
Ablation Studies on Deep Q Learning

Vraj Patel

261022581

vraj.patel@mail.mcgill.ca

Abstract

1 This project aims to conduct ablation studies for the DQN algorithm in two distinct
2 environments: Cartpole from Classic Control and Lunar Lander from Box2D. By
3 selectively removing key components of the DQN algorithm such as the Target
4 network, memory, and exploration rate. The aim is to analyze their individual
5 contributions to the algorithm's overall performance. Experiments employ the
6 same neural network architecture and hyperparameters across both environments,
7 allowed to compare and contrast the training convergence for each scenario. The
8 results of our study could provide valuable insights into the inner workings of
9 the DQN algorithm and inform the design of more efficient RL algorithms in the
10 future.

11 **1 Introduction**

12 Understanding the contribution of individual components or modules of an RL algorithm to its
13 overall performance can be challenging. Ablation studies, a type of experimental analysis used in RL
14 research, help to address this issue by selectively removing particular components of the algorithm
15 and analyzing the resulting impact on the agent's performance. This project focuses on conducting
16 ablation studies for DQN algorithm in two different environments, Cartpole and Lunar Lander. By
17 removing specific components such as the target network, buffer memory, and exploration policy, to
18 investigate their individual contributions to the algorithm's overall performance. To ensure our results
19 are generalizable and applicable across various environments, we trained the agent with the same
20 architecture and hyper-parameters. We also ran independent training sessions with different random
21 seeds to analyze the generalization issues in deep RL. By analyzing the results of our experiments,
22 we hope to gain insights into the strengths and weaknesses of the DQN algorithm and inform the
23 development of more efficient RL algorithms.

24 **2 Background**

25 **2.1 Cart Pole Environment**

26 Cartpole is one of the Classic Control environments used in reinforcement learning (RL). In this
27 environment, a pole is attached to a cart by an un-actuated joint, and the cart moves along a frictionless
28 track. The goal is to balance the pole by applying forces to the cart in the left and right direction. The
29 rewards for each step taken in Cartpole are +1, and to simplify the learning process and improve
30 convergence, we have capped the upper rewards to 500.

31 **2.2 Lunar Lander Environment**

32 This environment is a classic rocket trajectory optimization problem. The goal of this environment is
33 to control a lunar lander and guide it to a safe landing on the moon's surface. The lander is subject to

34 gravitational and thrust forces, and the agent must learn to apply the right amount of thrust to achieve
 35 a successful landing. A successful landing yields a high reward, while crashing or landing too fast
 36 results in a negative reward. For faster convergence the upper limit for rewards is capped at 250.

37 **2.3 Neural Network**

38 The neural network consists of a first layer that takes input of size equal to the state size of the
 39 environment, a hidden layer with a size of 64, and an output layer with the same size as the number
 40 of actions in the environment, with ReLU activation function. For optimization, I've used Adam
 41 Optimizer with Mean Squared Error to calculate the loss between target and expected Q - values.

42 **2.4 Soft updates of weights of Neural Network**

43 It is used to update the weights of the target network in a gradual manner, rather than updating them
 44 completely at each iteration. This helps to stabilize the learning process by reducing the variance
 45 of the target Q - values, which in turn helps to reduce the variance of the updates to the policy network.

$$Q_{target} = \tau Q_{policy} + (1 - \tau)Q_{target}, \quad \tau \ll 1 \quad (1)$$

46 Here, Q_{target} and Q_{policy} are the weights of the target network and policy network.

48 **3 Methodology**

49 **3.1 Exploration**

50 To enable my agent to achieve the maximum rewards in both environments, I have implemented
 51 the Epsilon Decay function, which gradually reduces the exploration rate with each time step. This
 52 function balances the trade-off between exploration and exploitation during the training phase.

53 **3.2 Training**

54 The agent was trained through 10 independent runs, each with 10 independent seeds, and for a total
 55 of 3000 episodes. At each time step, the state-action-reward-next state transition was stored in the
 56 experience replay buffer. The training process involved two neural networks: the policy network and
 57 the target network. The agent was trained on the policy network using mini-batches of transitions
 58 sampled from the experience replay buffer. The weights from the policy network were copied to the
 59 target network at a predefined update frequency. The target network was used to set the target values
 60 for the Q values, while the policy network was used to obtain the estimate of the state-action values.

61 **3.3 Hyperparameters**

62 Deep Reinforcement learning algorithms are very sensitive to hyperparameters, To have a better
 63 convergence it is very important to chose the values wisely. In the Table 1 is shown the used
 hyperparameter values.

Table 1: Values of hyperparameters

Best hyperparameters	
Batch size	128
Episodes	3000
γ	0.99
τ	0.01
Learn step	5
Learning rate	0.001
Memory size	10^4
Fixed exploration	0.2

65 **4 Experimental results**

66 **4.1 Training**

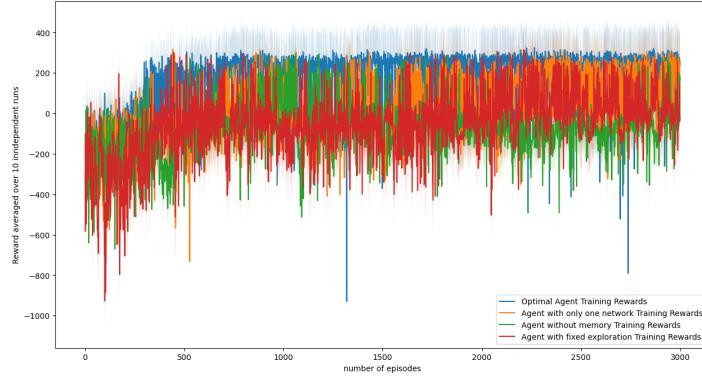


Figure 1: Lunar Lander training rewards

67 The training results for the Lunar Lander environment are displayed in Figure 1. Similar trends were
68 observed in the Cart Pole environment, where the optimal agent, with both policy and target networks,
69 converged significantly faster than other variants. The optimal agent also exhibited lower variance in
70 rewards, in contrast to the agent without memory.

71 Memory plays a crucial role in the DQN algorithm, as the agent is trained on batches of transitions
72 stored in experience replay. Interestingly, using only one network resulted in different outcomes, as
73 the target network is required to set the target values for the Q-values of the current state. Training the
74 agent with only one network requires more episodes to achieve the same performance as the optimal
75 agent with both policy and target networks.

76 As shown in Figure 1, the policy network alone took longer to converge, but continued training
77 ultimately led to better performance. In a fixed exploration rate was used, whereas others employed
78 epsilon decay exploration policies. Initially, the agent requires more exploration to discover the
79 optimal policy, but as it trains more, the exploration rate should decrease to avoid random actions
80 and choose more optimal ones. Results show that the agent with a fixed exploration rate had higher
81 variance than the optimal and one-network agents. Despite this, it still converged faster and more
82 efficiently than the agent without memory.

83 **4.2 Test**

Table 2: Average test rewards for 500 episodes

Algorithm	CartPole	Lunar Lander
Optimal	500	239.31
Only policy Network	265.858	196.15
Without memory	23.078	44.73
Fixed exploration	152.218	84.53

84 The test results for both environments are presented in Table 2. The results clearly show that the
85 optimal policy performs much better than the other types of policies. Even though the agent with
86 only a policy network has a single neural network, its performance is better than that of the agent
87 without memory and the one with fixed exploration. The results indicate that training an agent for
88 more episodes may lead to optimal policy-level performance, thereby reducing the computational
89 cost for calculating the loss and gradients of the neural network. The agent without memory is not
90 able to estimate the Q-values for the state, resulting in the agent taking random actions at each time
91 step. From Fig. 2 and 3, it can be observed that the agent with a fixed exploration rate performs
92 worse than the agent with only one policy network. This highlights the importance of having a good
93 exploration policy in reinforcement learning agents, as better exploration leads to better performance.

94 One important observation is that using the same neural network, hyperparameters, and exploration
 95 policy on different environments leads to different variances in the rewards. This suggests that having
 96 the same agent for different environments might not yield good results, as different environments
 97 require different sets of hyperparameters, neural network architecture, and exploration policies to
 98 achieve the best performance. Since the state space, reward structure, and actions are different for
 different environments, it is necessary to have different approaches for different environments.

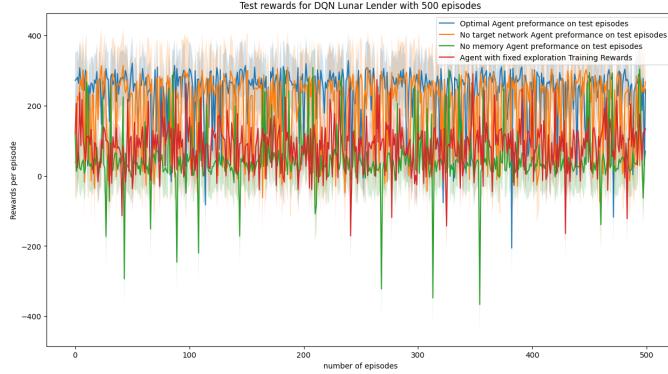


Figure 2: Lunar Lander test rewards

99

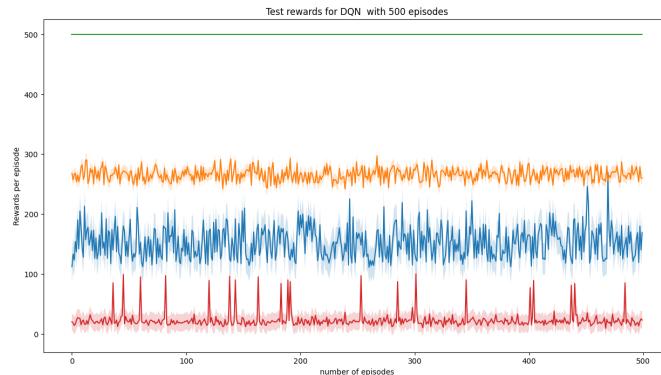


Figure 3: Cart Pole test rewards

100 5 Conclusion

101 The findings from the various Deep Q-learning algorithm variants emphasize the significance of
 102 each component for attaining optimal performance. Notably, using only a policy network can yield
 103 noteworthy outcomes, as long as there are adequate training episodes. Additionally, the exploration
 104 policy has a vital role in the agent's performance in the environment. Employing soft updates of the
 105 weight in the target network improves the stability and consistency of the results. The batch size
 106 and memory also have a substantial impact on the convergence of the Q-function, leading to faster
 107 and optimal performance. However, reproducing similar results in Deep Reinforcement Learning is
 108 challenging due to the reliance on critical parameters such as batch size, random seeds, learning rate,
 109 neural network architecture, exploration rate, learning step, and memory size. Therefore, selecting
 110 the appropriate parameter set is critical for achieving optimal performance in various environments.

111 6 Future work

112 Possible improvements can be done in many areas such as using Softmax policy for exploration, hyper
 113 parameters tuning, changing the reward structure. We can also improve the stability and performance
 114 of DQN algorithm by using Double Deep Q learning.

115 7 Appendix

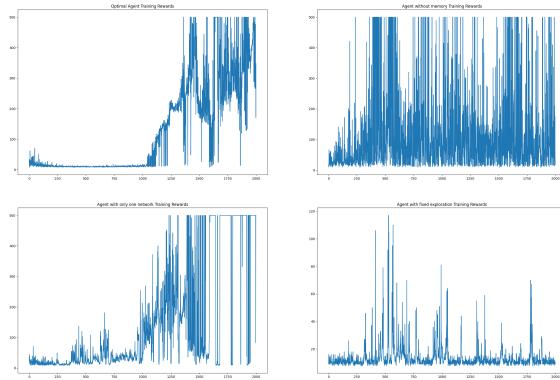


Figure 4: First 2000 episode rewards for Car Pole training for 5000 episodes

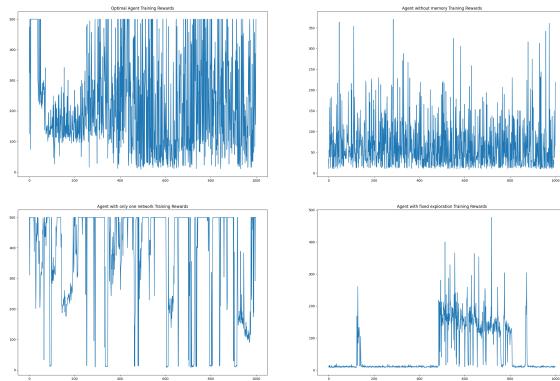


Figure 5: Overfitting during the last episodes while training Cart Pole environment for 5000 episodes

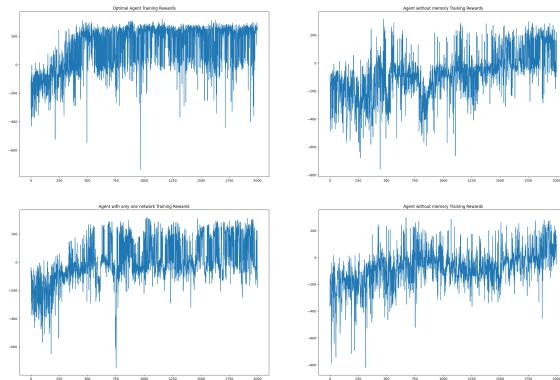


Figure 6: First 2000 episode rewards for Lunar Lander training for 5000 episodes

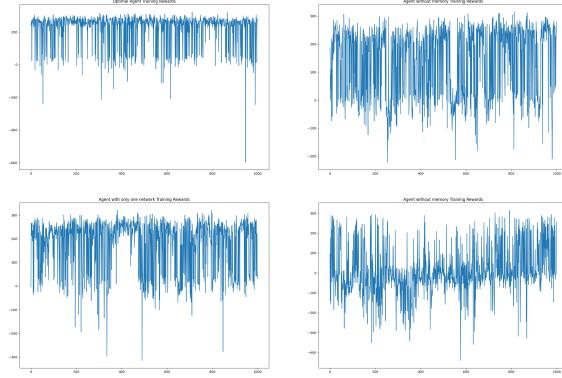


Figure 7: Overfitting during the last episodes while training Lunar Lander environment for 5000 episodes

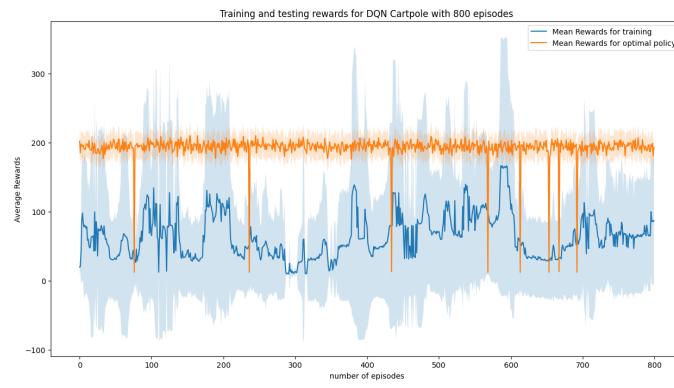


Figure 8: Without soft-updates of target network weights Cart Pole

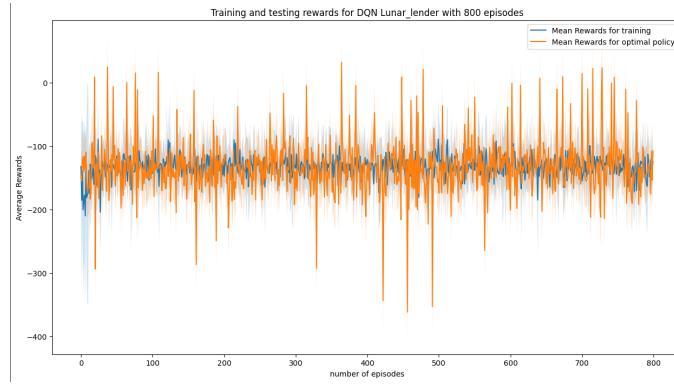


Figure 9: Without soft-updates of target network weights Lunar Lander

116 References

- 117 [1] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. and Riedmiller, M., 2013.
118 Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602.
- 119 [2] Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M.
120 and Silver, D., 2018, April. Rainbow: Combining improvements in deep reinforcement learning. In Proceedings
121 of the AAAI conference on artificial intelligence (Vol. 32, No. 1).

122 [3] Jomaa, H.S., Grabocka, J. and Schmidt-Thieme, L., 2019. Hyp-rl: Hyperparameter optimization by
123 reinforcement learning. arXiv preprint arXiv:1906.11527.