edX

**Audit Access Expires May 25, 2020**
You lose all access to this course, including your progress, on May 25, 2020.

# 8. Linear Q-Learning

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

Recall the linear approximation we chose.

$$Q\left(s, c, \theta\right) = \phi(s, c)^{T} \theta$$

with

Generating Speech Output

$$\phi\left(s, c\right) = \begin{bmatrix} \mathbf{0} \\ \vdots \\ \mathbf{0} \\ \psi_R\left(s\right) \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}$$

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

$$\theta = \begin{bmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_i \\ \vdots \\ \hat{\theta}_{d_C} \end{bmatrix}$$

With this notation, we get:

$$Q\left(s, c, \theta\right) = \psi_R(s)^T \hat{\theta}_c$$

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

$$\begin{bmatrix} Q\left(s, 1, \theta\right) \\ \vdots \\ Q\left(s, d_C, \theta\right) \end{bmatrix} = \begin{bmatrix} \hat{\theta}_1^T \\ \vdots \\ \hat{\theta}_{d_C}^T \end{bmatrix} \cdot \psi_R\left(s\right)$$

Generating Speech Output

## Epsilon-greedy exploration

1.0/1 point (graded)

Now you will write a function `epsilon\_ greedy` that implements the $\varepsilon$-greedy exploration policy using the current Q-function.

**Hint:** You can access $Q\left(s, c, \theta\right)$ using

`q_value = (theta @ state_vector)[tuple2index(action_index, object_index)]`

**Available Functions:** You have access to the NumPy python library as `np` and functions `tuple2index` and `index2tuple` . Your code should also use constants `NUM_ACTIONS` and `NUM_OBJECTS`

```
1 def epsilon_greedy(state_vector, theta, epsilon):
2     """Returns an action selected by an epsilon-greedy exploration poli
3
4     Args:
5         state_vector (np.ndarray): extracted vector representation
6         theta (np.ndarray): current weight matrix
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    if np.random.random() < epsilon:
13        action_index, object_index = np.random.randint(0, NUM_ACTIONS),
14                                     np.random.randint(0, NUM_OBJECTS)
15    else:
```

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

Correct

```python
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:
        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 2 of 25 attempts

ℹ 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 $(s, a, R(s, a), s')$.

Generating Speech Output

**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 $Q\left(s, a, \theta\right)$ using

```
q_value = (theta @ state_vector)[tuple2index(action_index, object_index)]
```

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

```
 1 def linear_q_learning(theta, current_state_vector, action_index, object
 2                       reward, next_state_vector, terminal):
 3     """Update theta for a given transition
 4
 5     Args:
 6         theta (np.ndarray): current weight matrix
 7         current_state_vector (np.ndarray): vector representation of cur
 8         action_index (int): index of the current action
 9         object_index (int): index of the current object
10         reward (float): the immediate reward the agent recieves from pl
11         next_state_vector (np.ndarray): vector representation of next s
12         terminal (bool): True if this epsiode is over
13
14     Returns:
15         None
```

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

Correct

Generating Speech Output

```
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 sta
        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 cu
        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

**CORRECT**

See full output

Submit      You have used 15 of 25 attempts

ⓘ  Answers are displayed within the problem

Generating Speech Output

# 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` $= 0.5$, `TESTING_EP` $= 0.05$ and the **learning rate** $\alpha = 0.01$: .

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

| 0.39 |
|---|

✔ **Answer:** 0.37

| Submit | You have used 2 of 6 attempts |
|---|---|

---

ⓘ   Answers are displayed within the problem

---

# Discussion

| | **Hide Discussion** |
|---|---|

**Topic:** Unit 5 Reinforcement Learning (2 weeks) :Project 5: Text-Based Game / 8. Linear Q-Learning

**Add a Post**

| Show all posts | by recent activity |
|---|---|

💬  [Tips from staff:] Evaluate linear Q-learning on Home World game
📌 Pinned   👤 Staff                                                          2

💬  updating theta
a small guidance which helped me: each update should not be for a full theta matrix, but for ...          12 new

☑  What is a dimension of theta and how to understand the hints equation          1 new   5

💬  plot
There is a plt.show() missing from the end of the agent_linear skeleton code          5

Generating Speech Output

**?** [Staff] Solver Problem

Hi, I got this error message "There was a problem running your solution (Staff debug: L379)." ...

2

**?** Epsilon-greedy

I think I'm overthinking the epsilon-greedy implementation... I understand that the q-value is ...

2

💬 My converged reward is not liked by the grader..

👤 Community TA

4

💬 Using the default hyperparameters won't let you pass the last question

7

💬 linear_q_learning ( theta .... - what is the point to pass it as a parameter if it is global?

Does anybody have and idea?

6

☑ How do you find phi(s, c)?

I'm feeling really dense, but how do you come up with phi (s, c) for updating theta? Do I need ...

9

☑ [Staff] Should theta be returned by linear_q_learning?

Currently, the comment says that linear_q_learning() returns None; It is not clear to me how t...

2

☑ [Staff or Student] q_value = (theta @ state_vector)[tuple2index(action_index, object_index)]

Greetings, I'm having a hard time interpreting this for some reason. 1) theta @ state_vector yi...

3

Generating Speech Output

2020-05-16, 11:06 p.m.