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 Privacy Policy.





Unit 5 Reinforcement Learning (2

Course > weeks)

> Project 5: Text-Based Game > 7. Linear Q-Learning

7. Linear Q-Learning

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

In this tab, you will implement the Q-learning algorithm with linear function approximation.

Recall the linear approximation we chose.

$$Q\left(s,c, heta
ight)=\phi(s,c)^T heta$$

with

$$\phi\left(s,c
ight)=\left[egin{array}{c} \mathbf{0}\ dots\ \mathbf{0}\ \psi_{R}\left(s
ight)\ \mathbf{0}\ dots\ \mathbf{0} \end{array}
ight]$$

Now, define $\hat{ heta}_i$ for i in range $1, d_C$ so that:

$$heta = \left[egin{array}{c} \hat{ heta}_1 \ dots \ \hat{ heta}_i \ dots \ \hat{ heta}_{d_C} \end{array}
ight]$$

With this notation, we get:

$$Q\left(s,c, heta
ight)=\psi_{R}(s)^{T}\hat{ heta}_{c}$$

In practice, we can implement $\hat{ heta}$ as a 2D array, so that

$$\left[egin{array}{c} Q\left(s,1, heta
ight) \ dots \ Q\left(s,d_{C}, heta
ight) \end{array}
ight] = \left[egin{array}{c} \hat{ heta}_{1}^{T} \ dots \ \hat{ heta}_{d_{C}}^{T} \end{array}
ight] \cdot \psi_{R}\left(s
ight)$$

Epsilon-greedy exploration

1.0/1 point (graded)

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

Hint: You can access $Q(s,c,\theta)$ using <code>q_value = (theta @ state_vector)[tuple2index(action_index, object_index)]</code>

Available Functions: You have access to the NumPy python library as <code>np</code> and functions <code>tuple2index</code> and <code>index2tuple</code>. Your code should also use constants <code>NUM ACTIONS</code> and <code>NUM OBJECTS</code>

```
def epsilon_greedy(state_vector, theta, epsilon):
    """Returns an action selected by an epsilon-greedy exploration policy

Args:
    state_vector (np.ndarray): extracted vector representation
    theta (np.ndarray): current weight matrix
    epsilon (float): the probability of choosing a random command

Returns:
```

```
(int, int): the indices describing the action/object to take

"""

# TODO Your code here
explore = np.random.random() <= epsilon
if explore:
    action index. object index = np.random.choice(NUM ACTIONS, 1)[0], np.random.choice(NUM OBJECTS, 1)[0]</pre>
```

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

Correct

```
def epsilon greedy(state vector, theta, epsilon):
    """Returns an action selected by an epsilon-greedy exploration policy
    Args:
        state vector (np.ndarray): extracted vector representation
        theta (np.ndarray): current weight matrix
        epsilon (float): the probability of choosing a random command
    Returns:
        (int, int): the indices describing the action/object to take
    coin = np.random.random_sample()
   if coin < epsilon:</pre>
        action index = np.random.randint(NUM ACTIONS)
        object_index = np.random.randint(NUM_OBJECTS)
    else:
        q values = theta @ state vector
        index = np.argmax(q values)
        action_index, object_index = index2tuple(index)
    return (action_index, object_index)
```

Test results

See full output
CORRECT
See full output

Submit

You have used 1 of 25 attempts

1 Answers are displayed within the problem

Linear Q-learning

1.0/1 point (graded)

Write a function [linear_q_learning] that updates the theta weight matrix, given the transition date (s, a, R(s, a), s').

Reminder: You should implement this function locally first. You should test this function along with the next one and make sure you achieve reasonable performance

Hint: You can access $Q\left(s,a, heta
ight)$ using <code>q_value = (theta @ state_vector)[tuple2index(action_index, object_index)]</code>

Available Functions: You have access to the NumPy python library as <code>np</code> . You should also use constants <code>ALPHA</code> and <code>GAMMA</code> in your code

```
theta (np.ndarray): current weight matrix
 7
          current_state_vector (np.ndarray): vector representation of current state
          action index (int): index of the current action
          object index (int): index of the current object
 9
10
          reward (float): the immediate reward the agent recieves from playing current command
          next_state_vector (np.ndarray): vector representation of next state
11
          terminal (bool): True if this epsiode is over
12
13
14
      Returns:
15
          None
```

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

Correct

```
def linear q learning(theta, current state vector, action index, object index,
                      reward, next_state_vector, terminal):
    """Update theta for a given transition
   Args:
        theta (np.ndarray): current weight matrix
        current state vector (np.ndarray): vector representation of current state
        action index (int): index of the current action
        object index (int): index of the current object
        reward (float): the immediate reward the agent recieves from playing current command
        next state vector (np.ndarray): vector representation of next state
        terminal (bool): True if this epsiode is over
   Returns:
        None
    .....
   q_values_next = theta @ next_state_vector
   maxq_next = np.max(q_values_next)
    q_values = theta @ current_state_vector
   cur index = tuple2index(action index, object index)
   q_value_cur = q_values[cur_index]
   target = reward + GAMMA * maxq_next * (1 - terminal)
   theta[cur_index] = theta[cur_index] + ALPHA * (
       target - q_value_cur) * current_state_vector
```

Test results

See full output

CUKKECI

See full output

Submit

You have used 4 of 25 attempts

• Answers are displayed within the problem

Evaluate linear Q-learning on Home World game

1/1 point (graded)

Adapt your run_episode function to call linear_Q_learning and evaluate your performance using hyperparmeters:

Set $[NUM_RUNS] = 5$, $[NUM_EPIS_TRAIN] = 25$, $[NUM_EPIS_TEST] = 50$, $[\gamma = 0.5]$, $[TRAINING_EP] = 0.5$, $[TRAINING_EP] = 0.05$ and the learning rate $\alpha = 0.001$.

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

0.3262413375213623

✓ Answer: 0.37

Submit

You have used 2 of 6 attempts

1 Answers are displayed within the problem

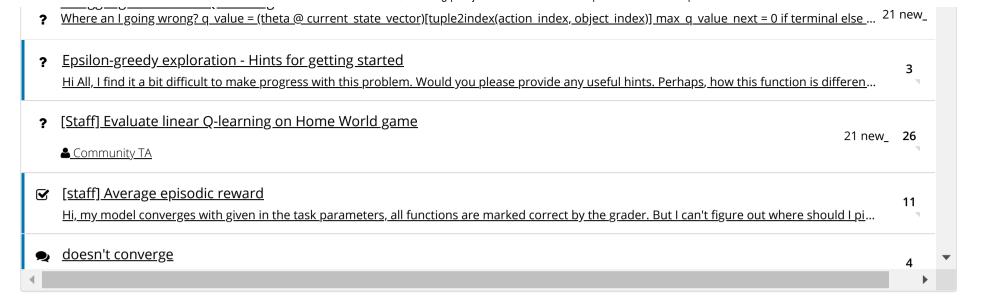
Discussion

Hide Discussion

Topic: Unit 5 Reinforcement Learning (2 weeks): Project 5: Text-Based Game / 7. Linear Q-Learning

Add a Post

Show all posts ▼ by recent a	activity
Typo but it is true terminal (bool): True if this epsiode is over IT IS OVER!	1
<u>Updating other rows of theta</u> <u>I only update one row of theta per the "Struggling with Q-Learning Update" discussion posted by Rohit, and one row of my output matches the c</u>	13
? How did you do epsilon-greedy exploration, linear Q-learning, and Evaluate? Here are my thoughts and what I have done. I'm getting confused by feature engineering. Is this similar to epsilon greedy() from three tabs ago? But at that time, we were given q func matrix	6
? [STAFF] Reset Greedy Tries (Staff debug: L364) Could staff please reset my Epsilon greedy tries. I have been having some issues that contributed to me wasting a lot of tries.	7
? [Staff] Epsilon-greedy exploration Hi, Over all this course i got along all subjects with more or less commodity but with this topic (Linear Q-learning) i can't get the big picture. I follo	6
? [Staff]There was a problem running the staff solution (Staff debug: L364) Liget the error again when I submit epsilon-greedy	9
Why don't we achieve optimal results? OK, same question as before. I received a green tick for my average episodic rewards after convergence. However, this number falls a fair bit sho	1
? [STAFF] Linear Q-learning or anyone! if terminal is True, then max q value nxt = 0 else max q value nxt is max of theta time time next state vector y = reward + (GAMMA * max q va	4



© All Rights Reserved