



# Adaptive Dynamic Stabilization of a Self-Balancing Scooter Using Epsilon-Greedy Q-learning

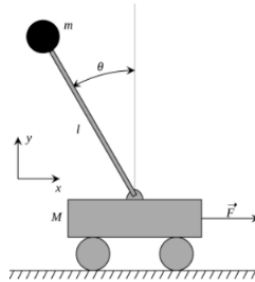
Hongyi Yang, Muhan Lin, Xuanyang Xu, Zixing Jiang

## 1 Introduction

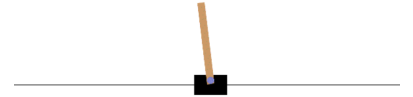
The self-balancing scooter (Figure 1a) is a common urban commuter tool. One of the most critical features in self-balancing scooter is dynamic stabilization, which allows the scooter to maintain its balance while in motion. Without this feature, it would be very difficult for riders to stay upright and control the scooter, and even cause safety accidents in severe cases. In order to achieve dynamic stabilization, the traditional automatic control methods, including Proportional-Integral-Derivative (PID) control and Model Predictive Control (MPC), model the self-balancing scooter as a cart-pole system (Figure 1b) and designs the controller by deriving its system dynamics. However, the dynamics-based control is sensitive to the rider's weight and center of mass (CoM). For riders of different weights and CoMs, the performance of stabilization varies. In this project, we aim to utilize the learning ability of artificial intelligence to design a dynamic stabilizer that adapts to different rider weights and CoMs.



(a) Man on a self-balancing scooter



(b) Cart-pole system dynamics



(c) OpenAI Cart Pole Environment

Figure 1: Illustrations of (a) the self-balancing scooter, (b) the cart-pole system dynamics, and (c) the OpenAI Cart Pole environment.

## 2 Problem Statement

In this project, we model the scooter as the cart-pole system (Figure 1b). The wheel of scooter is the cart the rider-scooter mass is the pole. Then the adaptive stabilization problem can be formulated as follows: control the inclination ( $\theta$ ) of a pole with **varying mass and CoM** within a certain range by controlling the movement ( $x$ ) of the cart (leftward or rightward). We intend to apply reinforcement learning (RL) to solve this problem. Namely, given the observation of the environment ( $\theta, \dot{\theta}, x, \dot{x}$ ), we want our agent (cart) to learn a policy (move leftward or move rightward) to stable the varying pole through trail and reward. We plan to implement our RL algorithm based on the [OpenAI Gym environment for cart-pole](#) (Figure 1c).

## 3 Epsilon-Greedy Q-Learning

In this project, we built our dynamic stabilization controller with Epsilon-Greedy Q-Learning, which is a well-developed reinforcement learning algorithm that balances exploration and exploitation in sequential decision-making problems. In this section, we will introduce the Epsilon-Greedy Q-Learning algorithm in detail and explain how we link this algorithm to the dynamic balancing problem we aim to solve.

States \ Actions	$\mathbf{a}_0$	$\mathbf{a}_1$	$\mathbf{a}_2$	$\dots$
$\mathbf{s}_0$	$Q(\mathbf{s}_0, \mathbf{a}_0)$	$Q(\mathbf{s}_0, \mathbf{a}_1)$	$Q(\mathbf{s}_0, \mathbf{a}_2)$	$\dots$
$\mathbf{s}_1$	$Q(\mathbf{s}_1, \mathbf{a}_0)$	$Q(\mathbf{s}_1, \mathbf{a}_1)$	$Q(\mathbf{s}_1, \mathbf{a}_2)$	$\dots$
$\mathbf{s}_2$	$Q(\mathbf{s}_2, \mathbf{a}_0)$	$Q(\mathbf{s}_2, \mathbf{a}_1)$	$Q(\mathbf{s}_2, \mathbf{a}_2)$	$\dots$
$\mathbf{s}_3$	$Q(\mathbf{s}_3, \mathbf{a}_0)$	$Q(\mathbf{s}_3, \mathbf{a}_1)$	$Q(\mathbf{s}_3, \mathbf{a}_2)$	$\dots$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$

Table 1: An example of Q-Table.  $Q(\mathbf{s}_i, \mathbf{a}_j)$  denotes the Q-Value of a state-action pair  $(\mathbf{s}_i, \mathbf{a}_j)$ . The higher  $Q(\mathbf{s}_i, \mathbf{a}_j)$  is, the more worthy to take action  $\mathbf{a}_j$  at state  $\mathbf{s}_i$ .

### 3.1 Q-Learning

Q-Learning is a classic value-based reinforcement learning algorithm that allow an agent to learn from its own actions and rewards in an environment. Unlike some other model-based control and reinforcement learning methods, it does not require a model of the environment. This feature helps to solve the traditional model-based control method's dependence on the rider's center of gravity mentioned in previous sections.

The goal of Q-Learning is to find an optimal policy that maximizes the expected value of the total reward over any sequence of actions. The main idea of Q-Learning is to use a table, called Q-Table, that stores the estimated value  $Q(\mathbf{s}, \mathbf{a})$  (Q-Value) of each action  $\mathbf{a}$  in each state  $\mathbf{s}$ , as shown in Table 1. The Q-Value, which reflects the value of taking a certain action in a certain state, is defined as the expectation of overall received rewards after taking action  $\mathbf{a}$  at state  $\mathbf{s}$ , as shown in (1).

$$Q(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left( \sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right) \quad (1)$$

where  $r(\mathbf{s}_t, \mathbf{a}_t)$  is the reward received from the environment for state-action pair  $(\mathbf{s}_t, \mathbf{a}_t)$  and  $\gamma$  is the discount factor.

In Q-Learning, in order to obtain the most accurate estimate of the Q-Value, the Bellman Equation is applied to update  $Q(\mathbf{s}, \mathbf{a})$  in every iteration that state-action pair  $(\mathbf{s}, \mathbf{a})$  is chosen, as shown in (2)

$$\text{New } Q(\mathbf{s}, \mathbf{a}) \leftarrow (1 - \lambda) \text{Current } Q(\mathbf{s}, \mathbf{a}) + \lambda \left[ r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}') \right] \quad (2)$$

where  $r(\mathbf{s}, \mathbf{a})$  is the reward received in this iteration,  $\mathbf{s}'$  is the transient state after taking  $\mathbf{a}$  at  $\mathbf{s}$ ,  $\mathbf{a}'$  is the legal action at  $\mathbf{s}'$ ,  $\gamma$  is the discount factor, and  $\lambda$  is the learning rate.

In the Q-Learning process, the agent updates the Q-Table as it explores the environment and receives feedback. The agent chooses the action  $\mathbf{a}^*(\mathbf{s})$  that has the highest value in the current state  $\mathbf{s}$ , according to the Q-Table, as shown in (3).

$$\mathbf{a}^*(\mathbf{s}) = \arg \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}) \quad (3)$$

This way, the agent learns to act optimally over time.

In our scenario, the state vector  $\mathbf{s}$  of the cart-pole system is continuous. In order to prevent infinitely many state-action pairs, we discretized the system's state space by intervals. The length of each interval is fine-tuned to ensure the effectiveness of learning.

### 3.2 Epsilon-Greedy Action Selection

In order to learn an optimal Q-table, we want our algorithm to explore all possible state-action pairs. Thus, in addition to the action selection criterion (3), an epsilon-greedy policy is applied. Given an epsilon-greedy factor  $0 \leq \varepsilon \leq 1$ , in each iteration, there is a probability of  $\varepsilon$  that the agent takes random actions to explore new possible state-action pairs, and there is a probability of  $1 - \varepsilon$  that the agent exploits investigated state-action pairs via the Q-Table. The overall logic flow of the Epsilon-Greedy Q-Learning is shown in algorithm 1.

---

**Algorithm 1** Epsilon-Greedy Q-Learning

---

**Require:** observation of state  $\mathbf{s}$ , epsilon-greedy factor  $\varepsilon \in [0, 1]$

**while** learning not terminate **do**

    Generate a random number  $i \in [0, 1]$

**if**  $i \leq \varepsilon$  **then**

        Take random action  $\mathbf{a}$

**else**

        Select action  $\mathbf{a}$  according to the Q-Table

▷ Equation (3)

**end if**

    Update Q-Table entry for  $(\mathbf{s}, \mathbf{a})$  according to the received reward  $r(\mathbf{s}, \mathbf{a})$

▷ Equation (2)

**end while**

---