

Reinforcement learning based control of an underactuated double pendulum system

Master's Thesis Nr. xxx

Scientific Thesis for Acquiring the Master of Science Degree
at the School of Engineering and Design
of the Technical University of Munich.

Thesis Advisor Laboratory for Product Development and Lightweight Design
Prof. Dr. Markus Zimmermann

Supervisor Laboratory for Product Development and Lightweight Design
Akhil Sathuluri
Hans Zweitkorrektor (Second corrector)

Submitted by Chi Zhang
Karl Köglspurger Straße 9, 80939, München
Matriculation number: 03735807
chi97.zhang@mytum.de

Submitted on Garching, 15.11.2023

Declaration

I assure that I have written this work autonomously and with the aid of no other than the sources and additives indicated.

Garching, 15.11.2023

Chi Zhang

Project Definition (1/2)

Initial Situation

Add your Project Brief here. If you don't need it, comment out the creation of this Project Brief in the main document `Thesis.tex`.

Goals

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Project Definition (2/2)

Contents of this Thesis

- Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
- Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
- Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
- Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet
 - Lorem ipsum dolor sit amet

Project Note

Master's Thesis	Nr. xxx
Supervisor	Akhil Sathuluri
Partners in industry/research	DFKI GmbH, Robotics Innovation Center
Time period	15.05.2023 - 15.11.2023

The dissertation project of Akhil Sathuluri set the context for the work presented. My supervisor Akhil Sathuluri mentored me during the compilation of the work and gave continuous input. We exchanged and coordinated approaches and results monthly.

An accurate elaboration, a comprehensible and complete documentation of all steps and applied methods, and a good collaboration with industrial partners are of particular importance.

Publication

I consent to the laboratory and its staff members using content from my thesis for publications, project reports, lectures, seminars, dissertations and postdoctoral lecture qualifications.

The work remains a property of the Laboratory for Product Development and Lightweight Design.

Garching, 15.11.2023

Chi Zhang

Akhil Sathuluri

Contents

1	Introduction	1
1.1	Motivation.....	1
1.1.1	Trajectory planning and tracking.....	3
1.1.2	Reinforcement learning based control.....	4
1.2	Problem setup.....	5
1.3	Contribution	6
1.4	Content	7
2	State of the art	9
2.1	Theory	9
2.1.1	Underactuated system.....	9
2.1.2	Dynamics of underactuated double pendulum system.....	10
2.2	Related work	13
2.2.1	Feedback cancellation	13
2.2.2	PFL:Partial Feedback Linearization	14
2.2.3	Energy shaping	14
2.2.4	Reinforcement-Learning based control	14
3	Methodology	15
3.1	Soft actor critic	15
3.2	Linear quadratic regulator	17
3.3	Combining SAC and LQR with region of attraction.....	18
3.4	Reward shaping	18
4	Experiment: training and simulation	21
4.1	Training setup.....	21
4.2	Training process.....	21
4.3	Simulation results	23
4.3.1	pendubot.....	23
4.3.2	acrobot	23
5	Experiment: hardware system	25
5.1	Hardware setup	25
5.2	system identification.....	26
5.3	sim2real problem	26
5.4	real hardware results.....	26

6 Discussion	29
6.1 introduction to leaderboard results.....	29
6.2 interpretation of simulation results	30
6.3 interpretation of hardware results	30
6.4 future work.....	31
Appendix	33
Appendix A An appendix.....	34

1 Introduction

Nonlinear systems, as their name suggests, don't adhere to linear relationships between inputs and outputs. This lack of linearity means the system's response to input changes is intricate and often unpredictable. In the real world, nearly all systems exhibit some form of nonlinearity. This nonlinearity manifests in various phenomena. For instance, in systems with multiple inputs and multiple outputs, interdependencies between variables become complex, posing a coupling problem. Chaotic behavior, often referred to as the butterfly effect, is another common issue. Even a slight alteration in initial conditions can drastically alter the system's outcome. The classic butterfly effect example illustrates this sensitivity — a butterfly's wings in Brazil potentially triggering a tornado in Texas. When dealing with control tasks, especially in many real-world systems, accounting for this nonlinearity is essential.

Gaining proficient control over nonlinear systems has been a primary objective in the field of control theory for a substantial period. Throughout history, control engineers have developed a diverse range of approaches to handle these intricate systems. The advent of robotics in recent times has brought forth new methodologies specifically designed to tackle the challenges posed by nonlinearities.

1.1 Motivation

Robots, which are programmable mechanical entities, are purposefully designed for autonomous or semi-autonomous task execution, showcasing mobility, manipulation, and interaction with their surroundings.

In the realm of modern robotics, these machines exemplify intricate, highly nonlinear mechanical systems. Some well-known instances include quadruped robotics, autonomous vehicles, quadcopters, and humanoid robots.

1 Introduction



Figure 1.1: Caption for Image 1



Figure 1.2: Caption for Image 2



Figure 1.3: Caption for Image 3



Figure 1.4: Caption for Image 4

1.1.1 Trajectory planning and tracking

To enhance their motion capabilities, reliable nonlinear control methods are essential. Traditionally, for planning complex motion in systems like those mentioned above, a two-step approach is followed, namely trajectory planning and trajectory tracking.

Trajectory planning involves computing a smooth and feasible path that a robot should follow to reach a specific target position or work point. A common method employed is trajectory optimization, aiming to minimize a cost function that factors in travel time, energy consumption, and smoothness of motion. Techniques such as gradient descent or genetic algorithms are often used for this optimization. Adhering to constraints such as maximum velocity and accelerations is crucial during this process.

Once the trajectory is successfully planned, trajectory tracking involves implementing control algorithms to guide the robot along the planned trajectory. Typically, a feedback control approach is utilized, continuously monitoring the robot's position and adjusting control inputs. However, in real-world systems, even with a well-planned trajectory, external disturbances, uncertainties, or system limitations may cause significant deviations. Hence, accurate state estimation and robust control are crucial in the realm of feedback control.

[maybe there a better example is to explain how atlas works] In the domain of industrial robot control, it's typical to separate the process into distinct planning and execution phases. This division arises from the fact that real-time responsiveness is not a stringent necessity in this context. When a task is defined, the planning phase kicks in, utilizing algorithms like linear interpolation and A* for trajectory generation. Following this, a steady and reliable control policy is implemented to ensure precise tracking of the generated trajectory during the execution phase.

Yet, in dynamic domains like automotive and flight control, the demand for real-time responsiveness takes center stage. Model Predictive Control (MPC) stands as a prime illustration of this critical requirement. MPC seamlessly integrates trajectory planning and execution into a unified framework. The process commences by projecting a sequence of control actions into the future as part of the planning stage. Subsequently, the calculated control inputs are meticulously fine-tuned to minimize the deviation between the actual system state and the planned trajectory. This strategic implementation effectively guides the system to closely track the intended trajectory.

MPC's brilliance lies in its ability to concurrently devise and optimize trajectories while swiftly adapting in real time to stay closely aligned with the planned trajectory, even when facing

various disturbances and uncertainties. This amalgamation plays a pivotal role in achieving precise control and adaptability, especially in rapidly changing and intricate environments.

Though the traditional approach of trajectory generation and tracking has proven effective for many systems, it has the following drawbacks:

- **Limited Adaptability:** Trajectory planning typically relies on predefined paths or trajectories, limiting adaptability to unforeseen changes or dynamic environments. If the environment changes significantly, the planned trajectory may no longer be optimal or even feasible.
- **Difficulty in Complex Environments:** In highly complex and cluttered environments, planning a feasible trajectory that avoids obstacles while reaching the goal can be challenging. The complexity increases with the number of obstacles and the intricacy of the environment.
- **Difficulty with Nonlinear Systems:** Trajectory planning struggles with highly nonlinear systems where the dynamics are hard to model accurately. Linearizing the system for planning purposes may lead to suboptimal or infeasible trajectories.
- **Static Planning:** Traditional trajectory planning is often static, assuming a stationary environment. It does not readily adapt to changing circumstances or dynamic obstacles, which limits its applicability in real-world scenarios.
- **High Computational Demands:** Some trajectory planning algorithms can be computationally intensive, especially for high-dimensional or complex robotic systems. This computational demand becomes a drawback, particularly in real-time or time-critical applications.

1.1.2 Reinforcement learning based control

In contrast to trajectory-based control, reinforcement learning (RL)-based control extracts an optimal policy through interactions with the environment, offering several advantages:

- **Adaptability and Flexibility:** RL enables systems to adapt and learn optimal behavior in environments that are dynamic and changing. The control policy can continuously evolve based on new experiences and acquired knowledge, ensuring adaptability to varying circumstances.
- **Less Model Information Required:** In contrast to traditional control methods that often necessitate a precise mathematical model of the system, RL can directly learn from inter-

actions with the environment without relying on an explicit model. This characteristic is particularly valuable in scenarios where system dynamics are complex or unknown.

- **Effective Handling of Nonlinearities and Complex Systems:** RL proves highly effective in dealing with highly nonlinear systems and complex control tasks that might pose challenges for traditional control methods. The use of neural network-based function approximation allows for the capture of intricate relationships between states and actions.
- **Efficient Handling of High-Dimensional Input Spaces:** RL demonstrates an ability to efficiently manage high-dimensional and continuous input spaces, a crucial feature in many real-world applications such as robotics, finance, and game playing.

Reinforcement learning, by learning directly from the environment, offers a dynamic and adaptable approach to control, making it particularly suitable for complex and nonlinear systems.

1.2 Problem setup

Within the realm of nonlinear systems, a particularly challenging class is underactuated systems. These systems are characterized by having fewer control inputs than degrees of freedom. This makes them notably harder to control compared to fully actuated systems. Interestingly, a majority of robots and even living beings in nature fall into the category of underactuated systems. Consequently, studying the control of underactuated mechanical systems holds significant universal relevance.

The double pendulum, a simple setup comprising two links connected by two rotational joints, is a prime example. The joints involved are the shoulder joint, directly connected to the world frame, and the elbow joint, situated between the two links. The end effector is positioned at the tip of the second link. Active control is achieved by attaching motors to the shoulder and elbow joints. In the domain of underactuated control, if the shoulder joint is actuated, the setup is known as a pendubot. On the other hand, if the elbow joint is actuated, it's referred to as an acrobot.

Despite its simple configuration, the system exhibits highly nonlinear and chaotic behavior. The double pendulum setup presents two classic tasks: swing-up and stabilization around the highest point. Research on swing-up and stabilization of the double pendulum can be traced back to the 1990s, and it continues to be a crucial testbed for validating the effectiveness of newly designed control algorithms.

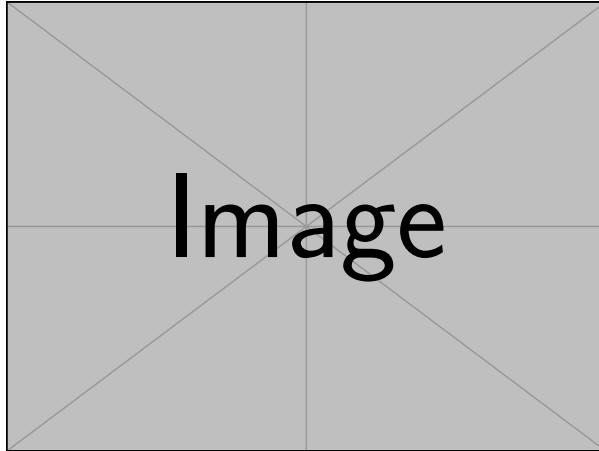


Figure 1.5: This is a sample image.

Our project's motivation is to develop a reinforcement learning-based control method suitable for underactuated control of the double pendulum system, specifically addressing swing-up and stabilization tasks. To evaluate the efficacy of this control method, we conduct both simulations and real system experiments.

1.3 Contribution

In this paper, our main contribution is as follows: developing an effective control strategy to achieve two key objectives with the double pendulum. The first task involves swinging the double pendulum from its lowest point to its highest point. The second task is to maintain stability at the highest point.

To tackle the swing-up task, we utilized a well-known model-free reinforcement learning algorithm called soft actor-critic. This algorithm allowed us to train a policy capable of reaching the region of attraction (RoA) of a continuous-time linear quadratic regulator (LQR) controller. Once the system enters the RoA, we seamlessly transition to the LQR controller to maintain stability around the highest point.

1.4 Content

The paper is structured as follows:

- **Chapter 2: State-of-the-Art**
 - This chapter provides an overview of current advancements in the field, summarizing essential theories, including those related to nonlinear and underactuated control.
- **Chapter 3: Methodology**
 - This chapter delves into the methodology, encompassing fundamental aspects of reinforcement learning, with a specific focus on the SAC algorithm. It explains the reward function used for training, the training procedure, and covers the concept of the LQR controller and how the combined controller was composed.
- **Chapter 4: Simulation Results**
 - In this chapter, we present the results obtained from simulations, showcasing the performance and behavior of the designed control strategy.
- **Chapter 5: Hardware Results**
 - This chapter reports the outcomes of experiments conducted on the hardware, providing insights into how we addressed the sim2real transfer problem.
- **Chapter 6: Discussion and Future Work**
 - The final chapter engages in a discussion about the obtained results and explores potential future directions for research and development.

2 State of the art

This chapter is about the state of the art.

2.1 Theory

This section is about the theory, for example dynamics and underactuated control.

2.1.1 Underactuated system

According to newtons second law($F = ma$), The dynamics of mechanical systems can be described as follows:

$$\ddot{q} = f(q, \dot{q}, u, t) \quad (2.1)$$

Where the state is given by a vector of positions(also known as the configuration vector), and a vector of velocities, \dot{q} .

For control, the second order differential equation can be rewritten as below:

$$\ddot{q} = f_1(q, \dot{q}, t) + f_2(q, \dot{q}, t)u \quad (2.2)$$

For a controlled dynamical system described by equation(2.2), if we have:

$$\text{rank}[f_2(q, \dot{q}, t)] < \dim[q] \quad (2.3)$$

then the system is underactuated at (q, \dot{q}, t) . There is another case for underactuation is even when f_2 is full rank, but additional constraints like $|\mathbf{u}| \leq 1$ can also make a system underactuated.

2.1.2 Dynamics of underactuated double pendulum system

As shown in figure 2.1, we model the dynamics of the double pendulum with 15 parameters which include 8 link parameters namely masses (m_1, m_2) , lengths (l_1, l_2) , center of masses (r_1, r_2) , inertias (I_1, I_2) for the two links, and 6 actuator parameters namely motor inertia I_r , gear ratio g_r , coulomb friction (c_{f1}, c_{f2}) , viscous friction (b_1, b_2) for the two joints and gravity.

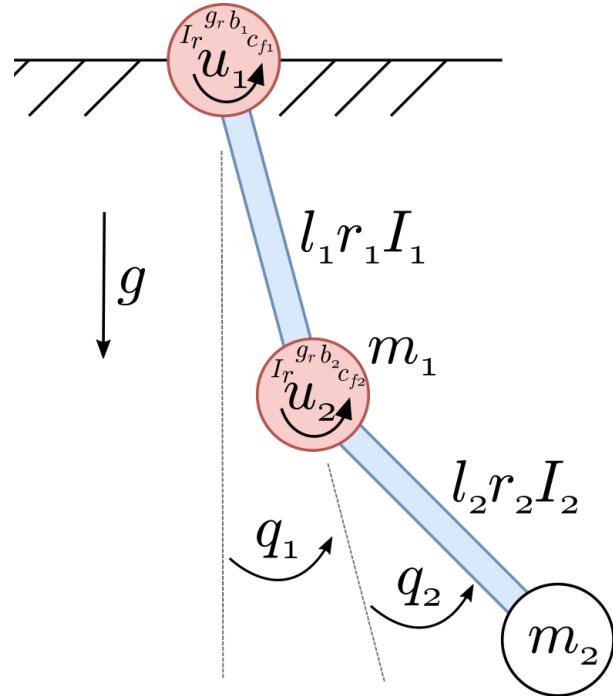


Figure 2.1: Double pendulum dynamics

The generalized coordinates $\mathbf{q} = (q_1, q_2)^T$ are the joint angles measured from the free hanging position. The state vector of the systems contains the position coordinates and their time derivatives: $\mathbf{x} = (\mathbf{q}, \dot{\mathbf{q}})^T$. The torque applied by the actuators are $\mathbf{u} = (u_1, u_2)$. The equation of motion for the dynamics of a dynamical system can be derived following the blow steps:

Step 1. Define the Lagrangian (L):

The Lagrangian (L) is defined as the difference between the kinetic energy (T) and the potential energy (U) of the system:

$$L = T - U \quad (2.4)$$

Step 2. Express the Kinetic Energy (T):

The kinetic energy (T) of the double pendulum is the sum of the kinetic energies of both pendulums. The kinetic energy for a pendulum is given by:

$$T = \frac{1}{2}m_1(\dot{x}_1^2 + \dot{y}_1^2) + \frac{1}{2}m_2(\dot{x}_2^2 + \dot{y}_2^2) \quad (2.5)$$

where m_1 and m_2 are the masses of the pendulums, (x_1, y_1) and (x_2, y_2) are their positions, and $\dot{x}_1, \dot{y}_1, \dot{x}_2, \dot{y}_2$ are their respective velocities.

Step 3. Express the Potential Energy (U):

The potential energy (U) of the double pendulum is the sum of the potential energies of both pendulums. The potential energy for a pendulum in a gravitational field is given by:

$$U = m_1gy_1 + m_2gy_2 \quad (2.6)$$

where g is the acceleration due to gravity.

Step 4. Formulate the Lagrange's Equation:

Use Lagrange's equation to derive the equations of motion for the generalized coordinates x_1, y_1, x_2, y_2 .

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_i} \right) - \frac{\partial L}{\partial q_i} = 0 \quad (2.7)$$

Step 5. Solve the Equations of Motion:

Solve the obtained set of second-order differential equations to determine the equations of motion for the system. The system dynamics with friction is:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} = Du + G(q) - F(\dot{q}) \quad (2.8)$$

Because the state vector is $\mathbf{x} = (\mathbf{q}, \dot{\mathbf{q}})^T$, the equation of motion can also be expressed as:

$$\begin{aligned} \dot{\mathbf{x}} &= f(x, u) \\ &= \begin{bmatrix} \dot{q} \\ M^{-1}(Du - C(q, \dot{q})\dot{q} + G(q) - F(\dot{q})) \end{bmatrix} \end{aligned} \quad (2.9)$$

Consider the forward kinematics of double pendulum system, the coordinate of the joint between the first link and the second link is $P_1 = (x_1, y_1)$, the coordinate of the end effector is $P_2 = (x_2, y_2)$.

$$\begin{cases} x_1 = l_1 \sin(q_1) \\ y_1 = -l_1 \cos(q_1) \end{cases} \quad (2.10)$$

$$\begin{cases} x_2 = l_1 \sin(q_1) + l_2 \sin(q_1 + q_2) \\ y_2 = -l_1 \cos(q_1) - l_2 \cos(q_1 + q_2) \end{cases} \quad (2.11)$$

Put equation(2.12) and (2.13) into (2.7)(2.8)(2.9)(2.10), we can get the mass matrix (with $s_1 = \sin(q_1), c_1 = \cos(q_1), \dots$)

$$\mathbf{M} = \begin{bmatrix} I_1 + I_2 + l_1^2 m_2 + 2l_1 m_2 r_2 c_2 + g_r^2 I_r + I_r & I_2 + l_1 m_2 r_2 c_2 - g_r I_r \\ I_2 + l_1 m_2 r_2 c_2 - g_r I_r & I_2 + g_r^2 I_r \end{bmatrix} \quad (2.12)$$

the Coriolis matrix:

$$\mathbf{C} = \begin{bmatrix} -2\dot{q}_2 l_1 m_2 r_2 \sin(q_2) & -\dot{q}_2 l_1 m_2 r_2 \sin(q_2) \\ \dot{q}_1 l_1 m_2 r_2 \sin(q_2) & 0 \end{bmatrix}, \quad (2.13)$$

The gravity vector:

$$\mathbf{G} = \begin{bmatrix} -gm_1 r_1 \sin(q_1) - gm_2 (l_1 \sin(q_1) + r_2 \sin(q_{1+2})) \\ -gm_2 r_2 \sin(q_{1+2}) \end{bmatrix}, \quad (2.14)$$

The friction vector:

$$\mathbf{F} = \begin{bmatrix} b_1 \dot{q}_1 + c_{f1} \arctan(100 \dot{q}_1) \\ b_2 \dot{q}_2 + c_{f2} \arctan(100 \dot{q}_2) \end{bmatrix} \quad (2.15)$$

(the $\arctan(100 \dot{q}_i)$ function is used to approximate the discrete step function for the coulomb friction)

and the actuator selection matrix \mathbf{D} :

$$\mathbf{D}_{full} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{D}_{pendu} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{D}_{acro} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad (2.16)$$

for the fully actuated system, the pendubot or the acrobot.

2.2 Related work

In this section, we want to introduce the methods of underactuated control and the works related.

2.2.1 Feedback cancellation

Fully actuated systems are dramatically easier to control than underactuated systems. For fully actuated systems with accurate, known dynamics(e.g. f_1 and f_2 are known for a second order system), it is possible to use feedback to effectively change an arbitrary control problem into the problem of controlling a simple linear system.

When f_2 has full row rank, it is invertible. Consider the potential nonlinear feedback control:

$$u = \pi(q, \dot{q}, t) = f_2^{-1}(q, \dot{q}, t)[u' - f_1(q, \dot{q}, t)] \quad (2.17)$$

where u' is the new control input(can be viewed as an internal input into your controller). Applying this feedback controller to equation(2.2) can lead to a linear, decoupled, second order system.

$$\ddot{q} = u' \quad (2.18)$$

In other words, if f_1 and f_2 are known and f_2 is invertible, the system is "feedback equivalent" to $\ddot{q} = u'$. Then we can use simple control method for linear systems to handle is equivalent system.

2.2.2 PFL:Partial Feedback Linearization

2.2.3 Energy shaping

example: torque limited simple pendulum

2.2.4 Reinforcement-Learning based control

Something about model free reinforcement learning and model based reinforcement learning.

3 Methodology

To achieve the swing up task, we employ the classical model-free reinforcement learning algorithm known as Soft actor critic to train a policy which is able to reach the region of attraction(RoA) of a continuous time linear quadratic regulator(LQR) controller. As soon as the system enters the RoA, we transition to the LQR controller to stabilize the entire system.

3.1 Soft actor critic

Soft Actor Critic (SAC) is a popular algorithm used in the field of reinforcement learning. It is originally designed for continuous action spaces, where the agent has an infinite choice of actions to take. In our problem scenario, the actuators of the double pendulum can be set to any value within the torque limit range. The position and velocity measurements obtained from the motors are also represented as continuous real numbers. Therefore, we choose the SAC algorithm to train the agent.

SAC optimizes a policy by maximizing the expected cumulative reward obtained by the agent over time. This is achieved through an actor and critic structure.

The actor is responsible for selecting actions based on the current policy in response to the observed state of the environment. It is typically represented by a shallow neural network that approximates the mapping between the input state and the output probability distribution over actions. SAC incorporates a stochastic policy in its actor part, which encourages exploration and helps the agent improve policies.

The critic, on the other hand, evaluates the value of state-action pairs. It estimates the expected cumulative reward that the agent can obtain by following a certain policy. Typically, the critic is also represented by a neural network that takes state-action pairs as inputs and outputs estimated value.

In addition to the actor and critic, a central feature of SAC is entropy regularization. The policy is trained to maximize a trade-off between expected return and entropy, which is a measure of randomness in the action selection. If x is a random variable with a probability density function P , the entropy H of x is defined as:

3 Methodology

$$H(P) = \mathbb{E}_{x \sim P} [-\log P(x)]$$

By maximizing entropy, SAC encourages exploration and accelerates learning. It also prevents the policy from prematurely converging to a suboptimal solution. The trade-off between maximizing reward and maximizing entropy is controlled through a parameter, α . This parameter serves to balance the importance of exploration and exploitation within the optimization problem. The optimal policy π^* can be defined as follows:

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(R(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot | s_t)) \right) \right]$$

During training, SAC learns a policy π_θ and two Q-functions Q_{ϕ_1}, Q_{ϕ_2} concurrently. The loss functions for the two Q-networks are ($i \in 1, 2$):

$$L(\phi_i, D) = \mathbb{E}_{(s, a, r, s', d) \sim D} \left[\left(Q_{\phi_i}(s, a) - y(r, s', d) \right)^2 \right],$$

where the temporal difference target y is given by:

$$\begin{aligned} y(r, s', d) &= r + \gamma(1 - d) \times \\ &\quad \left(\min_{j=1,2} Q_{\phi_{targ,j}}(s', \tilde{a}') - \alpha \log \pi_\theta(\tilde{a}' | s') \right), \\ &\quad \tilde{a}' \sim \pi_\theta(\cdot | s') \end{aligned}$$

In each state, the policy π_θ should act to maximize the expected future return Q while also considering the expected future entropy H . In other words, it should maximize $V^\pi(s)$:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] + \alpha H(\pi(\cdot | s)) \\ &= \mathbb{E}_{a \sim \pi} [Q^\pi(s, a)] - \alpha \log \pi(a | s) \end{aligned}$$

By employing an effective gradient-based optimization technique, the parameters of both the actor and critic neural networks undergo updates, subsequently leading to the adaptation of the policies themselves.

In conclusion, SAC's combination of stochastic policies, exploration through entropy regularization, value estimation, and gradient-based optimization make it a well-suited algorithm for addressing the challenges posed by continuous state and action spaces.

3.2 Linear quadratic regulator

We have the general form of a nonlinear system:

$$\dot{x}(t) = f(x(t), u(t)) \quad (3.1)$$

Linearize it around an operating point:

$$\dot{\bar{x}}(t) = A\bar{x}(t) + Bu(t) \quad (3.2)$$

The difference of the current state and the desired state is:

$$\bar{x} = x - x_{\text{op}}, \quad \dot{\bar{x}} = \dot{x} \quad (3.3)$$

where A and B can be calculated as:

$$A = \left. \frac{\partial f}{\partial x} \right|_{\text{op}}, \quad B = \left. \frac{\partial f}{\partial u} \right|_{\text{op}} \quad (3.4)$$

The cost function we are dealing with here is:

$$J = \int_0^{\infty} (x^T Q x + u^T R u) dt \quad (3.5)$$

Based on the Riccati equation:

$$A^T S + S A - S B R^{-1} B^T S + Q = 0 \quad (3.6)$$

the LQR control law for the linearized system:

$$u(t) = -K\bar{x}(t) \quad (3.7)$$

Where K can be calculated as:

$$K = R^{-1} B^T S \quad (3.8)$$

3.3 Combining SAC and LQR with region of attraction

This section is about how to use ROA to combine SAC and LQR

How to optimize the RoA of the LQR controller

First, construct Lyapunov function for stability check:

$$0 < V(\bar{x}) = \bar{x}^T S \bar{x} < \rho \quad (3.9)$$

Second, first derivative of Lyapunov function has to be negative

$$\dot{V}(\bar{x}) = 2\dot{\bar{x}}^T S \bar{x} < 0 \quad (3.10)$$

Choose suitable ρ and S that maximize the volume of an ellipsoid shaped RoA in 4D state space.

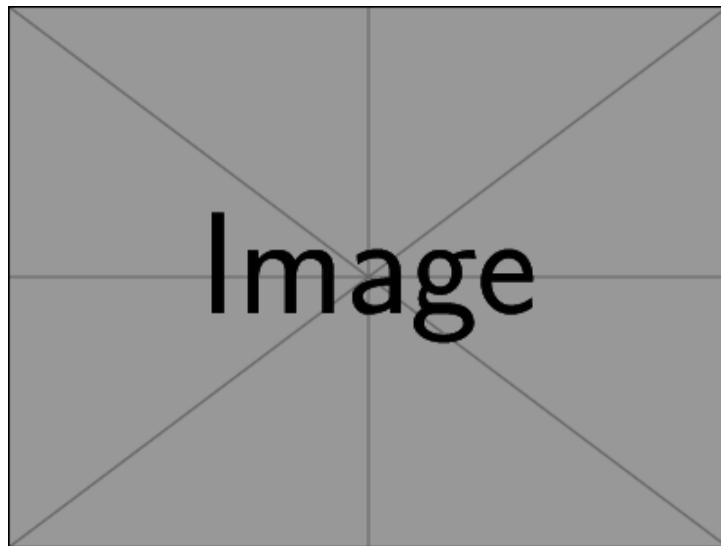


Figure 3.1: Region of Attraction

3.4 Reward shaping

This section is about the reward shaping problem of reinforcement learning training.

$$\begin{aligned}
 r(x, u) = & - (x - x_g)^T Q_{train} (x - x_g) - u^T R_{train} u \\
 & + \begin{cases} r_{line} & \text{if } h(p_1, p_2) \geq h_{line}, \\ 0 & \text{else} \end{cases} \\
 & + \begin{cases} r_{LQR} & \text{if } (x - x_g)^T S_{LQR} (x - x_g) \geq \rho, \\ 0 & \text{else} \end{cases} \\
 & - \begin{cases} r_{vel} & \text{if } |v_1| \geq v_{thresh}, \\ 0 & \text{else} \end{cases} \\
 & - \begin{cases} r_{vel} & \text{if } |v_2| \geq v_{thresh}, \\ 0 & \text{else} \end{cases}
 \end{aligned}$$

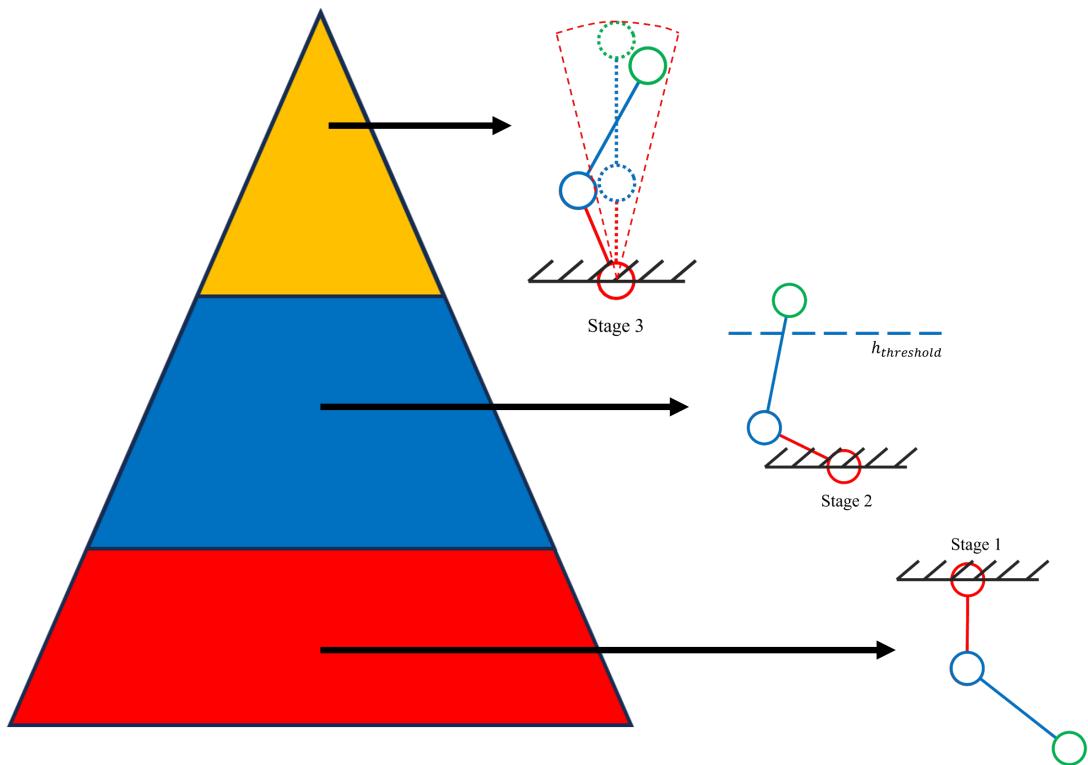


Figure 3.2: Interpretation of reward shaping

4 Experiment: training and simulation

This chapter is about the training and simulation.

ad;falknv.xzfvhlsakjdgnmdflk asdjgfmndfgb;lkjesf;mngfbl;kjfx.,mszedk.jfal;j

4.1 Training setup

This section is about training setup and environment building.

Robot	Quadratic Reward	Constant Reward	LQR
Pendubot	$Q_1 = 8.0$		$Q_1 = 1.92$
	$Q_2 = 5.0$	$r_{line} = 500$	$Q_2 = 1.92$
	$Q_3 = 0.1$	$r_{vel} = 0.0$	$Q_3 = 0.3$
	$Q_4 = 0.1$	$r_{LQR} = 1e4$	$Q_4 = 0.3$
	$R = 1e-4$		$R = 0.82$
Acrobot	$Q_1 = 10.0$		$Q_1 = 0.97$
	$Q_2 = 10.0$	$r_{line} = 500$	$Q_2 = 0.93$
	$Q_3 = 0.2$	$r_{vel} = 1e4$	$Q_3 = 0.39$
	$Q_4 = 0.2$	$r_{LQR} = 1e4$	$Q_4 = 0.26$
	$R = 1e-4$		$R = 0.11$

Table 4.1: Hyper parameters used for the SAC training and the LQR controller.

4.2 Training process

This section is about learning curve and tuning parameters

4 Experiment: training and simulation

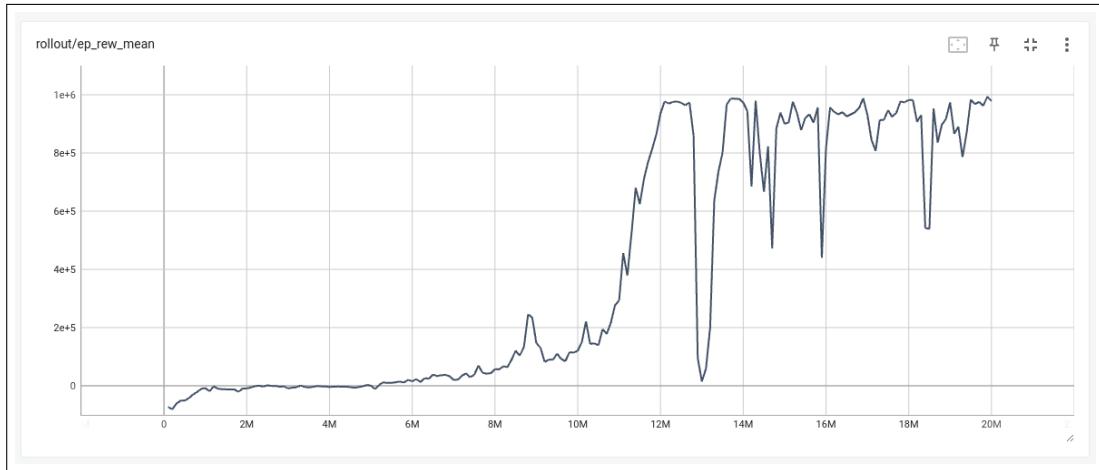


Figure 4.1: pendubot learning curve

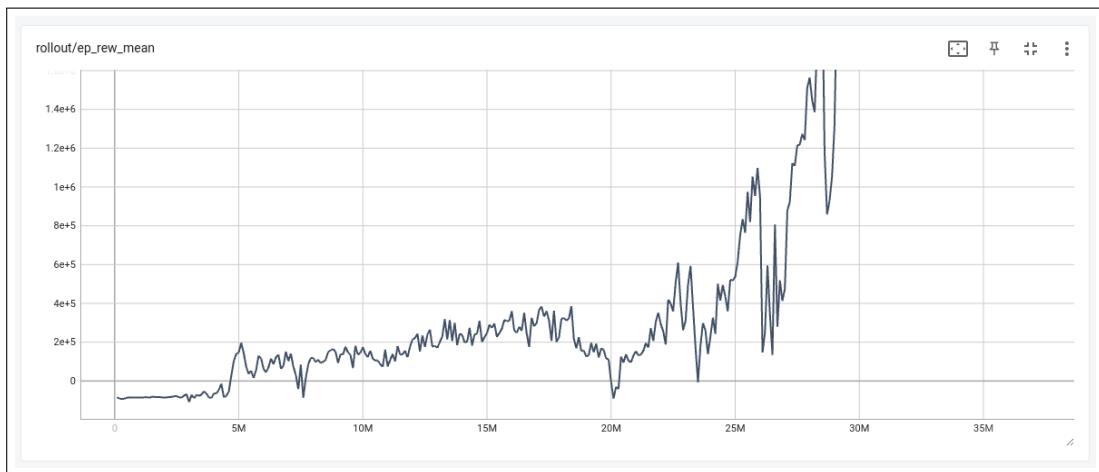


Figure 4.2: acrobot learning curve

4.3 Simulation results

This section is about simulation results in pendubot and acrobot.

4.3.1 pendubot

pendubot:

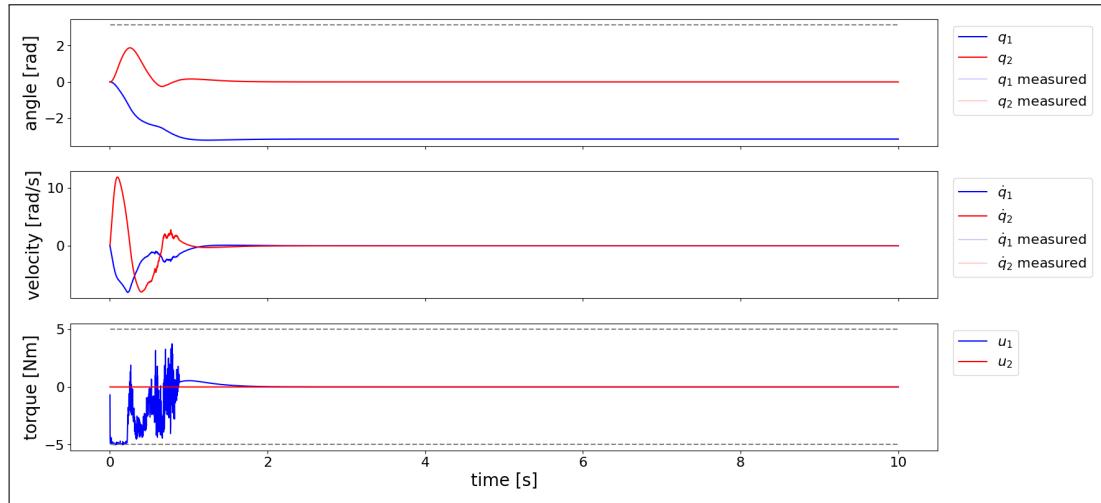


Figure 4.3: pendubot simulation result

4.3.2 acrobot

acrobot:

4 Experiment: training and simulation

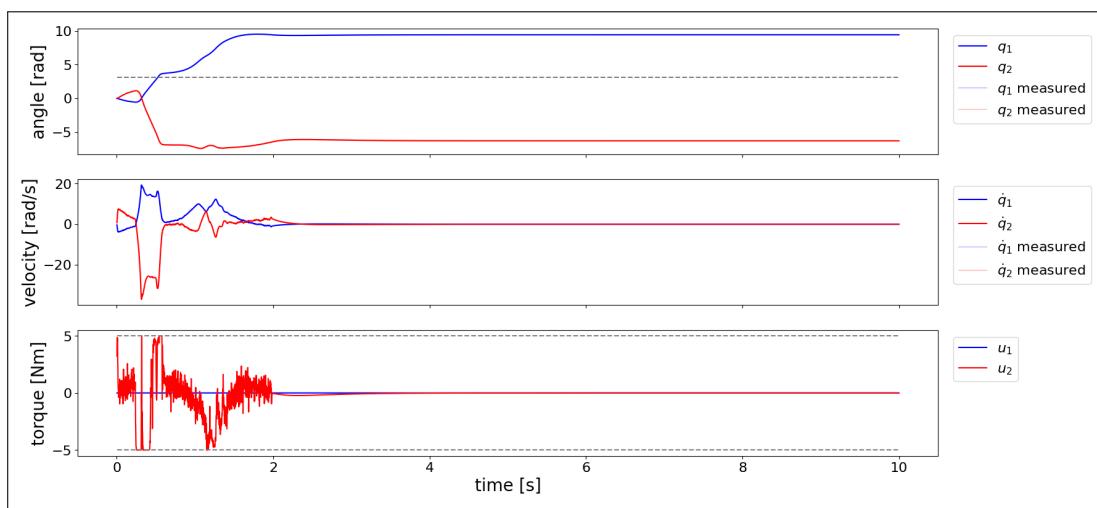


Figure 4.4: acrobot simulation result

5 Experiment: hardware system

This chapter is about the hardware experiment.

ad;falknv.xzfvhlsakjdgnmdflk asdjgfmndfgb;lkjesf;mngtbl;kjfx.,mszedk.jfal;j

5.1 Hardware setup

This section is about hardware setup and its features.

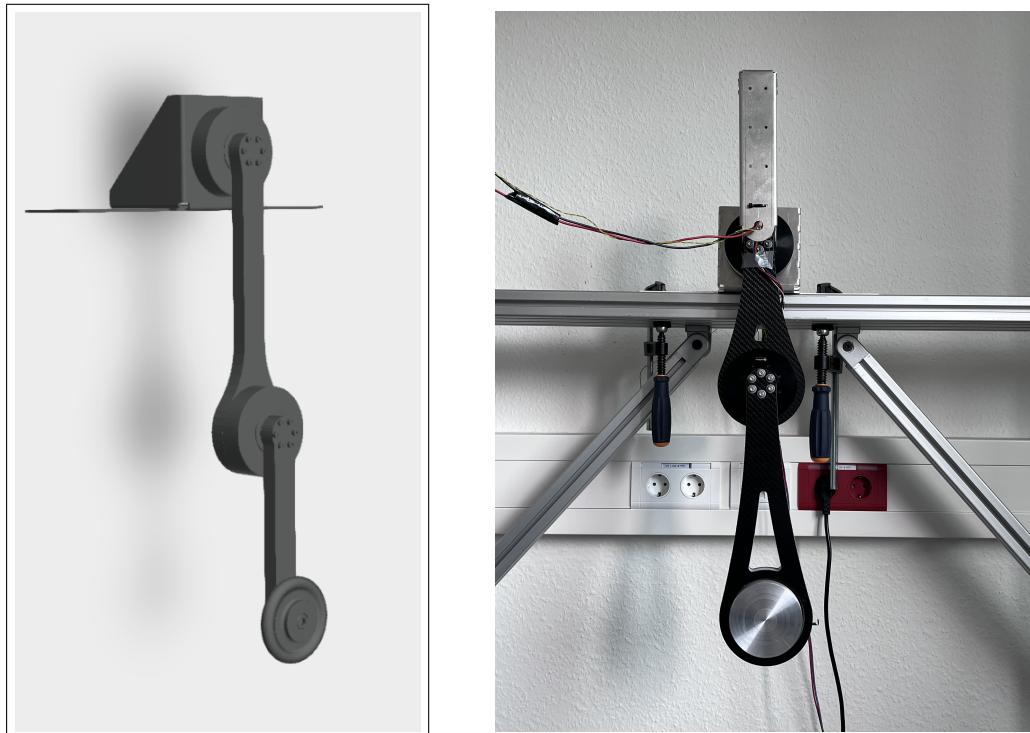
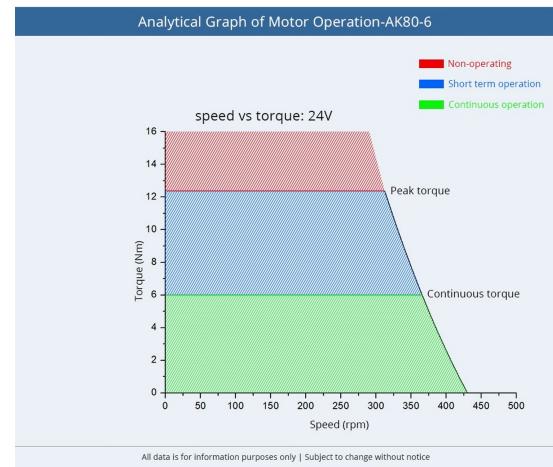


Figure 5.1: Double pendulum setup in CAD and in real world

5 Experiment: hardware system



(a) Figure A



(b) Figure B

Figure 5.2: Two figures side by side

5.2 system identification

This section is about the system identification problem when using hardware system.

5.3 sim2real problem

This section is about sim2real problem.

5.4 real hardware results

This section is about simulation results in pendubot and acrobot.

pendubot:

acrobot:

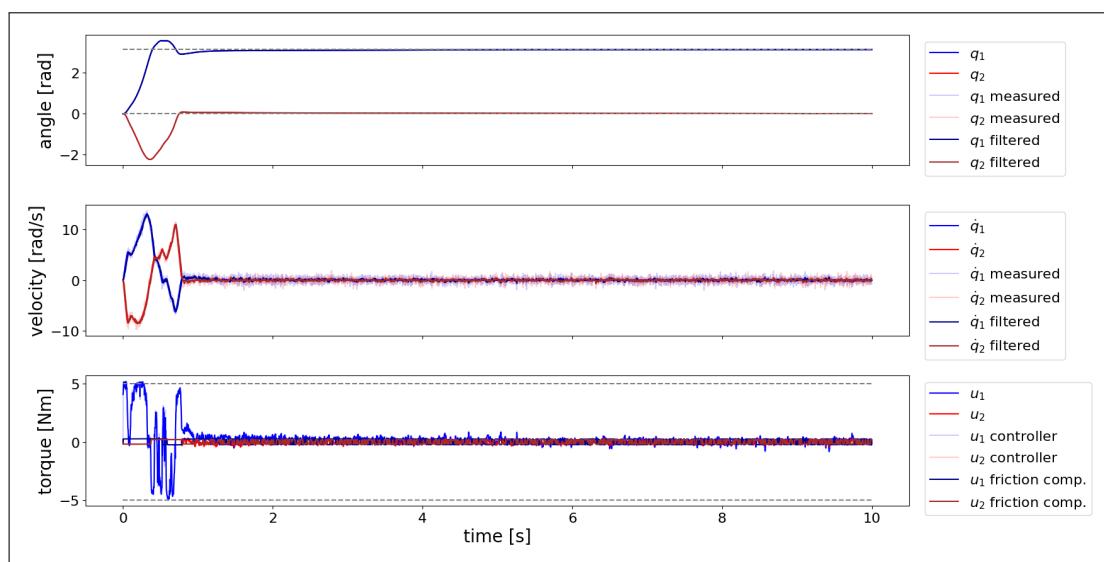


Figure 5.3: pendubot noisy simulation result

6 Discussion

This chapter is about the discussion of results.

ad;falknv.xzfvhlsakjdgnmdflk asdjgfmndfgb;lkjesf;mngfbl;kjfx.,mszedk.jfal;j

6.1 introduction to leaderboard results

This section is about simulation and real system leaderboard

- **Swingup Success** c_{success} : Whether the swingup was successful, i.e. if the end-effector is above the threshold line at the end of the simulation.
- **Swingup time** c_{time} : The time it takes for the pendubot to reach the goal region above the threshold line and stay there. If the end-effector enters the goal region but falls below the line before the simulation time is over, the swingup is not considered successful! The swingup time is the time when the end-effector enters the goal region and does not leave the region until the end.
- **Energy** c_{energy} : The mechanical energy used during the execution.
- **Max Torque** $c_{\tau,\text{max}}$: The peak torque that was used during the execution.
- **Integrated Torque** $c_{\tau,\text{integ}}$: The time integral over the used torque over the execution duration.
- **Torque Cost** $c_{\tau,\text{cost}}$: A quadratic cost on the used torques ($c_{\tau,\text{cost}} = \sum u^T R u$), with $R = 1$.
- **Torque Smoothness** $c_{\tau,\text{smooth}}$: The standard deviation of the changes in the torque signal.
- **Velocity Cost** $c_{\text{vel},\text{cost}}$: A quadratic cost on the joint velocities that were reached during the execution ($c_{\text{vel}} = \dot{q}^T Q \dot{q}$), with $Q = \text{identity}$.

6 Discussion

Criteria	Pendubot	Acrobot
Swingup Success	success	success
Swingup time [s]	0.65	2.06
Energy [J]	9.4	29.24
Max. Torque [Nm]	5.0	5.0
Integrated Torque [Nm]	2.21	4.57
Torque Cost [N ² m ²]	8.58	12.32
Torque Smoothness [Nm]	0.172	0.954
Velocity Cost [m ² /s ²]	44.98	193.78
RealAI Score	0.801	0.722

Table 6.1: Performance scores of our controller for pendubot and acrobot.

$$S = c_{\text{success}} \left(w_{\text{time}} \frac{c_{\text{time}}}{n_{\text{time}}} + w_{\text{energy}} \frac{c_{\text{energy}}}{n_{\text{energy}}} + w_{\tau,\max} \frac{c_{\tau,\max}}{n_{\tau,\max}} + w_{\tau,\text{integ}} \frac{c_{\tau,\text{integ}}}{n_{\tau,\text{integ}}} + w_{\tau,\text{cost}} \frac{c_{\tau,\text{cost}}}{n_{\tau,\text{cost}}} + w_{\tau,\text{smooth}} \frac{c_{\tau,\text{smooth}}}{n_{\tau,\text{smooth}}} + w_{\text{vel}, \text{cost}} \frac{c_{\text{vel}, \text{cost}}}{n_{\text{vel}, \text{cost}}} \right) \quad (6.1)$$

6.2 interpretation of simulation results

This section is about explaining the simulation results.

6.3 interpretation of hardware results

This section is about explaining the hardware results.

Criteria	Pendubot	Acrobot
Swingup Success	success	insuccess
Swingup time [s]	0.67	-
Energy [J]	37.12	-
Max. Torque [Nm]	5.0	-
Integrated Torque [Nm]	24.87	-
Torque Cost [N^2m^2]	78.7	-
Torque Smoothness [Nm]	0.774	-
Velocity Cost [m^2/s^2]	114.04	-
RealAI Score	0.298	-

Table 6.2: Real hardware performance scores of our controller for pendubot and acrobot.

6.4 future work

This section is to talk about things to be done.

Appendix

A An appendix

You can structure appendices, just like your thesis, with the \chapter, \section, and \subsection commands. Referencing also works as usual.

If your thesis does not contain an appendix, comment out the creation of the appendix at the appropriate place in the Thesis.tex file.