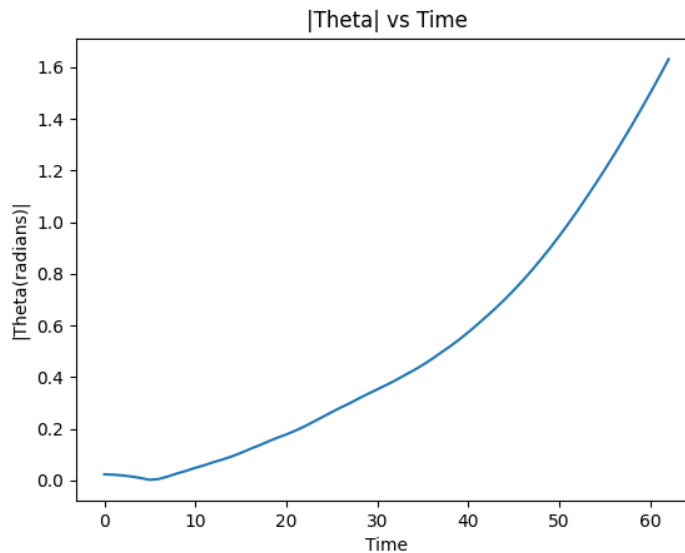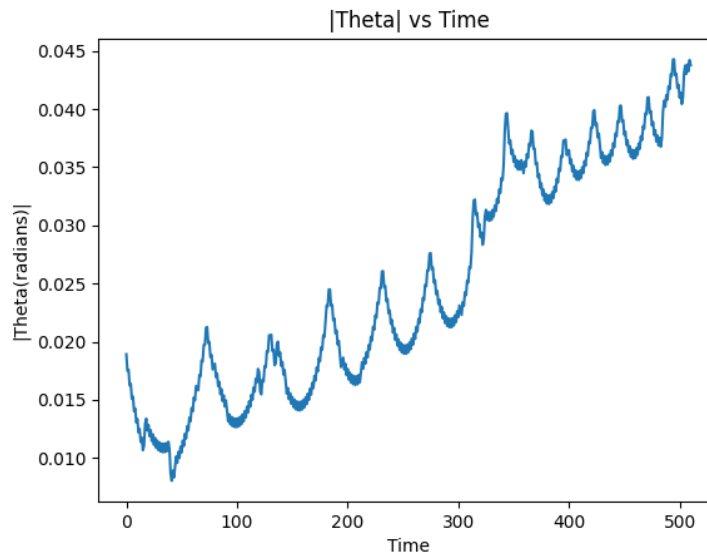# Homework 3 Part-1 Report

Jiahang Wang(261011319)

## Part A



theta vs time(initial simulation)



theta vs time(final simulation)

### Initial Simulation (Top Graph)

- In the initial simulation phase, the pendulum angle seems to deviate rapidly from the equilibrium position (θ=0) as time progresses. This suggests that the control system is not effectively stabilizing the pendulum, leading to an increasing magnitude of oscillation that likely results in system failure if the trend continues.

## Final Simulation (Bottom Graph)

- Compared to the initial simulation, this graph indicates that in the final simulation, the pendulum control appears to be more stable, with the angle variations contained within a relatively narrow band. This implies that the control algorithm might have learned how to maintain stability, although there are still fluctuations, the system does not go out of control as in the initial simulation.

In RL and Q-learning, random actions are crucial for exploring unknown strategies (exploration) and choosing the best-known actions for maximum reward (exploitation). Balancing these two aspects is key, with strategies like the epsilon-greedy method helping to find the right mix. As the algorithm learns, the reliance on random actions decreases, shifting towards more exploitation of learned knowledge.

# Part B

- Impose a negative reward, such as -1 or a larger negative number, when the car runs off-screen. The penalty should be significant enough to encourage the algorithm to avoid this behavior.

- Consider including the car's deviation from the screen's center as part of the penalty, increasing the negative reward as the car gets closer to the edge.

# Part C



final state value matrix

- Most of the matrix contains zeros, but there is a region with non-zero values. This suggests that the learning process has identified certain states as having higher expected returns than others.

# Part D

- The Discount Factor balances the importance of immediate versus future rewards, larger values favoring long-term strategies and smaller values focusing on immediate gains. It also ensures learning stability by preventing infinite reward accumulation.

- The Learning Rate determines how quickly new information overrides old, influencing the algorithm's adaptability and stability. A higher Learning Rate means new information has a greater impact on updating old knowledge. A low Learning Rate means the algorithm updates its knowledge more cautiously and might need more time to adapt to environmental changes.

# Other challenges

- Parameter Tuning: Identifying the optimal combination of parameters such as the Learning Rate, Discount Factor, and exploration strategies in Reinforcement Learning task like this requires extensive experimentation and adjustments.

- Dealing with the dynamics and uncertainties of the environment, such as noise. making it harder or longer for the algorithm to learn and make accurate predictions. This is crucial for real-world applications where conditions are rarely static or completely predictable.