
CSE 546: Reinforcement Learning

Assignment 1

Part 1

1
2
3 **Sumeet Milind Aher**
4 Department of Computer Science
5 UB ID: 50418795
6 sumeetmi@buffalo.edu
7

8 **Abstract**

9 **This is a submission for Assignment 1 for the course of Reinforcement**
10 **Learning.**

11 **1 Environment(Deterministic and Stochastic)**

12 **1.1 Set of Actions**

13 There are 4 allowed actions- a) Up, b) Down, c) Left, d) Right

14 **1.2 State**

16 There are 8 states in this environment:

17 a) Rivendell: This is where the elves will help our agent recover in the
18 journey.

19 b) Aragorn: This is where our heir to the throne named Aragorn will help our
20 agent reach the goal

21 c) Gollum: This is where the agent meets Gollum, a wild and deserted
22 creature living from ages in search for the ring and he will stop the agent.

23 d) Mountains of Mordor: These are very hard volcanic terrains filled with
24 goblins and orcs making it very hard and risky for our agent to reach the ring.

25 e) Door of Moria: This door leads the agent to an ancient hidden dwarf maze,
26 it's the door to go anywhere on the map. But, beware! Nobody knows these
27 tunnels, so the agent doesn't know where he will come out from once, he
28 arrives on this state. The agent is thrown on a random location in the
29 environment once it lands on this state.

30 f) Gandalf: This state is where the agent meets Gandalf and is blessed by his
31 wisdom!

32 g) Ring: This state is the terminal state where the agent will find the ring and
33 complete his quest.

34

35 **1.3 Rewards**

36 Each of the state have their own rewards:

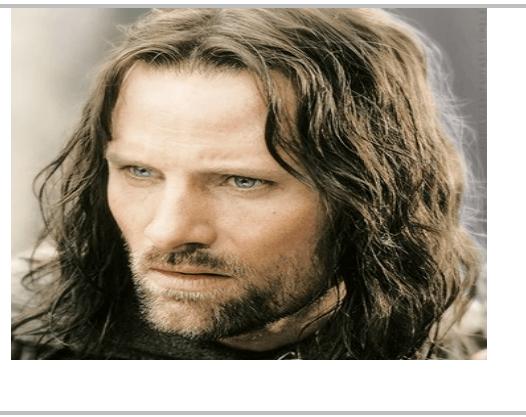
37 a) Rivendell: +2



38 Figure 1: Rivendell

39

40 b) Aragorn: +2



41 Figure 2: Aragorn

42

43 c) Gollum: -2



44 Figure 3: Golum

45

46 d) Mountains of Mordor: -4



47

Figure 4: Mordor

48

49 e) Door of Moria: 0



50

Figure 5: Moria

51

52 f) Gandalf: +4



53

Figure 6: Gandalf

54

55 g) Ring: 20



56 Figure 7: Ring

57

58

59

60 **2 Main Objective**

61 This is the environment of Lord of the Rings. Our agent (Frodo Baggins)
62 wants to reach to the Ring so he can destroy it. But the ring is far far away in
63 the land of Mordor. To reach, he must take help from the elves at Rivendell,
64 the heir to the throne Aragorn and Gandalf the White wizard. There is also a
65 mystical door which will randomly take the agent anywhere in the
66 environment. And our agent must stay away from Gollum and the mountains
67 of Mordor.

68

69 Agent: Frodo Baggins: Our agent is chosen to find the ring in the environment



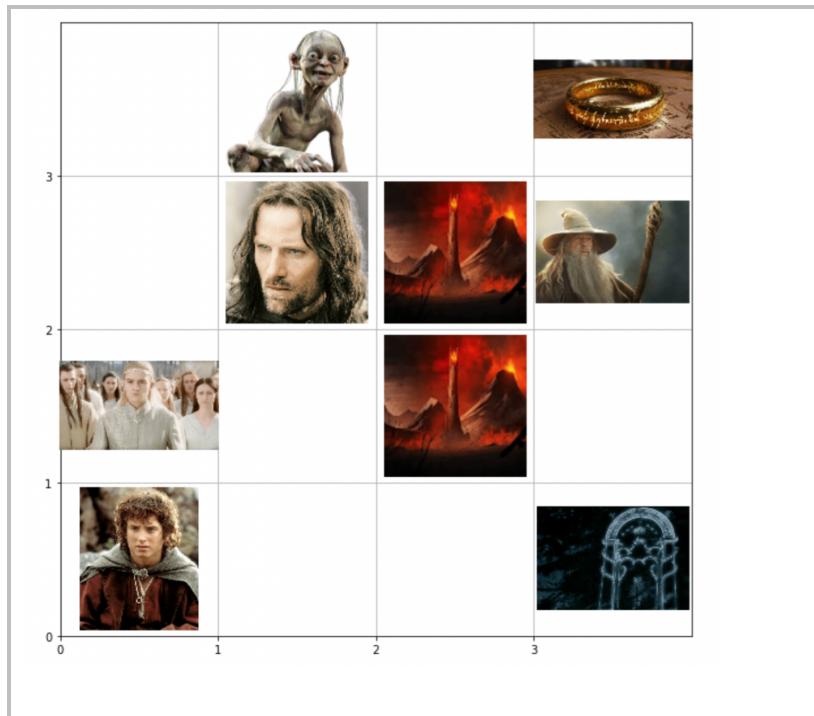
70 Figure 1: Agent

71

72

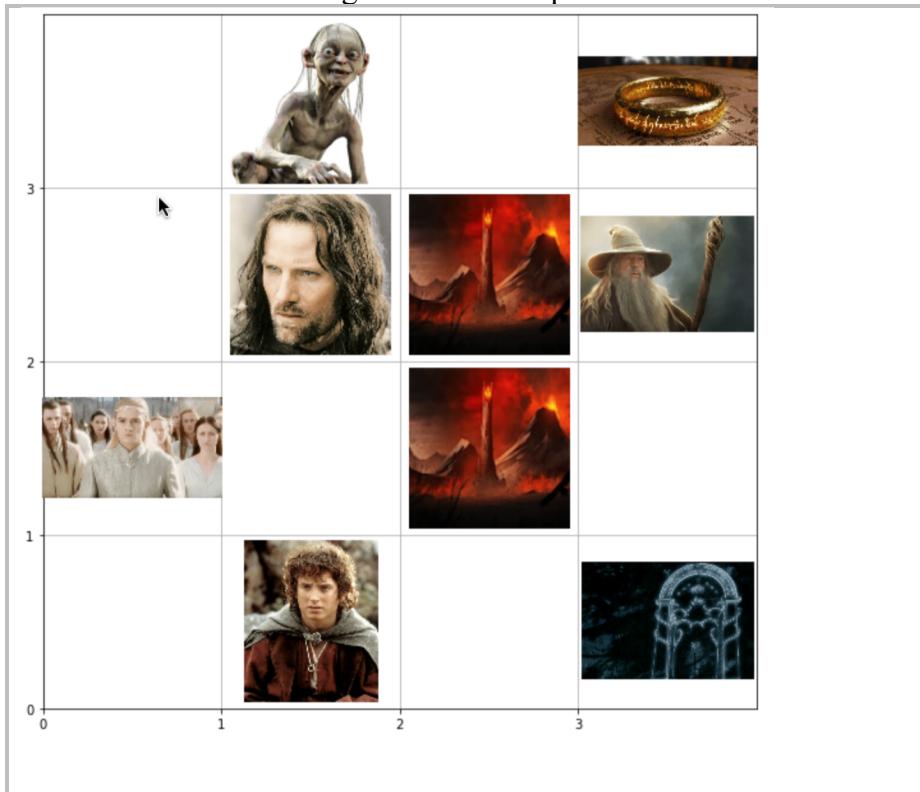
73 **3 Visualization**

74 These are some of the timesteps in the environment:



75

Figure 1: Timestep = 1



76

Figure 1: Timestep = 2



Figure 1: Timestep = 3

77

78

79

80 4 How did you define the stochastic environment?

81 In the stochastic environment, whenever the agent decides to take an action,
 82 the probability that the agent will be successful with that action is 95%, in the
 83 case of the remaining 5%, the agent will choose the UP action instead.

84

85 4.1 What is the difference between the deterministic and 86 stochastic environments?

87 In the deterministic environment, the agent would go right if the action
 88 decided is right! And similarly for other actions, the probability of execution
 89 of that action is 100%. On the other hand, in a stochastic environment, the
 90 probability of execution of a particular action is 95% and in the rest 5%, the
 91 agent will have to choose the UP action..

92

93 4.2 Safety in AI: Write a brief review explaining how you 94 99 ensure the safety of your environments.

95 These are the safety features of my environment:

96 4.2.1) Whenever the agent takes an action and ends up outside the
 97 environment, the code will clip the result between the limits of the grid length
 98 and width thus keeping the agent in the environment.

99 4.2.2) Because of 5.1, the agent might keep stepping outside the grid from a
100 rewarding location, and because of our safety feature, it will again land on the
101 same rewarding location, thus keep on collecting the reward, to stop this,
102 there is an additional safety feature which will give the reward only once.

103

Part 2

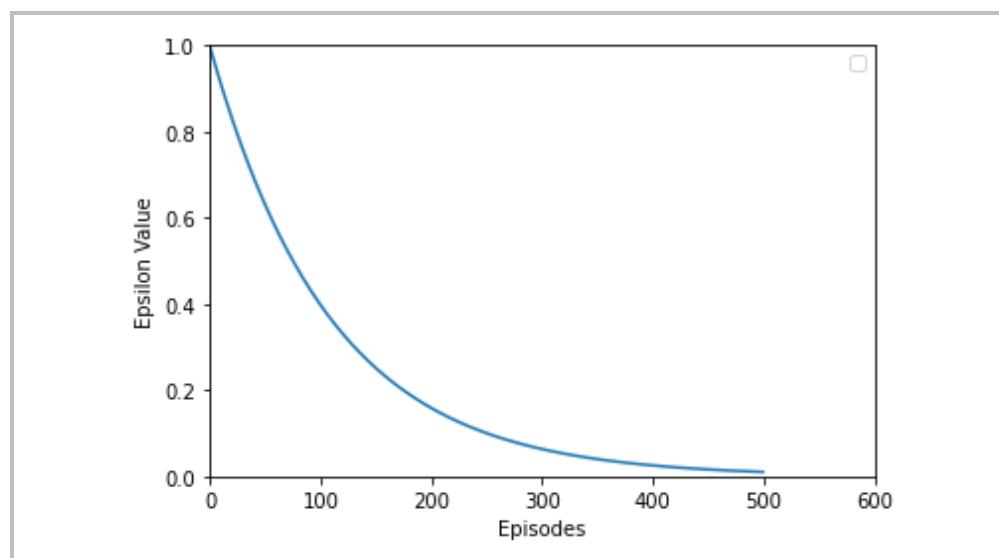
104

105

106 5 Results

107 All the results are for
108 alpha = 0.15
109 gamma = 0.95
110 episodes = 500
111 steps = 10
112 epsilon max = 10
113 epsilon min = 0.01
114

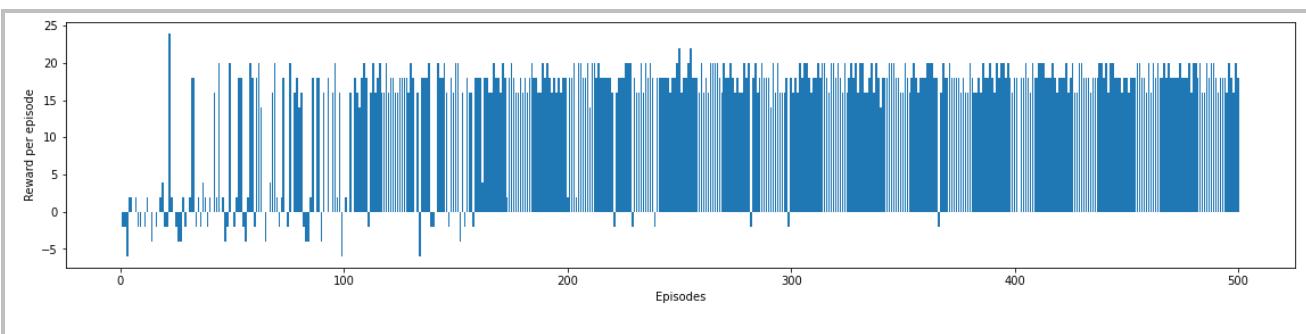
115 5.1) a) Applying Q-learning to solve the deterministic environment defined in Part 1.
116 Plots should include epsilon decay and total reward per episode.



117

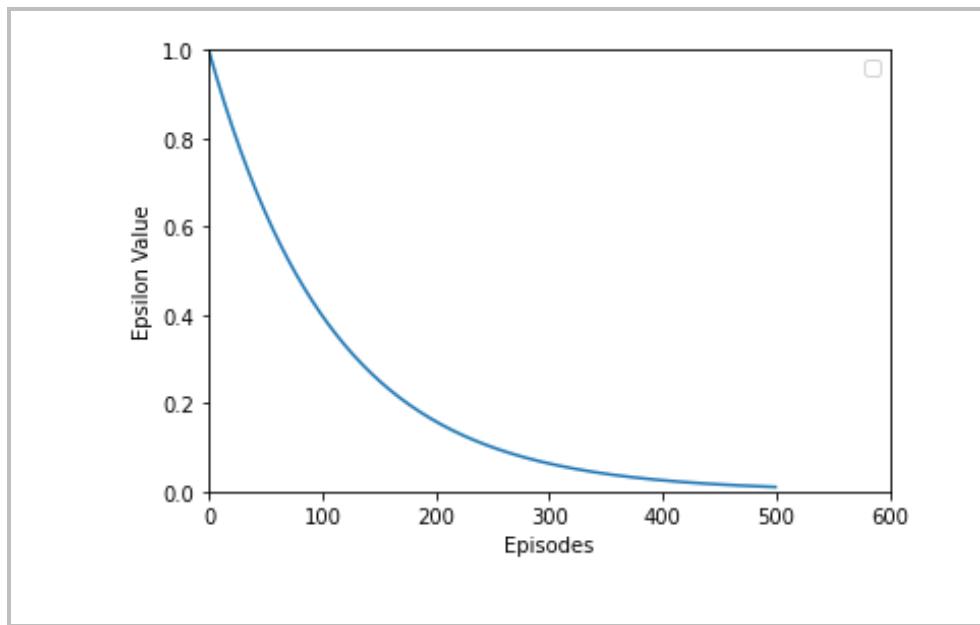
Figure 1: Epsilon Value vs Episode

118



119 Figure 2: Reward per episodes vs Epsiode for Q Learning Deterministic env

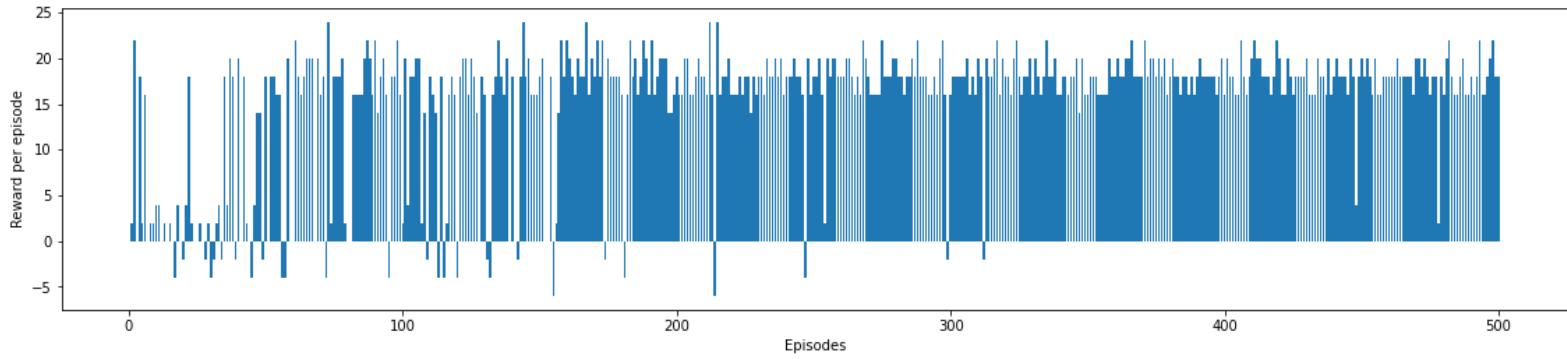
- 120 • b) Applying Q-learning to solve the stochastic environment defined in Part 1. Plots
121 should include epsilon decay and total reward per episode.



122 Figure 3: Epsilon decay vs Episode

123

124

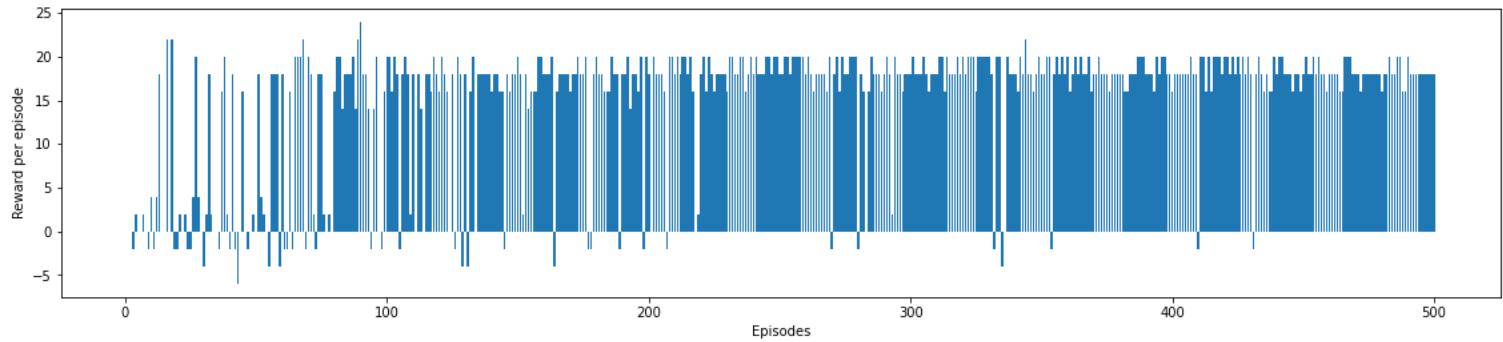


125
126

Figure 4: Reward per episode vs Episode For Q Learning Stochastic

127
128

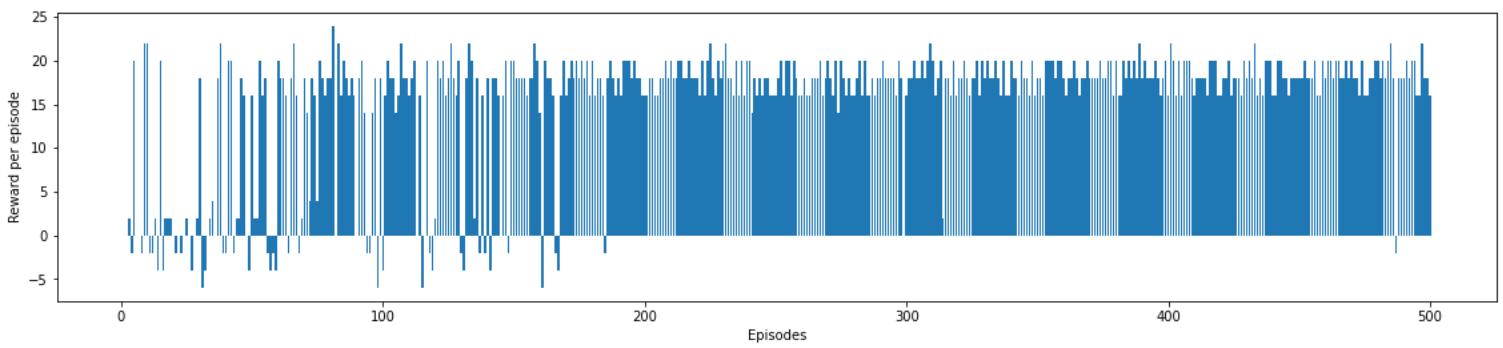
- c) Applying any other algorithm of your choice to solve the deterministic environment defined in Part 1. Plots should include total reward per episode.



129
130
131
132

Figure 5: Total Reward vs Episode For SARSA Deterministic env

- d) Applying any other algorithm of your choice to solve the stochastic environment defined in Part 1. Plots should include total reward per episode.



133

Figure 6: Reward per episodes vs episodes for SARSA Stochastic

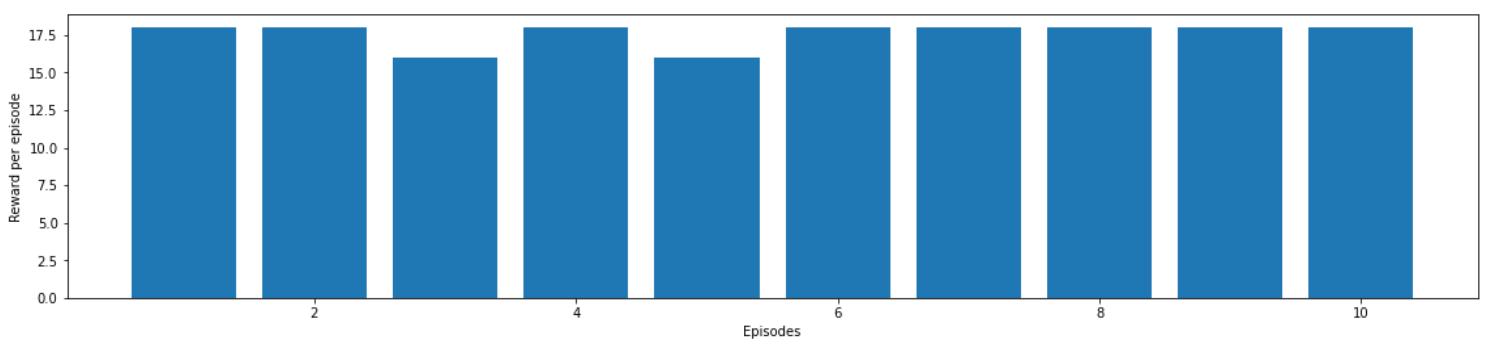
134

135

136

137

- e) Provide the evaluation results. Run your environment for at least 10 episodes, where the agent chooses only greedy actions from the learnt policy. Plot should include the total reward per episode.



138

Figure 7: Reward per episode vs Episodes (Deterministic env | Q learning)

139

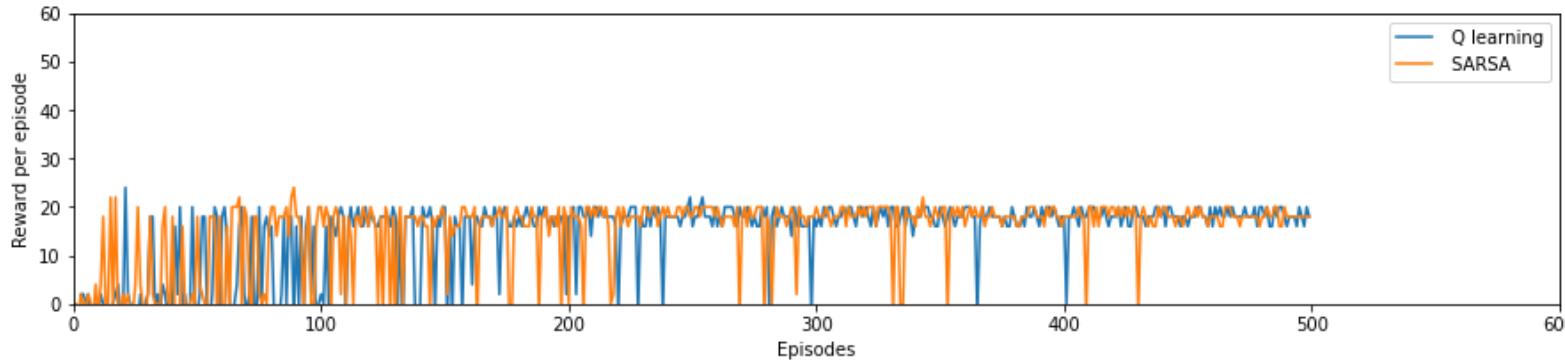
140 Discussion: Stochastic training takes more time both in SARSA as well as Q
141 learning compared to deterministic environment as the action determined
through the policy may or may not be executed.

142

143

144

- 5.2) Compare the performance of both algorithms on the same deterministic environment (e.g. show one graph with two reward dynamics) and give your interpretation of the results.



145

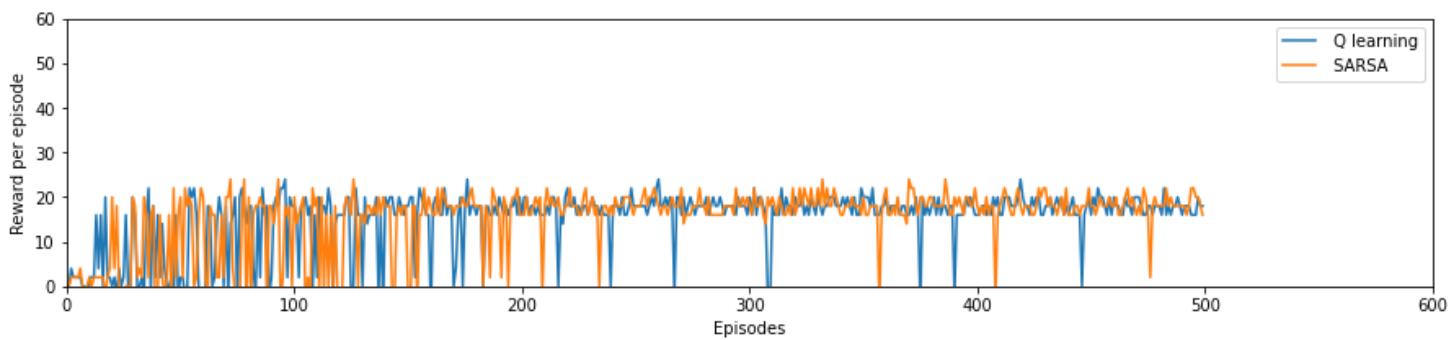
Figure 8: Reward per episode vs Episodes Deterministic env

146
147

Discussion: It looks like Q learning is learning faster in this scenario and combination of parameters a it is taking more reward per episode on average.

148
149
150

5.3) Compare how both algorithms perform in the same stochastic environment (e.g. show one graph with two reward dynamics) and give your interpretation of the results.



6.
7.

151

Figure 9: Reward per episode vs Episodes for Stochastic env

152
153
154

Discussion: It looks like SARSA is learning faster in this scenario and combination of parameters a it is taking more reward per episode on average than Q learning.

155

156
157

5.4) Briefly explain the tabular methods, including Q-learning, that were used to solve the problems. Provide their update functions and key features.

- 158 5.4)1. Q -learning – is a model free algorithm which update the value function based on
 159 Bellman equation
 160 5.4)2. Update Functions: To calculate the maximum expected future rewards for action at
 161 each state , we update the q-value for each state, action pair by adding the learning
 162 rate times the target minus the predicted q value.

1. Q-Learning:

163
$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

- 164 5.4)3. Key Features: Q learning is off policy so it learns the optimum value policy
 165 independent of agents actions.
 166 5.5) Sarsa – It's a slight variation of Q learning in such a way that it is an on policy
 167 algorithm such that the agent learns the value function according to action derived
 168 from current policy.
 169 5.5)1. Update Functions: To calculate the maximum expected future rewards for action at
 170 each state , we update the q-value for each state, action pair by adding the learning
 171 rate times the target minus the predicted q value.

2. SARSA:

172
$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha(r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$

- 173 5.5)2. Key Features:SARSA stands for STATE ACTION REWARD STATE ACTION, which is
 174 literally how the q value is update which is also seen from the equation where the
 175 present state, value pair is update based on the expected state action pair obtained
 176 from the current policy.

177 **5 Hyperparameter Tuning**

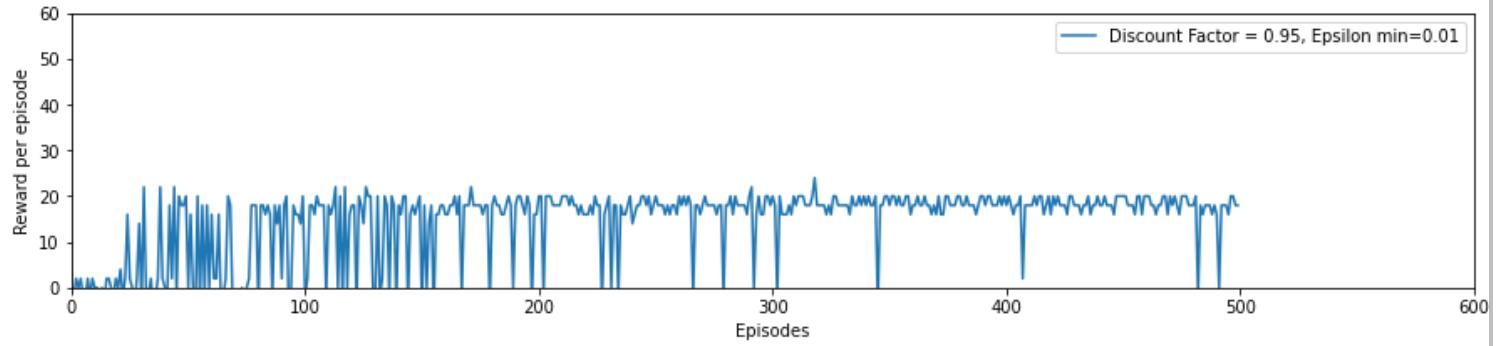
178 a) Discount factor [0.95, 0.9, 0.8]

179 c)Epsilon min values [0.01, 0.1, 0.2]

180 Lets compare the 9 combinations of these hyperparamters:

181 All the results are for Episodes = 500, steps =10

182 1. Discount factor 0.95 , Epsilon min values 0.01



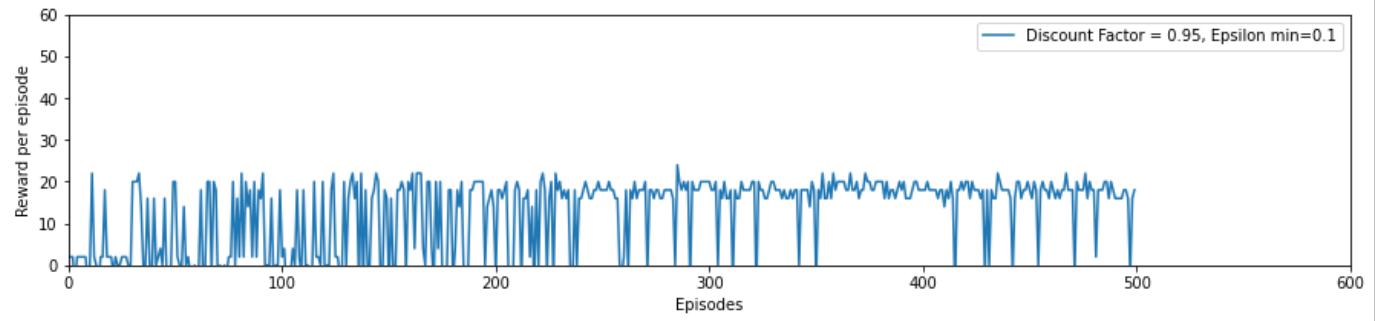
183

Figure 1: Reward per episode vs Episodes

184

185

2. Discount factor 0.95, Epsilon min values 0.1



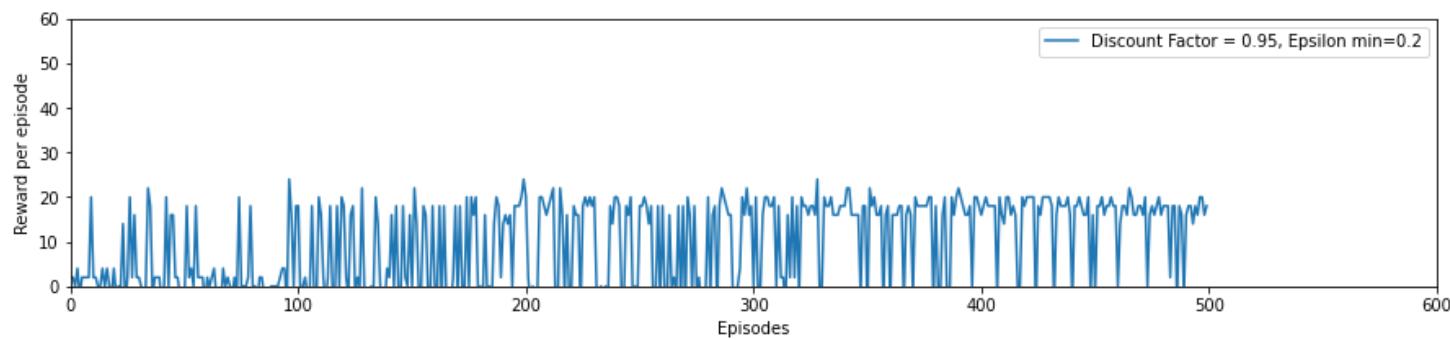
186

Figure 2: Reward per episode vs Episodes

187

188

3. Discount factor 0.95 Epsilon min values 0.2



189

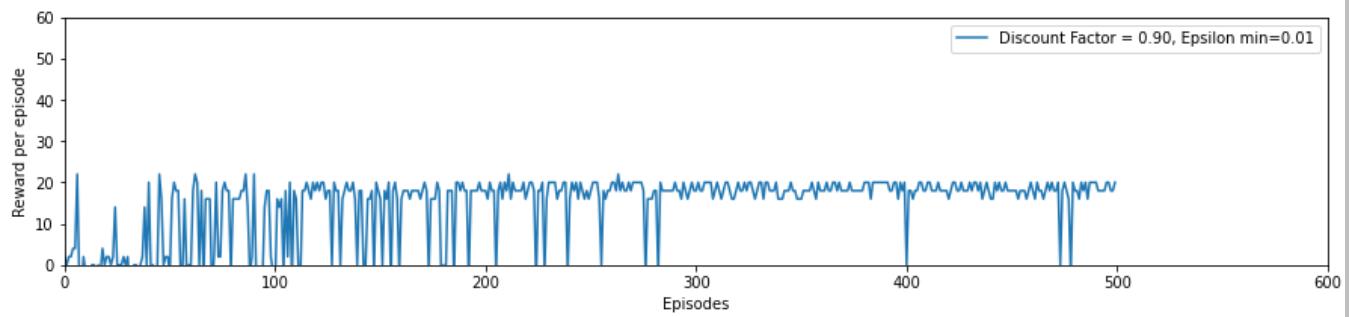
Figure 3: Reward per episode vs Episodes

190

Discussion: Comparing (1), (2), (3) , the best combination we have is
191 Discount factor = 0.95 and min Epsilon value = 0.01 as it clearly is learning much
192 faster to have the maximum reward per episode

193

4. Discount factor 0.9 , Epsilon min values 0.01

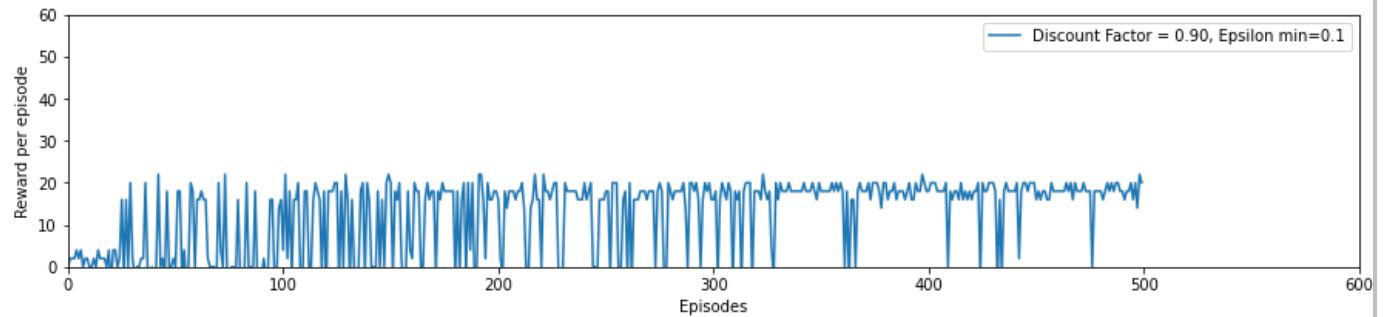


194

Figure 4: Reward per episode vs Episodes

195

196 5. Discount factor 0.9, Epsilon min values 0.1

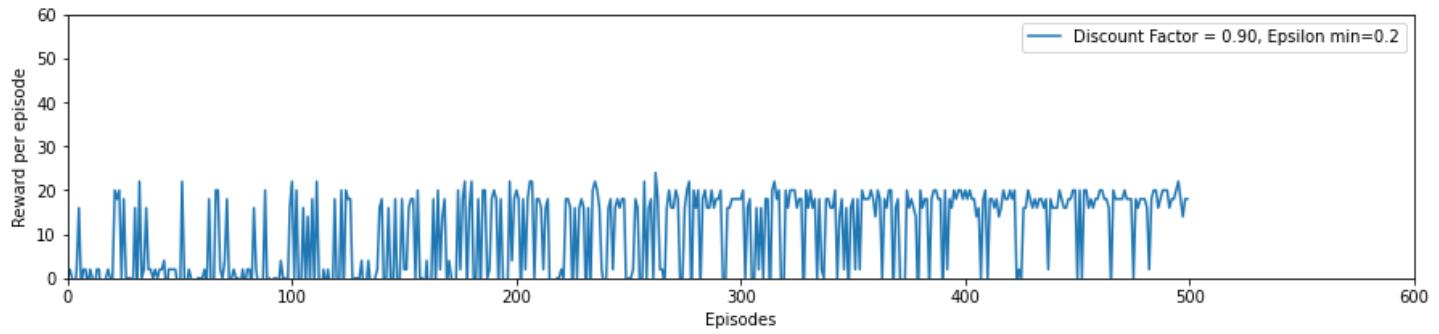


197

Figure 5: Reward per episode vs Episodes

198

6. Discount factor 0.9 Epsilon min values 0.2



199

Figure 6: Reward per episode vs Episodes

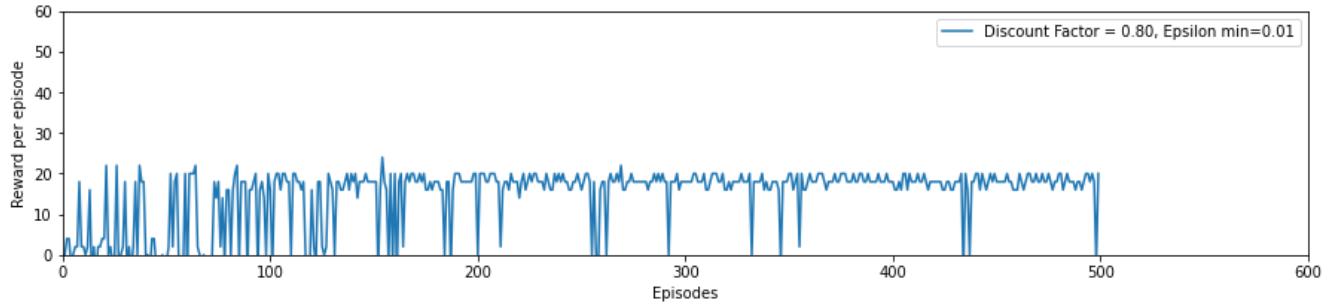
200

Discussion: Comparing (4), (5), (6) , the best combination we have is Discount factor = 0.90 and min Epsilon value = 0.01 as it clearly is learning much faster to have the maximum reward per episode

203

204

7. Discount factor 0.9 , Epsilon min values 0.01



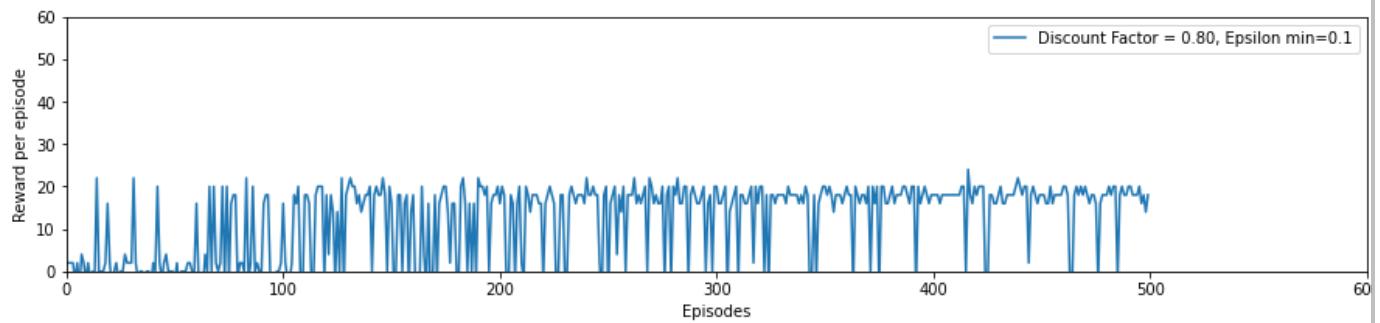
205

Figure 7: Reward per episode vs Episodes

206

207

8. Discount factor 0.9, Epsilon min values 0.1



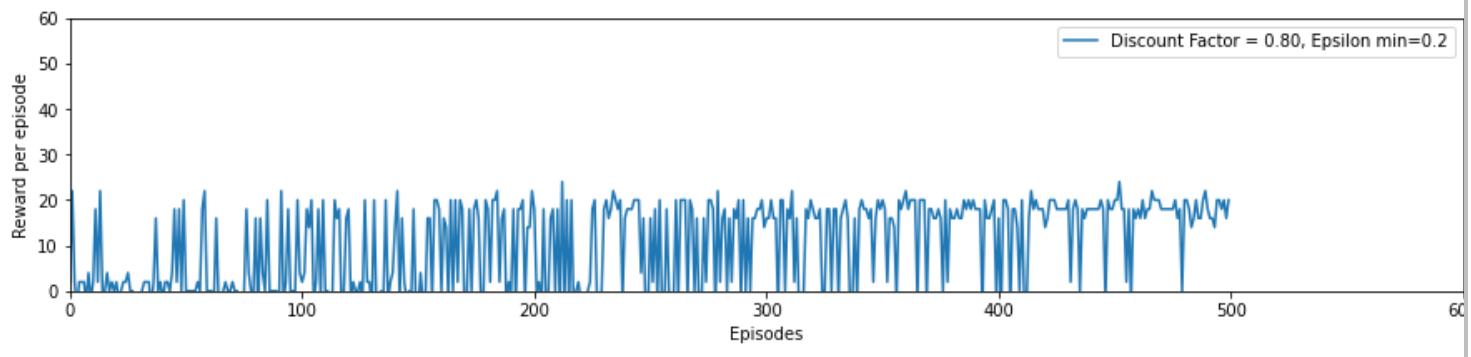
208

Figure 8: Reward per episode vs Episodes

209

210

9. Discount factor 0.9 Epsilon min values 0.2



211

Figure 9: Reward per episode vs Episodes

212 Discussion: Comparing (7), (8), (9) , the best combination we have is Discount
213 factor = 0.80 and min Epsilon value = 0.01 as it clearly is learning much faster to
214 have the maximum reward per episode

215 Conclusion: The best combination is Discount factor = 0.95 and min Epsilon
216 value = 0.01 which clearly shows a very fast learning of obtaining maximum
217 reward per episode aming all the 9 combinations.

218 **References**

219 <https://gym.openai.com/>