
Reinforcement Learning

Shageenth Sandrakumar

Department of Data Science and Analytics

Buffalo, NY 14226

shageent@buffalo.edu

Abstract

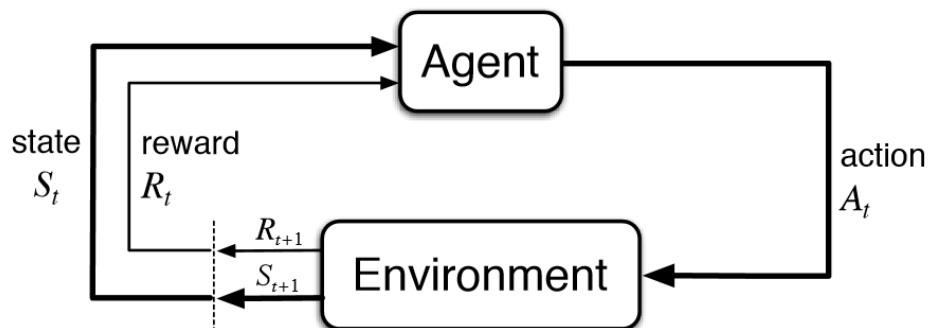
In this lab, we are required to work with Reinforcement Learning, a newer Machine Learning technique, that can train an agent in an environment. The agent will navigate the classic 4x4 grid-world environment to a specific goal. The agent will learn an optimal policy through Q-Learning which will allow it to take actions to reach the goal while avoiding the boundaries. We use a platform here called AI gym to facilitate the whole process of the construction of the agent and the environment.

Introduction

What is Reinforcement Learning?

Reinforcement Learning(RL) is one of the hottest research topics in the field of modern Artificial Intelligence and its popularity is only growing. Reinforcement Learning(RL) is a type of machine learning technique that enables an agent to learn in an interactive environment by trial and error using feedback from its own actions and experiences. Though both supervised and reinforcement learning use mapping between input and output, unlike supervised learning where the feedback provided to the agent is correct set of actions for performing a task, reinforcement learning uses rewards and punishments as signals for positive and negative behavior.

As compared to unsupervised learning, reinforcement learning is different in terms of goals. While the goal in unsupervised learning is to find similarities and differences between data points, in the case of reinforcement learning the goal is to find a suitable action model that would maximize the total cumulative reward of the agent. The figure below illustrates the action-reward feedback loop of a generic RL model.



49 **2. Open AI Gym**
50 Gym is a toolkit for developing and comparing reinforcement learning algorithms. It
51 makes no assumptions about the structure of your agent, and is compatible with any
52 numerical computation library, such as TensorFlow or Theano. The `gym` library is a
53 collection of test problems — environments — that you can use to work out your
54 reinforcement learning algorithms. These environments have a shared interface,
55 allowing you to write general algorithms.

56 **3. The Environment**
57 The Environment is a 4 by 4 grid environment described by two
58 things the grid environment and the agent. An observation space
59 which is defined as a vector of elements. This can be particularly
60 useful for environments which return measurements, such as in
61 robotic environments.
62 The core gym interface is `env`, which is the unified environment interface. The
63 following are the `env` methods that would be quite helpful to us:

64 `env.reset`: Resets the environment and returns a random initial state.

65 `env.step(action)`: Step the environment by one timestep.

66 Returns observation: Observations of the environment

67 reward: If your action was beneficial or not

68 done: Indicates if we have successfully picked up and dropped off a passenger, also
69 called one *episode*

70 info: Additional info such as performance and latency for debugging purposes

71 `env.render`: Renders one frame of the environment (helpful in visualizing the
72 environment)

73 **4. We have an Action Space of size 4**

74 0 = down

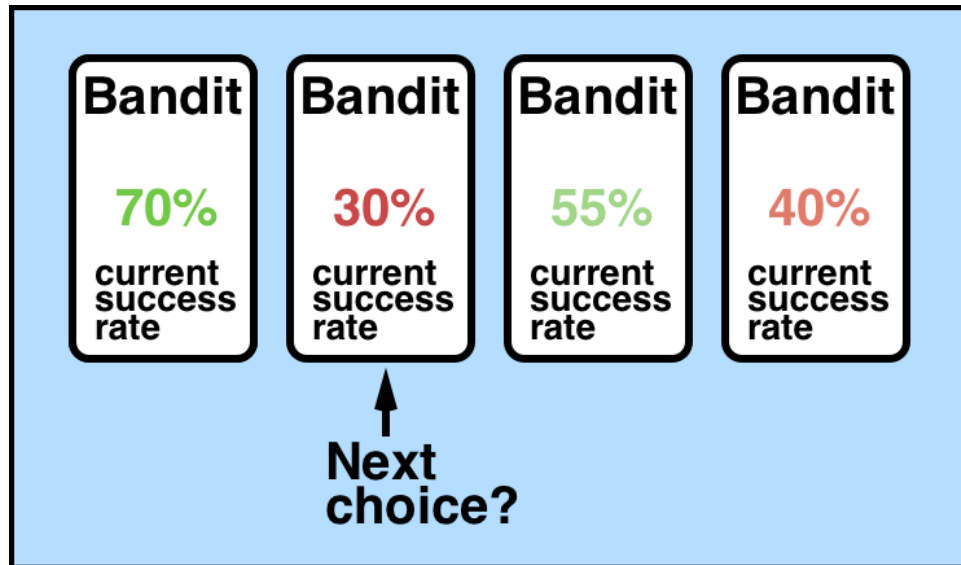
75 1 = up

76 2 = right

77 3 = left

78 **5. The MULti-Bandit Problem (E-Greedy Algorithm)**

79 The multi-armed bandit problem is a classic reinforcement learning example where
80 we are given a slot machine with n arms (bandits) with each arm having its own
81 rigged probability distribution of success. Pulling any one of the arms gives you a
82 stochastic reward of either $R=+1$ for success, or $R=0$ for failure. Our objective is to
83 pull the arms one-by-one in sequence such that we maximize our total reward
84 collected in the long run.
85



87

88

89 **The non-triviality of the multi-armed bandit problem lies in**
90 **the fact that we (the agent) cannot access the true bandit**
91 **probability distributions — all learning is carried out via the**
92 **means of trial-and-error and value estimation. So the question**
93 **is:**

94 This is our goal for the multi-armed bandit problem, and having such a strategy
95 would prove very useful in many real-world situations where one would like to select
96 the “best” bandit out of a group of bandits.

97 In this project, we approach the multi-armed bandit problem with a classical
98 reinforcement learning technique of an *epsilon-greedy agent* with a learning framework
99 of *reward-average sampling* to compute the action-value $Q(a)$ to help the agent improve
100 its future action decisions for long-term reward maximization.

101 In a nutshell, the epsilon-greedy agent is a hybrid of a (1) completely-exploratory agent
102 and a (2) completely-greedy agent. In the multi-armed bandit problem, a completely-
103 exploratory agent will sample all the bandits at a uniform rate and acquire knowledge
104 about every bandit over time; the caveat of such an agent is that this knowledge is never
105 utilized to help itself to make better future decisions! On the other extreme, a completely-
106 greedy agent will choose a bandit and stick with its choice for the rest of eternity; it will
107 not make an effort to try out other bandits in the system to see whether they have better
108 success rates to help it maximize its long-term rewards, thus it is very narrow-minded!

109 How do we perform this in our code?

110 We perform this by assigning a variable called epsilon. This epsilon switches between the
111 exploratory and the greedy agent. We choose a random number between 0 and 1, if this
112 number is less than epsilon, we tell the agent to explore if it is greater then we tell the
113 agent to be greedy. This tactic is used in our policy part of our code.

114Q-Learning Algorithm

115

116 Essentially, Q-learning lets the agent use the environment's rewards
117 to learn, over time, the best action to take in a given state.

118 In our environment, we have the reward table, that the agent will learn from. It does
119 something by looking at receiving a reward for taking an action in the current state, then
120 updating a *Q-value* to remember if that action was beneficial.

121 The values stored in the Q-table are called a *Q-values*, and they map to a (state, action)
122 combination.

123 A Q-value for a particular state-action combination is representative of the "quality" of
124 an action taken from that state. Better Q-values imply better chances of getting greater
125 rewards.

126

127 Q-values are initialized to an arbitrary value, and as the agent explores itself to the
128 environment and receives different rewards by executing different actions, the Q-
129 values are updated using the equation:

130
$$Q(\text{state}, \text{action}) \leftarrow (1 - \alpha)Q(\text{state}, \text{action}) + \alpha(\text{reward} + \gamma \max_a Q(\text{nextstate}, \text{allactions}))$$

131

132 Where:

133 - α (alpha) is the learning rate ($0 < \alpha \leq 1$) - Just like in supervised learning settings, α is
134 the extent to which our Q-values are being updated in every iteration.

135 - γ (gamma) is the discount factor ($0 \leq \gamma \leq 1$) - determines how much importance we
136 want to give to future rewards. A high value for the discount factor (close to 1)
137 captures the long-term effective reward, whereas, a discount factor of 0 makes our
138 agent consider only immediate reward, hence making it greedy.

139 What is this saying?

140 We are assigning (\leftarrow), or updating, the Q-value of the agent's current *state* and *action*
141 by first taking a weight $(1 - \alpha)$ of the old Q-value, then adding the learned value. The
142 learned value is a combination of the reward for taking the current action in the current
143 state, and the discounted maximum reward from the next state we will be in once we
144 take the current action.

145 Basically, we are learning the proper action to take in the current state by looking at
146 the reward for the current state/action combo, and the max rewards for the next state.
147 This will eventually cause our taxi to consider the route with the best rewards strung
148 together.

149 The Q-value of a state-action pair is the sum of the instant reward and the discounted
150 future reward (of the resulting state). The way we store the Q-values for each state and
151 action is through a Q-table

152 Q-Table

153 The Q-table is a matrix where we have a row for every state and a column for every
154 action. It's first initialized to 0, and then values are updated after training.

155

156 6. Epsilon

157 We want the odds of the agent exploring to decrease as time goes
158 on. One way to do this is by updated the epsilon every moment of
159 the agent during the training phase. Choosing a decay rate was
160 the difficult part. We want epsilon not to decay too soon, so the
161 agent can have time to explore. The way I went about this is I
162 want the agent the ability to explore the region for at least the
163 size of the grid movements. So in 25 movements the agent should
164 still have a good chance of being in exploratory phase. I chose
165 decay = $1/1.01$ because $(1/1.01)^{25} = 75\%$ which is a good chance
166 and still being in the exploratory phase within 0-25 moves. After
167 50 moves epsilon goes to 0.6 which is good odds of agent doing
168 both exploring and action orientated.

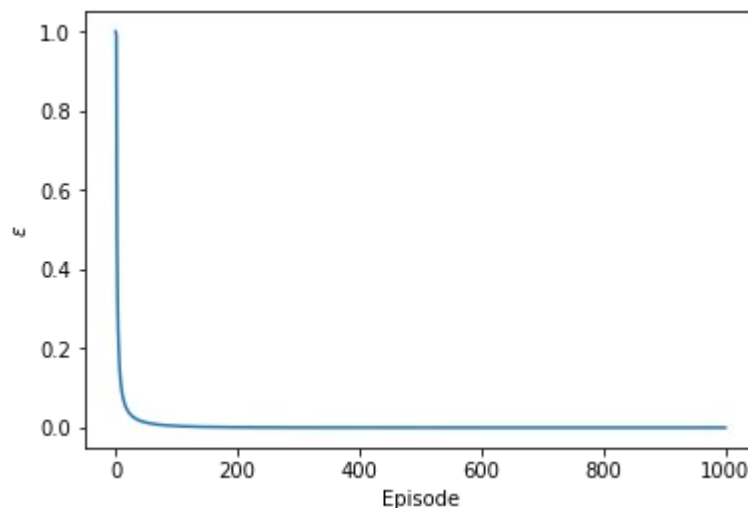
169

170 7. Results and Charts

171 Here we plot epsilon vs episode.

172 We see that we have a exponential decay going on here.

173



174
175
176

177 Then we plot rewards vs episode

178

179

180

181

182

183

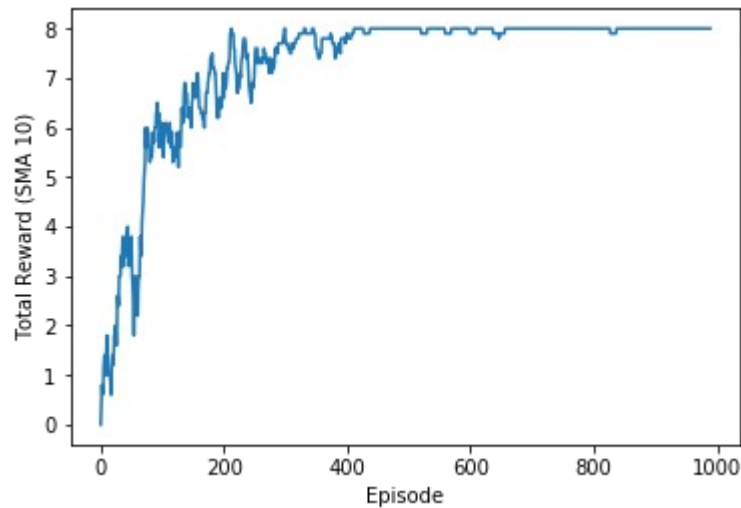
184

185

186

187

188



189 Here we see that the total rewards slowly goes up then plateaus after hitting 8. This is
190 because 8 is the optimal path for our algrotihm.

191

192 CONCLUSION

193 Our Reinforcement Algorithm had done fairly well. Our agent has been able to learn
194 the material in an effective manner. It has been able to reach its goal in 8 steps.

195

196

197

198

199

200

201

202

203

204

205

206