

---

# Performance of SARSA and Q-Learning on Cartpole and MountainCar

---

Rohan Akut  
260954888

## 1 MDP for Cartpole:

We have to balance a pole vertically on a cart in this problem. The entire description of the problem can be found in [1]

### 1.1 My understanding of the problem:

The problem supplies us with four states, namely: cart position, cart velocity, pole angle and pole angular velocity. These states are continuous, and the range of these values can be found in [2]. However, to formulate an MDP, we need discrete states. Hence I have converted the states from their continuous space to discrete space. The code to discretise the states has been referred from [3]

### 1.2 States:

As mentioned earlier, openAI supplies us with 4 states[2]. However, after referring to the explanation in [4], I understood that the two most important states in this problem are pole angle and pole angular velocity. As mentioned in [4] these two states are primarily responsible for maintaining a vertical position on the pole. Hence, I have just considered these two states to tackle this problem. Initially, I had just considered pole angle and ran the two algorithms. Post this, I have considered both the states (pole angle and pole angular velocity) and ran Sarsa and Q-learning on the two states.

### 1.3 Actions:

As mentioned in [2], the two possible actions are applying force to the left of the cart or right of the cart. I have not changed the actions, and the same actions are used for solving the problem.

### 1.4 Reward:

By default, openAI provides a reward of +1 for every timestep the pole stays upright[2]. However, I believe this reward function does not reflect the problem well. In my opinion, a reward should be provided for every correct action, and a penalty should be given for every wrong action. Hence I have used pole angle for my reward calculation. For every action taken, I check if the pole angle is greater than the pole angle in the previous timestamp. If the pole angle has increased, then it means that an incorrect action has been taken. Hence I punish the algorithm. The same applies when the pole angle decreases compared to the previous timestep. If the pole angle decreases, the algorithm is rewarded positively for taking the correct action. The idea for this custom reward has been referred from [5]

I have also added a huge penalty if the algorithm terminates before 30 timesteps. It is to be noted that the idea of adding this huge penalty was referred from [6].

## 30 1.5 Declarations to preserve academic integrity:

31 As mentioned previously, I have referred to multiple resources for formulating my MDP. In this  
32 section, I would like to summarise all the sources to ensure that the appropriate authors have been  
33 cited.

- 34 1. **States:** The idea for using just two states has been referred from [4] and the method for the  
35 discretisation of states has been referred from [3]
- 36 2. **Actions:** The possible actions have not been modified. Hence the default actions provided in  
37 [2] have been used.
- 38 3. **Reward:** The reward function has been designed by referring to two sources. The idea of  
39 penalising according to pole angle has been referred from [5]. Moreover, the idea of adding  
40 a huge negative reward when the algorithm finishes before 30 timesteps has been referred  
41 from [6].

42 I have also added the appropriate citations in the code file to preserve academic integrity.

## 43 2 MDP for MountainCar:

44 In this problem, we have to help the car reach the top of a mountain by using the laws of physics. The  
45 entire problem description can be found at [7]. In this section, I will explain the MDP formulation for  
46 solving this problem

### 47 2.1 States:

48 As mentioned in [8], the two states are: car position and car velocity. To solve this problem I have  
49 used both of these states. These states are available as continuous variables. Hence to discretise these  
50 states, I have referred to the code mentioned in [9].

### 51 2.2 Actions:

52 This problem has three possible actions: accelerate to the right, accelerate to the left, no accelera-  
53 tion[7]. I have used the default actions supplied by openAI for solving this problem.

### 54 2.3 Reward:

55 By default openAI gives a reward of -1 for every timestep where the car doesn't reach its goal. It  
56 provides a reward of 0 when the target position is achieved. However, I believe that this reward  
57 function would not help the car in reaching the target position. Hence I referred [5] to design my  
58 custom reward function. As mentioned in [5] I have given a positive reward every time the car gets  
59 closer to the goal. This means that instead of just setting a positive reward for the final target, the car  
60 would receive a reward if it climbs a certain portion of the hill. This idea of giving rewards according  
61 to car position is referred from [5]. I have given a pictorial representation of this custom reward  
62 function in Fig 1. The Fig 1 is not exactly to scale. However, it has been added to give a graphical  
63 representation of the custom reward chosen for this problem

## 64 2.4 Declarations to preserve academic integrity:

65 As mentioned previously, I have referred to multiple resources for formulating my MDP for Mountain  
66 Car Problem. In this section, I would like to summarise all the sources to ensure that the appropriate  
67 authors have been cited.

- 68 1. **States:** I have used the default states which are supplied by OpenAI[7]. However, these  
69 states have been discretised and the logic for discretisation has been referred from [9]

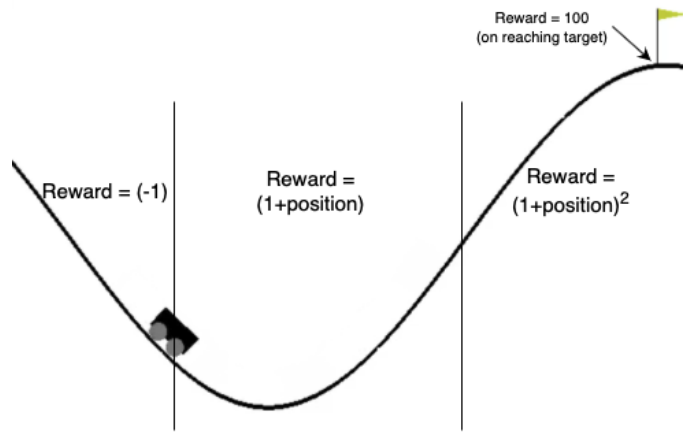


Figure 1: Custom Reward function for MountainCar[5]

2. **Actions:** The possible actions have not been modified. Hence the default actions provided in [7] have been used
3. **Reward:** The custom reward function has been referred from [5]. The pictorial representation of the reward function can be found in Fig 1.

### 3 Results for CartPole:

This section will describe the results that I obtained on running SARSA and Q-learning algorithms on CartPole problem. Before moving on to the MDP mentioned in Section 1.2 I have first created a simple MDP which just includes pole angle as a state. I have used this MDP to get a hang of openAI and understand the problem in a better sense. Hence initially, I will describe my findings on a Single State MDP for SARSA and Q-learning, and then I will move to the Two-State MDP mentioned in Section. 1.2

#### 3.1 Single State MDP:

I decided to run SARSA and Q-learning on a single state MDP because of the following reasons:

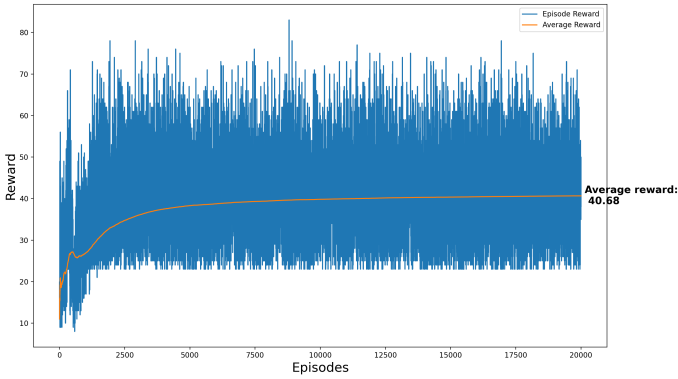
1. For this problem I was converting the continuous states into discrete values. However, I did not know the ideal values for discretisation. Hence I decided to test the discretisation values on a simple single-state MDP.
2. Even though the reward function has been referred from [5], I wanted to check if the reward function actually helps the models in converging faster

Hence I would like to mention that the purpose of Single-State MDP was **not to solve the problem, but to understand the system and the ideal discretisation values for the CartPole problem.**

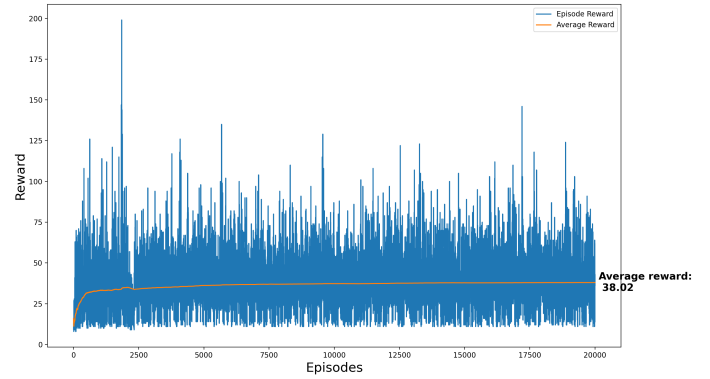
##### 3.1.1 Results for Single-State MDP:

In this section I will describe the results obtained after running SARSA and Q-learning on Single State Cartpole

**3.1.1.1. SARSA:** As mentioned in 3.1 the only state used in this MDP is pole angle. Using this as my state value I have tried to train SARSA. The SARSA model was trained for 20,000 episodes with a learning rate of 0.2 and discount factor of 0.7. The results of the training are as shown in Fig 2 The discussion for this figure has been provided in Section 3.1.2



(a) Single State SARSA



(b) Single State Q-learning

Figure 2: Performance of algorithms on Single State CartPole with learning rate = 0.2 and discount factor = 0.7

**3.1.1.2. Q-learning:** Similar to SARSA I have trained the Q-learning algorithm on CartPole problem for 20,000 episodes with a learning rate of 0.2 and discount factor of 0.7. The results for the training are as shown in Fig 2 Similar to SARSA, the discussion for Fig 2 has been provided in Section 3.1.2

### 3.1.2 Discussion for Single State MDP

In this section, I will provide an analysis of Fig 2. As mentioned in Section 3.1 the main reason to run a Single state MDP was to understand the problem properly. Hence after running both the algorithms, the following were my conclusions.

1. It is clear from Fig 2 that a single state MDP is not enough to solve the CartPole problem. However, it does show that both the models(SARSA and Q-learning) are learning and are converging to a maximum. This is encouraging as this proves that the implementation of these algorithms is, in fact, correct. With more information, these algorithms could solve the CartPole problem.
2. While designing the MDP, the reward function is crucial to help the model converge faster. If I had not added the custom reward function mentioned in Section 1.4 then both the algorithms do not converge properly. The average reward for both of the algorithms is around 20 as compared to 40 as visible in Fig 2. This shows that the implementation of my reward function is correct and it helps the algorithms in converging faster.
3. The discretisation level for this problem is very important. For example, if I do a coarse discretisation, I observed that SARSA and Q-learning never learnt anything worthwhile, and the average reward was less than 10. The same happened when the state discretisation was very fine. In this case these algorithms had too many options to choose from, and hence it got confused. Hence, it is imperative to strike a proper balance between the amount of discretisation to help these algorithms learn properly. It is to be noted that I have referred to [9], to find the ideal values for state discretisation. The same has been mentioned in the code as well.

Thus, after having an in-depth idea of the problem statement and understanding the ideal values of discretisation and reward functions, I shifted my focus from Single state MDP to a two-state MDP. In the following section, I will describe the performance of the 2-state MDP mentioned. The formulation of this MDP has been mentioned in Section 1.2

## 127 3.2 Two State MDP:

128 Building upon the information gained in Section 3.1.2, I have implemented a two-state MDP and  
 129 have run SARSA and Q-learning on this MDP. In this section, I would give a detailed analysis of the  
 130 performance of SARSA and Q-learning on this MDP

### 131 3.2.1 SARSA:

132 As mentioned in Section 3.1.2, the ideal reward function and the function approximation values  
 133 have been determined. However, we still need to find the ideal values of learning rate and discount  
 134 factor to solve CartPole problem successfully. Hence in this section, I would first present an analysis  
 135 of the performance of SARSA for different learning rates and discount factors. Then, once these  
 136 hyperparameters are set, I would run the SARSA algorithm on the CartPole problem.

137 **3.2.1.1. Determining Learning Rate and Discount Factor:** In order to determine these hyperpa-  
 138 rameters I have run SARSA on CartPole problem for 5000 episodes with the values of learning rate  
 139 varying from 0.1 to 1.0. The same has been done for discount factor. The results after running the  
 SARSA algorithm for different hyperparameters can be found in Fig 3. As visible in Fig 3 SARSA

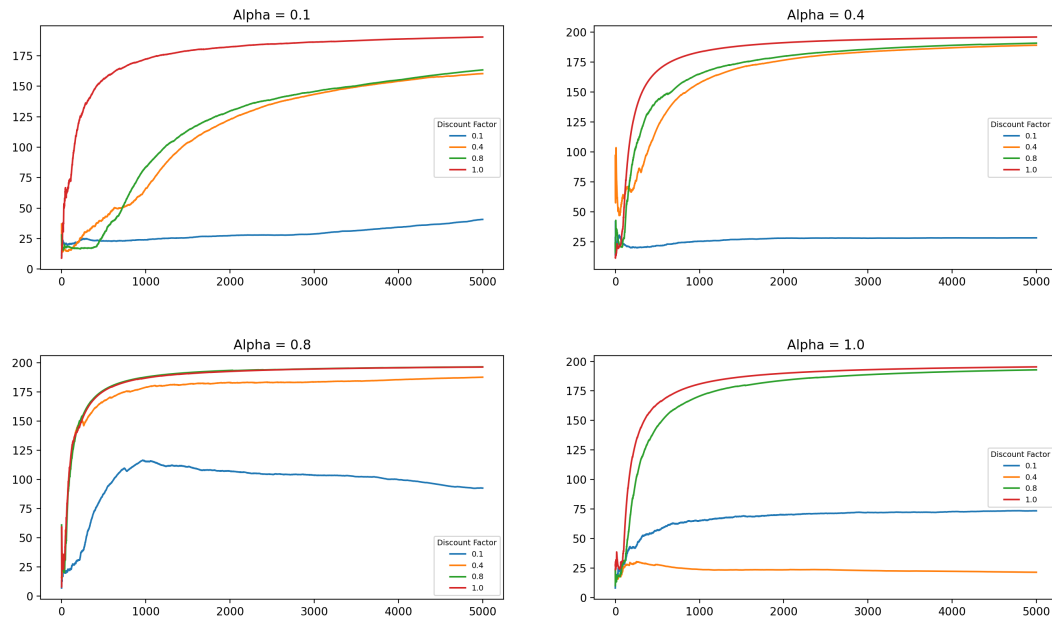


Figure 3: SARSA performance against different learning rate and discount factor values

140 gives a very high reward when learning rate = 0.8 and discount factor is either 0.8 or 1.0. Since the  
 141 discount factor of 1.0 could put an overdependence on future values, I have chosen value of 0.8 for  
 142 my learning rate and 0.8 for my discount factor.  
 143

144 **3.2.1.2. Performance of 2 State SARSA:** Since all the hyperparameters are fixed I have now run  
 145 SARSA on Cartpole problem. the algorithm is run for 3000 episodes with learning rate and discount  
 146 factor set to 0.8 as mentioned in Section 3.2.1.1. The performance of SARSA on Cartpole is as shown  
 147 in Fig 4. As visible in Fig 4 the SARSA reaches the maximum reward(200) in less than 500 episodes.  
 148 If I run this algorithm for a longer time, my algorithm would solve the problem. I will extend this  
 149 analysis to the Q-learning algorithm in the following sections.

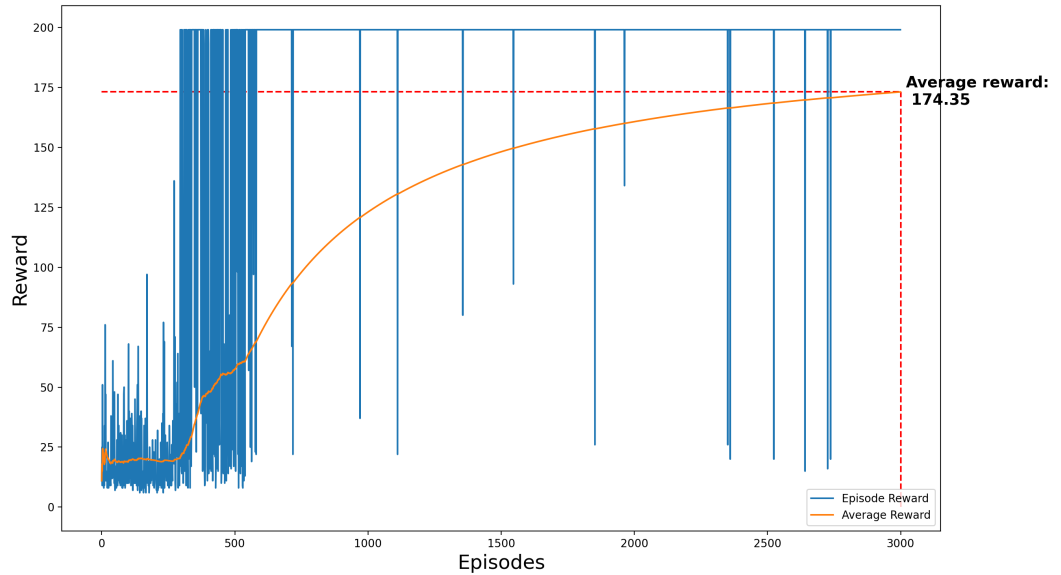


Figure 4: Two State SARSA with learning rate = 0.8 and discount factor = 0.8

### 150 3.2.2 Q-learning:

151 Similar to SARSA(Section 3.2.1) I have calculated the optimal values of learning rate and discount  
 152 factor for Q-learning algorithm. Post that I have run the Q-learning algorithm on Cartpole and  
 153 analysed the result.

154 **3.2.2.1. Determining Learning Rate and Discount Factor:** Similar to Section 3.2.1 I have varied  
 155 the hyperparameter values(learning rate and discount factor) between 0.1 to 1.0. The performance of  
 156 Q-learning on Cartpole is as mentioned in Fig 5. As visible in Fig 5 the learning rates of 0.4, 0.8 and  
 157 1.0 give a high reward. for the discount factor values of 0.4, 0.8 and 1.0. Since there is not much  
 158 difference in the reward for the above mentioned hyperparameter values, I have randomly chosen the  
 159 learning rate of 0.4 and discount factor of 0.8 for my problem.

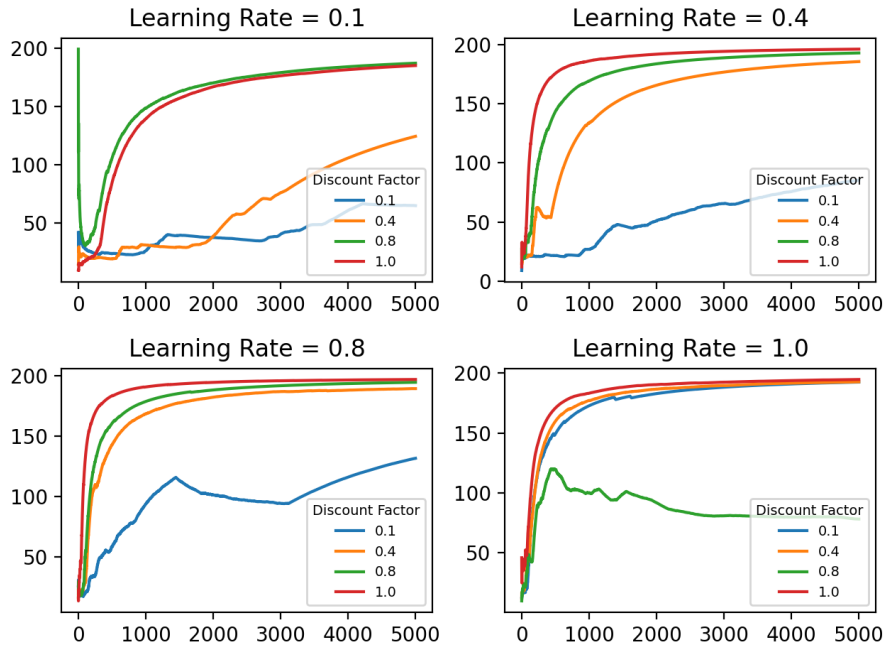


Figure 5: Q-learning performance against different learning rate and discount factor values

160 **3.2.2.2. Performance of Q-learning on 2State CartPole:** Since all the hyperparameters are fixed I  
 161 have now run Q-learning on Cartpole problem. The algorithm is run for 3000 episodes with learning  
 162 rate and discount factor set to 0.4 and 0.8 respectively. The performance of Q-learning on Cartpole is  
 as shown in Fig 6. As visible in Fig 6 the Q-learning reaches the maximum reward(200) in less than

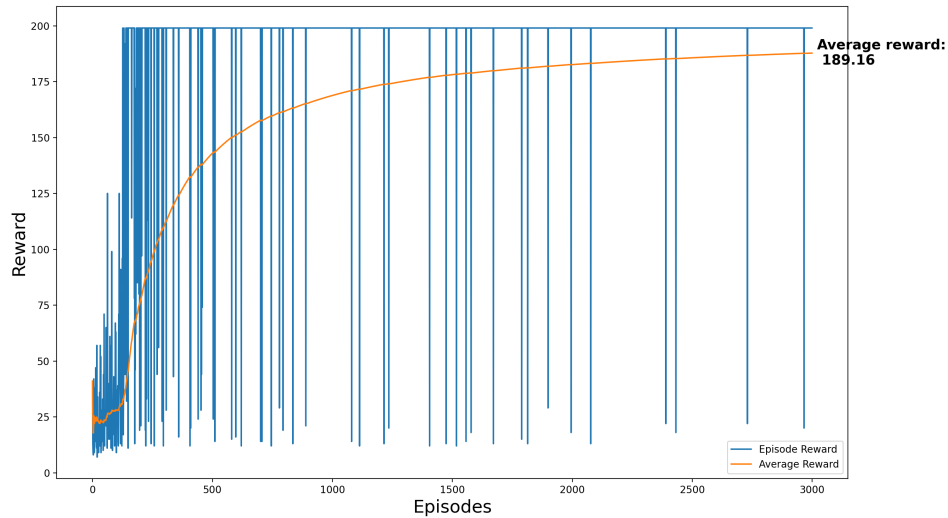


Figure 6: Two State Q-learning with learning rate = 0.4 and discount factor = 0.8

163  
 164 500 episodes. If I run this algorithm for a longer time, my algorithm would solve the problem.

165 Thus, in Section 3 I have initially determined the ideal reward function and the ideal discretisation  
166 values. Post that, I have determined the ideal values for hyperparameters and once all these values  
167 are set, I have run SARSA and Q-learning on the problem. I will perform a similar analysis on the  
168 MountainCar problem in the following section.

## 169 **4 Results for MountainCar:**

170 I have explained the MDP formulation for this problem in Section 2. In this section, I will discuss the  
171 performance of SARSA and Q-learning on the MountainCar problem using the MDP mentioned in  
172 Section 2.

173 Since I have not designed the reward function for this MDP and the discretisation has been referred  
174 from [5], I have not performed a Single State MDP analysis for the MountainCar problem. I have  
175 directly proceeded with a 2-state analysis.

176 In this section, I will analyse the performance of SARSA and Q-learning on the MountainCar  
177 problem. As mentioned in Section 2, I have used a two-state system to evaluate this problem. In  
178 the following sections, I will describe the performance of both the algorithms on the 2-state MDP  
179 mentioned in Section 2 .

### 180 **4.1 SARSA:**

181 Similar to Section 3, I have first determined the ideal values of hyperparameters and then tested the  
182 performance of SARSA on MountainCar. Hence in the following sections, I will first determine the  
183 optimal hyperparameters and then check the performance of SARSA on MountainCar.

#### 184 **4.1.1 Determining Learning Rate and Discount Factor:**

185 Similar to Section 3.2.1, the learning rate and discount factor are varied between 0.1 and 1.0. The  
186 algorithm has run for 5000 episodes, with every episode running for 1000 episodes. I have also  
187 performed an epsilon decay strategy to ensure that the agent explores the system before learning. The  
188 idea of epsilon decay has been referred from [9]. In order to ensure that I am choosing the optimal  
189 hyperparameter values, I have put an additional condition for termination. I have only terminated the  
190 training of the car if it reaches the terminal position more than five times. This ensures that only those  
191 hyperparameters are selected which learnt the optimal policy. The results can be found in Fig 7. From  
192 Fig 7 it is clear that when the discount factor value is high, the car reaches its goal five times and the  
193 training stops. These values are ideal values for training SARSA. In most cases, the car reaches the  
194 target(5 times) when the discount factor = 1.0. However, it is not common to set a discount factor of  
195 1.0. Hence, I have chosen a discount factor of 0.95. For learning rate, the ideal learning rate could be  
196 either 0.8 or 1.0. In this case I have chosen 0.8 as my learning rate.

### 197 **4.2 Performance of SARSA:**

198 I have used the same algorithm that I had used for CartPole problem. As mentioned in previous section  
199 the hyperparameter values are set to 0.95(discount factor) and 0.8(learning rate). The termination  
200 condition for this algorithm is when the algorithm reaches the target position ten times. The SARSA  
201 algorithm is run for 10,000 episodes with each episode running for 1000 timesteps. The results  
202 obtained after running the SARSA model on MountainCar are shown in Fig 8. As visible in Fig 8,  
203 the SARSA algorithm reaches the target position five times in just 5000 episodes. Thus the algorithm  
204 is terminated before the total assigned episodes(10,000) are completed. This proves that my SARSA  
205 implementation solves the MountainCar problem.



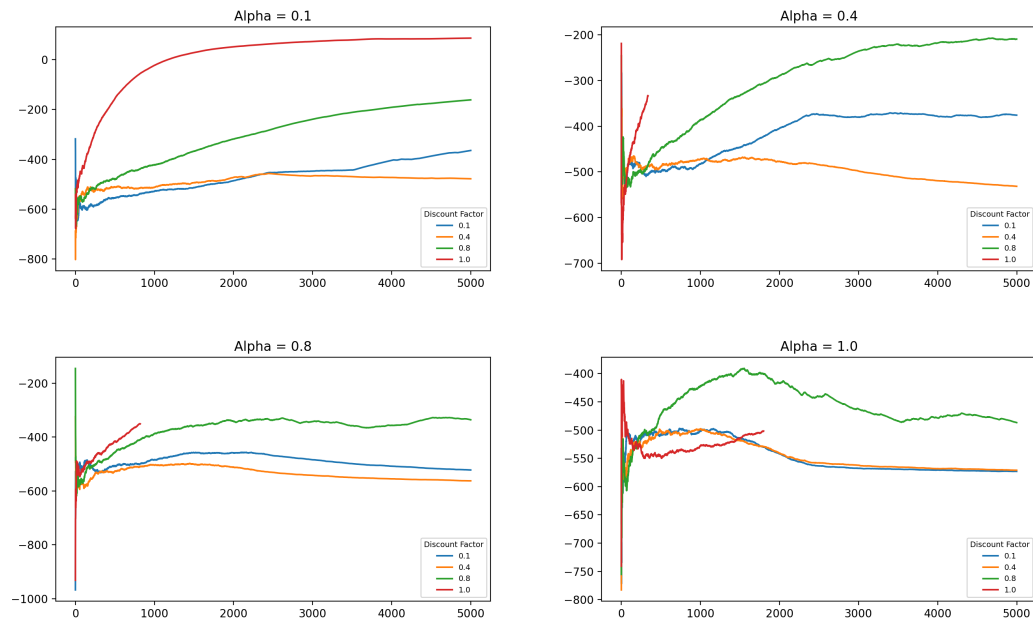


Figure 7: SARSA performance against different learning rate and discount factor values

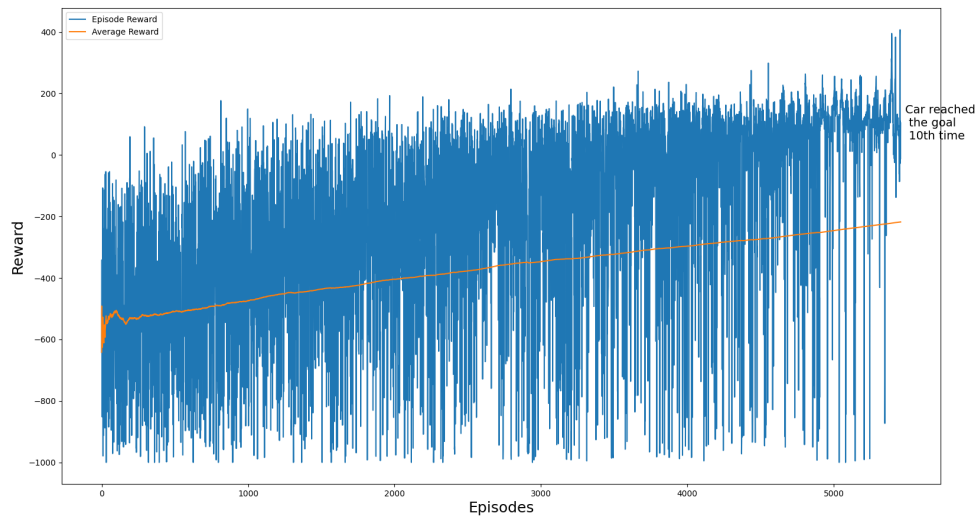


Figure 8: SARSA with learning rate = 0.8 and discount factor = 0.95

### 206 4.3 Q-learning:

207 In this section, I will discuss the performance of the Q-learning algorithm on the MountainCar  
 208 problem. Similar to the previous analysis, I would first find the optimal values for the hyperparameters,  
 209 following which I would test the performance of Q-learning on MountainCar.

### 4.3.1 Determining Learning Rate and Discount Factor:

Similar to Section 4.1.1 the learning rate and discount factor have been varied between 0.1 and 1.0. The algorithm has been run for 5000 episodes, with every episode running for 1000 timesteps. I have also performed an epsilon decay strategy to ensure that the agent explores the system before learning. The idea of epsilon decay has been referred from [9] .

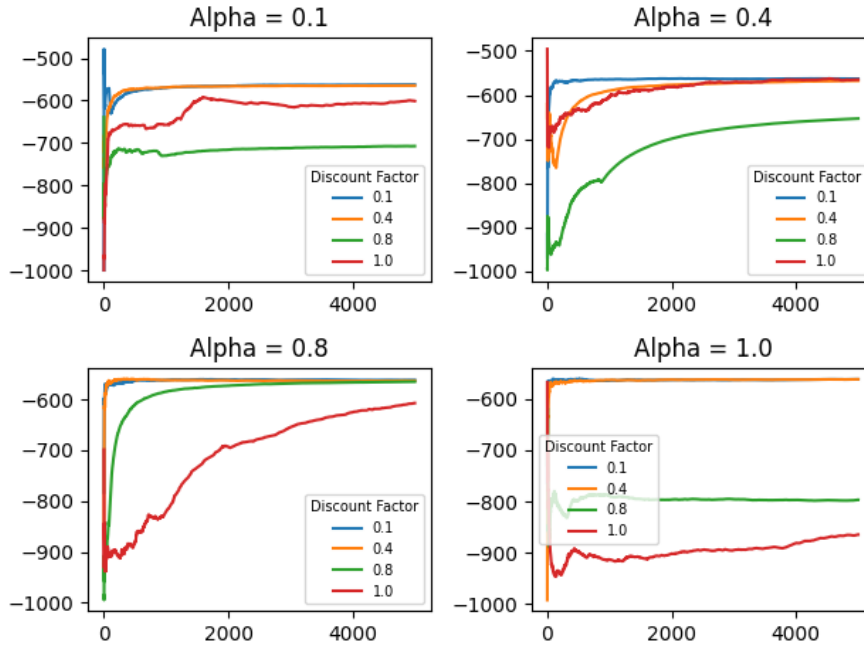


Figure 9: Q-learning performance against different learning rate and discount factor values

To ensure that I choose the optimal hyperparameter values, I have put an additional condition for termination. I have only terminated the training of the car, reaching the terminal position more than five times. This ensures that only those hyperparameters are selected to learn the optimal policy. The results can be found in Fig 9

From Fig 9 it is visible that my implementation of Q-learning was not able to solve this problem in 5000 episodes. I am not sure why the algorithm could not solve this problem as SARSA was able to solve it with similar settings. However, I have chosen the best hyperparameter values among the available choices to perform analysis. I believe the highest reward is obtained when the learning rate is 0.4. In that setting, the discount factors of 0.1, 0.4 and 1.0 give similar rewards. Hence I have randomly chosen the value of 0.4 as my discount factor. Thus the learning rate and discount factor both are set to 0.4. In the following section, I will test the performance of the Q-learning model on the MountainCar problem.

### 4.3.2 Performance of Q-learning:

In this section, I will test the performance of Q-learning on the MountainCar problem. The algorithm is run for 10,000 episodes, where each episode is run for 1000 timesteps. The algorithm is terminated when the car reaches the goal five times. The results of the performance are as shown in Fig 10. As visible in Fig 10 and as mentioned previously, the Q-learning model did not reach the goal in 10,000 episodes. If the Q-learning model had learnt properly, the algorithm would have terminated with a positive reward well before 10,000 episodes were completed. However, that is not the case as visible in Fig 10. Hence i believe my Q-learning implementation was **not** able to solve the MountainCar problem.

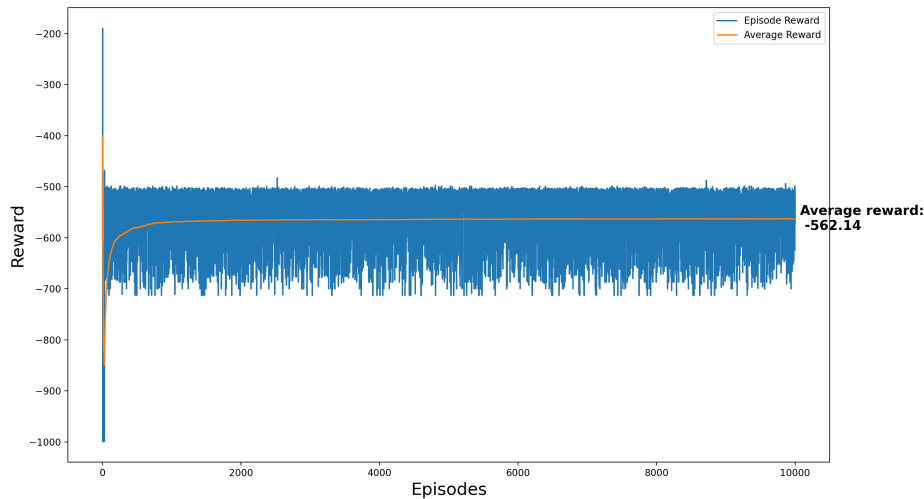


Figure 10: Q-learning with learning rate = 0.4 and discount factor = 0.4

## 5 Task 3:

### 5.1 On-policy and Off-policy:

In on-policy methods try to improve the same policy which is used by the agent to navigate the environment[15]. Similarly off-policy methods use two different policies. One policy is used by the agent to navigate through the environment while the second policy is the policy which actually learns from the environment.

SARSA is an example of on-policy model and Q-learning is an example of off-policy model. The reason SARSA is called as an on-policy model is because SARSA updates its current state Q-values by considering the **next action according to its current policy** [10]. This means that SARSA assumes that to move to the next state we have to follow the current policy and take the best action according to the policy[10]. Hence it is an on-policy model.

Q-learning is called an off-policy model because it follows 2 different policies. In Q-learning the agent is asked to take the next action greedily.[10] The agent takes the next action by not following the actual policy but by choosing the next action that has highest Q-value. However, when the update for the state-action value pairs happen, the updates happen according to the original policy and not the greedy policy. This means that in Q-learning, the agent follows a greedy policy while the update happens on the original policy[10]. Since two policies are followed in Q-learning, it is an off-policy method.

### 5.2 Model-based and Model-free Reinforcement Learning:

In model-based RL we need to have a model of the environment so that we can train the agent[15]. This means that in model-based RL the environment that we have modelled has an effect on the predictions/rewards of the agent [11]. Thus if we get the model wrong then our agent would never be able to successfully navigate through the environment. Dynamic Programming(DP) is an example of model-based RL[11]. In DP we need the transition probabilities among different states so that the agent can decide which state to choose that will maximise its reward. These transition probabilities are in a way a model of the environment[11]

262 In model-free RL, the agent has no predetermined environment model[15]. In fact the agent is  
263 initialised with some random values(also called as Q-values) and the agent learns to navigate through  
264 the environment on its own, by gaining experience over time[11]. Q-learning and SARSA are an  
265 example of model-free RL. These algorithms are initialised with random Q-values and the agent is  
266 run through the environment multiple times. As the agent gains more **experience** the Q-values are  
267 updated and eventually these algorithms converge to the optimal policy. Thus model-free methods  
268 learn from experience.

### 269 **5.3 Policy-based and Value-based RL:**

270 In Policy-based methods the agent chooses a random policy at the start. The initial policy does not  
271 have to be optimal. As the agent navigates through the environment, for every action the agent takes  
272 a gradient is generated and we follow a gradient ascent strategy to maximise the expected reward.  
273 Thus using the gradient ascent strategy we try to find the optimal policy[12]. In policy-based methods,  
274 the agent directly learns the policy and the actions are sampled from the policy directly[12]. This  
275 ensures that there is some stochasticity in the policy based methods. REINFORCE is an example of  
276 policy-based methods.

277 In REINFORCE, we select a random policy and take N steps using that policy. Once the steps  
278 are completed we calculate the discounted reward and using backpropagation we modify the action  
279 probabilities for the policy. We can then perform gradient ascent to maximise the expected reward.  
280 This procedure followed by REINFORCE is very similar to the procedure mentioned above. Hence  
281 REINFORCE is a policy-based RL algorithm

282 In value based methods we learn the state-action value pairs for an environment. Using these  
283 state-action pairs we can determine the optimal policy. However, in these methods our aim is not to  
284 learn the policy directly, but to learn the state-action values and then find the optimal policy from  
285 those values. Since we do not sample actions from a probability distribution, value-based methods  
286 are deterministic by nature. Q-learning and SARSA are examples of value based methods.

### 287 **5.4 Planning and Learning:**

288 Planning as the name suggests involves pre-deciding agent's actions before the agent navigates  
289 through the environment. For example in model-based methods we tend to determine the model of  
290 the environment and in that process we "plan" how the agent should navigate through the environment  
291 [15]. Thus pre-deciding how the agent should navigate through the environment is called planning[15].  
292 Dynamic Programming(DP) is an example of planning before the agent navigates through the system.  
293 In DP the transition probabilities are determined before the agent starts its navigation[15]. Thus  
294 determining the transition probabilities means determining the agent's actions before the agent starts  
295 navigating through the environment. This is known as planning.

296 Learning as the name suggests does not involve any planning. This means that we do not tell the  
297 agent any additional information and we let the agent explore the environment. The agent learns  
298 through experience and decides the optimal policy[15]. SARSA and Q-learning are examples of  
299 learning mechanism. These algorithms update the Q-values as they navigate through the environment  
300 and eventually they learn the optimal policy. These algorithms have no pre-assumption about the  
301 environment and all the experience is gained after navigating through the environment.

### 302 **5.5 Offline and Online RL:**

303 In RL the agent tries to navigate through the environment and during that process it gains experience.  
304 Using this experience the policy is modified and when repeated a large number of times we get an  
305 optimal policy. There are 2 categories through which an agent can gain this experience.

306 The first category is online learning, where the agent is made to run in an actual real world  
 307 environment. Thus using the real world experience the agent tries to learn the optimal policy[14].  
 308 Q-learning and SARSA could be given as examples of online RL as these models are generally  
 309 trained on actual environment. Even though, online RL is an ideal way to train an agent, this method  
 310 is very time consuming, expensive and sometimes even dangerous [14]. Hence, researchers are trying  
 311 to explore another alternative which is described below.

312 The second method is an offline method. In this method the agent tries to gain experience from  
 313 data which has been collected and stored. This setting is very similar to a supervised learning method  
 314 where a model tries to learn from a pre-collected and stored dataset[14]. This method is much more  
 315 viable because it is less time consuming and less expensive. However, this method suffers from  
 316 problems of its own. For example, this method is known to suffer from the problem of distributional  
 317 shifts[14] etc. Recently an algorithm called "Conservative Q-learning" was introduced in the field  
 318 of offline-RL [13]. This model tends to tackle the problem of distributional shift by conservatively  
 319 placing the importance on future unseen values.

## 320 References

- 321 [1] URL: <https://gym.openai.com/envs/CartPole-v0/>.
- 322 [2] URL: [https://github.com/openai/gym/blob/master/gym/envs/classic\\_](https://github.com/openai/gym/blob/master/gym/envs/classic_control/cartpole.py)  
 323 [control/cartpole.py](https://github.com/openai/gym/blob/master/gym/envs/classic_control/cartpole.py).
- 324 [3] URL: [https://github.com/philtabor/Youtube-Code-Repository/blob/master/](https://github.com/philtabor/Youtube-Code-Repository/blob/master/ReinforcementLearning/Fundamentals/sarsa.py)  
 325 [ReinforcementLearning/Fundamentals/sarsa.py](https://github.com/philtabor/Youtube-Code-Repository/blob/master/ReinforcementLearning/Fundamentals/sarsa.py).
- 326 [4] URL: [https://towardsdatascience.com/how-to-beat-the-cartpole-game-in-5-](https://towardsdatascience.com/how-to-beat-the-cartpole-game-in-5-lines-5ab4e738c93f)  
 327 [lines-5ab4e738c93f](https://towardsdatascience.com/how-to-beat-the-cartpole-game-in-5-lines-5ab4e738c93f).
- 328 [5] URL: [https://shiva-verma.medium.com/solving-reinforcement-learning-](https://shiva-verma.medium.com/solving-reinforcement-learning-classic-control-problems-openaigym-1b50413265dd)  
 329 [classic-control-problems-openaigym-1b50413265dd](https://shiva-verma.medium.com/solving-reinforcement-learning-classic-control-problems-openaigym-1b50413265dd).
- 330 [6] URL: <https://github.com/JackFurby/CartPole-v0/blob/master/cartPole.py>.
- 331 [7] URL: <https://gym.openai.com/envs/MountainCar-v0/>.
- 332 [8] URL: [https://github.com/openai/gym/blob/master/gym/envs/classic\\_](https://github.com/openai/gym/blob/master/gym/envs/classic_control/mountain_car.py)  
 333 [control/mountain\\_car.py](https://github.com/openai/gym/blob/master/gym/envs/classic_control/mountain_car.py).
- 334 [9] URL: [https://github.com/philtabor/Youtube-Code-Repository/blob/master/](https://github.com/philtabor/Youtube-Code-Repository/blob/master/ReinforcementLearning/Fundamentals/mountaincar.py)  
 335 [ReinforcementLearning/Fundamentals/mountaincar.py](https://github.com/philtabor/Youtube-Code-Repository/blob/master/ReinforcementLearning/Fundamentals/mountaincar.py).
- 336 [10] URL: [https://stats.stackexchange.com/questions/184657/what-is-the-](https://stats.stackexchange.com/questions/184657/what-is-the-difference-between-off-policy-and-on-policy-learning)  
 337 [difference-between-off-policy-and-on-policy-learning](https://stats.stackexchange.com/questions/184657/what-is-the-difference-between-off-policy-and-on-policy-learning).
- 338 [11] URL: [https://stats.stackexchange.com/questions/184657/what-is-the-](https://stats.stackexchange.com/questions/184657/what-is-the-difference-between-off-policy-and-on-policy-learning)  
 339 [difference-between-off-policy-and-on-policy-learning](https://stats.stackexchange.com/questions/184657/what-is-the-difference-between-off-policy-and-on-policy-learning).
- 340 [12] URL: <https://www.youtube.com/watch?v=cQf0QcpYRzE>.
- 341 [13] Aviral Kumar et al. "Conservative q-learning for offline reinforcement learning". In: *Advances*  
 342 *in Neural Information Processing Systems* 33 (2020), pp. 1179–1191.
- 343 [14] Sergey Levine et al. "Offline reinforcement learning: Tutorial, review, and perspectives on  
 344 open problems". In: *arXiv preprint arXiv:2005.01643* (2020).
- 345 [15] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,  
 346 2018.