

# Assignment 1

## Reinforcement Learning Programming - CSCN 8020

July 05, 2024

Student Name: Chen, Kun

Student ID: 8977010

### Problem 1: Off-policy Monte Carlo [30]

#### Problem Statement

Off-policy Monte Carlo with Importance Sampling. The environment can be described as follows:

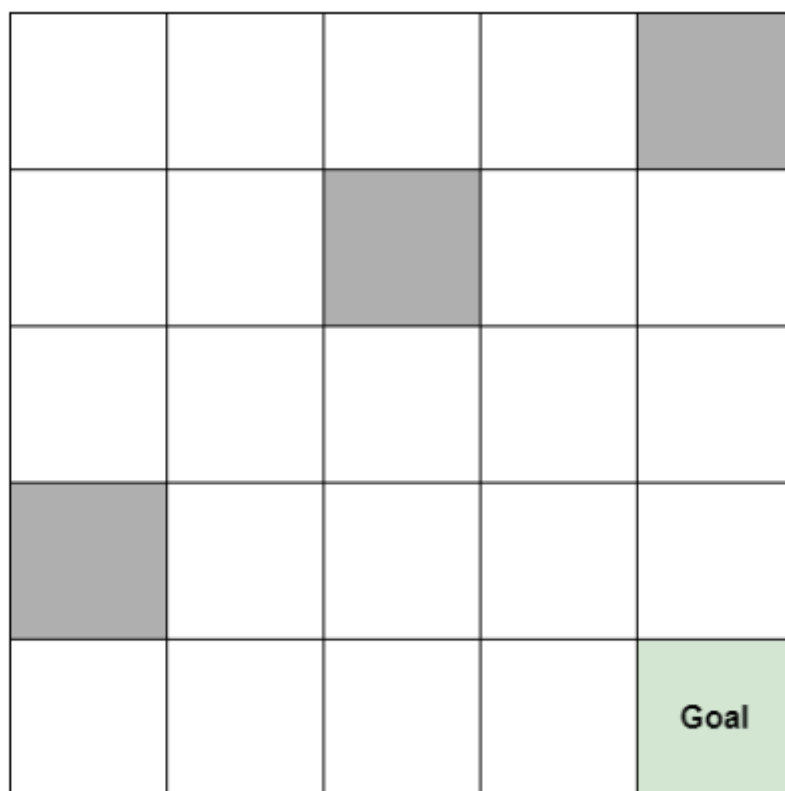


Figure 3: 5x5 Gridworld

- States: states are identified by their row and column, the same as a regular matrix. Ex: the state in row 0 and column 3 is  $s_{0,3}$  (Figure: 3)

Terminal/Goal state: The episode ends if the agent reached this state.  $s_{Goal} = s_{4,4}$

Grey states:  $\{s_{1,2}, s_{3,0}, s_{0,4}\}$ , these are valid but non-favourable states, as will be seen in the reward function.

- Actions:  $a_1$  = right,  $a_2$  = down,  $a_3$  = left,  $a_4$  = up for all states.
- Transitions: If an action is valid, the transition is deterministic, otherwise  $s' = s$
- Rewards  $R(s)$ :

$R(s) = +10$  if  $s = s_{4,4}$   
 $= -5$  if  $s \in S_{\text{grey}} = \{s_{0,4}, s_{1,2}, s_{3,0}\}$   
 $= -1$  otherwise

### Task

Implement the off-policy Monte Carlo with Importance sampling algorithm to estimate the value function for the given gridworld. Use a fixed behavior policy  $b(a|s)$  (e.g., a random policy) to generate episodes and a greedy target policy.

### Suggested steps

1. You can assume a specific discount factor (e.g.,  $\gamma = 0.9$ ) for this problem.
2. Use the same main algorithm implemented in lecture 4 in class.
3. Generate multiple episodes using the behavior policy  $b(a|s)$ .
4. For each episode, calculate the returns (sum of discounted rewards) for each state.
5. Use importance sampling to estimate the value function and update the target policy  $\pi(a|s)$

### Deliverables

- Full code with comments to explain key steps and calculations.
- Provide the estimated value function for each state.
- Important Compare the estimated value function obtained from Monte Carlo with the one obtained from Value Iteration in terms of optimization time, number of episodes, computational complexity, and any other aspects you notice.

### Answer:

#### 1. Define Gridworld Environment

States: Identify each state by its row and column.

Actions: Define actions (right, down, left, up).

Rewards: Define the rewards based on the given conditions.

#### 2. Define Off-Policy algorithm

**Input:** Arbitrary target policy  $\pi$

**Initialization:** For all  $s \in S$ ,  $a \in A(s)$ :

$Q(s,a) \leftarrow$  arbitrary value

$C(s,a) \leftarrow 0$

**Repeat forever:**

1.  $b \leftarrow$  behavior policy used to generate the episodes
2. Generate an episode using  $b$ :  $S_0, A_0, R_0, \dots, S_{T-1}, A_{T-1}, R_T, S_T$
3.  $G \leftarrow 0$

4.  $W \leftarrow 1$

5. **for**  $t = T-1, T-2, \dots, 0$ :

- $G \leftarrow \gamma G + R_{t+1}$
- $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
- $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$
- $W \leftarrow W \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$
- **if**  $W=0$  **then exit for loop**

Full code in mc\_frozen\_lake.py & mc\_maze\_agent.py

Provide the estimated value function for each state.

```
Estimated Action-Value Function (MC):
[[0.005281    0.00724844 0.0062977  0.00553023]
 [0.00584316 0.00965833 0.00617489 0.00649182]
 [0.00641625 0.         0.0147634  0.00565148]
 [0.00582435 0.04326686 0.         0.01343268]
 [0.         0.         0.         0.         ]
 [0.00741066 0.01077724 0.00907043 0.00548273]
 [0.00721103 0.02786813 0.         0.00610703]
 [0.         0.         0.         0.         ]
 [0.         0.1163666  0.05391322 0.01215465]
 [0.03588563 0.16461341 0.05526832 0.         ]
 [0.0099385  0.         0.02829157 0.00798514]
 [0.00994807 0.04778915 0.05711794 0.00855135]
 [0.02990819 0.11763607 0.1053016  0.         ]
 [0.05539837 0.23926065 0.16296364 0.03839555]
 [0.10813902 0.39805207 0.16843513 0.06542689]
 [0.         0.         0.         0.         ]
 [0.         0.07176263 0.11132978 0.0224493  ]
 [0.04861111 0.16833181 0.24397815 0.05353179]
 [0.11687859 0.41166927 0.40017191 0.11821972]
 [0.22700763 1.         0.40218119 0.16749923]
 [0.02869191 0.0245522  0.06685414 0.         ]
 [0.02669659 0.07084661 0.17368334 0.04009628]
 [0.0675012  0.1715602  0.3967523  0.11859705]
 [0.17171634 0.41023879 1.         0.22349161]
 [0.         0.         0.         0.         ]]
```

```
Target Policy (MC):
[['|' '|' '|' '>' '|' '<']
 ['|' '|' '|' '<' '|' '|' '|']
 ['>' '>' '|' '|' '|' '|']
 ['<' '>' '>' '|' '|' '|']
 ['>' '>' '>' '>' '<']]
```

```
Estimated Action-Value Function (Value Iteration):
[[0.43046721 0.4782969 0.4782969 0.43046721]
 [0.43046721 0.531441 0.531441 0.4782969 ]
 [0.4782969 0. 0.59049 0.531441 ]
 [0.531441 0.6561 0. 0.59049 ]
 [0. 0. 0. 0. ]
 [0.4782969 0.531441 0.531441 0.43046721]
 [0.4782969 0.59049 0. 0.4782969 ]
 [0. 0. 0. 0. ]
 [0. 0.729 0.729 0.59049 ]
 [0.6561 0.81 0.729 0. ]
 [0.531441 0. 0.59049 0.4782969 ]
 [0.531441 0.6561 0.6561 0.531441 ]
 [0.59049 0.729 0.729 0. ]
 [0.6561 0.81 0.81 0.6561 ]
 [0.729 0.9 0.81 0.729 ]
 [0. 0. 0. 0. ]
 [0. 0.729 0.729 0.59049 ]
 [0.6561 0.81 0.81 0.6561 ]
 [0.729 0.9 0.9 0.729 ]
 [0.81 1. 0.9 0.81 ]
 [0.6561 0.6561 0.729 0. ]
 [0.6561 0.729 0.81 0.6561 ]
 [0.729 0.81 0.9 0.729 ]
 [0.81 0.9 1. 0.81 ]
 [0. 0. 0. 0. ]]
```

```
Target Policy (Value Iteration):
[['|' '|' '|' '>' '|' '<']
 ['|' '|' '|' '<' '|' '|' '|']
 ['>' '|' '|' '|' '|' '|']
 ['<' '|' '|' '|' '|' '|']
 ['>' '>' '>' '>' '<']]
```

Compare the estimated value functions obtained from Monte Carlo (MC) and Value Iteration (VI)

#### 1. Optimization time:

- MC: Generally slower, as it requires sampling many episodes to converge.
- VI: Typically faster, as it directly computes the value function without sampling.

#### 2. Number of episodes:

- MC: Requires multiple episodes to estimate the value function. The number isn't explicitly stated, but it's likely in the hundreds or thousands.
- VI: Doesn't use episodes. It iteratively updates the value function until convergence.

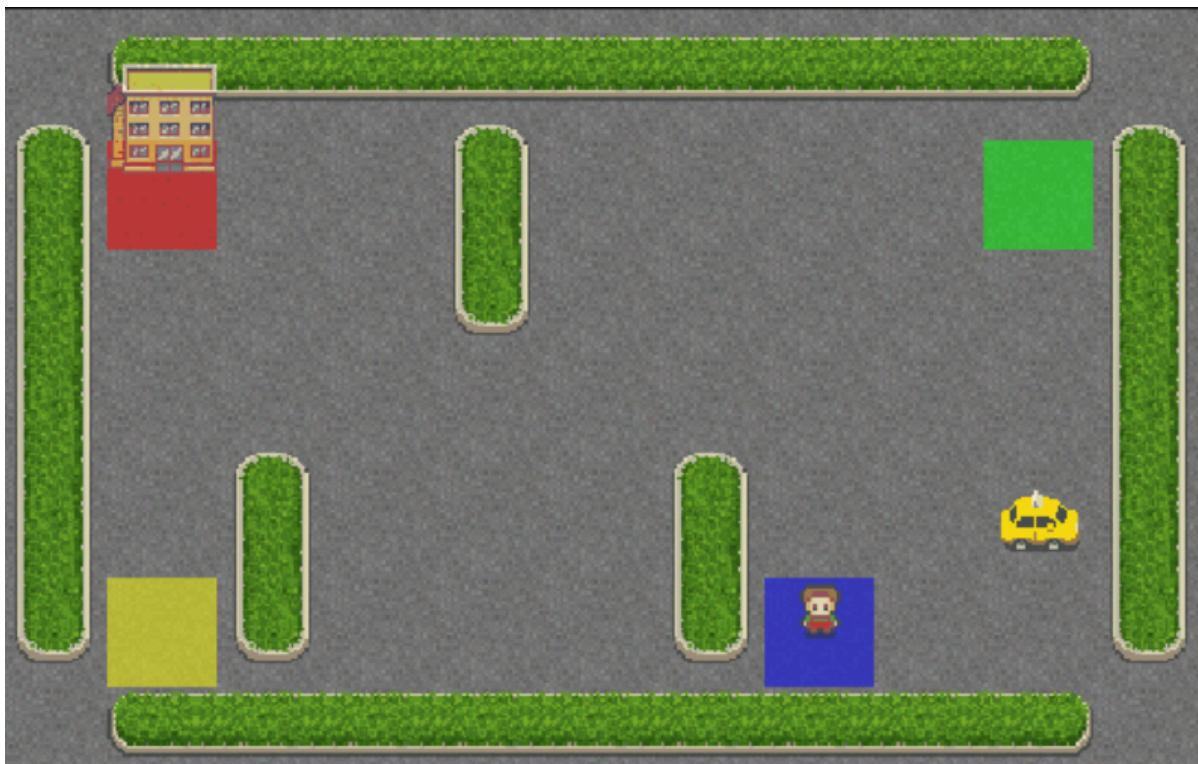
### 3. Computational complexity:

- MC:  $O(\text{number of episodes} * \text{average episode length})$
- VI:  $O(\text{number of iterations} * \text{number of states} * \text{number of actions})$

## Problem 2: Q-Learning [70]

### Introduction

You will work with the Taxi environment and implement QLearning. The way you interact with the environment will be very similar to the gym environments described in class. Therefore, most of the code we discussed is directly applicable.



The environment has 6 discrete actions, 500 discrete states, and the following rewards:

- -1 per step unless other reward is triggered.
- +20 delivering passenger.
- -10 executing "pickup" and "drop-off" actions illegally.

### Action space

- 0: Move south (down)
- 1: Move north (up)
- 2: Move east (right)
- 3: Move west (left)
- 4: Pickup passenger
- 5: Drop off passenger

### Destinations:

- 0: Red
- 1: Green
- 2: Yellow
- 3: Blue

An observation is returned as an `int()` that encodes the corresponding state, calculated by  $((\text{taxi row} * 5 + \text{taxi col}) * 5 + \text{passenger location}) * 4 + \text{destination}$

### Helper Utility

To assist you in understanding the environment further, you were provided with a file (assignment2 utils.py) that has a few methods to allow loading the environment and printing some basic information about it. It also allows you to switch to the detailed observation description using the state scalar value.

### Problem Statement

Implement the Q-Learning algorithm on the Taxi environment from OpenAI Gym. Train an agent to efficiently navigate and pick up passengers. Use the following hyperparameters:

- Learning Rate  $\alpha$ : 0.1
- Exploration Factor  $\epsilon$ : 0.1
- Discount Factor  $\gamma$ : 0.9

### Tasks

- Implement the Q-Learning algorithm and train an agent on the Taxi environment.
- Report the following metrics after training:

1. Total episodes
2. Total steps taken per episode
3. Return per episode
4. Average cumulative reward per episode (you can smooth your plot by using average cumulative reward every  $n$  episodes, keeping  $n$  constant across your runs).

Note: Average cumulative reward =  $\frac{1}{n} \sum_{i=1}^n R_i$

- Make a deliberate change to the following parameters (separately) and use each value once.
  - Learning Rate  $\alpha$  = [0.01, 0.001, 0.2]
  - Exploration Factor  $\epsilon$  = [0.2, 0.3]

### Deliverables

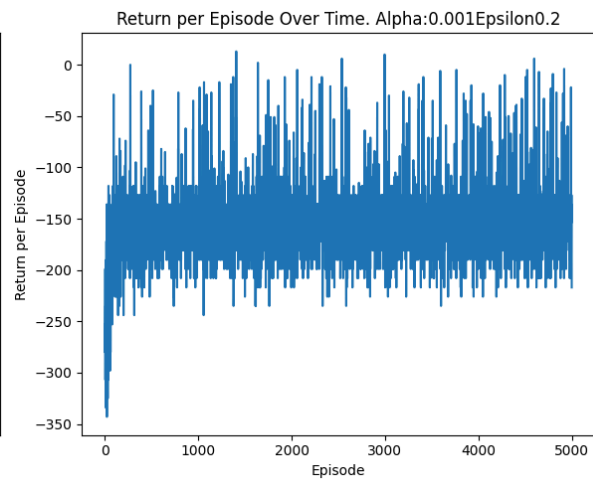
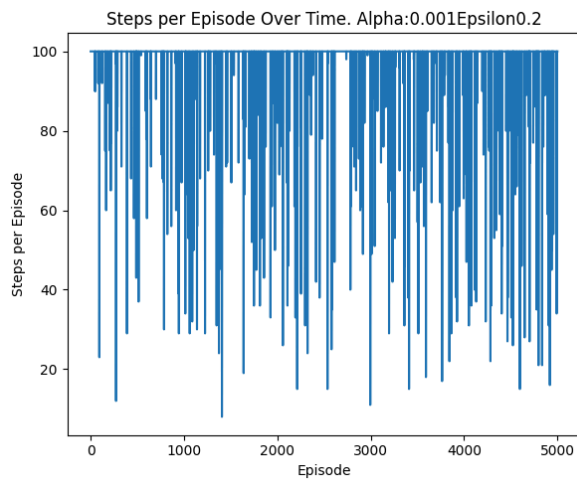
- Python code implementing Q-Learning and running it for the different hyperparameters. [20]
- A report (in PDF) containing the metrics, observations, and comments on the parameter change. Make sure you include appropriate plots for the previously stated metrics to support your arguments. [30]

- Based on your findings, choose what you think would be the best combination of learning rate and exploration factor and re-run the training. Report and comment on the differences observed. [20]

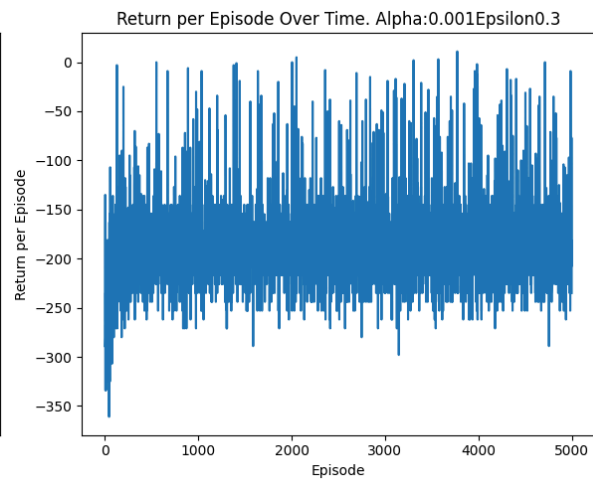
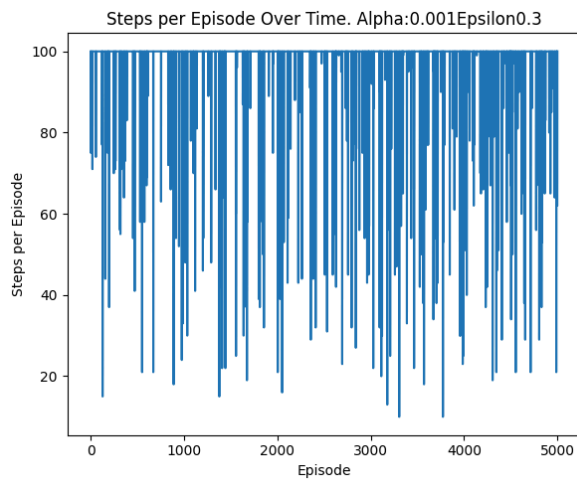
**Answer:**

1. Python code: in assignment2\_utils.py
2. metrics:

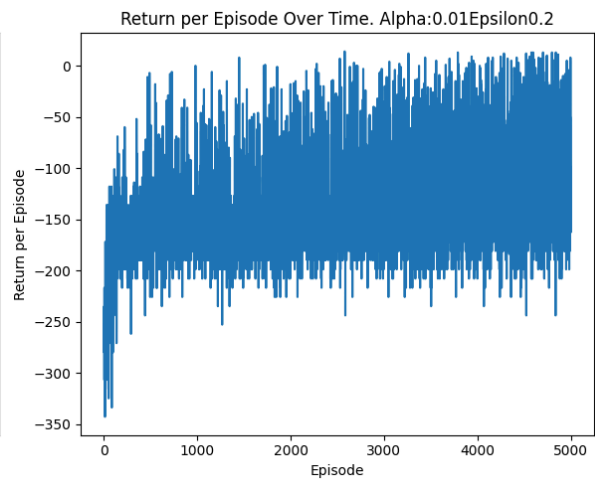
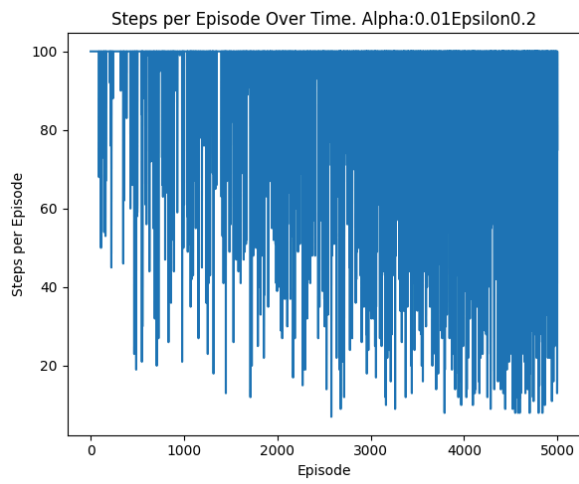
Alpha 0.001 Epsilon: 0.2



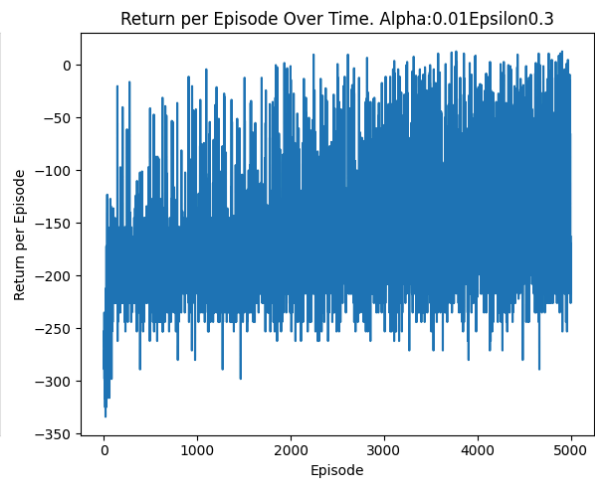
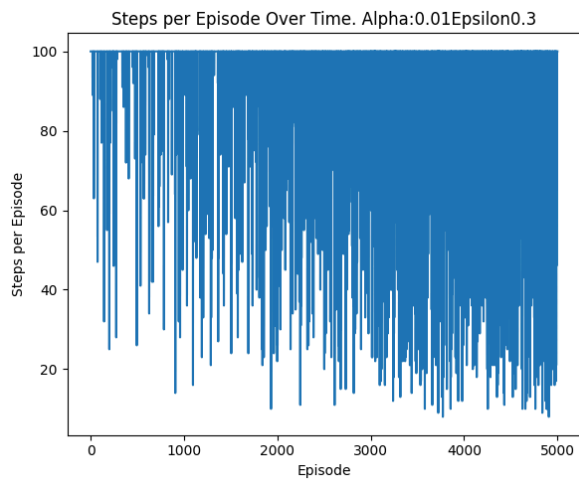
Alpha 0.001 Epsilon: 0.3



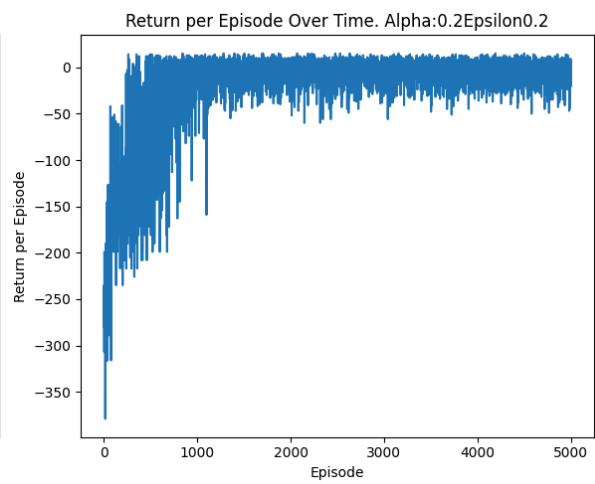
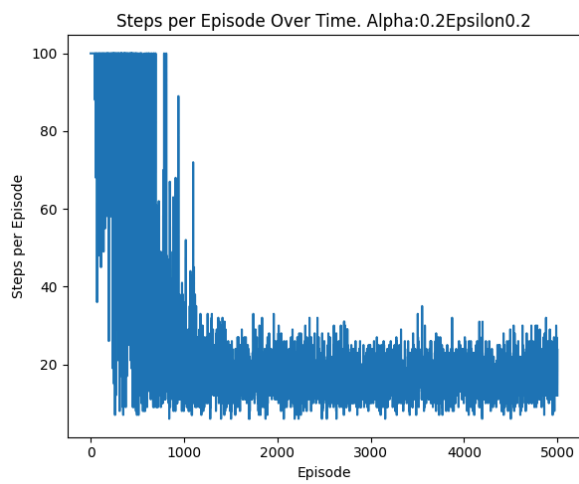
Alpha 0.01 Epsilon: 0.2



Alpha 0.01 Epsilon: 0.3

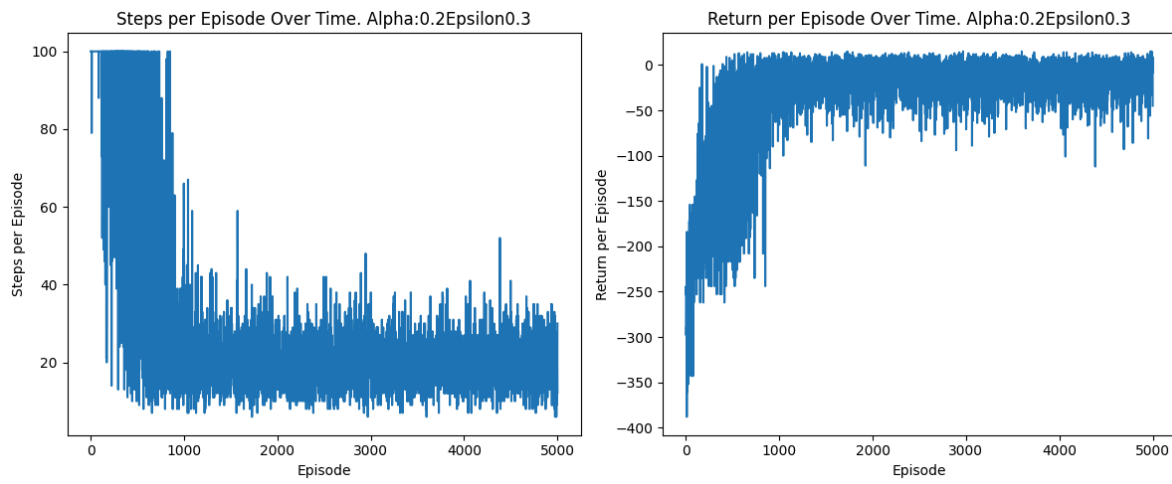


Alpha 0.2 Epsilon: 0.2



Alpha 0.2 Epsilon: 0.3





Looks like Alpha 0.2 Epsilon: 0.2 is the best parameter, fast convergence. Reducing the total steps per episode and increasing the return per episode.

By choosing the best learning rate and exploration factor based on average cumulative reward, the agent's performance should improve, resulting in more efficient and effective navigation and passenger pick-up/drop-off.