able to process around 200000000 frames. The optimizer used was RMSprop and the loss function was Huber loss. I chose the Huber loss function because for a double DQN there could be the occasional huge TD error. Even though it is rare this could make huge misestimates. Therefore using Huber loss counteracts this because the TD error is linear when it is large. Thus avoiding exploding gradients. I chose RMSprop as the optimizer since the gradients is non-stationary. RMSprop maintains the exponential moving average of squared gradients and rescales updates per parameter. Therefore being ideal in the Double DQN since it changes with the parameters.

As we can see from the results, the mean is 355.35. And all the average rewards are well above ≥ 40 . Something nice to also highlight is the graph results comparison of the prior experiment

and the final results. We can see in the prior experiment that it was trending up however could not have the trend we see with the final results. This is because we did not have as much training time as for the final model. Thus this model successfully hit the requirements. From the development of the 20% strategy, we were able to come up with a successful breakout model!

Sources

- <https://arxiv.org/pdf/1509.06461>
- <https://www.nature.com/articles/nature14236>
- <https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>
- <https://jonathan-hui.medium.com/rl-dqn-deep-q-network-e207751f7ae4>
- Gpt for debugging and graph setup.