EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the <a href="Privacy Policy">Privacy Policy</a>.





<u>Unit 5 Reinforcement Learning (2</u>

4. Tabular Q-learning for Home

Course > weeks)

> <u>Project 5: Text-Based Game</u> > World game

# 4. Tabular Q-learning for Home World game

Extension Note: Project 5 due date has been extended by 1 more day to September 6 23:59UTC.

In this section you will evaluate the tabular Q-learning algorithms for the *Home world* game. Recall that the state observable to the player is described in text. Therefore we have to choose a mechanism that maps text descriptions into vector representations.

In this section you will consider a simple approach that assigns a unique index for each text description. In particular, we will build two dictionaries:

- dict room desc that takes the room description text as the key and returns a unique scalar index
- dict quest desc that takes the quest description text as the key and returns a unique scalar index.

For instance, consider an observable state  $s=(s_r,s_q)$ , where  $s_r$  and  $s_q$  are the text descriptions for the current room and the current request, respectively. Then  $i_r=$  dict\_room\_desc[ $s_r$ ] gives the scalar index for  $s_r$  and  $i_q=$  dict\_quest\_desc[ $s_q$ ] gives the scalar index for  $s_q$ . That is, the textual state  $s=(s_r,s_q)$  is mapped to a tuple  $I=(i_r,i_q)$ .

Normally, we would build these dictionaries as we train our agent, collecting descriptions and adding them to the list of known descriptions. For the purpose of this project, these dictionaries will be provided to you.

## Evaluating Tabular Q-learning on Home World

1.0/1 point (graded)

The following python files are provided:

- framework.py contains various functions for the text-based game environment that the staff has implemented for you. Some functions that you can call to train and testing your reinforcement learning algorithms:
  - newGame()
    - Args: None
    - Return: A tuple where the first element is a description of the initial room, the second element is a description of the quest for this new game episode, and the last element is a Boolean variable with value *False* implying that the game is not over.
  - step\_game()
    - Args:
      - current\_room\_desc : An description of the current room

- current\_quest\_desc : A description of the current quest state
- action\_index : An integer used to represent the index of the selected action
- object\_index: An integer used to indicate the index of the selected object
- Return: the system next state when the selected command is applied at the current state.
  - next\_room\_desc : The description of the room of the next state
  - next\_quest\_desc : The description of the next quest
  - reward: A real valued number representing the **one-step** reward obtained at this step
  - terminal: A boolean valued number indicating whether this episode is over (either quest is finished, or the number of steps reaches the maximum number of steps for each episode).
- agent\_tabular\_QL.py contains various function templates that you will use to implement your learning algorithm.

In this section, you will evaluate your learning algorithm for the Home World game. The metric we use to measure an agent's performance is the cumulative discounted reward obtained per episode averaged over the episodes.

The evaluation procedure is as follows. Each experiment (or run) consists of multiple epochs (the number of epochs is NUM\_EPOCHS). In each epoch:

- 1. You first train the agent on NUM\_EPIS\_TRAIN episodes, following an  $\varepsilon$ -greedy policy with  $\varepsilon=$  TRAINING\_EP and updating the Q values.
- 2. Then, you have a testing phase of running <code>NUM\_EPIS\_TEST</code> episodes of the game, following an  $\varepsilon$ -greedy policy with  $\varepsilon=$  <code>TESTING\_EP</code>, which makes the agent choose the best action according to its current Q-values 95% of the time. At the testing phase of each epoch, you will compute the cumulative discounted reward for each episode and then obtain

the average reward over the <code>NUM\_EPIS\_TEST</code> episodes.

Finally, at the end of the experiment, you will get a sequence of data (of size NUM\_EPOCHS) that represents the testing performance at each epoch.

Note that there is randomness in both the training and testing phase. You will run the experiment NUM\_RUNS times and then compute the averaged reward performance over NUM\_RUNS experiments.

Most of these operations are handled by the boilerplate code provided in the <code>agent\_tabular\_QL.py</code> file by functions <code>run\_epoch</code> and <code>main</code>, but you will need to complete the <code>run\_episode</code> function.

Write a run\_episode function that takes a boolean argument (whether the epsiode is a training episode or not) and runs one episode.

**Reminder:** You should implement this function locally first. Make sure you can achieve reasonable performance on the Home World game before submitting your code

**Available Functions:** You have access to the NumPy python library as <code>np</code>, framework methods <code>framework.newGame()</code> and <code>framework.step\_game()</code>, constants <code>TRAINING\_EP</code> and <code>TESTING\_EP</code>, <code>GAMMA</code>, dictionaries <code>dict\_room\_desc</code> and <code>dict\_quest\_desc\_</code> and <code>previously implemented functions <code>epsilon\_greedy</code> and <code>tabular\_QLearning</code></code>

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

Correct

```
def run episode(for training):
    """ Runs one episode
    If for training, update 0 function
   If for testing, computes and return cumulative discounted reward
    Args:
        for training (bool): True if for training
    Returns:
        None
    .....
    epsilon = TRAINING EP if for training else TESTING EP
    gamma step = 1
   epi reward = 0
    (current room desc, current quest desc, terminal) = framework.newGame()
    while not terminal:
        # Choose next action and execute
        cur_room_desc_id = dict_room_desc[current_room_desc]
        cur_quest_desc_id = dict_quest_desc[current_quest_desc]
        (action index, object index) = epsilon greedy(cur room desc id,
                                                      cur_quest_desc_id,
                                                      q func, epsilon)
        (next_room_desc, next_quest_desc, reward,
        terminal) = framework.step game(current room desc, current quest desc,
                                         action index, object index)
        if for_training:
            # update Q-function.
            next_room_desc_id = dict_room_desc[next_room_desc]
            next_quest_desc_id = dict_quest_desc[next_quest_desc]
            tabular_q_learning(q_func, cur_room_desc_id, cur_quest_desc_id,
                               action index, object index, reward,
```

```
next_room_desc_id, next_quest_desc_id, terminal)

if not for_training:
    # update reward
    epi_reward = epi_reward + gamma_step * reward
    gamma_step = gamma_step * GAMMA

# prepare next step
    current_room_desc = next_room_desc
    current_quest_desc = next_quest_desc

if not for_training:
    return epi_reward
```

### Test results

```
Test: correctly computes reward

Testing computes correct reward for testing episode

Output:

eps greedy / e
```

Test: updates qvalue

Testing computes correct reward for testing episode

Output:

eps greedy / ql / eps greedy / ql /

<u>Hide output</u>

Submit

You have used 0 of 25 attempts

• Answers are displayed within the problem

# Report performance

2/2 points (graded)

In your Q-learning algorithm, initialize Q at zero. Set  $[num\_runs]=10$ ,  $[num\_epis\_train]=25$ ,  $[num\_epis\_test]=50$ ,  $\gamma=0.5$ ,  $[training\_ep]=0.5$ ,  $[training\_ep]=0.5$ ,  $[training\_ep]=0.05$  and the learning rate  $\alpha=0.1$ .

Please enter the number of epochs when the learning algorithm converges. That is, the testing performance become stable.

15

**✓ Answer:** 15

Please enter the average episodic rewards of your Q-learning algorithm when it converges.

Submit

You have used 3 of 6 attempts

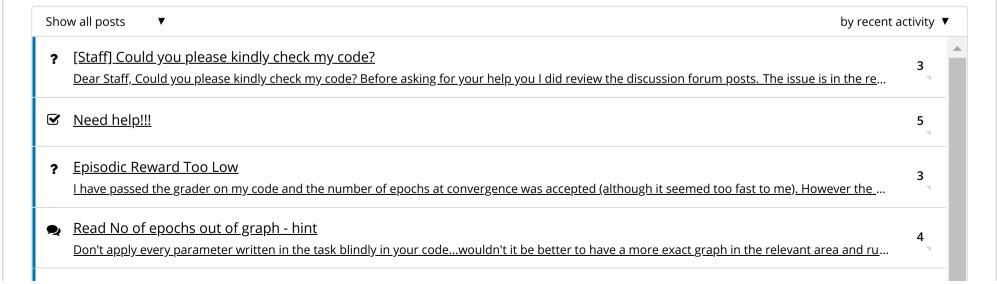
• Answers are displayed within the problem

#### Discussion

**Hide Discussion** 

**Topic:** Unit 5 Reinforcement Learning (2 weeks): Project 5: Text-Based Game / 4. Tabular Q-learning for Home World game

#### Add a Post



? [STAFF] The last question is super frustrating. Please add more attempts.  My algorithm is converging very fast, the first value got accepted. So I calculated the mean value of all the rewards from all episodes for this valu	5	
? [STAFF] Convergence velocity (can someone look my code?)  I read all the discussion about the convergence, I was getting a low convergence number and my code passed the grader even so. I found an err	4	
? run_episode(for_training)	4	
? Grader accepted but locally code has TypeError While trying to debug my code I decided to run it through the grader as I sometimes learn from the expected outputs it is giving and to my surpri	5	
? [staff] Tabular Q-learning Please check my code  I am almost there, geting the update qvalue correct but still struggling with the reward. I thought I have upadte the reward correctly, but it seems	5	
[staff] Report Performace - Please, have a look at my code.	5	
? [Staff] Check my run_episode code, pls  I have passed grader but I am curious why I am getting very low cumulative reward compared to expected reward? thanks.	5	
? [staff] why do I get staff debug error?  Hi, I receive the following message: There was a problem running your solution (Staff debug; L379). We couldn't run your solution (Staff debug; L	8	
run episode() half correct  I am stuck on deciding the initial value of epi reward in the run episode() function. Any hint, please?	12	
? Trying to understand the comment "initialize for each episode" in the beginning part of run_episode()	4	•

© All Rights Reserved