

Task 7.1D: Function approximation implementation

GitHub Link:

Objective: To Implement Task 1.1P with the following methods:

- Semi-Gradient Sarsa(0) (From Slide 14)
- Semi-Gradient TD(λ) (From Slide 9)

```
In [1]: #Loading all of our libraries...
import numpy as np
import matplotlib.pyplot as plt
# import gym
import sys

In [2]: #Connecting our Google Drive...
from google.colab import drive
drive.mount('/content/drive')
sys.path.insert(0, '/content/drive/MyDrive/Colab Notebooks/')

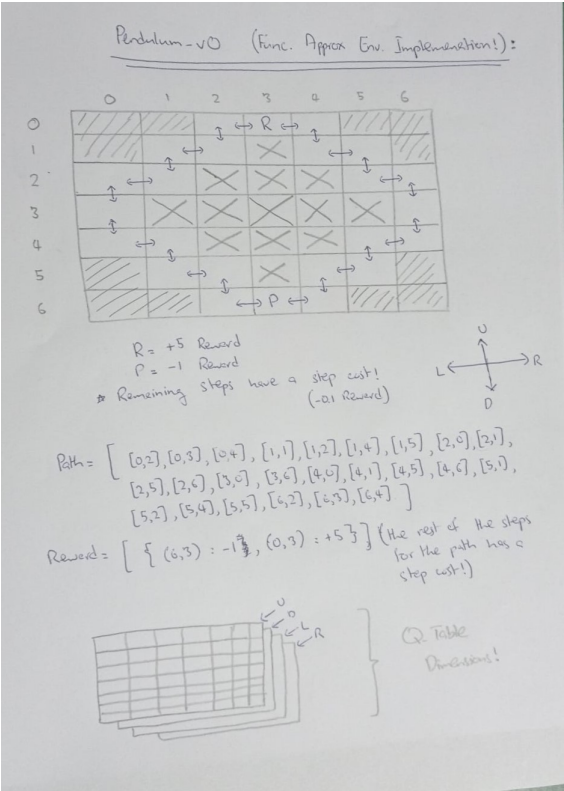
#Importing our GridWorld Module after connection...
from GW import Grid, print_values, print_policy

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
```

Creating our Environment

```
In [3]: #All the Constants...
ALL_POSSIBLE_ACTIONS = ('U', 'D', 'L', 'R')
num_episodes = 100
GAMMA = 0.9
ALPHA = 0.1
eps = 0.1
```

This is environment we have specified to the current model. This time we discretized the model based on the GridWorld. We were able to create our own Custom GridWorld based on the circular path of the Pendulum. Here's how we have developed our Environment on Sketch:



```
In [4]: #Creating the Pendulum Environment...

pendulum = Grid(7, 7, (6, 3))
step_cost = 0.1

#Dictionary of the rewards assigned at every step of the path...
rewards = {
    (0, 3): 5,
    (0, 2): step_cost,
    (0, 4): step_cost,
    (1, 1): step_cost,
    (1, 2): step_cost,
    (1, 4): step_cost,
    (1, 5): step_cost,
    (2, 0): step_cost,
    (2, 1): step_cost,
    (2, 5): step_cost,
    (2, 6): step_cost,
    (3, 0): step_cost,
    (3, 1): step_cost,
    (4, 0): step_cost,
    (4, 1): step_cost,
    (4, 5): step_cost,
    (4, 6): step_cost,
    (5, 1): step_cost,
    (5, 2): step_cost,
    (5, 4): step_cost,
    (5, 5): step_cost,
    (6, 2): step_cost,
    (6, 4): step_cost,
    (6, 3): 1
}

#Dictionary of the actions assigned at every step of the path...
actions = {
    (0, 3): ('D', 'R'),
    (0, 3): ('L', 'R'),
    (0, 4): ('L', 'D'),
    (1, 1): ('R', 'D'),
    (1, 2): ('L', 'U'),
    (1, 4): ('R', 'U'),
    (1, 5): ('L', 'D'),
    (2, 0): ('D', 'R'),
    (2, 1): ('L', 'U'),
    (2, 5): ('R', 'U'),
    (2, 6): ('R', 'U'),
    (3, 0): ('U', 'D'),
    (3, 1): ('U', 'D'),
    (4, 0): ('U', 'R'),
    (4, 1): ('L', 'D'),
    (4, 5): ('R', 'D'),
    (4, 6): ('L', 'U'),
    (5, 1): ('U', 'R'),
    (5, 2): ('L', 'D'),
    (5, 4): ('D', 'R'),
    (5, 5): ('L', 'U'),
    (6, 2): ('U', 'R'),
    (6, 3): ('L', 'R'),
    (6, 4): ('L', 'U'),
}

#Setting our Grid with the rewards and actions...
pendulum.set(rewards, pendulum)

In [5]: print(".....Rewards per state in the Environment.....\n")
print_values(rewards, pendulum)
```

-----Rewards per state in the Environment-----

```
0.00 |
0.00 |
-0.10 |
5.00 |
-0.10 |
0.00 |
0.00 |
```

```
0.00 |
-0.10 |
-0.10 |
0.00 |
-0.10 |
-0.10 |
0.00 |
```

```
-0.10 |
-0.10 |
0.00 |
0.00 |
0.00 |
-0.10 |
-0.10 |
```

```
-0.10 |
0.00 |
0.00 |
0.00 |
0.00 |
-0.10 |
```

```
-0.10 |
-0.10 |
0.00 |
0.00 |
0.00 |
-0.10 |
-0.10 |
```

```
0.00 |
-0.10 |
-0.10 |
0.00 |
-0.10 |
-0.10 |
0.00 |
```

```
0.00 |
0.00 |
-0.10 |
-1.00 |
-0.10 |
0.00 |
0.00 |
```

In [6]: #Ensiting all the possible actions per state...

```
state = list(actions.keys())
possible_actions = list(actions.values())

print("-----Possible Actions at every step in the Environment-----\n")
for i in range(len(state)):
    print("{} | {}".format(state[i], possible_actions[i]))
```

-----Possible Actions at every step in the Environment-----

```
(0, 2) | ('D', 'R')
(0, 3) | ('L', 'R')
(4, 4) | ('L', 'D')
(1, 3) | ('R', 'D')
(1, 2) | ('L', 'U')
(1, 4) | ('R', 'U')
(1, 5) | ('L', 'D')
(2, 4) | ('D', 'R')
(2, 3) | ('L', 'U')
(2, 5) | ('R', 'U')
(2, 0) | ('L', 'D')
(3, 0) | ('U', 'D')
(3, 4) | ('U', 'D')
(4, 0) | ('U', 'R')
(4, 1) | ('L', 'D')
(4, 3) | ('R', 'D')
(4, 6) | ('L', 'U')
(5, 1) | ('U', 'R')
(5, 2) | ('L', 'D')
(5, 4) | ('D', 'R')
(5, 5) | ('L', 'U')
(6, 2) | ('U', 'R')
(6, 3) | ('L', 'R')
(6, 4) | ('L', 'U')
```

Semi-Gradient Sarsa(0)

We will create some basic functionalities for our algorithms while the agent learns in the environment.

In [14]: # #function for getting the actions of our Optimal Policy...

```
# def max_dict(d):
#     # returns the argmax (key) and max (value) from a dictionary
#     max_key = None
#     max_val = float('-inf')
#     for k, v in d.items():
#         if v > max_val:
#             max_val = v
#             max_key = k
#     return max_key, max_val
```

In [15]: #Greedy/Exploration Function...

```
def epsilon_greedy_action(state, epsilon):
    #Exploration
    if np.random.uniform() < epsilon:
        return np.random.choice(ALL_POSSIBLE_ACTIONS) #Returns letter, not the index...
    #Exploitation
    else:
        _, maxarg_a = actions[(state[0], state[1])]
        return maxarg_a #Returns letter, not the index...
```

In [16]: #Function for mapping actions into index...

```
def action_map(a):
    if a == 'U':
        i = 0
    elif a == 'D':
        i = 1
    elif a == 'R':
        i = 2
    elif a == 'L':
        i = 3
    return i
```

In [17]: #Function to return reward with from a certain state with Grid Properties...

```
def step_function(s,a):
    # r = pendulum.move(a)
    i,j = s

    if s in actions.keys():
        if a == 'U':
            i -= 1
        elif a == 'D':
            i += 1
        elif a == 'R':
            j -= 1
        elif a == 'L':
            j += 1

        next_s = (i,j)
        if next_s in actions.keys():
            r = rewards.get(next_s, 0)
            #print(next_s,"-----",r)
            return next_s, r
        else:
            r = 0
            #print(next_s,"-----",r)

    #Undo the move...
    if a == 'U':
        i += 1
    elif a == 'D':
        i -= 1
    elif a == 'R':
        j += 1
    elif a == 'L':
        j -= 1
    #print("\nOut of bounds. Move Undone...")
    next_s = s
    return next_s, r

    else:
        r = 0
        #print("\nOut of bounds. Move Undone...")
        return s, r

#Running a small test for the function...
next_s, r = step_function((5,3),'D')
print(next_s, r)

(5, 3) 0
```

In [18]: #function for computing the gradient of the model...

```
def grad(Q, state, action):
    ci, cj = state
    gradient = np.zeros_like(Q)
    gradient = 1
    return gradient
```

For our Semi-Gradient SARSA(0), this is the pseudocode for us to implement:

Input:
A differentiable state-action value function $\hat{q}: \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$
A policy π if predicting or q_π if estimating (e.g. using $\varepsilon - greedy$)

Algorithm Parameter:
Step size $\alpha \in (0,1]$

Initialise:
 $\mathbf{w} \in \mathbb{R}^d$ arbitrarily e.g. $\mathbf{w} = 0$

Loop forever (for each episode):
 $S, A \leftarrow$ initial state and action of episode (e.g. using $\varepsilon - greedy$)
Loop for each step of the episode until $S \in \mathcal{S}^{(Terminal)}$:
Take action A , observe R, S'
If $S' \in \mathcal{S}^{(terminal)}$ then:
 $\mathbf{w} = \mathbf{w} + \alpha[R + \gamma \hat{q}(S, A, \mathbf{w}) - \nabla \hat{q}(S, A, \mathbf{w})]$, special case for terminal state can't include future state
else:
Choose A' as a function of $\hat{q}(S'; \mathbf{w})$ (e.g. using $\varepsilon - greedy$)
 $\mathbf{w} = \mathbf{w} + \alpha[R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$
 $S \leftarrow S'$
 $A \leftarrow A'$

```
In [19]: def semi_gradient_sarsa(num_episodes, alpha, gamma, epsilon):
total_reward_per_episode = []
average_reward_per_episode = []

for i in range(num_episodes):
    #pendulum.set_state(s) / state = env.reset()
    state = (6,3) # starting point of the Agent in the Environment...
    num_actions = len(ALL_POSSIBLE_ACTIONS) # Number of Actions
    Q = np.zeros((7, 7, num_actions)) # Q-table Initialized
    total_Q = np.zeros((7, 7, num_actions))
    total_reward = 0

    action = epsilon_greedy_action(state, epsilon)
    for t in range(200):

        # Getting next state and reward...
        next_state, reward = step_function(state, action)

        #Getting dimensions of the current and next state in order to update the Q-Table...
        ci, cj = state
        ni, nj = next_state

        #Get the next action...
        next_action = epsilon_greedy_action(next_state, epsilon)

        #Mapping the current and next action...
        a = action_map(action)
        next_a = action_map(next_action)

        #update the Q-Table...
        td_err = reward + gamma * Q[ni][nj][next_a] - Q[ci][cj][a]
        Q[ci][cj][a] += alpha * td_err * grad(Q[ci][cj][a], state, a)

        #creating the total sum of the reward...
        total_reward += reward

        #Assign the new state and action and repeat...
        state = next_state
        action = next_action

    #Assigning the final state values in the final step...
    if i + 1 == 100:
        total_Q = Q

    average_reward_per_episode.append(total_reward/200)
    total_reward_per_episode.append(total_reward)
    print("Episode --- [(i)/100] | Reward ---- {}".format(i + 1, total_reward))

return total_reward_per_episode, average_reward_per_episode, total_Q
```

Semi-Gradient TD(λ)

Again, we are going to implement some of the basic functionalities for this algorithm as well.

```
In [7]: # function for predicting the weights...
def predict(state, weights):
    i,j = state
    return np.dot(weights[i,j,:], np.ones(4))

In [8]: #function to return reward with from a certain state with Grid Properties...

def step_function(s,a):
    # r = pendulum.move(a)
    i,j = s

    if s in actions.keys():
        if a == 0: # 'U'
            i += 1
        elif a == 1: # 'D'
            i -= 1
        elif a == 2: # 'L'
            j += 1
        elif a == 3: # 'R'
            j -= 1

        next_s = (i,j)
        if next_s in actions.keys():
            r = rewards.get(next_s, 0)
            # print(next_s,"-----",r)
            return next_s, r
        else:
            r = 0
            #print(next_s,"-----",r)

    #Undo the move...
    if a == 0: # 'U'
        i += 1
    elif a == 1: # 'D'
        i -= 1
    elif a == 2: # 'L'
        j += 1
    elif a == 3: # 'R'
        j -= 1
    #print("\nOut of bounds. Move Undone...")
    next_s = s
    return next_s, r

    else:
        r = 0
        #print("\nOut of bounds. Move Undone...")
        return s, r

In [9]: #Running a small test for the function...
next_s, r = step_function((0,2),3)
print(next_s, r)

(0, 3) 5
```

Likewise, for our second algorithm, this is how we will implement from the below pseudocode:

Input:
The policy π to be evaluate
A differentiable function $\hat{v}: \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm Parameter:
Step size $\alpha \in (0,1]$
Trace decay rate $\lambda \in [0,1]$

Initialise:
 $\mathbf{w} \in \mathbb{R}^d$ arbitrarily e.g. $\mathbf{w} = 0$

Loop forever (for each episode):
Initialize S
Reset $\mathbf{z} = 0$
Loop for each step of the episode until $S \in \mathcal{S}^{(Terminal)}$:
Choose $A \sim \pi(\cdot | S)$
Take action A , observe R, S'
 $\mathbf{z} \leftarrow \gamma \lambda \mathbf{z} + \nabla \hat{v}(S, \mathbf{w})$
 $\hat{\delta} \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$
 $\mathbf{w} = \mathbf{w} + \alpha \hat{\delta} \mathbf{z}$
 $S \leftarrow S'$

```
In [10]: def semi_gradient_td_lambda(num_episodes, alpha, gamma, epsilon, lambda_):
total_reward_per_episode = []
average_reward_per_episode = []

for i in range(num_episodes):
    #pendulum.set_state(s) / state = env.reset()
    state = (6,3) # starting point of the Agent in the Environment...
    num_actions = len(ALL_POSSIBLE_ACTIONS) # Number of Actions
    weights = np.zeros((7, 7, num_actions)) # Q-table/weights Initialized
    total_Q = np.zeros((7, 7, num_actions))
    eligibility_trace = np.zeros((7,7,4)) # Resetting Eligibility Trace

    total_reward = 0
    # s = np.argmax(predict(state, weights))
    a = np.random.choice([2,3]) # To kick-start the algorithm...

    for t in range(200):

        # Getting next state and reward...
        next_state, reward = step_function(state, a) #env.step(action)

        #Getting dimensions of the current and next state in order to update the weights...
```

```

    ci, cj = state
    # Ri, Rj = next_state

    #Get the next action...
    next_a = predict(next_state, weights)

    #Computing TD Error (delta)...
    delta = reward + gamma * predict(next_state, weights) - predict(state, weights)

    #Updating the Q-Table/Weights...
    eligibility_trace *= lambda
    eligibility_trace[ci, ci, int(a)] += 1
    weights += alpha * delta * eligibility_trace

    #Creating the total sum of the reward...
    total_reward += reward

    #Assign the new state and action and repeat...
    state = next_state
    a = next_a

    #Assigning the final state values in the final step...
    if i + 1 == 100:
        Total_Q += weights

        average_reward_per_episode.append(total_reward/200)
        total_reward_per_episode.append(total_reward)
        print("Episode ---- ({} / 100) | Reward ---- {}".format(i + 1, total_reward))

    return total_reward_per_episode, average_reward_per_episode, Total_Q

```

Comparison of Results

SARSA(0) Results

```

In [20]: type, sarsa_args, Total_Q = semi_gradient_sarsa(num_episodes, ALPHA, GAMMA, eps)
print("Average Reward after 100 Episodes: ", np.mean(trype))

Episode --- (1/100) | Reward ---- -5.3
Episode --- (2/100) | Reward ---- -5.899999999999998
Episode --- (3/100) | Reward ---- -1.6000000000000003
Episode --- (4/100) | Reward ---- -1.6000000000000003
Episode --- (5/100) | Reward ---- -13.299999999999997
Episode --- (6/100) | Reward ---- -6.999999999999998
Episode --- (7/100) | Reward ---- -6.8999999999999955
Episode --- (8/100) | Reward ---- -1.8999999999999999
Episode --- (9/100) | Reward ---- -4.6
Episode --- (10/100) | Reward ---- -6.399999999999999
Episode --- (11/100) | Reward ---- -3.9000000000000017
Episode --- (12/100) | Reward ---- -6.599999999999999
Episode --- (13/100) | Reward ---- -9.699999999999998
Episode --- (14/100) | Reward ---- -4.4
Episode --- (15/100) | Reward ---- -7.899999999999998
Episode --- (16/100) | Reward ---- -18.899999999999998
Episode --- (17/100) | Reward ---- -5.299999999999999
Episode --- (18/100) | Reward ---- -4.4
Episode --- (19/100) | Reward ---- -3.8000000000000007
Episode --- (20/100) | Reward ---- -6.999999999999998
Episode --- (21/100) | Reward ---- -3.3000000000000004
Episode --- (22/100) | Reward ---- -4.799999999999999
Episode --- (23/100) | Reward ---- -7.699999999999998
Episode --- (24/100) | Reward ---- -7.899999999999997
Episode --- (25/100) | Reward ---- -6.7999999999999952
Episode --- (26/100) | Reward ---- -14.8999999999999963
Episode --- (27/100) | Reward ---- -5.3000000000000005
Episode --- (28/100) | Reward ---- -4.2
Episode --- (29/100) | Reward ---- -4.3999999999999995
Episode --- (30/100) | Reward ---- -5.3
Episode --- (31/100) | Reward ---- -5.5
Episode --- (32/100) | Reward ---- -6.399999999999998
Episode --- (33/100) | Reward ---- -4.3999999999999995
Episode --- (34/100) | Reward ---- -12.299999999999997
Episode --- (35/100) | Reward ---- -7.6999999999999975
Episode --- (36/100) | Reward ---- -6.799999999999998
Episode --- (37/100) | Reward ---- -6.599999999999999
Episode --- (38/100) | Reward ---- -5.3
Episode --- (39/100) | Reward ---- -12.999999999999997
Episode --- (40/100) | Reward ---- -13.199999999999996
Episode --- (41/100) | Reward ---- -4.8
Episode --- (42/100) | Reward ---- -8.6
Episode --- (43/100) | Reward ---- -6.799999999999999
Episode --- (44/100) | Reward ---- -4.9
Episode --- (45/100) | Reward ---- -3.1000000000000005
Episode --- (46/100) | Reward ---- -8.799999999999997
Episode --- (47/100) | Reward ---- -5.4999999999999994
Episode --- (48/100) | Reward ---- -5.3999999999999995
Episode --- (49/100) | Reward ---- -7.6999999999999975
Episode --- (50/100) | Reward ---- -1.6
Episode --- (51/100) | Reward ---- -5.699999999999998
Episode --- (52/100) | Reward ---- -12.899999999999997
Episode --- (53/100) | Reward ---- -4.8999999999999995
Episode --- (54/100) | Reward ---- -8.8999999999999986
Episode --- (55/100) | Reward ---- -7.6999999999999975
Episode --- (56/100) | Reward ---- -4.199999999999999
Episode --- (57/100) | Reward ---- -18.299999999999997
Episode --- (58/100) | Reward ---- -5.899999999999998
Episode --- (59/100) | Reward ---- -4.8
Episode --- (60/100) | Reward ---- -6.499999999999999
Episode --- (61/100) | Reward ---- -2.6
Episode --- (62/100) | Reward ---- -8.899999999999998
Episode --- (63/100) | Reward ---- -3.3000000000000003
Episode --- (64/100) | Reward ---- -10.599999999999998
Episode --- (65/100) | Reward ---- -12.599999999999996
Episode --- (66/100) | Reward ---- -5.899999999999999
Episode --- (67/100) | Reward ---- -6.699999999999999
Episode --- (68/100) | Reward ---- -6.399999999999999
Episode --- (69/100) | Reward ---- -10.699999999999998
Episode --- (70/100) | Reward ---- -9.699999999999998
Episode --- (71/100) | Reward ---- -6.399999999999998
Episode --- (72/100) | Reward ---- -4.2
Episode --- (73/100) | Reward ---- -4.8
Episode --- (74/100) | Reward ---- -9.599999999999996
Episode --- (75/100) | Reward ---- -3.3000000000000001
Episode --- (76/100) | Reward ---- -4.2
Episode --- (77/100) | Reward ---- -4.2
Episode --- (78/100) | Reward ---- -4.8000000000000001
Episode --- (79/100) | Reward ---- -4.4
Episode --- (80/100) | Reward ---- -5.899999999999999
Episode --- (81/100) | Reward ---- -5.3000000000000001
Episode --- (82/100) | Reward ---- -3.3000000000000001
Episode --- (83/100) | Reward ---- -6.999999999999998
Episode --- (84/100) | Reward ---- -13.199999999999974
Episode --- (85/100) | Reward ---- -7.899999999999998
Episode --- (86/100) | Reward ---- -3.1000000000000001
Episode --- (87/100) | Reward ---- -9.899999999999997
Episode --- (88/100) | Reward ---- -6.1
Episode --- (89/100) | Reward ---- -7.6999999999999975
Episode --- (90/100) | Reward ---- -8.399999999999997
Episode --- (91/100) | Reward ---- -6.399999999999999
Episode --- (92/100) | Reward ---- -10.899999999999998
Episode --- (93/100) | Reward ---- -6.8999999999999995
Episode --- (94/100) | Reward ---- -8.399999999999998
Episode --- (95/100) | Reward ---- -4.6
Episode --- (96/100) | Reward ---- -5.6999999999999975
Episode --- (97/100) | Reward ---- -3.1000000000000005
Episode --- (98/100) | Reward ---- -3.3000000000000003
Episode --- (99/100) | Reward ---- -5.6999999999999975
Episode --- (100/100) | Reward ---- -17.899999999999984
Average Reward after 100 Episodes: -6.579999999999998

In [21]: print("Average Reward per Episode: {}".format(sarsa_args))

Average Reward per Episode:
[-0.8265, -0.8264999999999999, -0.8080000000000002, -0.8064999999999995, -0.8249999999999999, -0.8344999999999996, -0.8055, -0.821, -0.8119999999999994, -0.8150000000000007, -0.8129999999999995, -0.8484999999999999, -0.8220000000000002, -0.8349999999999987, -0.8584999999999999, -0.8
2649999999999996, -0.8220000000000002, -0.8130000000000001, -0.8149999999999997, -0.8165000000000004, -0.8239999999999994, -0.8184999999999999, -0.8179999999999996, -0.8745999999999982, -0.8265000000000003, -0.821, -0.822, -0.8265, -0.8275, -0.8309999999999991, -0.822, -0.861499999999999
985, -0.8384999999999996, -0.8339999999999999, -0.8129999999999997, -0.8265, -0.864999999999998, -0.8659999999999998, -0.82, -0.841, -0.8339999999999996, -0.8245, -0.8150000000000003, -0.8149999999999998, -0.8214999999999981, -0.8219999999999996, -0.8349999999999986, -0.808, -0.8284999999999999, -0.865499
9999999985, -0.8244999999999991, -0.8349999999999986, -0.8384999999999988, -0.8284999999999998, -0.8514999999999998, -0.8254999999999988, -0.82, -0.8124999999999994, -0.8130000000000001, -0.8484999999999999, -0.8265, -0.8529999999999984, -0.8579999999999998, -0.8254999999999999, -0.8314999999999995, -0.8
3199999999999994, -0.8134999999999995, -0.8484999999999999, -0.8349999999999999, -0.8199999999999999, -0.821, -0.82, -0.8479999999999998, -0.8165000000000004, -0.821, -0.821, -0.8220000000000004, -0.8220000000000002, -0.8254999999999999, -0.8265000000000003, -0.8165000000000004, -0.8349999999999999, -0.8659999999999998, -0.8394
999999999987, -0.8155000000000002, -0.8484999999999998, -0.8385, -0.8384999999999988, -0.8413999999999998, -0.8313999999999994, -0.8649999999999998, -0.8344999999999996, -0.8413999999999994, -0.82, -0.8284999999999987, -0.8150000000000001, -0.8165, -0.8264999999999987, -0.8084999999999991]

In [22]: print("Q-Table in the 100th Episode: {}".format(Total_Q))

```

TD(Lambda) Results

Episode	Time [s]	Reward	Loss
Episode 1	2.1300	Reward -0.2	
Episode 2	2.3800	Reward -0.2	
Episode 3	2.4300	Reward -0.2	
Episode 4	2.7300	Reward -0.2	
Episode 5	2.7800	Reward -0.2	
Episode 6	2.9300	Reward -0.2	
Episode 7	3.0800	Reward -0.2	
Episode 8	3.1300	Reward -0.2	
Episode 9	3.2300	Reward -0.2	
Episode 10	3.3800	Reward -0.2	
Episode 11	3.4300	Reward -0.2	
Episode 12	3.5300	Reward -0.2	
Episode 13	3.6300	Reward -0.2	
Episode 14	3.7300	Reward -0.2	
Episode 15	3.8300	Reward -0.2	
Episode 16	3.9300	Reward -0.2	
Episode 17	4.0300	Reward -0.2	
Episode 18	4.1300	Reward -0.2	
Episode 19	4.2300	Reward -0.2	
Episode 20	4.3300	Reward -0.2	
Episode 21	4.4300	Reward -0.2	
Episode 22	4.5300	Reward -0.2	
Episode 23	4.6300	Reward -0.2	
Episode 24	4.7300	Reward -0.2	
Episode 25	4.8300	Reward -0.2	
Episode 26	4.9300	Reward -0.2	
Episode 27	5.0300	Reward -0.2	
Episode 28	5.1300	Reward -0.2	
Episode 29	5.2300	Reward -0.2	
Episode 30	5.3300	Reward -0.2	
Episode 31	5.4300	Reward -0.2	
Episode 32	5.5300	Reward -0.2	
Episode 33	5.6300	Reward -0.2	
Episode 34	5.7300	Reward -0.2	
Episode 35	5.8300	Reward -0.2	
Episode 36	5.9300	Reward -0.2	
Episode 37	6.0300	Reward -0.2	
Episode 38	6.1300	Reward -0.2	
Episode 39	6.2300	Reward -0.2	
Episode 40	6.3300	Reward -0.2	
Episode 41	6.4300	Reward -0.2	
Episode 42	6.5300	Reward -0.2	
Episode 43	6.6300	Reward -0.2	
Episode 44	6.7300	Reward -0.2	
Episode 45	6.8300	Reward -0.2	
Episode 46	6.9300	Reward -0.2	
Episode 47	7.0300	Reward -0.2	
Episode 48	7.1300	Reward -0.2	
Episode 49	7.2300	Reward -0.2	
Episode 50	7.3300	Reward -0.2	
Episode 51	7.4300	Reward -0.2	
Episode 52	7.5300	Reward -0.2	
Episode 53	7.6300	Reward -0.2	
Episode 54	7.7300	Reward -0.2	
Episode 55	7.8300	Reward -0.2	
Episode 56	7.9300	Reward -0.2	
Episode 57	8.0300	Reward -0.2	
Episode 58	8.1300	Reward -0.2	
Episode 59	8.2300	Reward -0.2	
Episode 60	8.3300	Reward -0.2	
Episode 61	8.4300	Reward -0.2	
Episode 62	8.5300	Reward -0.2	
Episode 63	8.6300	Reward -0.2	
Episode 64	8.7300	Reward -0.2	
Episode 65	8.8300	Reward -0.2	
Episode 66	8.9300	Reward -0.2	
Episode 67	9.0300	Reward -0.2	
Episode 68	9.1300	Reward -0.2	
Episode 69	9.2300	Reward -0.2	
Episode 70	9.3300	Reward -0.2	
Episode 71	9.4300	Reward -0.2	
Episode 72	9.5300	Reward -0.2	
Episode 73	9.6300	Reward -0.2	
Episode 74	9.7300	Reward -0.2	
Episode 75	9.8300	Reward -0.2	
Episode 76	9.9300	Reward -0.2	
Episode 77	10.0300	Reward -0.2	
Episode 78	10.1300	Reward -0.2	
Episode 79	10.2300	Reward -0.2	
Episode 80	10.3300	Reward -0.2	
Episode 81	10.4300	Reward -0.2	
Episode 82	10.5300	Reward -0.2	
Episode 83	10.6300	Reward -0.2	
Episode 84	10.7300	Reward -0.2	
Episode 85	10.8300	Reward -0.2	
Episode 86	10.9300	Reward -0.2	
Episode 87	11.0300	Reward -0.2	
Episode 88	11.1300	Reward -0.2	
Episode 89	11.2300	Reward -0.2	
Episode 90	11.3300	Reward -0.2	
Episode 91	11.4300	Reward -0.2	
Episode 92	11.5300	Reward -0.2	
Episode 93	11.6300	Reward -0.2	
Episode 94	11.7300	Reward -0.2	
Episode 95	11.8300	Reward -0.2	
Episode 96	11.9300	Reward -0.2	
Episode 97	12.0300	Reward -0.2	
Episode 98	12.1300	Reward -0.2	
Episode 99	12.2300	Reward -0.2	
Episode 100	12.3300	Reward -0.2	
Average Reward	12.3300	Reward -0.2	

```
In [13]: print("Q-Table in the 100th Episode: \n", Total_Q_)
```

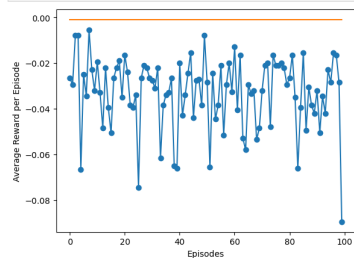
[illegible]

Semi-Gradient SARSA(0) vs Semi-Gradient TD(Lambda)

```
In [28]: #Plotting the average rewards per episode...
```

```
#Plotting the average rewards per episode...
x = [x for x in range(100)]
y1 = sarsa_arpe
y2 = td_arpe

plt.gca() # shorthand for "get current axis"
plt.plot(x, y1, marker = 'o')
plt.plot(x, y2)
plt.xlabel("Episodes")
plt.ylabel("Average Reward per Episode")
plt.show()
```



References

- https://www.gymnasium.dev/environments/classic_control/pendulum/
- <https://numpy.org/doc/stable/reference/>
- [https://www.learnataci.com/tutorials/reinforcement-q-learning-scratch-python-openai-gym/](https://www.learndatasci.com/tutorials/reinforcement-q-learning-scratch-python-openai-gym/)
- Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.
- Semi-Gradient SARSA: <https://web.stanford.edu/class/csych290/readings/SuttonBartoIPRLbook2nded.pdf> (pp. 152–154)
- Class Slides (Week 7, Weekof_07_FunctionApproximation1_Slides 9,14)
-
- <https://www.schools.com/python/matplotlib/markers.asp>