



MITx 6.86x

Machine Learning with Python-From Linear Models to Deep Learning[Course](#)[Progress](#)[Dates](#)[Discussion](#)[Resources](#)[Course](#) / [Unit 5. Reinforcement Learning \(2 weeks\)](#) / [Project 5: Text-Based Game](#)[< Previous](#)

3. Q-learning Algorithm

[Bookmark this page](#)

Project due May 10, 2023 08:59 -03 Completed

In this section, you will implement the Q-learning algorithm, which is a model-free algorithm for finding an optimal Q-function. In the tabular setting, the algorithm maintains the Q-value for all possible state-action pairs. Starting from a random Q-function, the agent continuously collects experiences and updates its Q-function.

From now on, we will refer to $\mathbf{c} = (a, b)$ as "an action" although it really is an action with parameters.

Q-learning Algorithm

- The agent plays an action \mathbf{c} at state \mathbf{s} , getting a reward $R(\mathbf{s}, \mathbf{c})$ and observing the next state \mathbf{s}' .
- Update the single Q-value corresponding to each such transition:

$$Q(\mathbf{s}, \mathbf{c}) \leftarrow (1 - \alpha) Q(\mathbf{s}, \mathbf{c}) + \alpha [R(\mathbf{s}, \mathbf{c}) + \gamma \max_{\mathbf{c}' \in C} Q(\mathbf{s}', \mathbf{c}')]]$$

Tip: We recommend you implement all functions from this tab and the next one before moving on to the next tab online. Make sure you achieve reasonable performance on the *Home World* game.

Single step update

1.0/1 point (graded)

Write a function `tabular_q_learning` that updates the single Q-value, given the transition $(\mathbf{s}, \mathbf{c}, R(\mathbf{s}, \mathbf{c}), \mathbf{s}')$.

Reminder: You should implement this function locally first. You can read through the next tab to see the context in which this function is called.

Available Functions: You have access to the NumPy python library as `np`. You should use `ALPHA` and `GAMMA` in your code.

```

1 def tabular_q_learning(q_func, current_state_1, current_state_2, action_index,
2                         object_index, reward, next_state_1, next_state_2,
3                         terminal):
4     """Update q_func for a given transition
5
6     Args:
7         q_func (np.ndarray): current Q-function
8         current_state_1, current_state_2 (int, int): two indices describing
9         action_index (int): index of the current action
10        object_index (int): index of the current object

```

Epsilon-greedy exploration

1.0/1 point (graded)

Note that the Q-learning algorithm does not specify how we should interact in the world. It merely updates the values based on the experience collected. If we explore randomly, we would most likely not get anywhere. A better option is to exploit what we have learned, as summarized by current Q-values. We can always act greedily with respect to the current Q-values, i.e., take an action that maximizes the Q-value. Of course, early on, these are not necessarily the best actions. For this reason, a typical exploration strategy is to follow a so-called ϵ -greedy policy. In this policy, we take a random action out of the set of actions with probability ϵ , and otherwise follow the greedy action. This policy balances exploration vs exploitation. A large value of ϵ means exploring more (randomly), while a small ϵ means exploiting what we have learned. A small ϵ , on the other hand, will generate experience consistent with the current estimates of Q-values.

Now you will write a function `epsilon_greedy` that implements the ϵ -greedy exploration policy using the current Q-function.

Reminder: You should implement this function locally first. You can read through the notebook to get the context in which this function is called.

Available Functions: You have access to the NumPy python library as `np`. Your code should use the constants `NUM_ACTIONS` and `NUM_OBJECTS`.

```
1 def epsilon_greedy(state_1, state_2, q_func, epsilon):
2     """Returns an action selected by an epsilon-Greedy exploration policy
3
4     Args:
5         state_1, state_2 (int, int): two indices describing the current state
6         q_func (np.ndarray): current Q-function
7         epsilon (float): the probability of choosing a random command
8
9     Returns:
10        (int, int): the indices describing the action/object to take
11    """
12    coin = np.random.random_sample()
13
14    if coin < epsilon:
15        action_index = np.random.randint(NUM_ACTIONS)
```

Press ESC then TAB or click outside of the code editor to exit

Correct

Test results

? Random index value of Epsilon-greedy exploration

I passed the first test case (epsilon==0) of Epsilon-greedy exploration and failed on the second test case(eps

💬 "if random.random() < epsilon" for calculating <epsilon probability does not work for epsilon-

Hi with same implementation for rest of code, but changing the if condition from what is mentioned in the titl

? Question on the "terminal" flag for Single-step update

I'm not clear what is meant to happen in the single step update when the terminal parameter is True. I figured

? The grader on the greedy_epsilon question picks a Q-value that is not the maximum

After a few tries and misses on the epsilon greedy, I decided to add some prints to see what the Q matrix loc

💬 Meaning of (s,c,R(s,c),s')

💬 q_func

Is q_func a 4-D matrix? If it is a 4-D matrix why is it called q_func instead of q_values?

? Q_func and states

I have two questions. 1. Why are there two indexes for current_state? I thought state is just a room index out

? Basic question on q_func input

A very basic (or even silly) question: I can see for q_func input: q_func.shape is (2, 5, 3, 7) I understand there

💬 Single step update : tabular_q_learning

For this function are we setting q_func[current_state_1, current_state_2, action_index, object_index] = 0 but i

💬 Single Step answers off by < 0.01

My answers for the single step are off by less than 0.01. I seem to have entered in the formula correctly. Rath

? Absolutely no idea what "terminal" means

Getting correct output for false terminal, can't get the correct output for true terminal. According to framewo

💬 Epsilon-greedy exploration

edX

[About](#)

[Affiliates](#)

[edX for Business](#)

[Open edX](#)

[Careers](#)

[News](#)

Legal

[Terms of Service & Honor Code](#)

[Privacy Policy](#)



© 2023 edX LLC. All rights reserved.

深圳市恒宇博科技有限公司 [粤ICP备17044299号-2](#)