
CSE 574: Introduction to Machine Learning

Project 4: Tom and Jerry in Reinforcement Learning

Anunay Rao

anunayra@buffalo.edu

1 Introduction

The project is to apply the deep reinforcement learning algorithm to make the agent learn to navigate in the grid-world environment. Here, the Tom is an agent and Jerry is the goal. The main task for the agent is to find the shortest path to the goal.

2 Implementation

2.1 Build a 3-layer neural network using Keras library

```
model = Sequential()

### START CODE HERE ### (= 3 lines of code)
model.add(Dense(128, input_shape=(self.state_dim,), activation='relu'))
model.add(Dense(self.action_dim, input_shape=(128,), activation='relu'))
model.add(Dense(self.action_dim))
### END CODE HERE ###
```

3-layer Neural network is created to train the agent as we are using deep reinforcement learning algorithm which implies applying deep learning to Reinforcement learning.

2.2 Implement exponential decay formula for epsilon

$$\epsilon = \epsilon_{min} + (\epsilon_{max} - \epsilon_{min}) * e^{-\lambda|S|}$$

where $\epsilon_{min}, \epsilon_{max} \in [0, 1]$

λ is hyperparameter for epsilon

$|S|$ is total number of steps

```
### START CODE HERE ### (= 1 line of code)
self.epsilon = self.min_epsilon + (self.max_epsilon - self.min_epsilon)*np.exp(-1*self.lamb*np.absolute(self.steps))
### END CODE HERE ###
```

Here ϵ is the exploration rate which makes the agent to take random action. Higher the value of ϵ more the percentage chance that agent will select random action. Initially it is better for the agent to take random action so that it can try all kinds of possibilities before finding any pattern. When the agent does not take the action randomly it selects the action based on the current state and the action that

will maximize the reward. During the process we want to reduce the random actions of the agent as it will tend to find the optimal policy after exploring so we need exponentially decaying epsilon. When epsilon eventually decays to 0 the policy becomes greedy i.e it do not take random actions anymore. Higher the values of λ smaller is the exploration rate.

2.3 Implementing Q-function

$$Q_t = \begin{cases} r_t, & \text{if episode terminates at step } t + 1 \\ r_t + \gamma * \max_a Q(s_{t+1}, a), & \text{otherwise} \end{cases}$$

```
### START CODE HERE ### (~ 4 line of code)
if st_next is None:
    t[act] = rew
else:
    t[act] = rew + self.gamma*np.amax(t[act])
### END CODE HERE ###
```

Q-Learning is a value-based reinforcement learning algorithm which is used to find the optimal action-selection policy using a Q function. The goal is to maximize the value function Q. Q(state, action) returns the expected future reward of that action at that state. Here, γ is a discount factor which is multiplied by future rewards discovered by the agent to make them of less worth than the immediate rewards.

3 Results:

Configuration:

$\epsilon_{min} = 0.05$

$\epsilon_{max} = 1$

$\gamma = 0.99$

Number of episodes = 10000

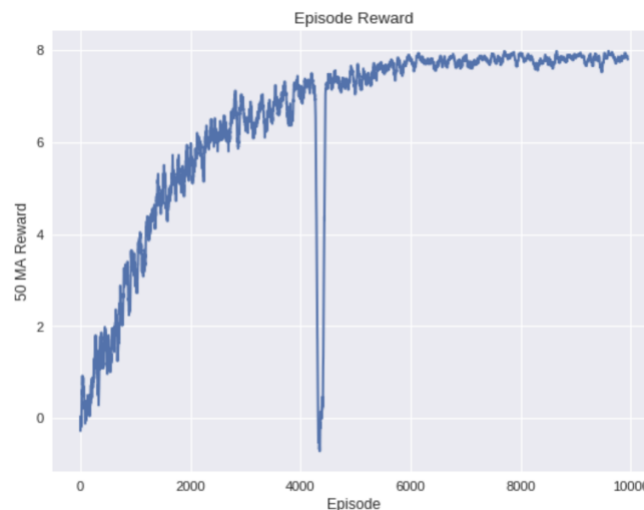


Figure 1: Variation of Rewards with number of Episodes

63 For the above configuration, it seen quite evident from the graph that the agent
 64 converge to the optimal policy in about closer to 6000 episodes. It took around
 65 760 seconds to run for 10000 episodes.

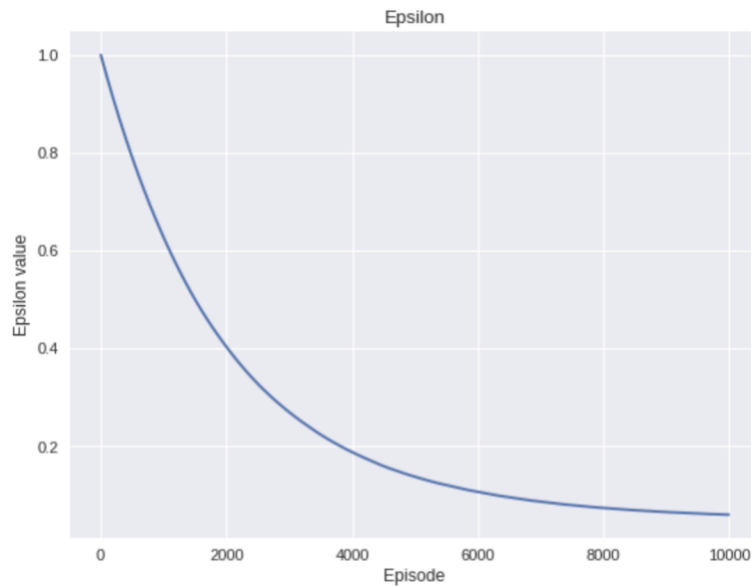


Figure 2: Epsilon is decaying exponentially with number of episodes

4 Hyperparameter Tuning:

4.1 Variaton of total mean-reward with different number of episodes

Configuration:

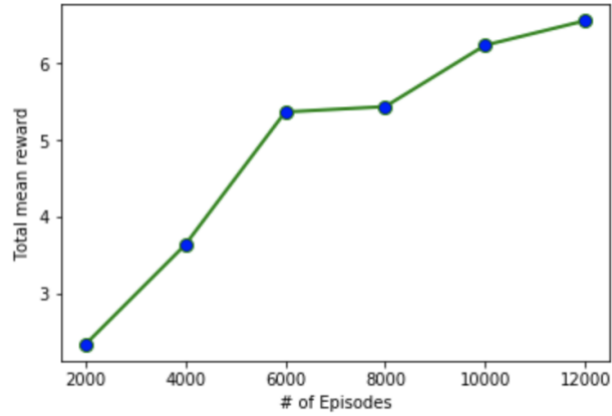
$$\epsilon_{min} = 0.05$$

$$\epsilon_{max} = 1$$

$$\gamma = 0.99$$

Number of episodes = 10000

# of Episodes	2000	4000	6000	8000	10000	12000
Total Mean Reward	2.34	3.64	5.37	5.44	6.24	0.21



80

81 4.2 Variation of total mean epsilon value with different 82 number of episodes

83 Configuration:

84 $\epsilon_{min} = 0.05$

85 $\epsilon_{max} = 1$

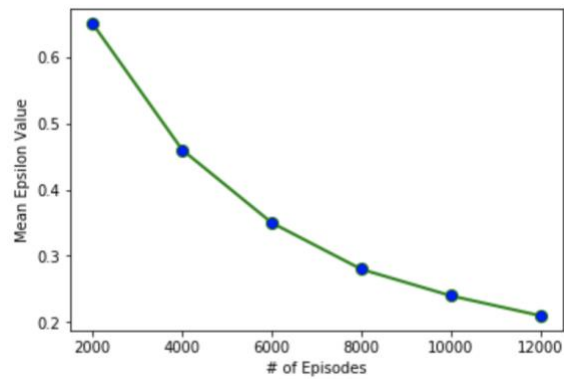
86 $\gamma = 0.99$

87 Number of episodes = 10000

88

# of Episodes	2000	4000	6000	8000	10000	12000
Total Mean Epsilon	0.65	0.46	0.35	0.28	0.24	0.21

89



90

91 4.3 Variaton of total mean-reward with ϵ_{max}

92 Configuration:

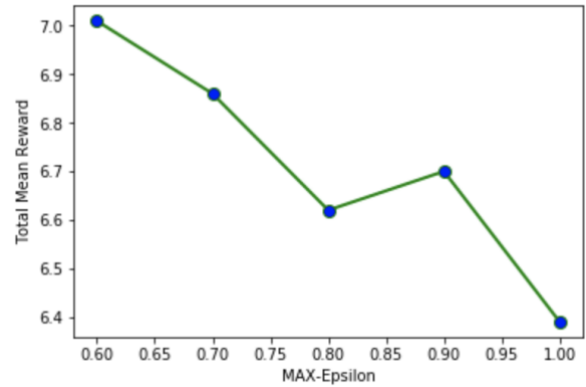
93 $\epsilon_{min} = 0.05$

94 $\gamma = 0.99$

95 Number of episodes = 10000

96

ϵ_{max}	0.6	0.7	0.8	0.9	1
Total Mean reward	7.01	6.86	6.62	6.70	6.39



4.4 Variaton of total mean-reward with ϵ_{min}

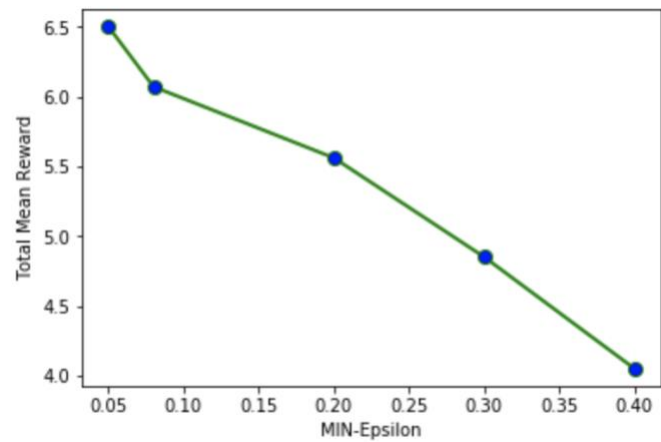
Configuration:

$\epsilon_{max} = 1$

$\gamma = 0.99$

Number of episodes = 10000

ϵ_{min}	0.05	0.08	0.2	0.3	0.4
Total Mean reward	6.50	6.07	5.56	4.85	4.05



114 5 Writing Tasks

115 **5.1 Explain what happens in reinforcement learning if the agent always**
 116 **chooses the action that maximizes the Q-value. Suggest two ways to force the**
 117 **agent to explore?**

119 If the agent always chooses the action that maximizes the Q-value and if it has not
 120 found the optimal policy to select which actions to perform then the agent will
 121 perform non-optimal actions repeatedly. But if the agent has not found the optimal
 122 policy then exploring actions which do not have the highest Q-value might
 123 allow it to find a better policy.

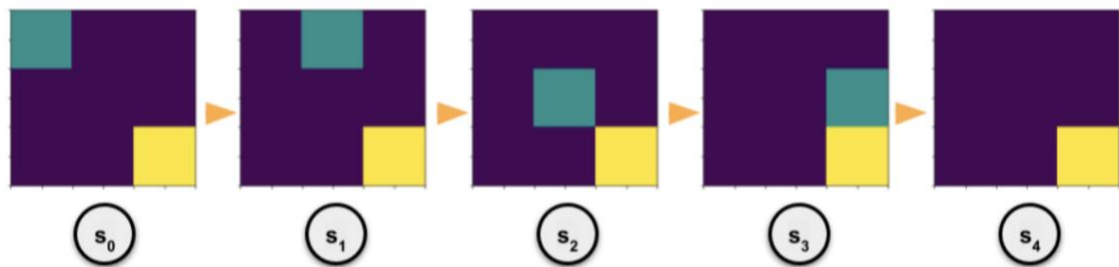
124 The two ways to force the agent to explore are:

- 126 • Allow it to pick random actions occasionally.
- 127 • Set the initial Q-values high so that the unexplored region good look to explore.

128 **5.2 Calculate Q-value for the given states and provide all the calculation**
 129 **steps.**

130

131



132

133 For State 4, it is a terminal state which means that agent has reached the goal so
 134 Q-value will be zero for any action corresponding to state 4.

135

STATE	ACTION			
	UP	DOWN	LEFT	RIGHT
0				
1				
2				
3				
4	0	0	0	0

136

137 For State 3,

138 $Q(S3, \text{Down}) = 1 + \gamma * \max_a Q(s_{t+1}, a) = 1 +$

139 $0.99 * \max_a (Q(S4, \text{UP}), Q(S4, \text{DOWN}), Q(S4, \text{LEFT}), Q(S4, \text{RIGHT}))$

140 $= 1 + 0.99 * 0 = 1$

141

142

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0				
1				
2				
3		1		
4	0	0	0	0

143

144 $Q(S3, \text{RIGHT}) = 0 + 0.99 * \max_a Q(S3, a) = 0 + 0.99 * 1 = 0.99$

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0				
1				
2				
3		1		0.99
4	0	0	0	0

145

146 $Q(S3, \text{LEFT}) = Q(S3, \text{UP})$ as these actions are symmetric.

147 $Q(S3, \text{LEFT}) = -1 + 0.99 * \max_a Q(S2, a)$

148

149 For State 2,

150 $Q(S2, \text{RIGHT}) = 1 + 0.99 * \max_a Q(S3, a) = 1 + 0.99 * 1 = 1.99$

151 $Q(S2, \text{DOWN}) = Q(S2, \text{RIGHT})$ as they are symmetric actions.

152 $Q(S2, \text{UP}) = -1 + 0.99 * \max_a Q(S1, a)$

153 $Q(S2, \text{LEFT}) = Q(S2, \text{UP})$ as these are symmetric actions.

154

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0				
1				
2		1.99		1.99
3		1		0.99
4	0	0	0	0

155

156

157

158 For State 1,

159 $Q(S1, \text{DOWN}) = 1 + 0.99 * \max_a Q(S2, a) = 1 + 0.99 * 1.99 = 2.97$

160 $Q(S1, \text{RIGHT}) = Q(S1, \text{DOWN})$ as these are symmetric actions

161 $Q(S1, \text{UP}) = 0 + 0.99 * 2.97 = 2.94$

162 $Q(S1, \text{LEFT}) = -1 + 0.99 * \max_a Q(S1, a)$

163

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0				
1	2.94	2.97		2.97
2		1.99		1.99
3		1		0.99
4	0	0	0	0

164

165 For State 0,

166 $Q(S0, \text{RIGHT}) = 1 + 0.99 * \max_a Q(S1, a) = 1 + 0.99 * 2.97 = 3.94$

167 $Q(S0, \text{DOWN}) = Q(S0, \text{RIGHT})$ as these are symmetric actions

168 $Q(S0, \text{UP}) = 0 + 0.99 * \max_a Q(S0, a) = 0.99 * 3.94 = 3.90$

169 $Q(S0, \text{LEFT}) = Q(S0, \text{UP})$ as these are symmetric actions.

170

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0	3.90	3.94	3.90	3.94
1	2.94	2.97		2.97
2		1.99		1.99
3		1		0.99
4	0	0	0	0

171

172 $Q(S1, \text{LEFT}) = -1 + 0.99 * \max_a Q(S1, a) = -1 + 0.99 * 3.94 = 2.90$

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0	3.90	3.94	3.90	3.94
1	2.94	2.97	2.90	2.97
2		1.99		1.99
3		1		0.99
4	0	0	0	0

173

174 $Q(S2, \text{UP}) = -1 + 0.99 * \max_a Q(S1, a) = -1 + 0.99 * 2.97 = 1.94$

175

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0	3.90	3.94	3.90	3.94
1	2.94	2.97	2.90	2.97
2	1.94	1.99	1.94	1.99
3		1		0.99
4	0	0	0	0

176

177 $Q(S3, \text{LEFT}) = -1 + 0.99 * \max_a Q(S2, a) = -1 + 0.99 * 1.99 = 0.97$
178

STATE	<u>ACTION</u>			
	UP	DOWN	LEFT	RIGHT
0	3.90	3.94	3.90	3.94
1	2.94	2.97	2.90	2.97
2	1.94	1.99	1.94	1.99
3	0.97	1	0.97	0.99
4	0	0	0	0

179
180
181
182