
Q Learning and Deep Q Network

Ling Fei Zhang

Department of Computer Science
McGill University
Montreal, QC
lzhang133@gmail.com

Abstract

In this paper, we compare two widely used RL algorithms, Q Learning and Deep Q Network (DQN), in the context of two games, CartPole and Lunar Lander. Our results demonstrate that DQN learns faster than Q Learning, but can be more unstable during the learning process. In contrast, Q Learning learns more steadily but at a slower rate. These findings suggest that the choice of the algorithm depends on the specific application and trade-offs between speed and stability.

1 Introduction

Reinforcement learning is a subset of artificial intelligence that concerns itself with developing agents that can learn to take actions in a stochastic environment with the goal to maximize a reward function. In recent years, reinforcement learning has gained considerable attention as it is a powerful approach that can be used in a wide range of applications, from robotics and gaming to finance and healthcare. Among the many algorithms in reinforcement learning, Q Learning has been one of the most popular and most used. However, Q Learning also has its limitations when it comes to handling large state spaces and continuous action spaces. To solve this problem, significant progress has been made by exploiting the advantages of deep learning in reinforcement learning, resulting in the “Deep Q Network (DQN)” algorithm [4]. DQN combines Q Learning with deep neural networks to learn policies that can handle complex state-action spaces. In this project, the two models will be tested against each other in the context of two games: CartPole and Lunar Lander.

Q Learning is a model free reinforcement learning algorithm that uses a table to store the Q values for each state-action pair. The Q value represents the expected return that the agent can obtain by taking a specific action in a given state. To learn, Q Learning updates the Q values based on the Bellman equation. To be more specific, the algorithm learns by updating the expected return in terms of the immediate reward and the expected return in the next state. In simple environments, Q Learning is very effective. However, the model struggles in more complex environments, usually due to the high dimensionality. This is known as the curse of dimensionality, where the number of actions increases exponentially with the number of degrees of freedom [3], making it impractical to store and update the Q values for all the state-action pairs.

To overcome the limitations of Q Learning, DQN was introduced. DQN combines Q Learning with deep neural networks to learn policies that can handle state-action spaces. The deep neural network approximates the Q values, which enables the agent to learn a function that maps the states to the Q values. Using neural networks as our function approximator, we can handle large state spaces as well as continuous action spaces.

However, no algorithm is perfect, as DQN faces its own challenge of instability. In fact, DQN can be quite unstable during the learning process and also requires careful hyperparameter tuning. The instability arises mainly from the non-stationarity of the learning process, where the Q values change as the agent learns, which then affects the target values used in the update rule. The need of careful

tuning of hyperparameters arises from the complexity of the DQN algorithm, which involves the use of multiple layers of neural networks, different learning rates and replay buffers.

2 Background

We consider tasks where the agent interacts with two environments: CartPole and Lunar Lander. Both environments have some level of stochasticity, as they have random initial states. At each time step, the agent must learn to select an action $a_t \in \mathcal{A} = \{1, \dots, K\}$. The action is then executed and a reward and the next state is observed. In fact, any episode can be defined as a sequence of the form $s_1, a_1, s_2, a_2, \dots, a_{n-1}, s_n$, where s_n is a terminal state. All sequences are assumed to terminate in a finite number of time steps. As a result, we have that all sequences are finite Markov decision processes (MDP) [4].

The goal of the agent is to select actions in a way such that the sum of rewards is maximized. We make the assumption that future rewards are discounted by a factor γ at each time step, and we define the future discounted return at time t as

$$R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \quad (1)$$

where T is the time step at which the game terminates.

Then, we define the optimal action-value function $Q^*(s, a)$ to be the function that achieves the maximum expected return. Formally, this means that after observing some state s and taking some action a , we must have

$$Q^*(s, a) = \max_{\pi} \mathbb{E}(R_t | s_t = s, a_t = a, \pi) \quad (2)$$

where π is a policy mapping states to actions.

It's important to note that the optimal action-value function follows an important identity: the *Bellman Equation*. The basic intuition of the *Bellman Equation* is as follows: if we know the optimal value $Q^*(s', a')$ for all actions a' of the next time step, then the optimal strategy is simply to select the action a' such that we maximize the expected value of $(r + \gamma Q^*(s', a'))$.

3 Methodology

3.1 Q Learning

The idea of Q Learning is to select an action given a particular state from a function Q , which the agent in turns receives an reward and a next state. The agent then updates Q in respect to maximizing the long term expected sum of rewards. The core of the Q Learning algorithm is defined by the update rule as follows

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (3)$$

In this case, the learned action-value function Q directly approximates q_* , the optimal action-value function, independent of the policy being followed. As a result, this drastically simplifies the analysis of the algorithm. As it turns out, all that is required for proper convergence is that all state-action pairs continue to be updated[6].

Q Learning's algorithm uses a table to store the Q values for each of the state-action pairs, and its limitation lies on the dimensionality of the state space. In fact, since both the CartPole and Lunar Lander games have continuous observation space, it is practically impossible to store the Q values for all state-action pairs. Thus, we must first discretize the continuous states into discrete states [1]. This can be done by taking the range of values for each feature of the state, and dividing the range into a specified number of bins. In the case that a feature has unlimited range (i.e. $\min = -\infty$ and $\max = \infty$), we can set a minimum and maximum manually based on the observations. For example, in CartPole, the cart velocity and the pole angular velocity are unbounded. However, from the observations, we decided that a minimum of -3.5 and a maximum of 3.5 is a good range for the cart velocity. To train our agent, we simply pass the required arguments to the Agent_Q class, and

train it on a desired number of episodes. In our case, we number of training episodes is 500 for both games.

We used very similar hyperparameters for both games, with the exception of the learning rate. In particular, for both games, we used a discount factor of 0.99, a start epsilon of 1, a minimum epsilon of 0.01, a discretization done by dividing the range of a feature into 10 equal sized bins [2], and an epsilon linear decay with a slope of 5^{-4} . With some experimentation, we discovered that the best learning rates for CartPole and Lunar Lander are 0.25 and 0.005 respectively.

3.2 DQN

To build our agent for DQN, a few steps are necessary as we need a deep neural network to approximate our Q values. Since DQN is a off policy reinforcement learning algorithm, it uses a technique called experience replay. Thus, our first step is to implement `ReplayMemory`, which acts as a replay buffer of past experiences. Each memory consist of a 5-tuple, namely (state, action, reward, next state, done). We also implemented the `sample` method since the algorithm learns by sampling from the buffer, and updates the Q network accordingly.

Next, we implemented the deep neural network, which consists of two hidden layers, each with a default number of 128 units. We used ReLU activation function, Adam's optimizer and calculated our loss with `SmoothL1Loss`. We decided to use the Huber loss during our implementation to make our agent more robust to outliers when the estimates of Q values are noisy. In fact, the Huber loss acts like the mean squared error when the error is small, and acts like the mean absolute error when it's large [5]. We used In the forward method, our network takes in a state, and outputs a Tensor of shape `action_space`, each element corresponding to the Q value of the state-action pair.

With this, we can build our RL agent. We built the agent with both a policy as well as a target network in order to make learning more stable. For our experiment, we initialized a replay buffer of size 100000, and trained our agent for 500 episodes. The actions are sampled based on an epsilon-greedy exploration. The full algorithm of DQN is described in Algorithm 2.

Unlike Q Learning, DQN does not need to discretize the states of a continuous observation space. This is because the continuous features can be passed directly into the deep neural network, and obtain the state-action Q values as outputs. However, DQN is computationally expensive compared to Q Learning, since we are working with two neural networks and performing gradient descent.

We used the same hyperparameters for both games. Specifically, we used a learning rate of 0.003, a discount factor of 0.99, an initial epsilon of 1, a minimum epsilon of 0.01, an epsilon decay of 5^{-4} , a tau of 0.005 and a batch size of 64.

4 Algorithm

4.1 Q Learning

Algorithm 1 Q Learning(episodes, α, ϵ, γ)

```

1: Initialize  $Q(s, a)$  for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$  arbitrarily.
2: Set  $Q(\text{terminal}, \cdot) = 0$  for all terminal states.
3: for each episode in episodes do
4:   Initialize  $s$ 
5:   done  $\leftarrow$  False
6:   while not done do
7:     Choose  $a \in \mathcal{A}$  from  $s$  using policy derived from  $Q$   $\triangleright$   $\epsilon$ -greedy, Softmax
8:     Take action  $a$ 
9:     Observe reward  $r$  and next state  $s'$ 
10:     $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_a Q(s', a) - Q(s, a)]$   $\triangleright$  Eq. (3)
11:     $s \leftarrow s'$ 
12:   end while
13: end for

```

In Algorithm 1, α is the learning rate, $0 < \gamma \leq 1$ is the discount factor, and ϵ is the exploration factor for the agent.

4.2 Deep Q Network

Algorithm 2 DQN(episodes, $\alpha, \epsilon, \gamma, C$)

```

1: Initialize replay buffer  $\mathcal{D}$  with maximum capacity 100000
2: Initialize policy network  $Q$  with random weights  $\theta$ 
3: Initialize target network  $\tilde{Q}$  with random weights  $\tilde{\theta}$  ▷ more stable learning
4: for each episode in episodes do
5:   Initialize  $s$ 
6:   Set  $t \leftarrow 0$  ▷ time step
7:    $done \leftarrow \text{False}$ 
8:   while not  $done$  do
9:     Choose  $a \in \mathcal{A}$  from  $s$  using policy network  $Q$  ▷  $\epsilon$ -greedy, Softmax
10:    Take action  $a$ 
11:    Observe reward  $r$ , next state  $s'$  and  $done$ 
12:    Store transition  $(s, a, r, s', done)$ 
13:    Sample a minibatch of random transitions  $(s, a, r, s', done)$  from  $\mathcal{D}$ 
14:    Set  $\hat{y} = \begin{cases} r & \text{if } s \text{ is a terminal state} \\ r + \gamma \max_a (s', a; \tilde{\theta}) & \text{otherwise} \end{cases}$  ▷ use target network
15:    Perform a gradient descent step on  $(\hat{y} - Q(s, a; \theta))^2$  with respect to the policy network parameters  $\theta$ 
16:    if  $t \bmod C = 0$  then
17:       $\tilde{Q} \leftarrow Q$  ▷ resets target network
18:    end if
19:     $C \leftarrow C + 1$ 
20:  end while
21: end for

```

5 Results

After training both agents for 500 episodes, we can observe the training return for both agents in Fig. 1.

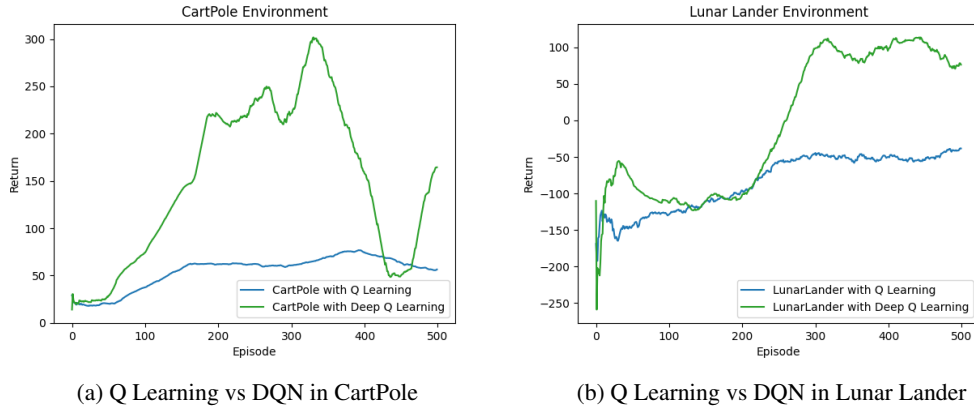


Figure 1: 500 training episodes for Q Learning and DQN. Returns are averaged over the last 100 episodes.

Immediately, we notice a significant learning improvement in DQN over Q learning in both games.

5.1 CartPole

In the CartPole environment, we can see that DQN easily reached an average return of 300 shortly after 300 training episodes. We also observe a big drop in average return after approximately 350 episodes. We suspect that this comes from the fact that the agent escaped a local minima, and tries to work its way to a global minima as the return peaks upwards again at around episode 450. We can see that DQN's return grows quickly, but also much more unstably to Q Learning. This is a direct consequence of the algorithm, as one of its limitations is the instability during the learning process.

As for Q Learning, we can see that the agent learns at a much slower but more steady rate compared to DQN. We can see that even after 500 episodes of Q Learning training, the agent seems to only obtain an average return of a little less than 100. This could be that the agent is stuck in a local minima. Due to the good returns, the agent might decide to stay in the local minima, rather than risking to escape it.

5.2 Lunar Lander

In the Lunar Lander environment, we can once more see that DQN outperforms Q Learning. Once more, we can see the properties of each algorithm. Clearly, DQN's learning is fast but unstable. The agent reaches a return of -50 at episode 50, but it decreases to approximately -125 at episode 150. Then, the agent seems to have learnt some important information, as the return spikes up to +100 shortly after. It then once more comes to a plateau for the remaining of the training episodes.

On the other hand, Q Learning has a slower but stable learning. The agent doesn't have any drastic change in average returns. Rather, the average return slowly increases.

6 Conclusion

In this paper, we have studied two popular reinforcement learning algorithms, Q Learning and DQN, on two different games, CartPole and Lunar Lander. Our results showed that DQN was able to learn faster in both games, but it exhibited some instability during the learning process. In contrast, Q Learning learned more steadily but at a slower pace. These findings suggest that the choice of algorithm depends on the specific application and trade-offs between speed and stability.

Future Work As a future direction, we suggest exploring modifications of Deep Q Learning, such as deep deterministic policy gradient (DDPG), which combines Q Learning with policy gradients to learn a deterministic policy directly. This approach has been shown to be effective in continuous action spaces and could potentially address some of the instability issues observed in DQN. Overall, our work highlights the importance of carefully selecting the appropriate reinforcement learning algorithm for a given task.

References

- [1] When should discretization of observations be considered?, 2022.
- [2] Rohan Gupta. An introduction to discretization techniques for data scientists, 2019.
- [3] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2019.
- [4] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.
- [5] George Seif. Understanding the 3 most common loss functions for machine learning regression, 2019.
- [6] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.