



PROJECT

Train a Smartcab to Drive

A part of the Machine Learning Engineer Nanodegree Program

PROJECT REVIEW

CODE REVIEW 6

NOTES

▼ agent.py 6

```

1 import random
2 import math
3 from environment import Agent, Environment
4 from planner import RoutePlanner
5 from simulator import Simulator
6
7 class LearningAgent(Agent):
8     """ An agent that learns to drive in the Smartcab world.
9         This is the object you will be modifying. """
10
11     def __init__(self, env, learning=False, epsilon=1.0, alpha=0.5):
12         super(LearningAgent, self).__init__(env) # Set the agent in the environment
13         self.planner = RoutePlanner(self.env, self) # Create a route planner
14         self.valid_actions = self.env.valid_actions # The set of valid actions
15
16         # Set parameters of the learning agent
17         self.learning = learning # Whether the agent is expected to learn
18         self.Q = dict() # Create a Q-table which will be a dictionary of tuples
19         self.epsilon = epsilon # Random exploration factor
20         self.alpha = alpha # Learning factor
21
22         #####
23         ## TO DO ##
24         #####
25         # Set any additional class parameters as needed
26         self.iter = 1
27         self.a = 0.05
28         self.b = 0.05
29
30
31     def reset(self, destination=None, testing=False):
32         """ The reset function is called at the beginning of each trial.
33             'testing' is set to True if testing trials are being used
34             once training trials have completed. """
35
36         # Select the destination as the new location to route to
37         self.planner.route_to(destination)
38
39         #####
40         ## TO DO ##
41         #####
42         # Update epsilon using a decay function of your choice
43         # Update additional class parameters as needed
44         # If 'testing' is True, set epsilon and alpha to 0
45         if testing:
46             self.epsilon = 0
47             self.alpha = 0
48         else:
49             self.epsilon -= 0.05
50             self.epsilon = math.cos(self.a*self.iter)
51             if self.alpha > 0.2:
52                 self.alpha = max( math.cos(self.b*self.iter), 0.2)
53             self.iter += 1
54
55         return None
56
57     def build_state(self):
58         """ The build_state function is called when the agent requests data from the

```

```

59     environment. The next waypoint, the intersection inputs, and the deadline
60     are all features available to the agent. """
61
62     # Collect data about the environment
63     waypoint = self.planner.next_waypoint() # The next waypoint
64     inputs = self.env.sense(self)           # Visual input - intersection light and traffic
65     deadline = self.env.get_deadline(self)   # Remaining deadline
66
67     is_red_light = True if inputs['light'] == 'red' else False
68     #left = inputs['left']
69     #right = inputs['right']
70     #oncoming = inputs['oncoming']
71     can_go_left = True if inputs['oncoming'] != 'forward' and inputs['oncoming'] != 'right' else False

```

SUGGESTION

We really should NOT be self creating states here. As the agent should be able to learn your self created states based solely on its process.

(<https://discussions.udacity.com/t/number-and-identity-of-states/186637>)

The goal of reinforcement learning is to learn these rules without them being hard-coded into the algorithm. Try to include these states directly and see that Q-matrix will converge such as your agent will be able to learn these rules on its own. But I will let this slide, since these are actually the correct way way to do this. But in the future, please don't do this in reinforcement learning.

```

72     can_go_right = True if inputs['left'] != 'forward' else False
73
74     #####
75     ## TO DO ##
76     #####
77     # Set 'state' as a tuple of relevant data for the agent
78     state = (waypoint, is_red_light, can_go_left, can_go_right) # left, right, oncoming
79
80     return state
81
82
83 def get_maxQ(self, state):
84     """ The get_max_Q function is called when the agent is asked to find the
85         maximum Q-value of all actions based on the 'state' the smartcab is in. """
86
87     #####
88     ## TO DO ##
89     #####
90     # Calculate the maximum Q-value of all actions for a given state
91
92     actions = self.Q[state]
93
94     maxQ = None
95     max = float('-inf')
96
97     for act in actions.keys():
98         if actions[act] > max:
99             maxQ = act
100             max = actions[act]
101
102     return maxQ

```

REQUIRED

You should be returning the maximum Q-value for the state, not necessarily the action associated for it. As your function right now return something like `right`, as we should be returning a number such as `0.837`. The reason for this is in the `choose_action()` function.

```

103
104
105 def createQ(self, state):
106     """ The createQ function is called when a state is generated by the agent. """
107
108     #####
109     ## TO DO ##
110     #####
111     # When learning, check if the 'state' is not in the Q-table

```

REQUIRED

Make sure you Q-table is only initialized when Learning == True

```

112     # If it is not, create a new dictionary for that state
113     # Then, for each action available, set the initial Q-value to 0.0
114     if self.learning:
115         if state not in self.Q:
116             self.Q[state] = dict()
117             for act in self.valid_actions:
118                 self.Q[state][act] = 0.0
119

```

```

120         return
121
122
123     def choose_action(self, state):
124         """ The choose_action function is called when the agent is asked to choose
125             which action to take, based on the 'state' the smartcab is in. """
126
127         # Set the agent state and default action
128         self.state = state
129         self.next_waypoint = self.planner.next_waypoint()
130         action = None
131
132         #####
133         ## TO DO ##
134         #####
135         # When not learning, choose a random action
136         if not self.learning:
137             action = self.valid_actions[random.randint(0, len(self.valid_actions)-1)]
138         # When learning, choose a random action with 'epsilon' probability
139         # Otherwise, choose an action with the highest Q-value for the current state
140         else:
141             if self.epsilon > random.random():
142                 action = random.choice(self.valid_actions)
143             else:
144                 action = self.get_maxQ(state)

```

REQUIRED

Your agent should be choosing a random action from a choice of actions that have the highest Q-value. For example, since all actions are initialized with a reward of zero, it's possible that all four actions are considered "optimal". Not having the agent choose a random action from this would imply that it always chooses, perhaps, the first available option. This is incorrect behavior:

```

STATE X:
-- 'forward' : 0.00
-- 'left'    : 0.00
-- 'right'   : -1.023
-- 'None'    : 0.00

```

The agent should choose one of 'forward', 'left', or 'None' with equal probability, since they are all considered optimal with the current learned policy. Using a list would be a good idea.

```

145         return action
146
147
148
149     def learn(self, state, action, reward):
150         """ The learn function is called after the agent completes an action and
151             receives an award. This function does not consider future rewards
152             when conducting learning. """
153
154         #####
155         ## TO DO ##
156         #####
157         # When learning, implement the value iteration update rule

```

REQUIRED

Make sure you Q-Learning algorithm is only updated when Learning == True

```

158         # Use only the learning rate 'alpha' (do not use the discount factor 'gamma')
159         self.Q[state][action] = (self.Q[state][action] * (1 - self.alpha)) + (reward * self.alpha)

```

AWESOME

Good work with your Bellman equation!

```

160         return
161
162
163
164     def update(self):
165         """ The update function is called when a time step is completed in the
166             environment for a given trial. This function will build the agent
167             state, choose an action, receive a reward, and learn if enabled. """
168
169         state = self.build_state()           # Get current state
170         self.createQ(state)                  # Create 'state' in Q-table
171         action = self.choose_action(state)   # Choose an action
172         reward = self.env.act(self, action) # Receive a reward
173         self.learn(state, action, reward)    # Q-learn
174
175         return
176
177

```

```

178 def run():
179     """ Driving function for running the simulation.
180         Press ESC to close the simulation, or [SPACE] to pause the simulation. """
181
182     #####
183     # Create the environment
184     # Flags:
185     #   verbose    - set to True to display additional output from the simulation
186     #   num_dummies - discrete number of dummy agents in the environment, default is 100
187     #   grid_size  - discrete number of intersections (columns, rows), default is (8, 6)
188     env = Environment()
189
190     #####
191     # Create the driving agent
192     # Flags:
193     #   learning    - set to True to force the driving agent to use Q-learning
194     #   * epsilon   - continuous value for the exploration factor, default is 1
195     #   * alpha     - continuous value for the learning rate, default is 0.5
196     agent = env.create_agent(LearningAgent, learning=True, alpha=1)
197
198     #####
199     # Follow the driving agent
200     # Flags:
201     #   enforce_deadline - set to True to enforce a deadline metric
202     env.set_primary_agent(agent, enforce_deadline=True)
203
204     #####
205     # Create the simulation
206     # Flags:
207     #   update_delay - continuous time (in seconds) between actions, default is 2.0 seconds
208     #   display      - set to False to disable the GUI if PyGame is enabled
209     #   log_metrics  - set to True to log trial and simulation results to /logs
210     #   optimized    - set to True to change the default log file name
211     sim = Simulator(env, update_delay=0.01, log_metrics=True, optimized=True)
212
213     #####
214     # Run the simulator
215     # Flags:
216     #   tolerance    - epsilon tolerance before beginning testing, default is 0.05
217     #   n_test       - discrete number of testing trials to perform, default is 0
218     sim.run(n_test=10)
219
220
221 if __name__ == '__main__':
222     run()
223

```

► report.html

► logs/sim_improved-learning.txt

► logs/sim_default-learning.txt

Learn the [best practices for revising and resubmitting your project](#).

RETURN TO PATH

Rate this review

[Student FAQ](#)