

# Report – Project 2

Name: Sandesh Kumar Srivastava

UBIT: Sandeshk

## Part 1 - Implementing and applying DQN

In this assignment, we are implementing value function approximation algorithm, deep Q learning (DQN) following DeepMind's paper. Following are some of the salient features of DQN algorithm:

1. Experience replay: In order to efficiently select data for training of large neural networks, DeepMind team introduced the technique of experience replay. Experience replay allows us to store the experience gained  $(S_t, a_t, r_t, S_{t+1})$  in a memory segment for various episodes and then sample a subset of these experiences for training the neural network. Experience replay provides various benefits such as better efficiency as each experience is used in many iterations as training data; random sampling breaks the correlation among training data and it also avoids negative feedback loops that might shift the model to an unwanted local minima.
2. Target network: We use the concept of target network in order to solve the moving target problem by using a separate model for generating target. We periodically sync the target model with the actual model weights. Target network provides more stability to the algorithm by introducing a delay between an update to action value function  $Q$  and its effect on training data thus avoiding oscillations.
3. Approximate action-value function  $Q \leftarrow q(s, w)$ : Since determining the action-value function optimally using Bellman equation,  $Q_{i+1}(s, a) = E [r + \gamma \max_{a'} Q_i(s', a') | s, a]$ , is impractical for large systems, we use a function approximator to estimate the action-value functions. Furthermore, instead of linear function approximator, DeepMind team used a non-linear function approximator in neural network.

### Environment:

I have used Open AI Gym's **CartPole-v1** environment for checking DQN performance. Some of the details of this environment are as follows:

- a) Observation/State: It is of the form of Box (4) where the four values correspond to cart position, cart velocity, pole angle and pole angular velocity.
- b) Actions: It is of the form of Discrete (2) and consists of Action set  $(A) = \{0, 1\}$  where 0 means push cart to left and 1 means push cart to right
- c) Reward: It is 1 for each step till the game terminates. Reward set  $\mathbb{R} = \{1\}$
- d) Episode termination: Episode terminates when pole angle is more than 12 degrees or cart reaches the edge of display (cart position more than 2.4) or episode length is more than 500.
- e) Starting state: It is randomly selected from all possible states.

## DQN Algorithm:

### DQN Algorithm

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$

**For** episode = 1,  $M$  **do**

Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

With probability  $\epsilon$  select a random action  $a_t$

otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$

Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$

Every  $C$  steps reset  $\hat{Q} = Q$

**End For**

**End For**

Alina Vereshchaka (UB)

CSE4/510 Reinforcement Learning, Lecture 7.3

June 16, 2020

8 / 16

### Implementation details:

1. I have implemented two classes:

- **DQNAgent:** It is the agent class where all the main logic to implement DQN algorithm resides.
- **Environment:** It is the class which prepares the CartPole-v1 environment using gym's environment and uses DQNAgent as agent to solve this environment.

2. In DQNAgent class, I have defined following methods:

- **\_\_init\_\_:** Is used to initialize the environment variables and hyper parameters. It also creates variables for storing rewards metrics and also creates two neural network models viz. model and target.
- **build\_model:** Function to create neural network used by the algorithm.
- **step:** Implements epsilon-greedy approach of selecting a random action with probability epsilon and optimal action with probability  $1 - \epsilon$ .
- **sync\_target:** Is used to update the target model weights with model weights after every  $C$  steps.
- **replay:** Is used to sample minibatch of experiences from memory and use it for performing gradient descent step.
- **update\_epsilon():** It updates the epsilon value after every episode.

3. In Environment class, I have defined following methods:

- **\_\_init\_\_:** Is used to create the environment from Open AI gym class and create an agent instance of DQNAgent. It is also used to initialize some hyper parameters like number of episodes for training etc.
- **train():** Here we perform different steps for various episodes. It calls different agent methods for performing various tasks.

## Results:

We record the award collected by the agent for each episode and observe that initially the agent performs poorly collecting very few rewards for each episode. However, after sufficient training, the agent starts performing better and towards the end it is able to achieve the maximum award quite a number of times.

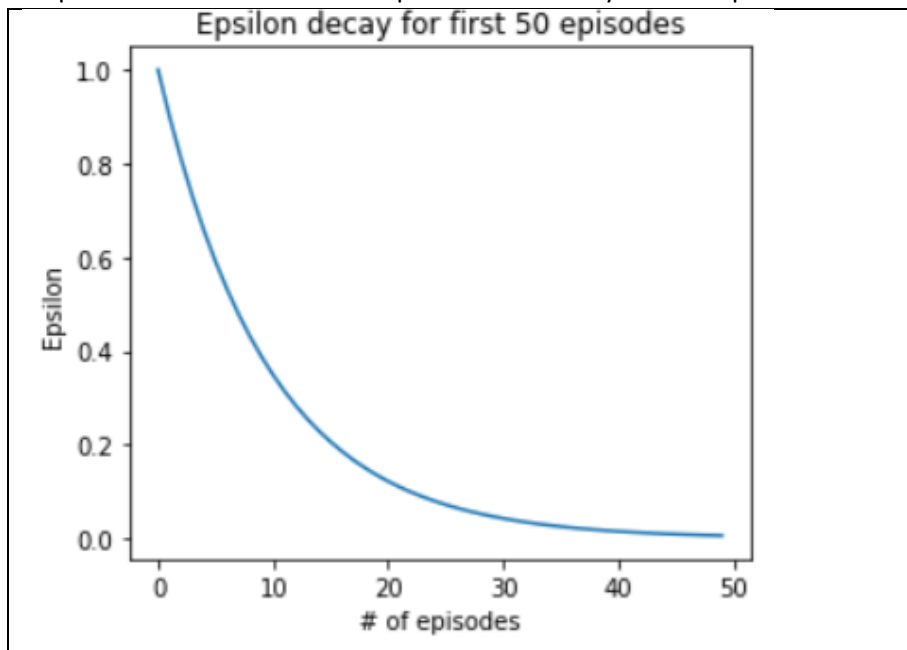
Hyper-parameters used for capturing below results are listed here:

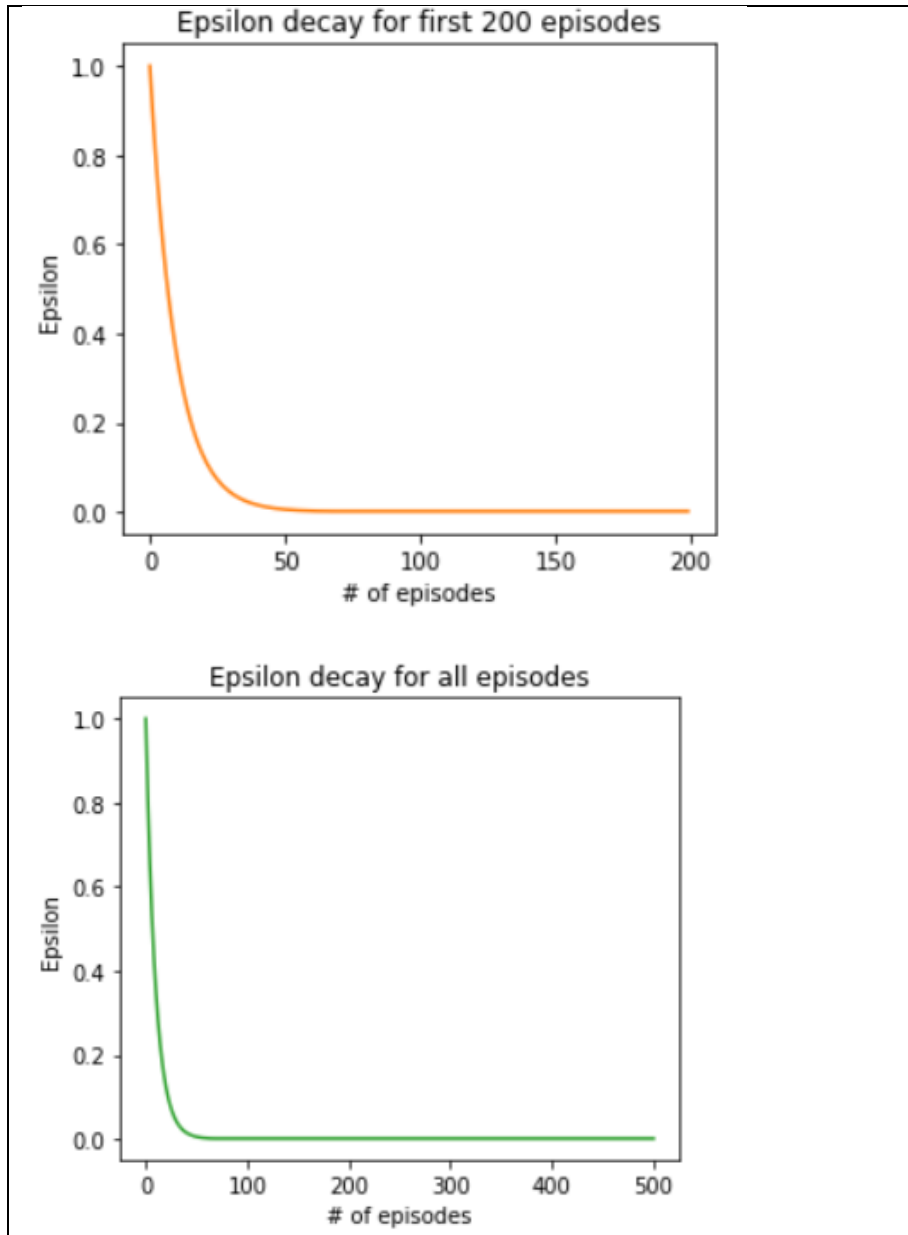
```
self.memory = deque(maxlen=50000)
self.gamma = 0.9
self.epsilon_max = 1.0
self.epsilon_min = 0.001
self.epsilon_decay = 0.9
self.batch_size = 32

self.number_of_episodes = 500
self.C = 1
self.reward_average_episodes = 10
```

The results are captured with following graphs:

1. **Epsilon decay graph:** Here I have shown the variation of epsilon value with the number of episodes encountered. To illustrate it better, I have captured 3 graphs of epsilon-decay with first 50 episodes then with first 200 episodes and finally with all episodes.





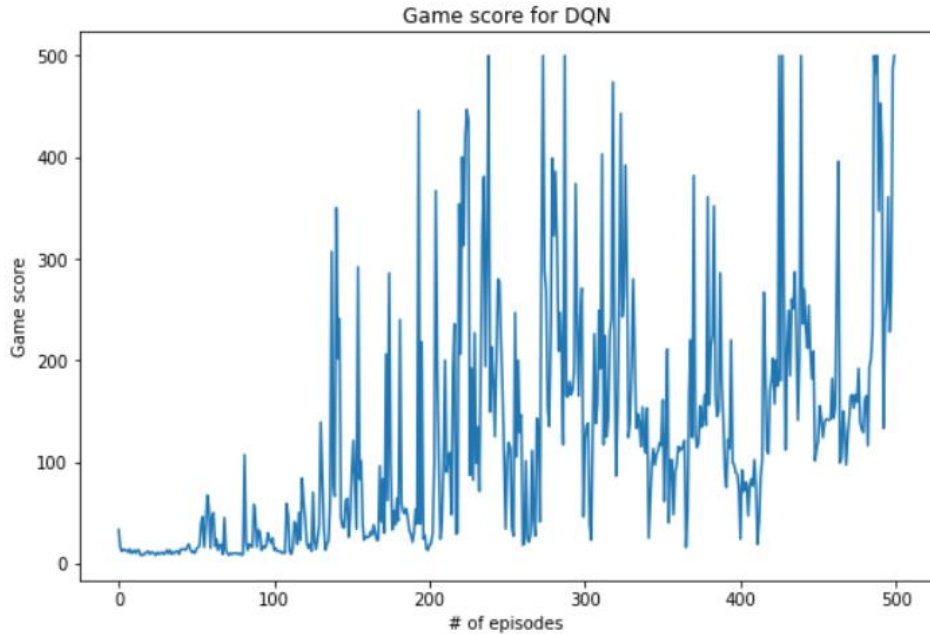
2. Game score graph: This graph captures the reward collected by agent for each episode.

```
Game score for DQN: [33.0, 16.0, 12.0, 14.0, 13.0, 13.0, 11.0, 14.0, 10.0, 11.0, 13.0, 10.0, 13.0, 13.0, 8.0, 9.0, 8.0, 10.0, 11.0, 12.0, 9.0, 11.0, 11.0, 9.0, 8.0, 11.0, 9.0, 10.0, 11.0, 9.0, 10.0, 13.0, 10.0, 13.0, 9.0, 11.0, 10.0, 12.0, 12.0, 9.0, 14.0, 14.0, 14.0, 13.0, 16.0, 19.0, 14.0, 11.0, 12.0, 10.0, 14.0, 16.0, 17.0, 39.0, 46.0, 16.0, 36.0, 67.0, 55.0, 15.0, 46.0, 50.0, 17.0, 24.0, 13.0, 18.0, 18.0, 9.0, 45.0, 13.0, 10.0, 8.0, 10.0, 9.0, 10.0, 10.0, 10.0, 10.0, 9.0, 10.0, 8.0, 9.0, 107.0, 21.0, 13.0, 19.0, 17.0, 16.0, 58.0, 43.0, 16.0, 33.0, 30.0, 13.0, 17.0, 15.0, 18.0, 30.0, 23.0, 20.0, 25.0, 13.0, 15.0, 12.0, 12.0, 12.0, 10.0, 11.0, 10.0, 59.0, 48.0, 12.0, 9.0, 15.0, 41.0, 36.0, 19.0, 50.0, 23.0, 84.0, 60.0, 44.0, 20.0, 15.0, 19.0, 12.0, 70.0, 21.0, 14.0, 25.0, 38.0, 139.0, 85.0, 36.0, 13.0, 19.0, 22.0, 60.0, 307.0, 79.0, 66.0, 350.0, 202.0, 241.0, 45.0, 36.0, 35.0, 62.0, 64.0, 26.0, 48.0, 91.0, 121.0, 77.0, 34.0, 292.0, 83.0, 101.0, 38.0, 23.0, 26.0, 26.0, 26.0, 31.0, 28.0, 38.0, 31.0, 23.0, 23.0, 96.0, 41.0, 83.0, 30.0, 206.0, 62.0, 286.0, 99.0, 33.0, 52.0, 38.0, 64.0, 42.0, 240.0, 59.0, 51.0, 49.0, 54.0, 45.0, 32.0, 29.0, 21.0, 30.0, 53.0, 39.0, 446.0, 39.0, 218.0, 24.0, 28.0, 14.0, 13.0, 18.0, 20.0, 30.0, 97.0, 367.0, 189.0, 100.0, 24.0, 41.0, 102.0, 200.0, 90.0, 105.0, 109.0, 48.0, 189.0, 236.0, 29.0, 30.0, 354.0, 206.0, 400.0, 313.0, 421.0, 447.0, 437.0, 87.0, 192.0, 82.0, 227.0, 99.0, 134.0, 71.0, 200.0, 331.0, 381.0, 194.0, 314.0, 500.0, 149.0, 213.0, 182.0, 125.0, 197.0, 280.0, 277.0, 217.0, 168.0, 98.0, 34.0, 103.0, 119.0, 113.0, 35.0, 27.0, 247.0, 105.0, 200.0, 129.0, 146.0, 18.0, 20.0, 101.0, 25.0, 21.0, 29.0, 111.0, 42.0, 27.0, 143.0, 132.0, 41.0, 308.0, 500.0, 288.0, 269.0, 153.0, 135.0, 230.0, 399.0, 323.0, 386.0, 318.0, 209.0,
```

```

247.0, 206.0, 117.0, 500.0, 165.0, 164.0, 179.0, 166.0, 170.0, 191.0, 374.0, 238.0, 165.0, 243.0,
271.0, 46.0, 122.0, 136.0, 138.0, 39.0, 23.0, 132.0, 226.0, 138.0, 156.0, 249.0, 192.0, 403.0,
117.0, 224.0, 125.0, 143.0, 227.0, 241.0, 474.0, 211.0, 86.0, 152.0, 268.0, 443.0, 243.0, 264.0,
392.0, 292.0, 124.0, 140.0, 195.0, 280.0, 202.0, 133.0, 147.0, 139.0, 115.0, 154.0, 124.0, 109.0,
153.0, 25.0, 60.0, 97.0, 113.0, 97.0, 107.0, 111.0, 119.0, 116.0, 161.0, 61.0, 121.0, 211.0, 40.0,
102.0, 102.0, 48.0, 90.0, 97.0, 115.0, 111.0, 114.0, 113.0, 121.0, 16.0, 50.0, 168.0, 220.0, 124.0,
382.0, 130.0, 114.0, 118.0, 155.0, 134.0, 143.0, 166.0, 136.0, 361.0, 156.0, 213.0, 225.0, 352.0,
166.0, 145.0, 150.0, 286.0, 178.0, 141.0, 91.0, 75.0, 122.0, 113.0, 220.0, 99.0, 97.0, 89.0, 88.0,
73.0, 24.0, 92.0, 71.0, 80.0, 79.0, 47.0, 76.0, 83.0, 76.0, 102.0, 76.0, 19.0, 41.0, 85.0, 102.0,
267.0, 194.0, 112.0, 108.0, 170.0, 181.0, 202.0, 157.0, 200.0, 175.0, 500.0, 180.0, 500.0, 268.0,
112.0, 208.0, 249.0, 185.0, 260.0, 251.0, 287.0, 223.0, 141.0, 206.0, 500.0, 236.0, 270.0, 233.0,
212.0, 254.0, 209.0, 182.0, 209.0, 101.0, 113.0, 120.0, 155.0, 144.0, 124.0, 137.0, 142.0, 142.0,
141.0, 142.0, 182.0, 143.0, 154.0, 283.0, 396.0, 99.0, 104.0, 150.0, 132.0, 97.0, 125.0, 143.0,
165.0, 166.0, 153.0, 166.0, 159.0, 192.0, 140.0, 132.0, 129.0, 161.0, 165.0, 116.0, 194.0, 201.0,
230.0, 500.0, 483.0, 500.0, 347.0, 453.0, 414.0, 133.0, 236.0, 258.0, 361.0, 228.0, 270.0, 487.0,
500.0]

```

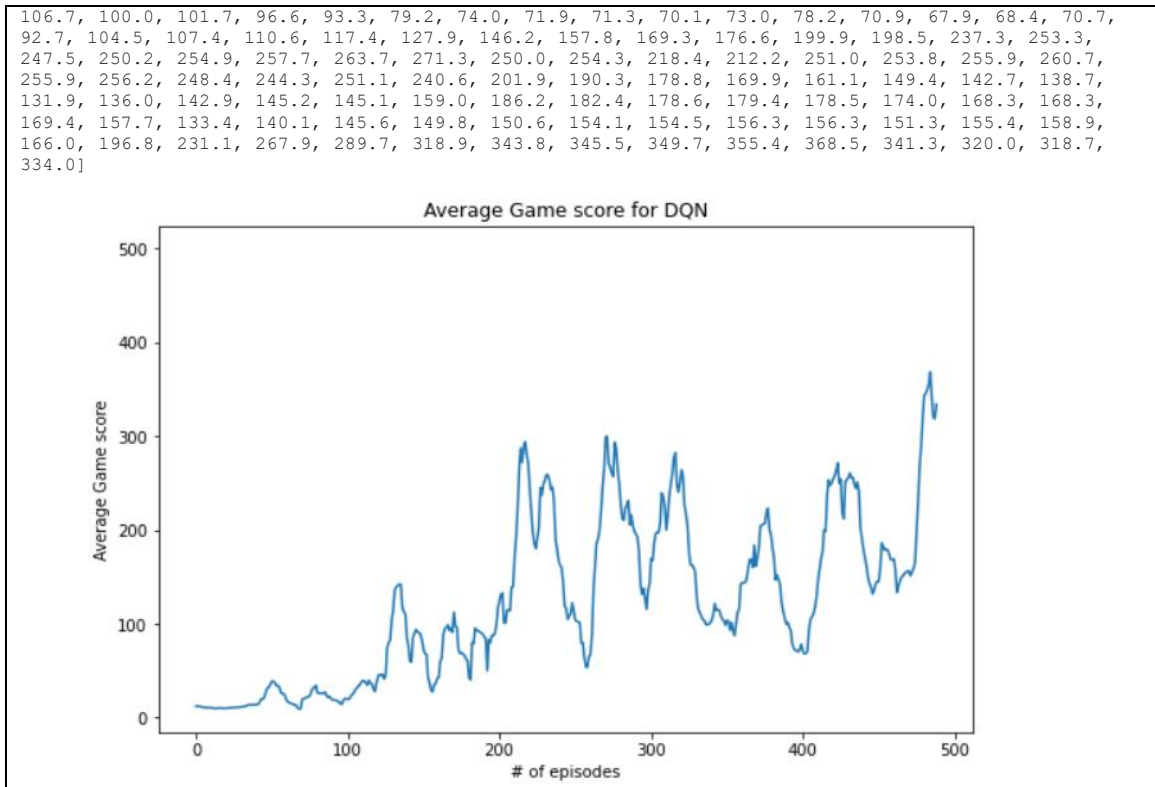


3. Average score graph: This graph captures the average reward collected by the agent for every `self.reward_average_episodes`.

```

Average Game score for DQN: [12.1, 12.2, 12.1, 11.6, 11.2, 10.9, 10.5, 10.6, 10.7, 10.3, 10.4,
10.2, 9.8, 9.8, 10.0, 10.1, 10.1, 10.1, 9.8, 9.9, 10.1, 10.0, 10.4, 10.5, 10.5, 10.6, 10.8, 10.9,
10.9, 11.3, 11.4, 11.8, 11.8, 12.5, 13.3, 13.7, 13.6, 13.6, 13.7, 13.7, 13.9, 14.2, 16.8, 19.8,
19.5, 21.7, 27.3, 31.6, 32.1, 35.3, 38.7, 38.7, 37.2, 33.9, 34.1, 32.3, 26.5, 25.5, 25.3, 21.7,
17.5, 16.8, 15.3, 15.0, 14.2, 13.4, 13.4, 9.9, 9.4, 9.3, 19.2, 20.3, 20.7, 21.6, 22.3, 22.9, 27.8,
31.1, 31.9, 34.3, 26.6, 25.8, 26.2, 25.8, 25.9, 27.3, 23.8, 21.5, 22.4, 20.4, 18.9, 18.8, 18.3,
18.0, 17.2, 15.3, 14.0, 17.9, 20.2, 20.1, 19.5, 19.8, 22.7, 25.1, 26.0, 29.9, 31.2, 33.7, 34.9,
38.1, 39.2, 39.2, 37.0, 34.6, 39.7, 36.8, 35.9, 30.0, 27.8, 37.3, 43.8, 45.9, 45.3, 46.0, 41.2,
45.1, 74.4, 79.8, 82.6, 103.7, 115.4, 135.9, 139.1, 140.8, 142.1, 142.3, 118.0, 112.7, 110.9, 85.0,
76.9, 60.5, 59.4, 85.0, 89.8, 93.7, 91.1, 90.8, 88.6, 82.1, 72.6, 68.0, 67.4, 42.0, 36.8, 29.0,
27.5, 34.8, 36.3, 42.0, 42.4, 59.9, 63.3, 88.1, 94.9, 95.9, 98.8, 93.0, 95.3, 91.2, 112.2, 97.5,
96.4, 72.7, 68.2, 69.4, 67.4, 66.5, 62.2, 61.0, 42.3, 40.3, 79.8, 78.8, 95.2, 93.1, 92.7, 91.2,
90.4, 89.2, 85.9, 85.0, 50.1, 82.9, 80.0, 87.6, 87.2, 89.9, 98.8, 117.0, 124.0, 131.5, 132.7, 100.8,
100.8, 114.4, 114.9, 113.8, 139.0, 139.6, 170.6, 191.4, 222.6, 262.5, 287.3, 272.4, 288.7, 293.9,
281.2, 270.5, 243.9, 219.7, 197.6, 186.0, 180.4, 191.1, 203.3, 245.1, 237.3, 248.7, 253.5, 258.9,
258.6, 253.5, 243.1, 245.4, 230.8, 190.6, 179.1, 168.1, 161.8, 160.6, 144.4, 119.1, 116.1, 104.9,
108.1, 111.2, 122.4, 113.9, 104.0, 102.8, 101.8, 101.2, 79.4, 80.0, 64.2, 54.0, 53.7, 65.1, 67.2,
87.9, 135.4, 162.1, 186.1, 190.3, 199.6, 219.9, 245.5, 264.6, 299.1, 300.1, 271.0, 266.9, 260.6,
257.0, 293.5, 287.0, 263.5, 249.1, 227.1, 212.3, 210.5, 223.2, 226.4, 231.2, 205.5, 216.1, 204.3,
198.6, 195.6, 192.4, 177.2, 142.1, 131.5, 137.6, 127.1, 115.6, 135.9, 142.9, 169.6, 167.5, 186.0,
196.2, 197.3, 197.4, 207.7, 239.5, 235.7, 225.1, 200.0, 215.1, 237.0, 248.8, 260.9, 277.4, 282.5,
247.5, 240.4, 251.3, 264.1, 257.5, 226.5, 216.9, 204.4, 176.7, 162.9, 162.9, 159.8, 155.6, 130.1,
115.9, 112.3, 108.9, 104.7, 103.9, 99.6, 99.1, 99.8, 100.6, 104.2, 110.3, 121.7, 114.4, 114.9,
114.4, 108.1, 105.2, 103.3, 98.7, 103.7, 103.0, 93.2, 101.3, 92.7, 87.5, 99.5, 112.5, 115.2, 141.9,
143.8, 143.8, 144.3, 147.7, 159.5, 168.8, 168.6, 160.2, 183.9, 161.3, 169.6, 180.7, 204.1, 205.2,
206.3, 207.0, 219.0, 223.2, 201.2, 194.7, 180.9, 170.6, 146.7, 152.1, 147.5, 142.2, 122.5, 113.5,

```



## Part 2 - Improving DQN:

DeepMind also observed that DQN suffers from overestimations for some games in the Atari 2600 domain so they suggested improvements in the form of Double DQN (DDQN) algorithm. For this section, I have implemented DDQN algorithm which is based on using two estimators – one for obtaining the best action and the other for evaluating state-action value  $Q$  for evaluating the selected action. The basic idea for DDQN is that it is less likely for both estimators to overestimate at the same time.

### DDQN Algorithm:

#### Implementation details:

I have implemented two classes:

- **DDQNAgent**: It is the agent class where all the main logic to implement DQN algorithm resides. and is similar to `DQNAgent` class.
- **dEnvironment**: It is the class which prepares the CartPole-v1 environment using gym's environment and uses `DDQNAgent` as agent to solve this environment. It is analogous to `Environment` class.

## Double Deep Q Network

### Algorithm 1 : Double Q-learning (Hasselt et al., 2015)

---

Initialize primary network  $Q_\theta$ , target network  $Q_{\theta'}$ , replay buffer  $\mathcal{D}$ ,  $\tau < 1$

**for** each iteration **do**

**for** each environment step **do**

        Observe state  $s_t$  and select  $a_t \sim \pi(a_t, s_t)$

        Execute  $a_t$  and observe next state  $s_{t+1}$  and reward  $r_t = R(s_t, a_t)$

        Store  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{D}$

**for** each update step **do**

        sample  $e_t = (s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}$

        Compute target Q value:

$Q^*(s_t, a_t) \approx r_t + \gamma Q_{\theta'}(s_{t+1}, \text{argmax}_{a'} Q_{\theta'}(s_{t+1}, a'))$

        Perform gradient descent step on  $(Q^*(s_t, a_t) - Q_\theta(s_t, a_t))^2$

        Update target network parameters:

$\theta' \leftarrow \tau * \theta + (1 - \tau) * \theta'$

---

The algorithm for DDQN is very similar to DQN and the only change is where we select the best action and evaluate the state-action value Q:

DQN	DDQN
<code>Q_prime = max(self.target_model.predict(state_prime)[0])</code>	<code>best_action = np.argmax(self.model.predict(state_prime)[0])</code>
<code>state_target[0][action] = reward + self.gamma*Q_prime</code>	<code>Q_prime = self.target_model.predict(state_prime)[0][best_action]</code>
	<code>state_target[0][action] = reward + self.gamma*Q_prime</code>
<code>self.target_model.set_weights(self.model.get_weights())</code>	<code>target_weights[i] = model_weight*self.tau + target_weight * (1-self.tau)</code>
	<code>self.target_model.set_weights(target_weights)</code>

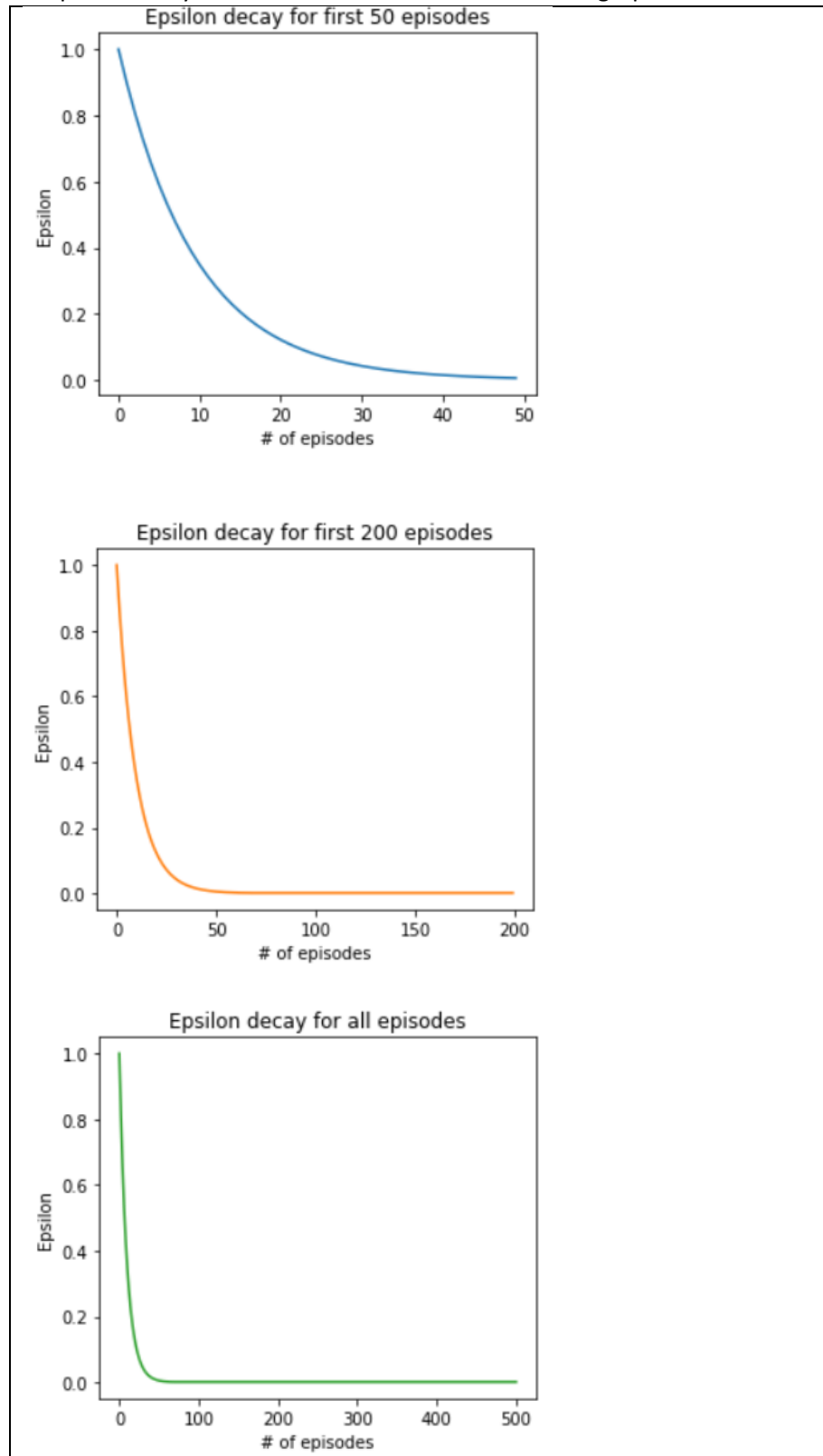
### Results:

The results are captured with following graphs by keeping all hyper parameters same as that for DQN except for tau:

```
self.tau = 0.99
```

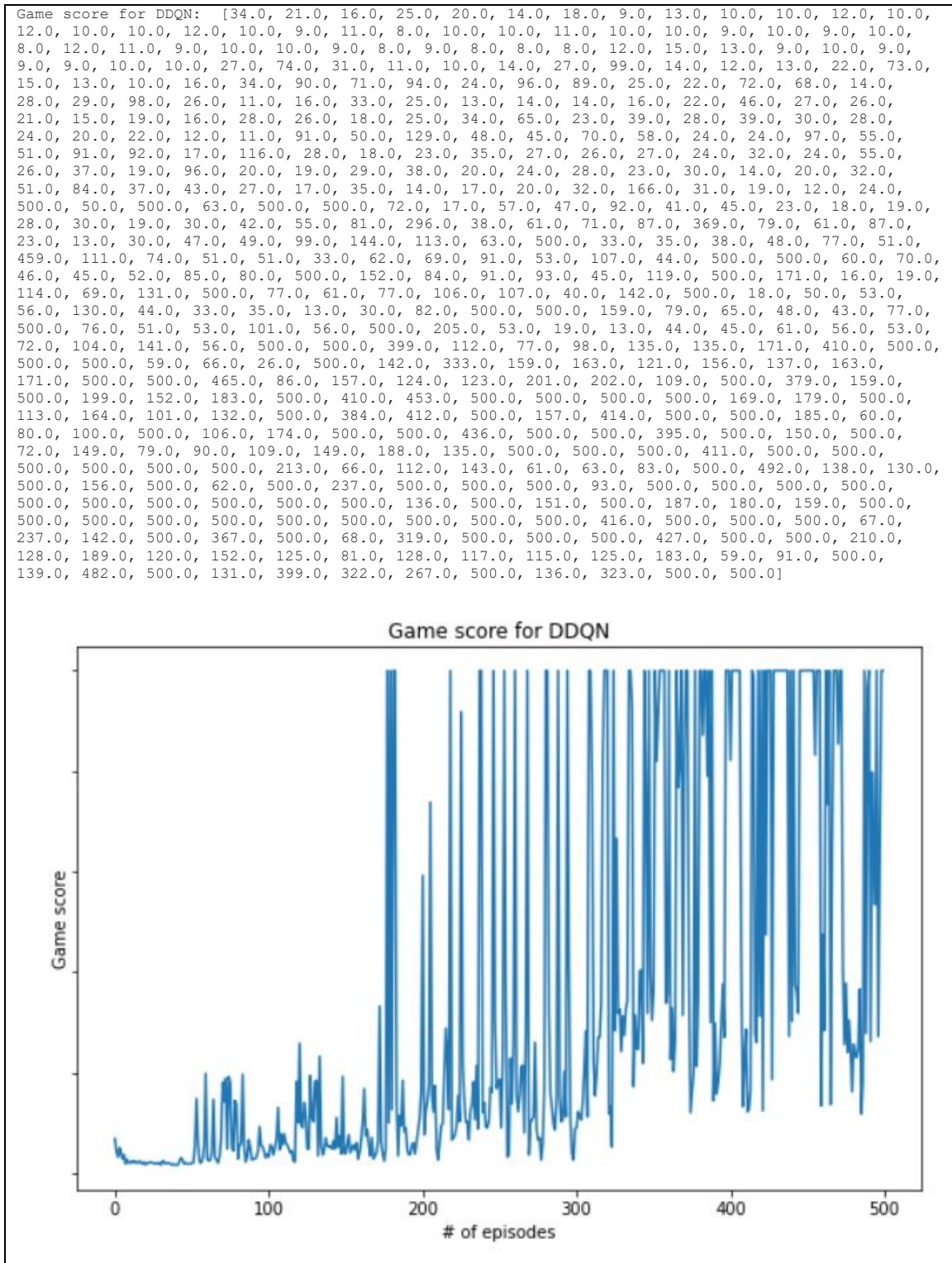
1. **Epsilon decay graph:** Here I have shown the variation of epsilon value with the number of episodes encountered. To illustrate it better, I have captured 3 graphs of epsilon-decay with first

50 episodes then with first 200 episodes and finally with all episodes. Since the parameters used for epsilon decay are same as that for DQN so the two graphs are identical.

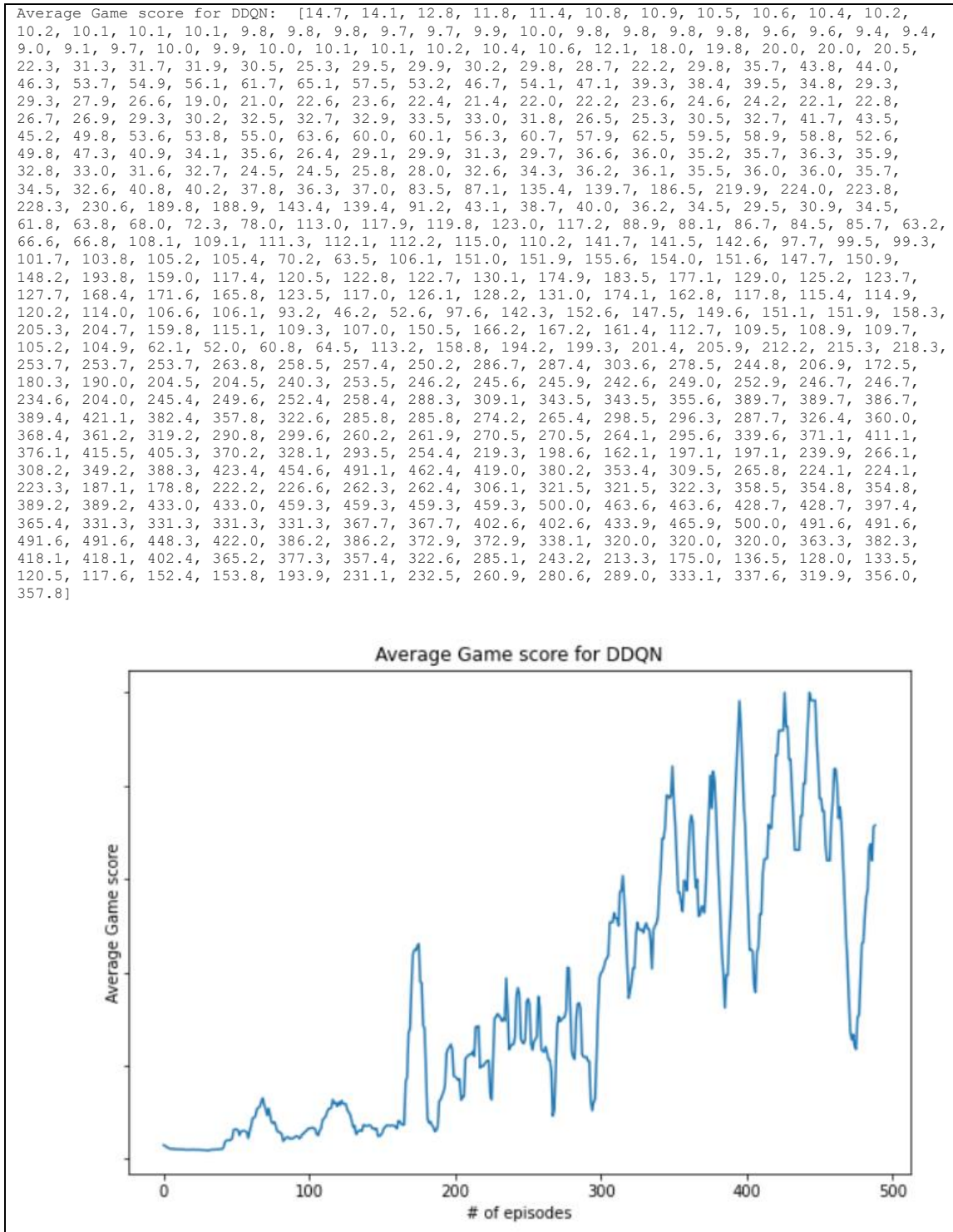




1. **Game score graph:** This graph captures the reward collected by agent for each episode.

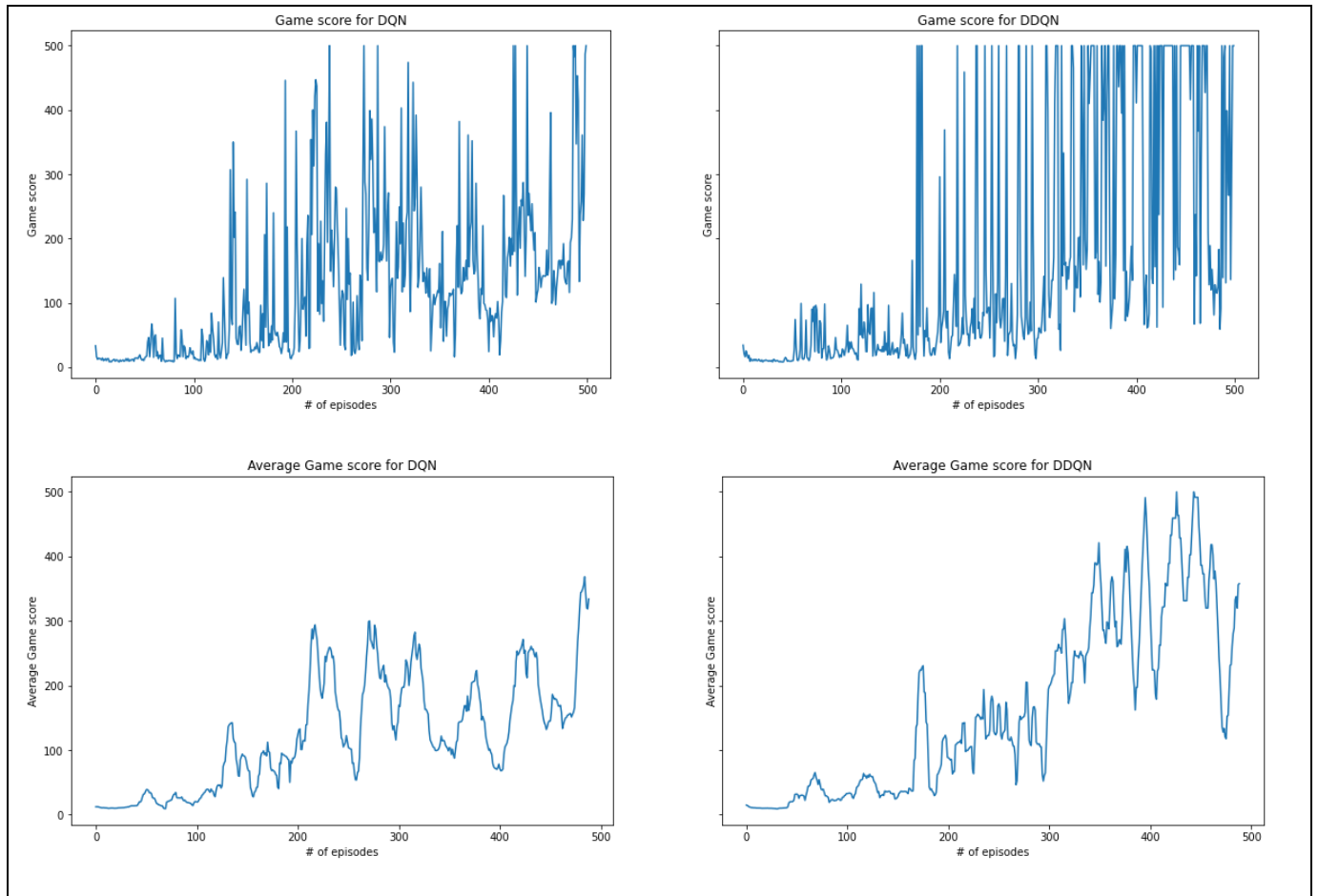


2. **Average score graph:** This graph captures the average reward collected by the agent for every `self.reward_average_episodes`.



### Comparison of DQN and DDQN performance:

Next, we study the comparative performance of DQN and DDQN for the same initial set of hyperparameters:



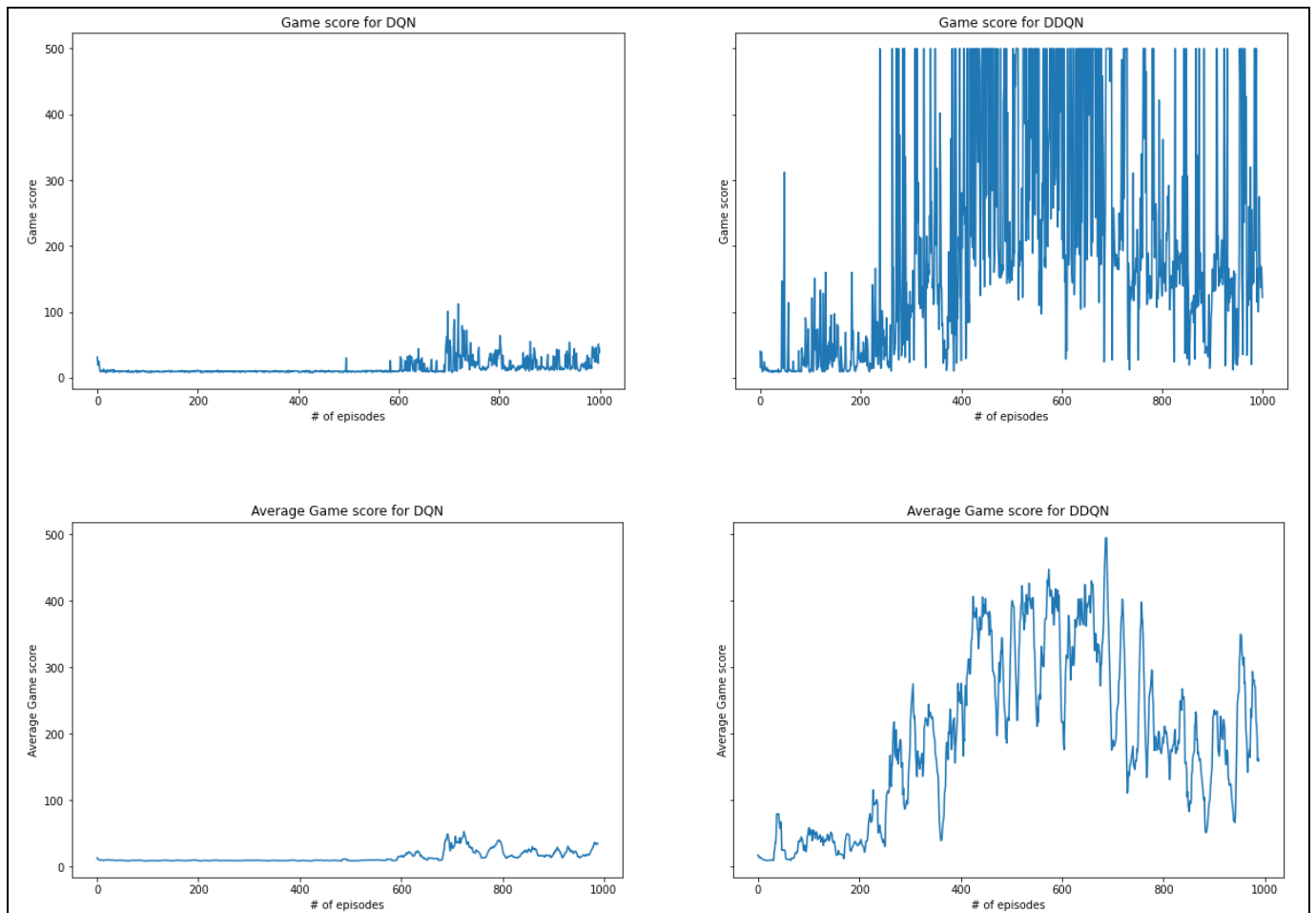
As we can see that keeping other hyper-parameters constant, normally DDQN performs better than DQN for the total and average scores collected. It can also be seen from the graphs that DDQN starts performing better earlier (fast learner) than DQN (slow learner) over the same set of episodes. So, we can say that DDQN learns better policy quickly as compared to DQN.

## Effect of other parameters on performance of DQN and DDQN:

In this section, I have tried to show the effect of modifying some of the hyperparameters on the game score and average score calculated for DQN and DDQN.

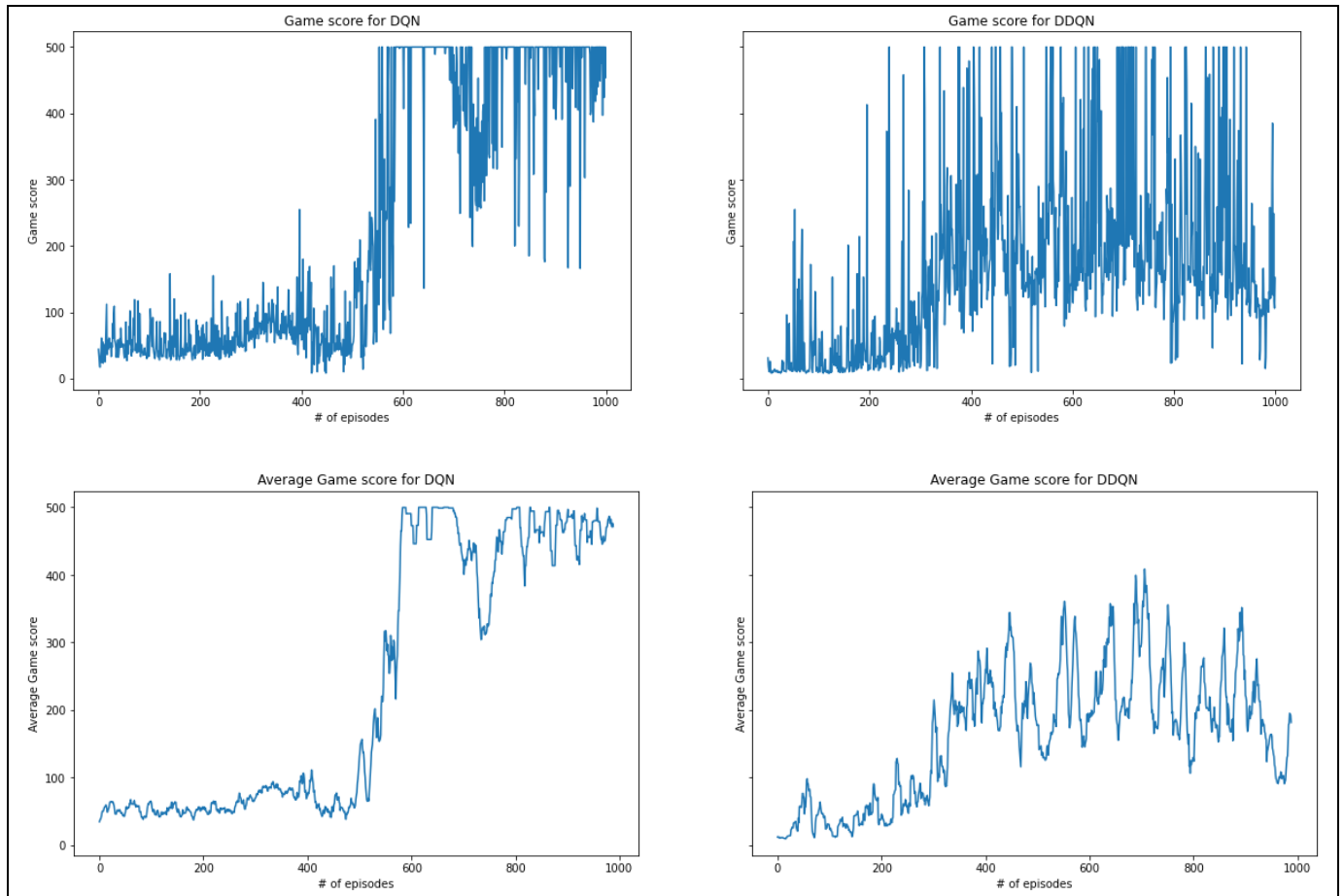
```
1. self.number_of_episodes = 1000  
   self.C = 10
```

Here I have kept the sync parameter C to 10 and observed the behavior for 1000 episodes. As we can see in this case, DQN does not perform very well but DDQN performs very well. It can be due to the slow learning of DQN as compared to DDQN as the number of weight sync is reduced.



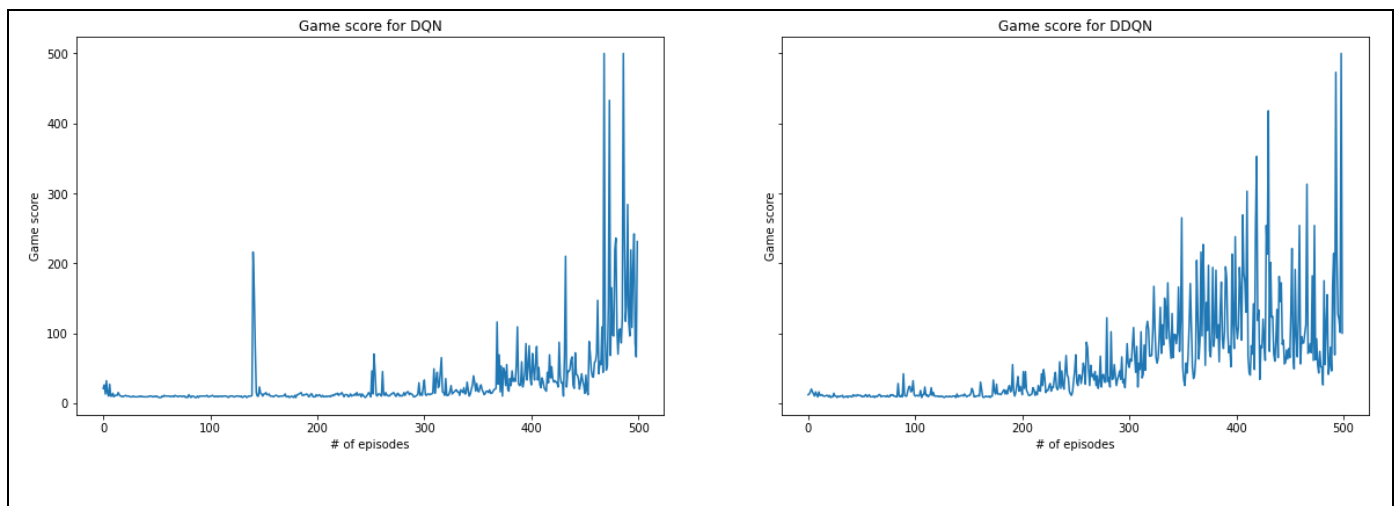
```
2. self.number_of_episodes = 1000  
   self.C = 5
```

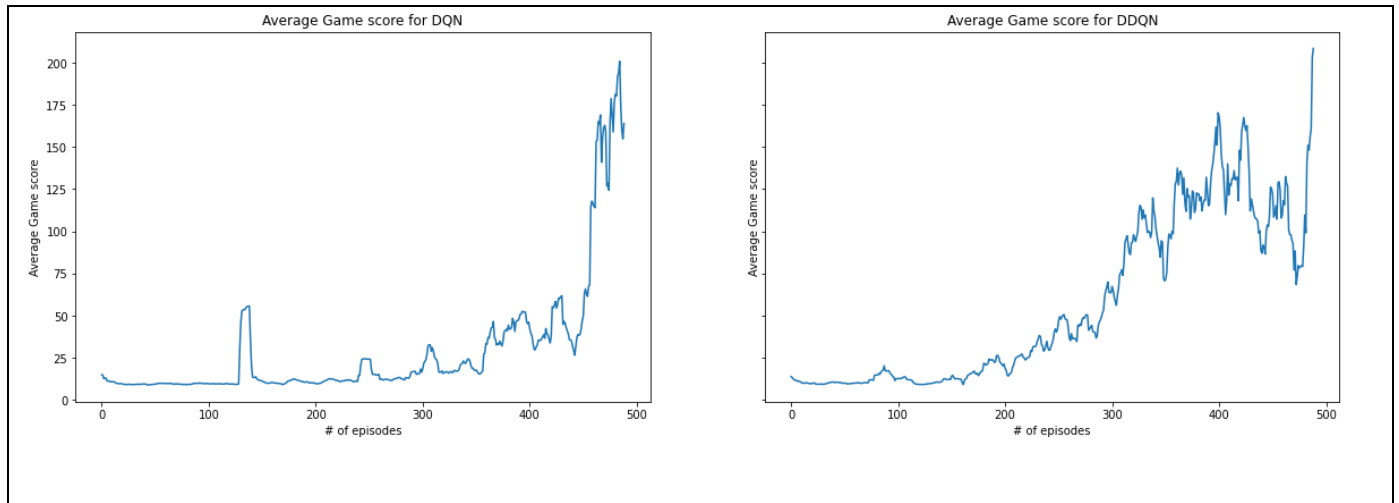
Here I have kept the sync parameter C to 5 and observed the behavior for 1000 episodes. As we can see in this case, DQN starts performing well later than DDQN but it picks up towards the later episodes.



```
3. self.number_of_episodes = 500
   self.C = 5
```

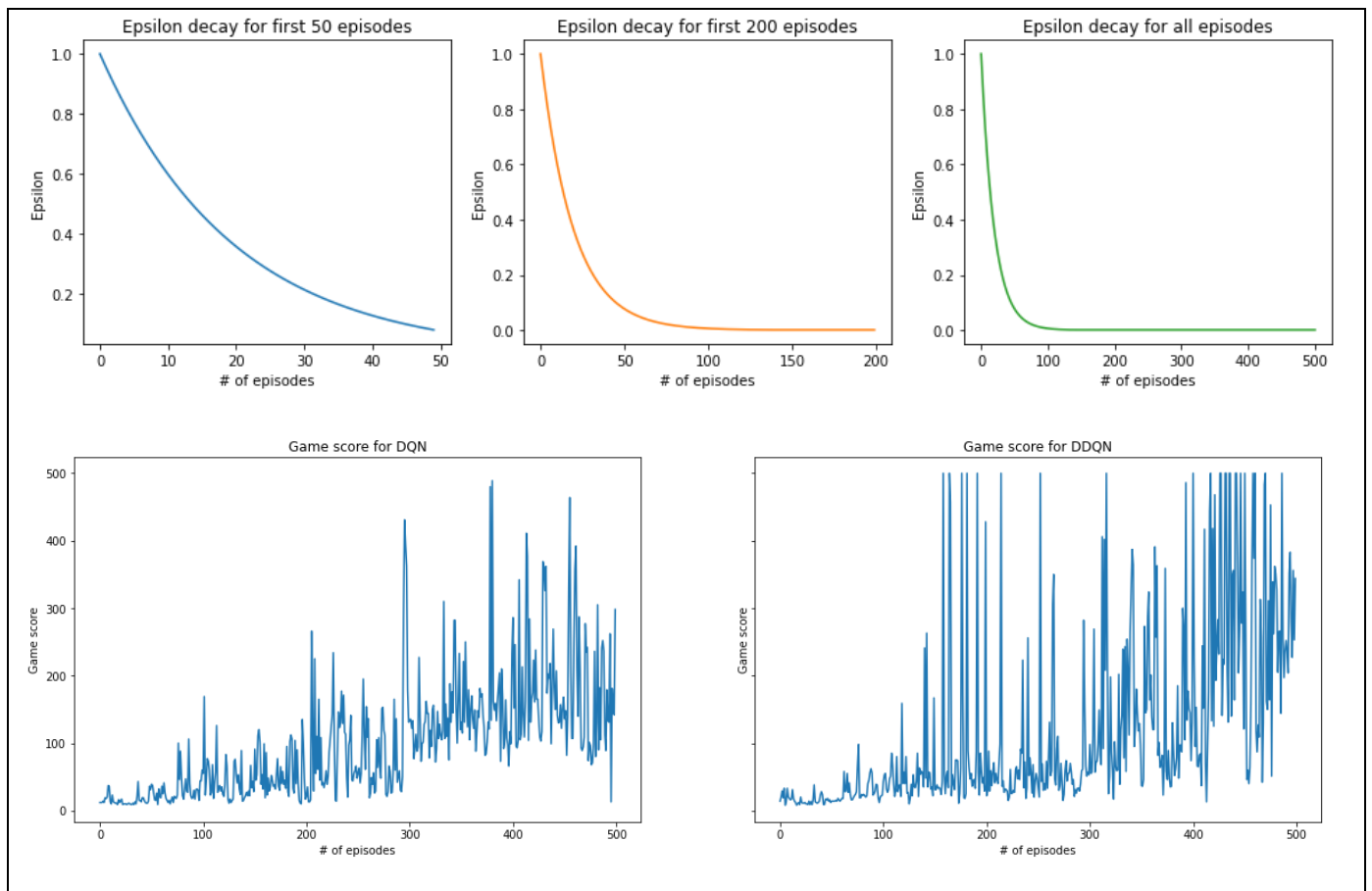
Here I have kept the sync parameter C to 5 and observed the behavior for 500 episodes. As we can see in this case, DQN does not perform very well but DDQN performs relatively well. We can again see the slow learning of DQN as compared to DDQN from these results.

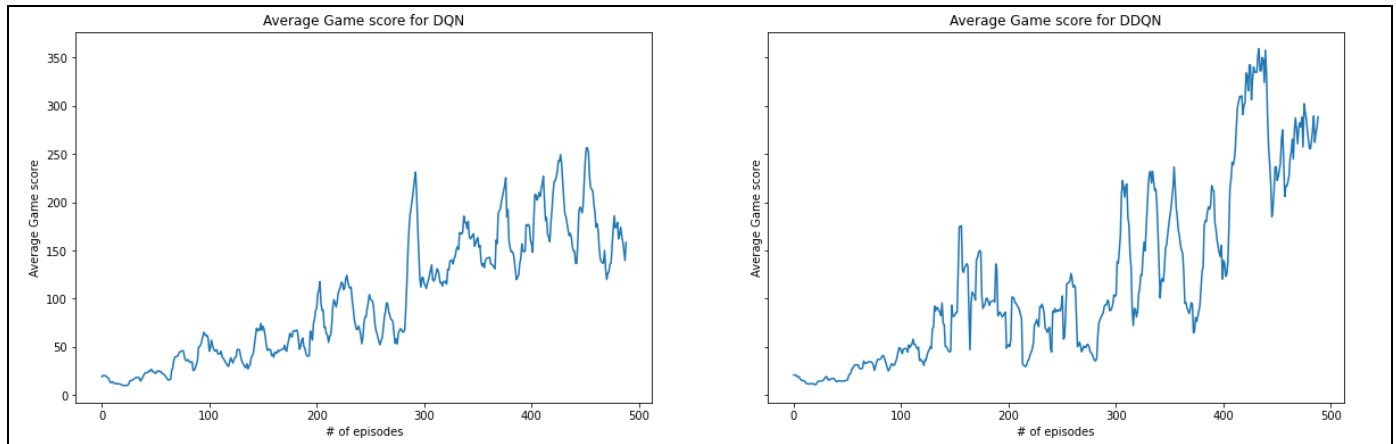




```
4. self.epsilon_decay = 0.95
   self.C = 1
```

Here I have kept the `epsilon_decay` parameter to 0.95 and observed the behavior for 500 episodes. As we can see in this case, epsilon decay happens relatively slower and DQN and DDQN both perform relatively well for later episodes. We can again see the slow learning of DQN as compared to DDQN from these results.





During different trials I have also observed that the results vary even with same set of parameters. It can be due to the neural network training, random minibatch selection and also different server configurations at runtime.

### **Conclusion:**

From the results obtained, we can see that both DQN and DDQN perform poor initially but as the number of episodes increases they learn from the experiences and start performing better as the number of episodes increases. In comparison, DDQN learns faster and starts performing better soon as compared to DQN. Towards the end, we also saw the effects of different hyper parameter on the performance of both DQN and DDQN.