# Optimal Input Excitation for Parameter Estimation in Electromechanical Brake Actuators

Studiengang Mechatronik
Master thesis in the department of Electrical Engineering and Information Technology
by Shengya Guo
Date of submission: 2. December 2024

Erstprüfer: Prof. Dr.-Ing. Rolf Findeisen
Betreuer: Dr.-Ing. Eric Lenz
Betreuer: M. Sc. Daniel Stümke
Darmstadt

TECHNISCHE
UNIVERSITÄT
DARMSTADT

CCPS
Control and Cyber-Physical Systems

Electrical Engineering and
Information Technology
Department

Optimal Input Excitation for Parameter Estimation in Electromechanical Brake Actuators
Studiengang Mechatronik

Master thesis in the department of Electrical Engineering and Information Technology by Shengya Guo

Date of submission: 2. December 2024

Darmstadt

Technische Universität Darmstadt
Institut für Automatisierungstechnik und Mechatronik
Fachgebiet Control and Cyberphysical Systems
Prof. Dr.-Ing. Rolf Findeisen

## Task definition

With the ongoing electrification of passenger vehicles, the transition from hydraulic to electromechanical brake (EMB) systems becomes attractive. This new kind of braking system promises advantages, as it removes bulky hydraulic components, enables fast wheel individual brake actions, and integrates seamlessly with the X-by-wire paradigm, but it also poses new challenges which need to be overcome. To enable the reliable operation of EMB systems, exact knowledge about system parameters e.g., actuator motor torque constant, friction coefficients or the brake pad stiffness, is essential [1]. Due to aging and wear or manufacturing imperfections these quantities are uncertain and may vary over the component's lifetime. Therefore, a method which reliably estimates system parameters is of great importance. An example of such a method in the context of EMB is presented in [2], where a genetic algorithm is used to obtain a control input which maximizes an information gain criterion (D-Optimality). In other engineering domains, the optimal excitation problem is tackled with Reinforcement Learning (RL) [3] and Chance Constraints [4]. The aim of the offered master's thesis is to investigate how optimal excitation strategies can be realized (open- or closed-loop, with RL or other optimization techniques, etc.) for a mechatronic system such as the electromechanical brake.

### Tasks:

- Literature research on parameter estimation for mechatronic systems
- Formulation of a simplified system model for analysis of excitation strategies
- Implementation of a suitable excitation strategy
- Evaluation based on simulation results and optionally also on prototype hardware

This thesis is offered in cooperation with Robert Bosch GmbH (Corporate Research).

Beginn: ——
Ende: ——
Seminar: ——

**Erklärung zur Abschlussarbeit gemäß § 22 Abs. 7 APB TU Darmstadt**

Hiermit erkläre ich, Shengya Guo, dass ich die vorliegende Arbeit gemäß § 22 Abs. 7 APB der TU Darmstadt selbstständig, ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt habe. Ich habe mit Ausnahme der zitierten Literatur und anderer in der Arbeit genannter Quellen keine fremden Hilfsmittel benutzt. Die von mir bei der Anfertigung dieser wissenschaftlichen Arbeit wörtlich oder inhaltlich benutzte Literatur und alle anderen Quellen habe ich im Text deutlich gekennzeichnet und gesondert aufgeführt. Dies gilt auch für Quellen oder Hilfsmittel aus dem Internet.

Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Mir ist bekannt, dass im Falle eines Plagiats (§ 38 Abs. 2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei einer Thesis des Fachbereichs Architektur entspricht die eingereichte elektronische Fassung dem vorgestellten Modell und den vorgelegten Plänen.


Darmstadt, 2. December 2024

_____

S. Guo

## Kurzfassung

Zusammenfassung entsprechend der Dokumentensprache. In diesem Fall Deutsch.

## Abstract

Additional abstract in English.

# Contents

# 1 Introduction

Nowadays, electrification of vehicles and autonomous driving have became the most two popular trends of the entire vehicle manufacture industry. Under the encouragement from policies in different country, the battery electric vehicle (BEV) market has exploded in recent years. More and more brand have start the turn their production into electrocution. Back in to the days, traditional motor vehicle use the 12-Voltage platform and the Electronic Control Unit (ECU) to control the whole system running, the main target of the system is to keep the vehicle functional in any circumstance. To enable new functions, such as Adaptive Cruise Control (ACC), the amount of networking within the domains of in traditional vehicle EE system architectures has been increased [5]. This leads to the result, that the need for the accomplish of new functions increase the requirement of electrical power in vehicles [6]. With the respect of this, the electric vehicle platform could brings a lot of advantages, for example the possibility to carry high performance computer, control all components on chassis with a single centralized system etc, which make the vehicle more suitable for the autonomous driving.

With the growing interest of autonomous driving and the support of electric vehicle platform, a new technology with the name of X-by-wire has show up. X-by-wire technology replace the traditional mechanical connection between operator and chassis components with electronic signal. By doing this, the advanced X-by-wire chassis could reduce the time of system reaction and improve the performance and safety of vehicles [7].

Brake-by-wire (BBW) is a part of the X-by-wire chassis. Different from the traditional hydraulic brake system, which the brake paddle connected to a hydraulic cylinder and achieve the control of the hydraulic brake system, the brake paddle in brake-by-wire system doesn't directly drive the brake system, but will generate a control signal to controller and then manipulates the actuator. brake-by-wire systems can be categorized in main types based on their sources and methods of regulating braking force: the electro-hydraulic brake (EHB) system and the electromechanical brake (EMB) system [8].

The EHB system use a electro-pump to drive hydraulic system and generate the clamping force for brakes. Due to the advantages of good compatibility and easy implementation, EHB has become a mature solution in the current market [9]. As the hydraulic fluid system is a main part in EHB systems, the drawback of hydraulic system such as slow braking response, possibility of leaks and the consequence of pollution still existed [10]. Moreover, the elimination of the high cost of hydraulic system always stands for a goal of manufactures.

In EMB systems, motor-generated torque is transmitted directly to the brake's friction components via mechanisms like gear reduction and ball screws, generating frictional braking torque [10]. Because the whole hydraulic system was saved, the disadvantages of EHB systems doesn't existed for EMB. Moreover, from a physics perspective, EMB systems can achieve faster braking responses, which are urgently needed for autonomous driving. In many dangerous scenarios, the autonomous vehicle need to made the decision in millisecond. With the help of EMB systems, it will accomplish higher driving performance and shorter braking distance than experienced human driver. EMS systems also have

several other advantages, including reduced system volume and weight, integrated with a parking brake, and compatibility with active safety control systems [11], [12].

With all the advantages above, the fact is that EMB have not show up in the mass production yet. Functional safety concepts is a main reason of it [13]. EMB actuators need to provide a maximum caliper clamping force of 40 kN to provide sufficient braking force [14]. EMB needs to work under all conditions like hot, cold, strong vibration etc, which the normal motor can't functional normally.

To guarantee rapid response in all circumstance, the EMB need to have the ability to be controlled precisely in very short time. Thus, build up an accurate model for EMB is the research hard point for researchers. In order to apply the control task and simulation, a precise model of EMB is necessary. Hence, the identification of model parameter is very important for the performance and accuracy of EMB models. Form this reason, plenty research of parameter estimation of EMB have been conducted, including nonlinear optimization, xxx and xxx **<empty citation>**

With the developing of artificial intelligence (AI) and the growing of computing performance of edge device, more and more AI approaches was used in the vehicle area. Not only the most popular topic in the vehicle industry task like end to end autonomous driving and large language models (LLMs) is involved in AI approaches, in more specific area, AI was also widely used to solve the industry problem including complex and non-linear control for autonomous vehicle, optimal design of motor for electric vehicle and optimization of battery management systems (BMS) [15]–[19]. This leads to the idea that, AI approaches could might also help with searching for the optimal excitation for identification of EMB parameters, which will be conducted in the following text in detail.

In this thesis, LLMs and translation tools are used to check grammar and make the writing more suitable for scientific purposes [20], [21].

## 1.1 Motivation

As the EMB could brings a lot of benefits which had already show in early text, there still existing various challenges to be solve. The most critical thing is that autonomous driving system with fast response capability puts high demands on the accuracy of EMB system model. Unlike general systems, the EMB is a dynamic system, and the system characteristics of the EMB will change accordingly with the passage of time or changes in temperature. This poses an even greater challenge to system modelling, as the system parameters determined once using conventional methods cannot be guaranteed to be correct in the long term. In addition to this, traditional methods require the assumption of an a priori value for the system parameters and the optimal design of experiments based on this a priori value [3]. The problem is that this a priori value cannot be determined in advance, and if the a priori value deviates significantly from the actual value, this significantly reduces the performance of the design results.

In contrast to traditional time-consuming parameter identification methods, this thesis attempts to solve this problem by applying artificial intelligence, more specifically reinforcement learning, to determine the optimal input excitation for parameter identification. The main aim of this thesis is to investigate two questions, firstly whether reinforcement learning can be used to optimise the best input excitation for parameter identification, and secondly how reinforcement learning can be designed in such a way that the agent is free from its dependence on the a priori values of the system parameters. For the first problem, an agent with reinforcement learning based on a priori values of system parameters is trained and its performance is tested. For the second problem, a reinforcement learning framework is designed

based on the improvement of existing methods and its robustness to non-deterministic parameter a priori values.

Another implication of this thesis is that its approach is not only limited to apply to EMB systems. More precisely, this thesis attempts to develop a more general approach for optimising input excitations for parameter discrimination. In this paper, EMB is used more as a platform to apply and test the method. Therefore another motivation for this topic is that, with simple modifications, the method in this thesis can also be used in other systems to optimise the input excitation for parameter discrimination.

To do the identification, the conventional approach is do different test for different parameter. In this thesis, approaches to do the identification of EMB model parameter was developed.

## 1.2 Structure of the Thesis

Chapter 2 describes the current state of development of optimal experimental design and parameter identification methods, comparing their application in the field of EMB and the application of reinforcement learning methods in optimal input exaction.

In Chapter 3, the modelling of the electromechanical braking system was described and a simplified model of the EMB used in this thesis is presented.

Chapter 4 presents the methods for applying reinforcement learning to the optimizing experimental design in EMB, with detailed explanations of reinforcement learning framework formulation, parameter estimation, and environment setup.

Chapter 5 describes the application of the reinforcement learning framework to single-parameter optimal input excitation, where training results for reinforcement learning are shown and the performance of two observation space structures is compared.

Chapter 6 extends the application to the multi-parameter domain, describing the multi-parameter approach and presenting the training results.

Chapter 7 summarises the content of this thesis and provides an outlook of future works.

# 2 State of the Art

Papers on general optimal design Papers on parameter estimation and its use in EMB Papers on RL and some RL approaches already used for optimal input exaction and parameter estimation.

Compare the pros and cons of existing approaches for optimal input exaction and parameter estimation and make a clear view in the table. Came to the clue of the necessities of this thesis and which problem or drawback form the existing research are overcome in this thesis.

# 3  Modelling of the Electromechanical Brake System

To have the idea of the EMB system, a mathematics model is necessary. Talking with the parameter estimation, the way of modelling have direct affect of the performance of the modelling. For a functional part with high quality control task and rapid response ability like EMB, the model of the system should be as precise as possible. As EMB is a complex system, this leads to the high level of parameters to be consider of during the modelling.

Nowadays, a lot of researches have been done in electromechanical brake, but mostly focus on the disc EMB system. Compare with the disc brakes, drum brakes is not wildly used for vehicles in recent years. The reason of it including, drum brakes have lower performance facing the cooling problem and the traditional respect from customer that considering drum brakes is cheaper etc. However, drum brakes still have advantages under the more and more restrict politic restriction, especially concern about the lower piratical emission because of the closed drum [22]. Thus manufacture like Bosch start to do research on electromechanical drum brake again and in this thesis, research was conducted based on a electromechanical drum brake system from Bosch.

In this chapter, the structure of the to be modelled EMB system will be introduced first. Then, different approaches to do the modelling of the EMB system will be introduced. Later the key parameters in different point of view for the modelling of EMB will be described. Finally the simplify of the EMB model for the later use will be conducted.

The subject of different degrees of accuracy when modelling different parts of hydraulic systems has been well described by several authors, where only a simplified modelling is usually applied [23].

## 3.1  Structure of the Electromechanical Drum Brake

As show in Figure 3.1, the electromechanical drum brake to analysed have the following components: the brake shoes, the brake lining, return springs and the brake drum. When the ECU receives a braking signal, it will activates the electro motor, which drives the actuator and push brake shoes and brake pads against the drum. When braking is released, the motor reverses, and a spring returns the brake shoes and brake pads to its original position [24], [25].

Consider of the structure of the electromechanical drum brake, several force needs to take into consideration during the modelling. First of all, the propose of the break is to make the wheel speed lower. When the actuator force the brake shoes move against the brake drum, the drake lining will have contact with the drum and the clamping force will show up, finally this clamping force will transformed in to brake moment. Secondly, the force form the return springs is also a main part in the model. Another big part of force which need take in to consideration is the force from the actuator, namely the electro motor.
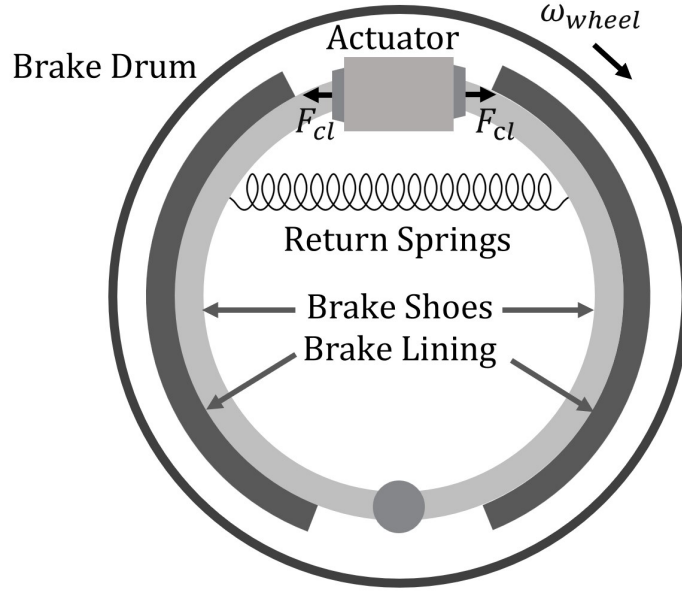
Figure 3.1: Mechanical structure of the electromechanical drum brake to be analyzed [26]

There is still a kind of torque which can't neglected, that is the friction torque. The friction torque is everywhere, from inside the motor to the actuator. In following text, these torque will be deal in detail.

## 3.2 Lumped parameter EMB Model

In this section, all the torque which have been mention in the above text will be treat carefully in the following text. As we can already seen, that torques happens all over the EMB, which make it difficult for the modelling. According to researches, the lumped parameter EMB model always been used [1]. In this thesis, three main torque are considered in the lumped parameter EMB Model: the motor torque $T_m$, the load torque $T_l$ and the friction torque $T_f$. Consider of the dynamic behavior of the motor, the EMB have been modeled as:

$$J\ddot{\theta} = T_m - T_f - T_l \tag{3.1}$$

The motor torque $T_m$ serves to accelerate the brake shoes, while the load torque $T_l$ and frictional torque $T_f$ act in opposition, resisting this acceleration. Assuming a linear relationship, the electromagnetic motor torque $T_m$ cloud be regard as directly proportional to the motor torque constant $k_m$ and the motor current $i_m$ with the equation

$$T_m = k_m i_m \tag{3.2}$$

The load torque $T_l$ is proportional to the clamping force $F_{cl}$ and the overall gear ratio $\gamma$
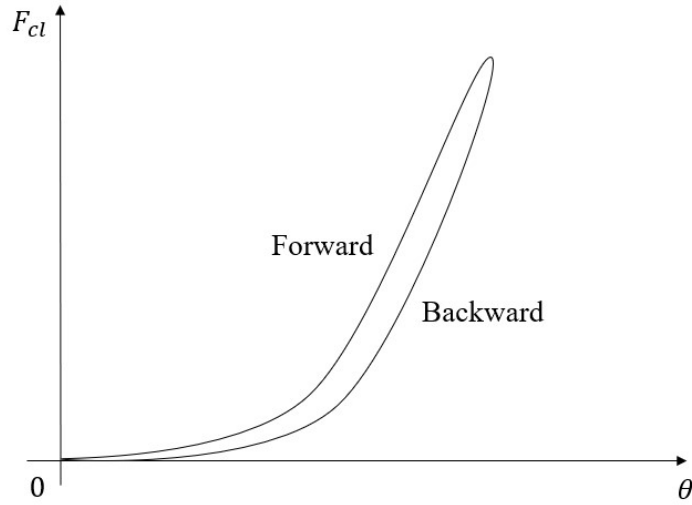
$$T_l = \gamma F_{cl} \tag{3.3}$$

Figure 3.2: Qualitative hysteresis characteristics between clamping force and motor angle

The gear ratio $\gamma$ represents the characteristics of the entire mechanical system, transforming the force from the braking shoe side to torque on the motor side. The friction torque $T_f$ could be regard as a equation which included the motor angle $\theta$ and the motor velocity $\dot{\theta}$. The detail of calculation of clamping force $F_{cl}$ and friction torque $T_l$ will be conducted in the following subsection.

### 3.2.1 Stiffness Model

The EMB is a complex nonlinear system, and its nonlinear properties are reflected in the stiffness model. The characteristic curve, also known as the stiffness of the EMB system, describes the relationship between the motor angle $\theta$ and the clamping force $F_{cl}$. The clamping force is quantified in relation to the drive motor through the utilisation of force sensors positioned on the EMB. As illustrated in Figure 3.2, the stiffness curve is nonlinear. It is assumed that 0 degrees represents the point of contact between the brake pad and the brake drum. The brake force stiffness curve can be divided into two intervals. During the initial range, the brake force grows at a slow rate. However, from later on range, the clamping force grows nearly exponentially with angle. The stiffness model of the EMB exhibits hysteresis curve characteristics, as outlined in the referenced literature [27]. In comparison to the process in which the brake pads leave the brake drum due to the action of the reset spring, when the brake pads are actuated by the actuator and move towards the brake drum, the curve can reflect a higher stiffness. Furthermore, the hysteresis characteristics are also dependent on factors such as speed and temperature. Given that the parameter identification process generally focuses more on the braking process than the release process, the hysteresis characteristics are ignored in this model to maintain a simplified representation.

### 3.2.2 Friction Model

For a long time, friction has been intensively studied in physical systems. The classic model of friction is static model, it including Coulomb friction $F_c$ and viscous friction $F_v$. By applying this classic model to
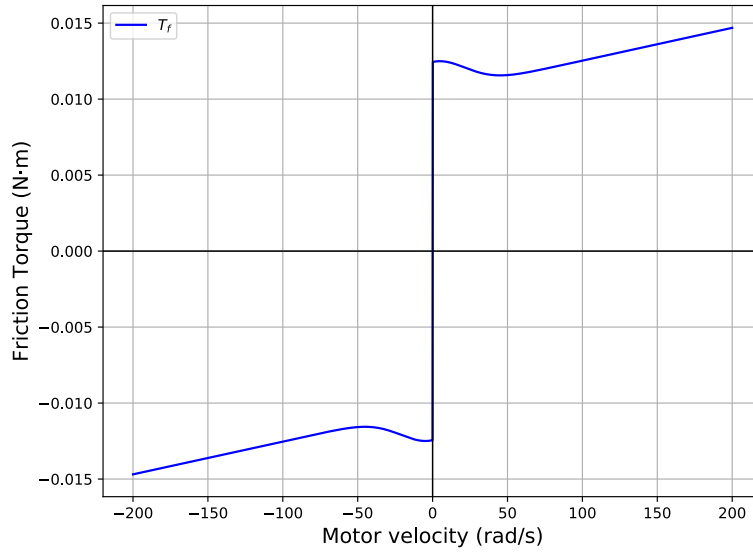
Figure 3.3: Illustration of the Stribeck friction model

the EMB system, the friction torque equal to the summation of these two frictions torque

$$T_f = T_c \operatorname{sign}(\dot{\theta}) + f_v \dot{\theta} \tag{3.4}$$

$T_c$ is the friction torque, $f_v$ is the viscous friction torque parameter and $\dot{\theta}$ is the motor velocity.

The drawback in equation (3.4) is, in this static model, the static friction which higher than the Coulomb friction is not included. Stribeck observed a continues dependence for static friction and dynamic friction [28]. One commonly used model for the non-linearity with Stribeck behavior considered in EMB is expressed as

$$T_f = T_c + (T_s - T_c)e^{-|\dot{\theta}/\dot{\theta}_\sigma|^{\delta_\sigma}} + f_v v \tag{3.5}$$

where $\dot{\theta}_\sigma$ represents the Stribeck velocity. The friction model with Stribeck effect is show in Figure 3.3 This type of model has been widely used for a long period. By holding the velocity, the friction function can be easily determined by measuring the friction force.

In EMB modelling, an additional friction component, which is dependent on the load, could also be considered. This is modeled as a linear function of the clamping force $F_{cl}$, with $f_l$ as the proportionality coefficient

$$T_s = T_{s0} + f_l F_{cl} \tag{3.6}$$

One of the limitations of the aforementioned model is the necessity to ascertain the point at which the motor speed is at its lowest. Karnopp discovered that this issue can be circumvented by incorporating a small zero-speed interval $\theta$ into the model and preventing the motor from repeatedly alternating between positive and negative equations of motion [29]. In accordance with the aforementioned model, the friction force may be either arbitrary or defined, contingent upon whether the speed is less than the

low-speed interval threshold. With this theory, the EMB system friction could be built as

$$
T_f = \begin{cases} T_c \, \text{sign}(\dot{\theta}) + f_v \dot{\theta} & \text{if } |\dot{\theta}| > \varepsilon \\[2mm] T_e & \text{if } |\dot{\theta}| < \varepsilon \text{ and } T_e < T_s \\[2mm] T_s & \text{otherwise} \end{cases} \tag{3.7}
$$

$T_e$ stands for the external torque and $T_f$ stands for the load-independent static friction torque. When the motor speed with in the zero-speed interval and the external torque smaller than the static friction torque, the friction torque will be equal to the external torque, as the system could be seen as don't move.

However, this model is not without its limitations. When the motor speed stay in the zero-speed interval and the external torque smaller than the static friction torque, the friction force is not explicitly given, rendering this model unsuitable for cases where the friction force must be derived.

Olsson et al.[30] said many friction phenomena do not show up in constant velocity experiments. Thus, in addition to static friction modeling systems, dynamic friction modeling has also been studied. Given the high precision of the EMB, the dynamic friction model is more complex and requires detailed modeling of the system to analyze the dynamic properties of the system. Consequently, the dynamic friction model is also considered in the modeling.

There is dynamic friction model as LuGre model [31]. The LuGre model interprets friction as the deflection of tiny bristles between surfaces. When a force is applied, these bristles bend like springs, and if they bend too much, they start to slip. In steady-state motion, the average deflection depends on speed, becoming smaller at higher velocities. This reduction in friction at higher speeds represents the Stribeck effect. The model also includes rate-dependent behaviors, capturing dynamic effects like frictional lag.

But this model is often too complex and need to be simplify. This so called steady-state approximation allows us to focus on the frictional behavior after the dynamic transients have subsided, effectively making the model act quasi-static. LuGre model in static could be described as

$$
\begin{aligned}
T_f &= \sigma_0 g(\dot{\theta}) \, \text{sign}(\dot{\theta}) + \sigma_2 \dot{\theta} \\
&= \left( T_c + (T_s - T_c) e^{-\left(\frac{\dot{\theta}}{\omega_s}\right)^2} \right) \text{sign}(\dot{\theta}) + f_v \dot{\theta}
\end{aligned} \tag{3.8}
$$

As the equation showed, the LuGre model in static form will be the same as it has been shown in equation (3.5), with the static friction model and the Stribeck effect been considered. This show the unify of different ways of friction modelling.

## 3.3 Simplify of the Model of Electromechanical brake System

The preceding section provides a comprehensive account of the various factors that must be taken into account when constructing a nonlinear EMB system. In theory, the more detailed the modeling of a system, the more accurately it reflects the actual performance of the physical system. Nevertheless, in practice, the most detailed model is not always necessary. In the majority of cases, the model must be
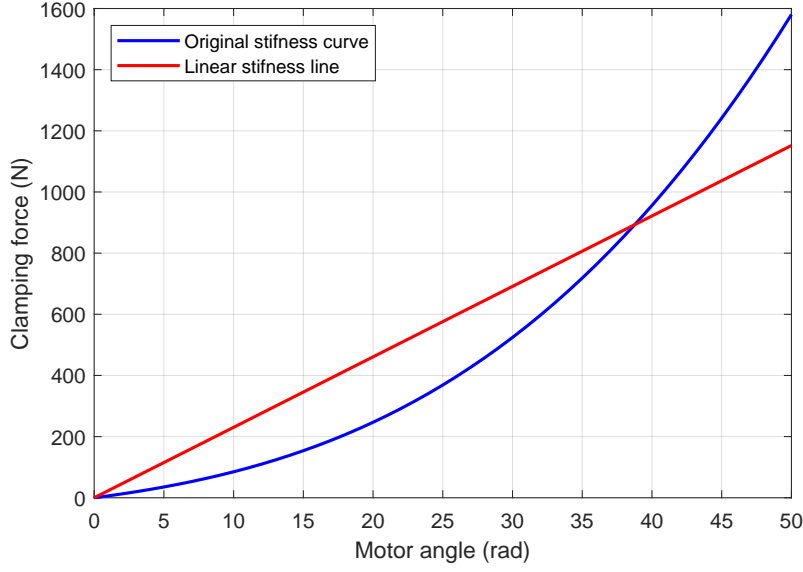
Figure 3.4: Comparison of original and linearized stiffness curves for clamping Force

simplified in order to facilitate computation. This thesis investigates the potential for optimal input excitation of reinforcement learning methods for parameter identification. In theory, each parameter has an optimal excitation, and overly complex EMB models will have a greater number of parameters that need to be recognized. This will complicate the problem that reinforcement learning needs to solve at the outset. Furthermore, this thesis is primarily a research project on EMB systems, and its primary objective is not to solve the engineering problems of actual EMB systems. Consequently, the corresponding simplifications of the models for each of the aforementioned parts were carried out.

### 3.3.1 Simplify of the Stiffness Model

For the stiffness model, a polynomial is generally fitted to the experimentally collected relationship between motor angle and clamping force, and the fitted result is mostly a cubic curve

$$F_{cl} = k_1\theta^3 + k_2\theta^2 + k_3\theta \tag{3.9}$$

Figure 3.4 illustrates the red curve representing the relationship between the motor angle and the clamping force, based on data collected from an established and well-validated electromechanical drum brake model experiment. Introducing three stiffness parameters, $k_1 - k_3$, would greatly increase the dimensional complexity of the reinforcement learning problem, necessitating parameter identification for each of these components. To mitigate this issue, a linear interpolation method is used in the modeling process. The linear clamping force is represented by

$$F_{cl} = k_1\theta \tag{3.10}$$

Although using linear interpolation increases the deviation between the system model and the real system behavior, it reduces the three stiffness parameter dimensions to one parameter dimension, which helps to reduce the learning difficulty of the reinforcement learning agent.
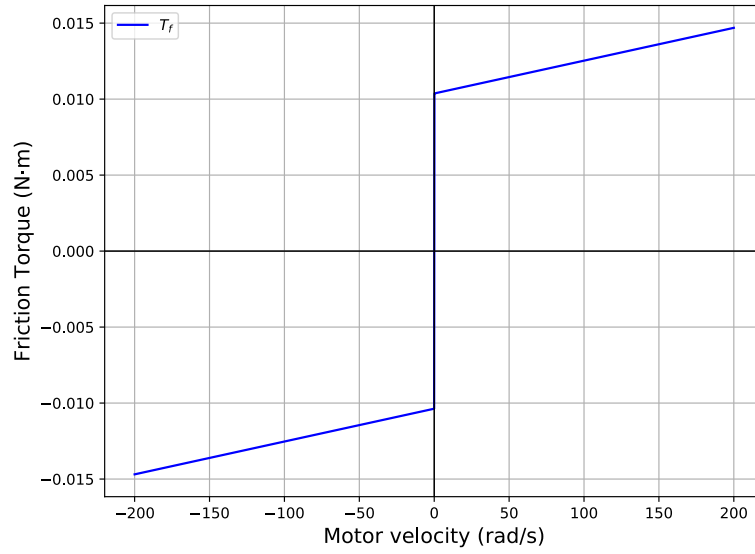
Figure 3.5: Illustration of the simplified friction model

To reduce the dimension level of the stiffness parameter, data from validated EMB model was imported into $Matlab$, and the $CurveFittingToolbox$ was used to derive the linear stiffness line (red line) through linear interpolation. The resulting red and blue curves in Figure 3.4 demonstrate the original stiffness response and the simplified linear approximation, respectively. The slope of the linear curve corresponds to the stiffness parameter $k_1$.

### 3.3.2 Simplify of the Friction Model

The modelling of friction model have talk a lot in the before pages. It can be observed that complex friction models bring high accuracy while likewise adding more parameter dimensions. Taking the Stribeck model as an example, the modeling session introduces more parameters such as Stribeck velocity according to equation (3.5). When the load-dependent static friction torque is considered, as showed in equation (3.6), it will not only brings other parameters like $T_{s0}$ and $f_l$, also make the friction model relevant to the stiffness model, which will make the analyse of the whole EMB model even more complicated.

To limit the dimensional of the parameters in the reinforcement learning process, the primary consideration when modeling friction is to ensure the simplicity of the model. Consequently, the classical friction model, as presented in equation (3.4), is adopted due to its straightforward representation. Figure 3.5 illustrates the simplified friction model used in this study.

### 3.3.3 Summary of the EMB Model

After all the simplify of the early pages, finally a simplified lumped parameter nonlinear EMB model was introduced and showed in the state space:

$$\dot{x} = f(x, u), \quad y = h(x) \tag{3.11}$$

where $x = [x_1, x_2]^T$ is the state with the motor position $x_1$ and the motor velocity $x_2$ and the $u$ is the motor current as the input of the system. The simplified nonlinear function is considered as:

$$f(x, u) = \begin{bmatrix} x_2 \\ \frac{1}{J} \left( k_m u - \gamma F_{cl} - T_f \right) \end{bmatrix} \tag{3.12}$$

where $J$ is the lumped inertia of the EMB system, $k_m$ is the motor torque constant, $\gamma$ is the gear ratio, $F_{cl}$ is the clamp force, $T_f$ is the friction torque. In this simplified lumped parameter EMB model, we firstly consider the stiffness characteristics of the clamp force is linear function:

$$F_{cl}(x_1) = \begin{cases} k_1 x_1 & \text{if } x_1 \geq 0, \\ 0 & \text{if } x_1 < 0. \end{cases} \tag{3.13}$$

where $k_1$ is the stiffness characteristics parameter, which is not totally true but it is simple for the calculation. The friction model here only take coulomb and viscous friction in to consideration, given by:

$$T_f(x_2) = T_c \operatorname{sign}(x_2) + f_v x_2 \tag{3.14}$$

where $T_c$ and $f_v$ are the coulomb friction torque and viscous friction coefficients.

The system output also the observation is the motor position $x_1$ and the motor velocity $x_2$, hence:

$$h(x) = [x_1, x_2]^T \tag{3.15}$$

# 4 Method

Following the development of the EMB model in Chapter 3, an investigation was initiated into the input excitation to this model. As demonstrated in Chapter 2, the optimization problem associated with identifying the input excitation of the system parameters has been extensively studied. In this process, a variety of methods, including MPC and RL, are employed. However, these methods are not without limitations. In particular, reinforcement learning methods have been applied only in localized areas and still have significant shortcomings. This chapter will highlight the methodology, propose a new reinforcement learning framework for solving optimal input excitation for system parameter identification, and describe the different theoretical methods used in this structure from the beginning to the end. It will also describe how to perform the optimal experimental design for the above EMB model, how to combine it with reinforcement learning, and the theoretical methods used in this process.

## 4.1 Formulation the EMB model for the optimization of input excitation

Consider this model will be used in later research with reinforcement learning, to prevent unrealistic situation from happening, constraints must also be set. The first consideration is the input of the EMB model, which is contingent upon the technical parameters of the motor. It can be deduced that the input voltage of the motor is limited to a range of -6V to 6V. Secondly, the range of motion of the motor is constrained by the physical limitations of the system. The motor is unable to rotate at an infinite angle; therefore, the rotation angle of the motor is set to be between -10 rad and 100 rad. Meanwhile, according to the The test data obtained from the EMB model in the laboratory indicates that the maximum rotational speed of the motor is 500 rad/s. Over speed of the motor will not only leads to greater friction, generating excessive heat, but also decrease the control precision of the system. The motor may behave too fast for the response controller, leading to braking force errors and potentially compromising vehicle safety and comfort. Consequently, the motor rotational speed has been set to a range between -500 rad/s and 500 rad/s.

As motor speed sensors are employed in the actual parameter identification process, the measurements of the sensors are discrete sequences. Consequently, the continuous state space equations mentioned above must be discredited. There are numerous methods of discrediting the system, including Euler's method, RK-4 method, and so forth. However, to ensure the optimal simplicity of the model, Euler's method is utilized in this thesis

$$x_d(k + 1) = x_d + t_s f(x_d(k), u(k)) \tag{4.1}$$

where Ts is the time step. Since the sampling frequency of the sensor used in the experimental platform is 1000 Hz, ts is determined to be 0.001s in this model accordingly.

Furthermore, the sensor noise utilized to quantify the system's input and output must be considered during the model's construction, as it will directly influence the accuracy of the parameter identification. In accordance with the specifications of the actual task, the EMB model investigated in this thesis is equipped with a high-precision angle sensor, but lacks a speed sensor. The motor speed is provided by a kind of state observer, which introduces a relative higher degree of error. The sensor noise is typically modeled as Gaussian noise, which obeys a Gaussian distribution $v_o \sim \mathcal{N}(0, R)$ with mean $0$ and variance $R$. For the angle sensor, the variance of the Gaussian noise is $1 \times 10^{-6}$. For the state observer used to observe the motor speed, the variance of the Gaussian noise is $1$. The input signal is presumed to be entirely free from nose, leading to the assumption that $v_i = 0$. After the discrediting the system and considering the input noise, the EMB model in the subsection 3.3.3 will turn to be like

$$
\begin{aligned}
x_d(k+1) &= f_d(x_d(k), u(k)) \\
y_d(k) &= h_d(x_d(k)) \\
y_m(k) &= h_d(x_d(k)) + v_o(k), \quad k = 1, ..., n
\end{aligned}
\tag{4.2}
$$

where the state variable $x_d(k) = [x_1(k), x_2(k)]^T$ represents the motor position and motor velocity at time step $k$. Here, denotes the system output $y_d(k)$ at time step $k$, while corresponds to the measured output $y_m(k)$ at the same time step.

In conclusion, a streamlined discrete EMB model has been developed to facilitate optimal input excitation for the investigation of reinforcement learning and parameter identification. Based on the conventional parameter identification approach and the simulation model outcomes, parameters that have been included in the EMB system has been identified in the table 4.1.

In these parameters, the lumped inertia of the EMB system $J$ and total gear ratio $\gamma$ could be seen as fixed and provide by the producer, so not included in the parameter estimation. The motor constant $k_m$ and the stiffness parameter $k_1$ will change over the use of EMB, so it is two parameter to be estimated. The friction parameters are also critical for the EMB system, as they directly affect braking performance. These parameters can be influenced by factors such as temperature, wear, and contamination, which can alter the system's response and reliability. So at the end, the to estimated EMB parameters $\theta$ have 4 dimensional

$$
\theta = [k_m, k_1, T_c, f_v]^T
\tag{4.3}
$$

## 4.2  Optimal experimental design and Fisher information matrix

Optimal experimental design is an approach to design a input trajectory to excited the system and inferring the parameter from the output trajectory. In this section, the formulation of optimal experimental design (OED) in the context of the EMB model was discussed, focusing on concepts such as Fisher Information Matrix (FIM) and sensitivity analysis.

Optimal Experimental Design is a methodology used to design an input trajectory that effectively excites the system, enabling accurate inference of system parameters from the resulting output trajectory. The

| Parameter | Value | Unit |
|:---:|:---:|:---:|
| **Motor** | | |
| $J$ | $4.624 \times 10^{-6}$ | kg·m$^2$ |
| $k_m$ | $2.170 \times 10^{-2}$ | m/rad |
| $\gamma$ | $1.889 \times 10^{-5}$ | m |
| **Stiffness** | | |
| $k_1$ | $2.304 \times 10^{1}$ | N/rad |
| **Friction** | | |
| $T_C$ | $1.037 \times 10^{-2}$ | N·m |
| $f_v$ | $2.160 \times 10^{-5}$ | N·s/rad |
| **Constraints** | | |
| $u_{max}$ | 6 | V |
| $u_{min}$ | -6 | V |
| $x_{1,max}$ | 100 | rad |
| $x_{1,min}$ | -10 | rad |
| $x_{2,max}$ | 500 | rad/s |
| $x_{2,min}$ | -500 | rad/s |
| **Noise** | | |
| $R_{x_1}$ | $1 \times 10^{-6}$ | - |
| $R_{x_2}$ | 1 | - |

Table 4.1: Parameters for the Simplified EMB Model

goal of OED is to maximize the information obtained from experiments, and a key metric for assessing this information is the Fisher Information Matrix (FIM), FIM takes the form [32]

$$
\begin{aligned}
M(t) &= \int_0^t \left( \frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}} \right)^{\top} R^{-1} \left( \frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}} \right) dt, \\
&= \int_0^t \left( \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(t)}{\partial \boldsymbol{\theta}} \right)^{\top} R^{-1} \left( \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(t)}{\partial \boldsymbol{\theta}} \right) dt
\end{aligned}
\tag{4.4}
$$

In equation (4.4), the part $\frac{\partial \mathbf{h}}{\partial \boldsymbol{\theta}}$ is called sensitivity. It reflects the relationship between system outputs and the unknown parameters. Sensitivity analysis determines how changes in the input trajectory affect the system output with respect to the parameters. If a particular input trajectory causes significant changes in the output with respect to the parameters, it implies that the parameters are highly sensitive to this input, suggesting that the input trajectory contains a high level of information about the parameters. Thus, sensitivity analysis is essential for evaluating how effective a given input is for parameter estimation.

Furthermore, FIM incorporates the effect of measurement noise. As previously stated, measurement noise was seen as Gaussian noise with a mean value of zero. In multiple-input and multiple-output (MIMO) systems, the greater the noise in the output measurements of a given dimension, the more inaccurate the measurements become. This results in a reduction in the information content of that dimension's outputs generated in the parameter identification process. The inverse matrix of the output

noise variance matrix $R^{-1}$ is included in equation (4.4). This form amplifies the amount of information contained in the output dimension with measurements variance less than 1 and reduces the amount of information contained in the output dimension with measurements variance greater than 1.

The Fisher Information Matrix (FIM) measures the amount of information available in the output for estimating the system parameters. It quantifies how much the system output changes in response to variations in the parameters, thereby indicating the quality of the parameter estimation. Under the assumption that parameter estimates are unbiased and that measurement noise is uncorrelated, the inverse of the FIM provides an approximation of the Cramér-Rao lower bound [32]. This bound represents the theoretical lower limit for the variance of unbiased estimators, closely related to the variance-covariance matrix of the estimated parameters. By maximizing some criteria of the FIM, the lower bound is minimized , thus let the system have the possibility to get a lower variance of unbiased estimators, which means improving the precision of parameter estimation. Therefore, when designing the input excitation trajectory, criteria for input excitation are often evaluated using scalar measures of the FIM $M$ and the inverse matrix of FIM $M^{-1}$, including [33]:

1) D-optimality: Minimizes the determinant of inverse of the FIM $|M^{-1}|$, ensuring that the volume of the confidence ellipsoid for parameter estimates is minimized. Also, the convex function $log|M^{-1}|$ impels the mutual information between system output and parameters [34].

2) A-optimality: Minimizes the trace of the inverse of the FIM $tr(M^{-1})$, reducing the average variance of parameter estimates.

3) E-optimality: Minimizes the biggest eigenvalue of inverse of the FIM $\lambda_{max}(M^{-1})$, minimizing the length of the largest uncertainty axis within the joint confidence region of the parameters, effectively reducing the largest possible parameter errors.

Suppose a system model has two parameters that need to be identified, according to Equation (4.4), FIM $M$ is a semi-positive definite symmetric matrix, so the inverse matrix of FIM also take the form

$$M^{-1} = \begin{bmatrix} m_{11} & m_{12} \\ m_{12} & m_{22} \end{bmatrix} \tag{4.5}$$

The estimates $\hat{\theta}_1$ and $\hat{\theta}_2$ are the unbiased estimators for the parameters. The elements $m_{11}$ and $m_{22}$ represent the diagonal entries of the matrix $M^{-1}$, while its eigenvalues are $\lambda_1$ and $\lambda_2$. Figure 4.1 intuitively shows an example for these three optimality with this parameters in two dimensions. The blue-shaded area in the figure represents the area of the confidence ellipse. Given that the two parameters constitute a two-dimensional planar space, the D-optimality is reflected in the minimization of the ellipse area. The semi-axis lengths of the two axes of the confidence ellipse correspond to the two eigenvectors of the FIM inverse matrix. Therefore, the E-optimality is reflected in the minimization of the scales of the axes of the confidence intervals. The A-optimality is aimed at minimizing the trace, i.e., $m_{11} + m_{22}$, of the FIM inverse matrix. This is equivalent to the overall reduction of the uncertainty in all directions.

In summary, OED involves designing inputs that maximize the information available for parameter estimation. The FIM is a fundamental tool in this process, the design of optimal input trajectory based on optimality criteria, while sensitivity analysis helps determine the effectiveness of these inputs in estimating the system parameters accurately.
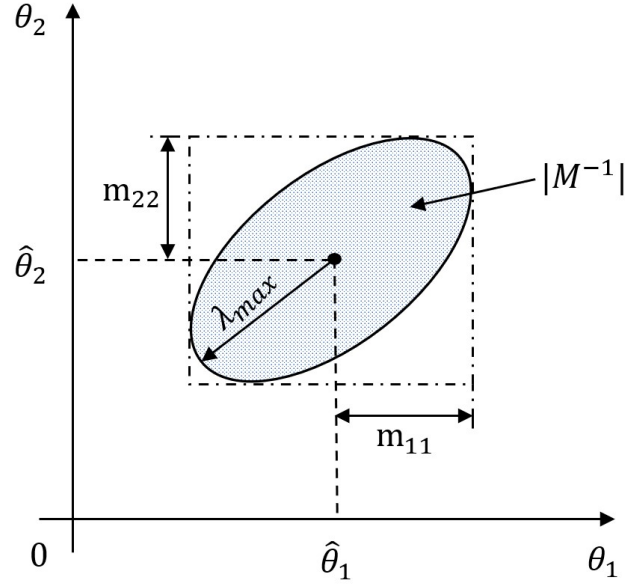
Figure 4.1: Confidence ellipsoid for parameters in two dimensions based on paper [33]

## 4.3 Formulation of the optimal experimental design in EMB system

In this section, the formulation of an optimal input excitation for a time-discrete EMB system is presented, incorporating the Fisher Information Matrix into an optimization problem. Assume the existence of a time-discrete EMB model, as described in Section 4.1 and Equation (4.2). Given a trajectory $U$ over $n$ steps, the objective is to determine the optimal input trajectory by solving the following optimization problem

$$\min_{U} \begin{cases} |M^{-1}| \text{ or } \log|M^{-1}|, & \text{for D-optimality,} \\ \text{tr}(M^{-1}), & \text{for A-optimality,} \\ \lambda_{\max}(M^{-1}), & \text{for E-optimality.} \end{cases} \tag{4.6}$$

Here, $M$ is the Fisher Information Matrix, and the optimization criteria (D-, A-, or E-optimality) are chosen depending on the desired objective: maximizing determinant-based informativeness, minimizing overall parameter uncertainty, or reducing the worst-case uncertainty.

The optimization is subject to the following boundary conditions and constraints

$$
\begin{aligned}
x_d(1) &= x_0, \\
x_{1,min} &\leq x_1 \leq x_{1,max}, \\
x_{2,min} &\leq x_2 \leq x_{2,max}, \\
u_{min} &\leq u \leq u_{max}, \\
k &\leq t_{max}
\end{aligned}
\tag{4.7}
$$

where $x_d$ is the vector of state variables, $u$ is the system input trajectory, $k$ is the time step and $t_{max}$ is the maximum number of steps defined for the parameter estimation process. The constraints ensure

that the system operates within physical and operational limits, such as state bounds, input saturation and time restriction.

In the optimal experimental design specific to the EMB system, further detailed constraints were considered. According to the requirement of experiment, the EMB motor will start from a zero position and velocity. The experiment will last for 500 ms, and at the end of this period, the motor should return to the zero position. These experimental setup requirements will be implemented later during the construction of the EMB environment using the Gymnasium framework and the design of the reward function for reinforcement learning.

### 4.3.1 Reformulating the optimization problem for reinforcement learning

In reinforcement learning, the goal is to train an agent to maximize the return over an episode. Therefore, the minimization problem in Equation (4.7) cannot directly serve as the return for the RL agent. To solve the excitation design task, this optimization problem is reformulated into a framework that can combine with reinforcement learning methods. One intuitive approach is to simply apply a negative sign to the objective function to convert it into a maximization problem. While this approach is feasible, there is still another problem that exists.

Since the FIM is changing over time during the parameter estimation experiment, as shown in Equation (4.4), the objective function still requires calculating the inverse of the FIM, $M^{-1}$, at each time step. If the dimension of the matrix $M$ is large, this significantly increases computational difficulty and causes longer computation times. Therefore, the objective function for the optimization problem needs to be modified to avoid matrix inversion operations when practically applied to reinforcement learning. The approach for this modification will be presented below.

Suppose that the FIM is an $n$-dimensional nonsingular matrix with eigenvalues $\lambda$ that satisfy the following relationship

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n \tag{4.8}$$

the determinant of the FIM can be expressed as:

$$|M| = \prod_{k=1}^{n} \lambda_k, \tag{4.9}$$

and the determinant of the inverse of the FIM is given by:

$$|M^{-1}| = \prod_{k=1}^{n} \frac{1}{\lambda_k} = \frac{1}{|M|}. \tag{4.10}$$

When the D-optimality criterion with logarithm is applied, the object function becomes

$$\log|M^{-1}| = \log|\frac{1}{M}| = -\log|M| \tag{4.11}$$

With this approach, the optimization problem for the D-optimality criterion is transformed into a maximization problem, avoid to calculate the inverse matrix at each time step. Similarly, under the E-optimality criterion, the objective becomes

$$\min_{U} \lambda_{\max}(M^{-1}) \rightarrow \max_{U} \lambda_{\min}(M). \tag{4.12}$$

However, for the A-optimality criterion, which involves the trace of the inverse FIM

$$\mathrm{tr}(M^{-1}) = \sum_{k=1}^{n} \frac{1}{\lambda_k} \neq \frac{1}{\mathrm{tr}(M)}, \tag{4.13}$$

Hence the calculation of the inverse FIM is still necessary for the A-optimality criterion. Therefore, due to the consideration of the simplicity of computation, D-optimality and E-optimality criteria are primarily considered in this thesis.

### 4.3.2 Normalization of parameters

The advantage of D-optimality over the other two optimality criteria is its independence from the numerical scales of the parameters. In the model studied in this thesis, as shown in Table 4.1, the magnitude of the a priori values of each parameter ranges from $10^{-5}$ to $10^{1}$. It can be concluded from the calculation of sensitivity that the system output is generally more sensitive to parameters with smaller numerical scales. Since significantly value difference could existed between different parameters, in order to keep the variance of the different parameters comparable, the parameters need to be normalized when applying the other two optimality criteria. The normalised parameter vector for a $n$ dimensional parameter vector could be written as:

$$\bar{\theta} = \left[ \frac{\theta_1}{\theta_{nom,1}}, \frac{\theta_2}{\theta_{nom,2}}, \ldots, \frac{\theta_n}{\theta_{nom,n}} \right]^T \tag{4.14}$$

### 4.3.3 Normalization of outputs

Since sensitivity depends on the magnitude of the change in outputs with respect to parameter variations, different output scales in a multi-output system can significantly influence the calculation of FIM and sensitivity. For instance, if some parameters in a system only affect specific output dimensions, and the range of these output dimensions is significantly larger than others, the corresponding dimension will contribute larger values to the FIM. This implies that the input excitation will yield higher information content for the parameters affecting those output dimensions. Therefore, to ensure balanced sensitivity and an accurate representation of information content, not only do the parameters need to be normalized, but the output scales should also be normalized.

For a system with $p$ output dimensions, represented as $y = h(x) = [y_1, \ldots, y_p]$, normalization can be applied as follows:

$$\bar{h} = D_y^{-1} h(x) \tag{4.15}$$

where $D_y$ is a diagonal matrix defined as:

$$D_y = diag([y_{1,max}, y_{2,max}, \ldots, y_{p,max}]) \tag{4.16}$$

This normalization ensures that all outputs are brought to a comparable scale, thereby preventing any single output dimension with a larger range from disproportionately influencing the FIM. By normalizing the outputs in this way, a more balanced sensitivity analysis is achieved, which helps in obtaining an accurate estimation of the parameters and ensures that the information content from all output dimensions is represented equitably.

### 4.3.4 Calculation of the Fisher information matrix <mark>with discrete system</mark>

As described in the Equation (4.4), the calculation of FIM involves the sensitivity $\frac{\partial h}{\partial \theta}$. Since the EMB model is discretized using Euler's method, the equation (4.4) must also be transferred into discretized version. Consider the simplified lumped parameter nonlinear time-discrete EMB model with measurement noise as showed in (4.2)

$$
\begin{aligned}
x_d(k+1) &= f_d(x_d(k), u(k)) \\
y_m(k) &= h_d(x_d(k)) + v_o(k), \quad k = 1, ..., n
\end{aligned}
\tag{4.17}
$$

The Fisher information matrix with n-dimensional parameters after the normalization of parameters and outputs, can be calculated using the following formulation [2]

$$
M = \sum_{k=1}^{n} \left( \frac{d\bar{h}_k}{d\bar{\theta}} \right)^T R^{-1} \left( \frac{d\bar{h}_k}{d\bar{\theta}} \right) = \sum_{k=1}^{n} (M_{\bar{\theta},k}),
$$
$$
\frac{d\bar{h}_k}{d\bar{\theta}} = \begin{bmatrix} \frac{d\bar{h}_k}{d\theta_1} & \frac{d\bar{h}_k}{d\theta_2} & \cdots & \frac{d\bar{h}_k}{d\theta_n} \end{bmatrix},
\tag{4.18}
$$

with the time steps $k$ written in the subscripts. The sensitivity of each output with respect to each parameter is obtained as

$$
\frac{d\bar{h}_k}{d\bar{\theta}_i} = \frac{\partial \bar{h}_k}{\partial x_d^T} \frac{\partial x_{d,k}^T}{\partial \bar{\theta}_i} + \frac{\partial \bar{h}_k}{\partial \bar{\theta}_i}
\tag{4.19}
$$

Here, $\frac{\partial x_{d,k}^T}{\partial \theta}$ represents the sensitivity of each state variable with respect to each parameter. In the continuous domain, this sensitivity can be calculated by integrating[4]

$$
\frac{d}{dt} \frac{\partial x}{\partial \theta}(t) = \frac{\partial f}{\partial x} \frac{\partial x}{\partial \theta}(t) + \frac{\partial f}{\partial \theta}
\tag{4.20}
$$

with initial condition

$$
\frac{\partial x}{\partial \theta}(0) = \frac{\partial x_0}{\partial \theta}
\tag{4.21}
$$

Transfer this function into the discrete domain, it turns int

$$
\frac{\partial x_{d,k+1}^T}{\partial \bar{\theta}_i} = \frac{\partial f_{d,k}}{\partial x_d^T} \frac{\partial x_{d,k}^T}{\partial \bar{\theta}_i} + \frac{\partial f_{d,k}}{\partial \bar{\theta}_i},
\tag{4.22}
$$

With the consideration of no sensitivity or information is present at the beginning of the parameter estimation experiment, the initial condition of the calculation is given by

$$
\frac{\partial x_{d,1}^T}{\partial \theta} = 0
\tag{4.23}
$$

Using this approach, the Fisher information matrix in the discrete system is computed iteratively at each time step, and the final FIM is obtained as the number of total time steps $n$ reaches.

### 4.3.5 Implementation of the Calculation of the Fisher information matrix in Python

After completing the theoretical formulation of the calculation of FIM, the corresponding model in Python is implemented. Since the equations in the before section basically always calculate the gradient, automatic differentiation tools are suitable for this task. This implementation utilizes PyTorch's Autograd feature for automatic differentiation [35], which simplifies calculations of the FIM.

In the Python script, equations from the previous subsection is encapsulated in a class named $FI\_matrix$, which simulates the EMB system using a discrete representation and estimates the FIM iteratively. Functions contained in the class including:

$f(x, u, theta)$: Defines the state update equations of the EMB model, including dynamic components such as motor torque, load torque, and friction, as showed in equation (4.18).

$h(x)$: Represents the output function that provides system outputs, such as motor position and velocity.

$jacobian(x, u, theta)$: Uses PyTorch's Autograd to compute the Jacobian of the state function f with respect to both the state variables and system parameters.

$sensitivity\_x()$ and $sensitivity\_y()$: Calculate the sensitivity of the states and outputs with respect to the system parameters recursively, as described in equation (4.19) and (4.22).

$fisher\_info\_matrix()$: Computes the FIM by accumulating the sensitivity information over all time steps.

The use of PyTorch's Autograd for automatic differentiation greatly simplify the calculation of Jacobians, ensuring accurate and efficient sensitivity computation throughout the simulation. Furthermore, as this script is incrementally updated, at each time step not only can the FIM be readily computed, but also its corresponding state variables and system outputs, as well as the sensitivity of each state variable and output with respect to each parameter. This enables the monitoring and evaluating of the trend of FIM and sensitivity at an episode, which is a valuable contribution to the next reinforcement learning framework.

## 4.4 Initialization and normalization of the Fisher Information Matrix

After the automated computational FIM was written, predetermined inputs were used to test the performance of the script. Different system input trajectories were selected during the testing process, such as step inputs, ramp inputs, and sinusoidal inputs. In comparison to the other two optimality criteria, the performance of the D-optimality criterion was initially evaluated, as it represents the optimality in the parameter space.

Figure 4.2 illustrates the progression of the system input $u$, state parameters, $x_1$ and $x_2$, the determinant of the FIM, denoted by $|M|$, the logarithm of the determinant, represented by $log(|M|)$, and the condition number of the FIM, $\kappa(M)$, over time in response to a constant input trajectory. However, unexpected behavior was noted when calculating the log determinant of the FIM. Figure 4.2(c) shows he log determinant curve is discontinuous. Further examination of the determinant values from Figure 4.2(b) revealed that this discontinuity was caused by negative values for the determinant of the FIM at certain time points, which caused the logarithm function to return a NaN (not-a-number) value. Since the FIM is a symmetric semi-positive definite matrix, its determinant should not be negative.
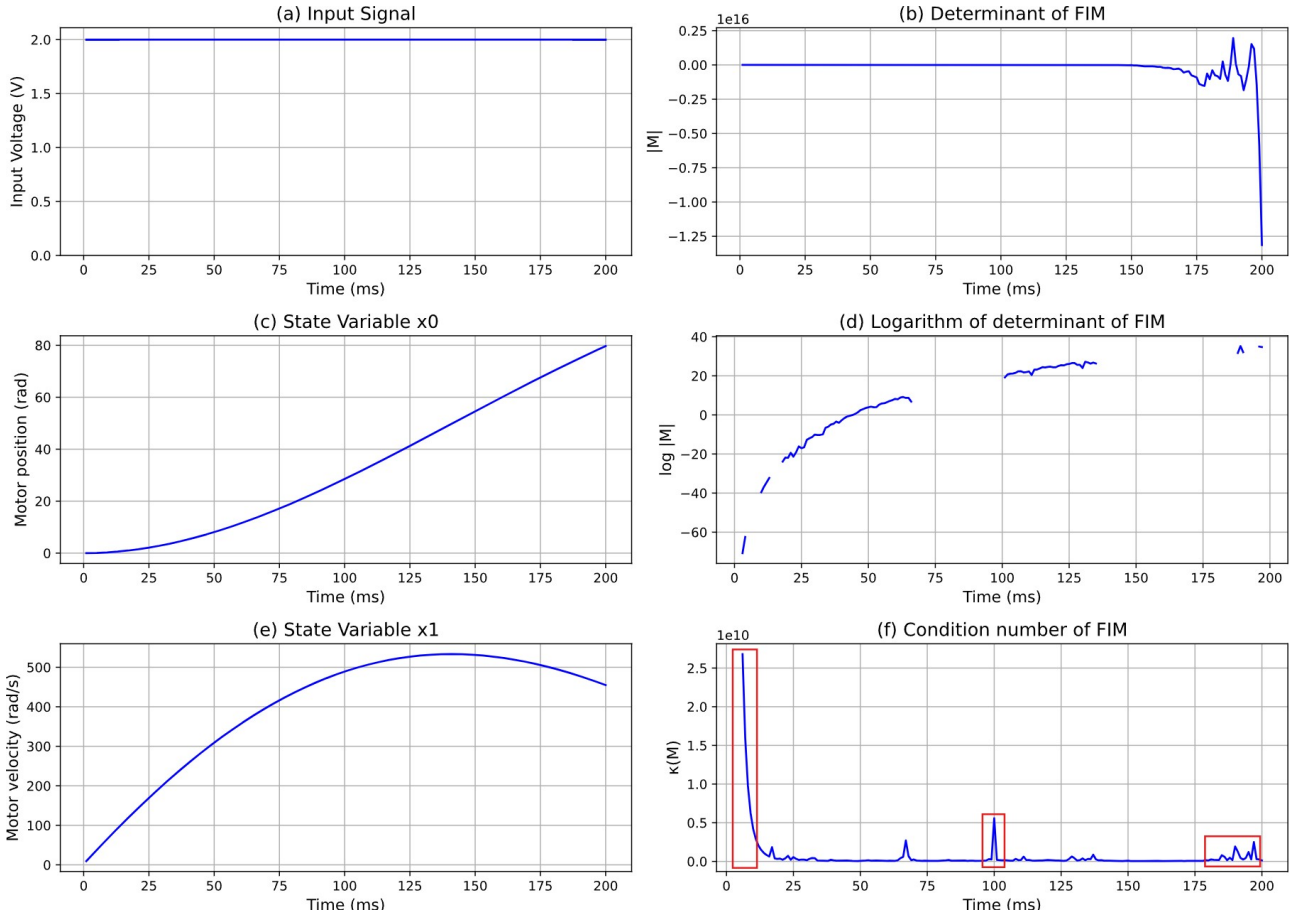
Figure 4.2: System Response Over Time to a Constant Input Trajectory: (a) System input $u$. (b,c) State parameters $x_1$ and $x_2$. (d) Determinant of the FIM $|M|$. (e) Logarithm of the determinant $log(|M|)$. (f) Condition number of the FIM, $\kappa(M)$

The erroneous negative determinant values were likely due to floating-point errors inherent in the determinant calculation, particularly when the matrix approaches singularity. The FIM was found to be approximately singular at certain time points, as evidenced by a sudden increase in the conditional number of the matrix at those points. As the red box showed in Figure 4.2(f), during the experiment, the conditional number of the FIM spiked unexpectedly at several intervals, which resulted in the matrix becoming ill-conditioned and leading to numerical instability during determinant calculation.

In reinforcement learning, where the input trajectory is determined by the agent, it is not feasible to anticipate when the FIM will become nearly singular and design the input to avoid this issue. Furthermore, the determinant of the FIM tends to increase rapidly over time, which could lead to potential overflow issues. To mitigate this problem, method was developed to scale the FIM matrix, ensuring numerical stability and preventing such issues during determinant calculation.

Furthermore, a step reward function was devised based on the characteristics of reinforcement learning in conjunction with the D-optimality criterion. The issue that needed to be addressed was that, due to the incremental nature of FIM calculation—where the FIM at each time step is equivalent to the FIM from the previous time step plus a rank-1 update—the determinant of the FIM cannot become

non-zero until at least $n$ steps have been completed. This behavior presents a challenge for calculating the reinforcement learning reward function in the initial phase. Consequently, an initialization method for the FIM was developed and demonstrated to have no adverse effect on the optimization results.

### 4.4.1 Initialization of the Fisher Information Matrix

In the design of the calculation of the FIM, an assumption is made

$$M_0 = 0 \tag{4.24}$$

indicating that no information is obtained at the beginning of the experiment. Since at each time step during the training, the FIM is updated by a rank-1 update, initializing the FIM as in Equation (4.24) poses an issue: at least during the initial $n$ steps, the step reward cannot be calculated because $rank(M_k)$ is less than $n$. The determinant of an $n$-dimensional matrix with rank less than $n$ is zero. To address this, the FIM is initialized with a sufficiently small value $\epsilon$ multiplied by an identity matrix $I$ of appropriate dimension:

$$M_{init} = \epsilon I \tag{4.25}$$

As add an identity matrix will change the FIM in the end, it could be prove the affection could be neglected with Courant-Fisher theorem. To proceed, the Courant-Fischer theorem is first introduced and used to prove the following inequality

$$\lambda_i(A) + \lambda_1(B) \leq \lambda_i(A + B) \leq \lambda_i(A) + \lambda_n(B) \tag{4.26}$$

with eigenvalues $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$

**Theorem 1** (Courant-Fischer Theorem). *[36] Let $\mathbf{M}$ be a symmetric matrix with eigenvalues $\mu_1 \geq \mu_2 \geq \cdots \geq \mu_n$. Then,*

$$\mu_k = \max_{\substack{S \subseteq \mathbb{R}^n \\ \dim(S)=k}} \min_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}^T \mathbf{M} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \min_{\substack{T \subseteq \mathbb{R}^n \\ \dim(T)=n-k+1}} \max_{\substack{\mathbf{x} \in T \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}^T \mathbf{M} \mathbf{x}}{\mathbf{x}^T \mathbf{x}},$$

*where the maximization and minimization are over subspace $S$ and $T$ of $\mathbb{R}^n$.*

Consider a matrix $M = A + B$ with a unit vector $\mathbf{x}$, by using the Courant-Fisher Theorem, the $i$-th eigenvalue of the matrix $M$, $\lambda_i$ can be characterized as

$$\lambda_i = \max_{\substack{S \subseteq \mathbb{R}^n \\ \dim(S)=k}} \min_{\substack{\mathbf{x} \in S \\ \mathbf{x} \neq 0}} (\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{B} \mathbf{x}) \tag{4.27}$$

because $\mathbf{x}^T \mathbf{x} = 1$. Could be noticed that

$$\mathbf{x}^T \mathbf{B} \mathbf{x} \geq \lambda_1(\mathbf{B}) \tag{4.28}$$

for any unit vector $\mathbf{x}$, because $\lambda_1(\mathbf{B})$ is the smallest eigenvalue of $\mathbf{B}$. Thus

$$\lambda_i(\mathbf{A} + \mathbf{B}) \geq \max_{\substack{\dim(S)=i \\ \mathbf{x} \in S, \|\mathbf{x}\|=1}} \min \left( \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda_1(\mathbf{B}) \right),$$
$$= \lambda_i(\mathbf{A}) + \lambda_1(\mathbf{B}) \tag{4.29}$$

Similarly could be proved that

$$
\lambda_i(\mathbf{A} + \mathbf{B}) \leq \max_{\substack{\dim(S)=i \\ \mathbf{x} \in S, \|\mathbf{x}\|=1}} \min \left( \mathbf{x}^T \mathbf{A} \mathbf{x} + \lambda_1(\mathbf{B}) \right),
$$

$$
= \lambda_i(\mathbf{A}) + \lambda_n(\mathbf{B})
$$

(4.30)

Thus the inequality (4.26) is proved.

When the FIM $M_n$ is initialized with a sufficiently small value $\epsilon$ multiplied by an identity matrix $I$, inequality (4.26) becomes the following equation

$$
\lambda_i(M_n + M_{init}) = \lambda_i(M_n) + \epsilon
$$

(4.31)

Thus, initializing the FIM with $M_0 = \epsilon I$ ensures that the matrix remains non-singular while having a negligible effect on its eigenvalues. This allows the objective function to be evaluated consistently from the start, ensuring stability in the early stages of the experiment without compromising the optimization process or the accuracy of the results.

### 4.4.2 Normalization of the Fisher Information Matrix

Another observation from Figure 4.2 is that during the experiment, the value of $det\,(M_k)$ can easily grow to an extremely large number due to the large values of the matrix elements. This can lead to a consequence where the condition number of the FIM also becomes large, making the calculation of the determinant unstable. To restrict the growth of $det\,(M_k)$, a normalization approach with normalization factor $\alpha_k$ is designed

$$
\bar{M}_k = \alpha_k M_k
$$

(4.32)

where $\bar{M}_k$ is the normalized FIM.

After the normalization, it is important to keep the step reward unchanged, otherwise it will change the object function of the optimization problem. This is possible due to the properties of logarithms and determinants, which allows Equation (5.12) to be rewritten in

$$
\begin{aligned}
r_k &= log|M_k| - logdet|M_{k-1}| \\
&= log\frac{|M_k|}{|M_{k-1}|} \\
&= log\frac{\alpha_k|M_k|}{\alpha_k|M_{k-1}|} \\
&= log|\bar{M}_k| - log|\frac{\alpha_k}{\alpha_{k-1}}\bar{M}_{k-1}|
\end{aligned}
$$

(4.33)

Thus, the original FIM at time step $k$ is not needed for the calculation of step reward $r_k$.

Due to Equation (4.2), the FIM is actualization by rank-1 update, FIM at time step $k$ is

$$
M_k = M_{k-1} + M_{\bar{\theta},k}
$$

(4.34)

Hence, Equation (4.32) changes into

$$
\begin{aligned}
\bar{M}_k &= \alpha_k M_k \\
&= \alpha_k (M_{k-1} + M_{\bar{\theta},k}) \\
&= \frac{\alpha_k}{\alpha_{k-1}} \bar{M}_{k-1} + \alpha_k M_{\bar{\theta},k}
\end{aligned}
\tag{4.35}
$$

Thus, the update for the normalized FIM in this approach also have no requirement for the calculation of original FIM, avoiding the occurrence of a potentially large determinant.

One reasonable way to choose the initial normalized FIM is make it same to the initial FIM in last subsection, and to design $\alpha_k$ is to always keep $det\left(\bar{M}_k\right)$ related to the determinate of the initial matrix

$$
\alpha_k = \left( \frac{det\left(M_{init}\right)}{det\left(\bar{M}_{k-1}\right)} \right)^{\frac{1}{n}}
\tag{4.36}
$$

An experiment with the same input signal after initialization and normalization was conducted and is shown in Figure 4.3. As observed from (b) and (d), the determinant of FIM is positive and logarithm is now continues. From (d), it can be seen that the condition number does not suddenly increase, but grows continuously as the experiment progresses.

## 4.5  Reinforcement Learning

In this thesis, reinforcement learning methods will be used to carry out input excitation optimisation used for EMB parameter identification. Reinforcement learning is a machine learning method. The content of it can be abstracted into two parts, the agent and the environment. The agent can gain experience and learn by interacting with the environment, giving the agent the ability to accomplish set tasks autonomously. In the application of this thesis, the environment is the EMB and the agent is responsible for giving the input excitation used for parameter identification.

Reinforcement learning based on the Markov Decision Process (MDPs), it is an abstraction of the real world. An MDPs is defined by a tuple $(S, A, P, R, \gamma)$, where $S$ is the set of states, $A$ is the set of actions, $P$ represents the transition probability between states, $R$ is the reward function, and $\gamma$ is the discount factor for the future rewards [37]. Figure 4.4 illustrates this process. Starting from the initial state $s_0$, the agent takes an action $a_1$ that leads to a new state $s_1$ and receives a reward $r_1$, and then this process will continue in the same way. The goal of the agent is to find an optimal policy $\pi(a|s)$ which could based on the state $s$, do the action $a$ and maximize the expected cumulative return

$$
G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}
\tag{4.37}
$$

And the expected cumulative return based on strategy $\pi_\theta$ is

$$
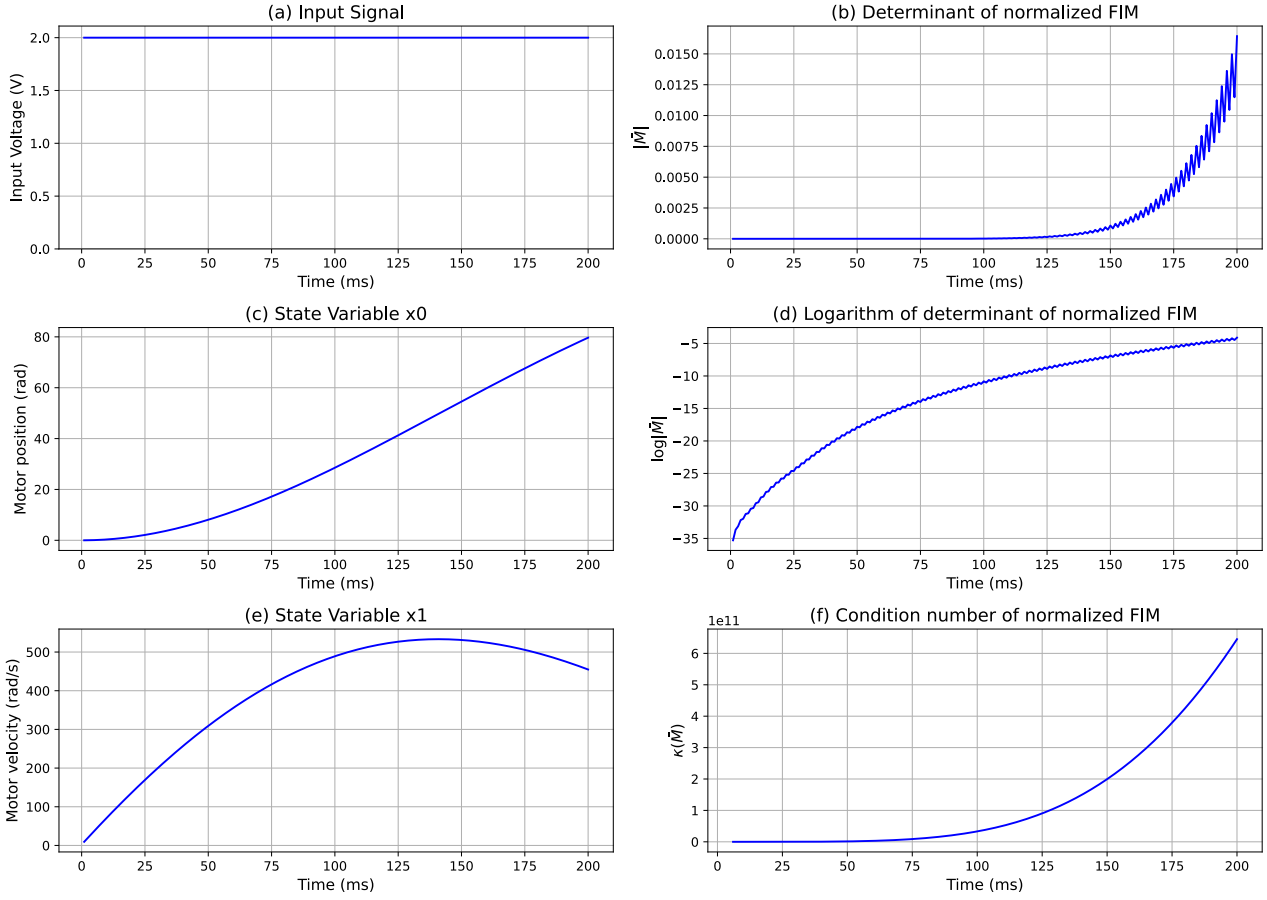J = \mathbb{E}_\pi(G_t)
\tag{4.38}
$$

Figure 4.3: System response over time to a constant input trajectory after initialization and normalization of FIM: (a) System input $u$. (b,c) State parameters $x_1$ and $x_2$. (d) Determinant of the FIM $|\bar{M}|$. (e) Logarithm of the determinant $log(|\bar{M}|)$. (f) Condition number of the FIM, $\kappa(\bar{M})$

Action-value function and state-value function build up the core part of the reinforcement learning. The action-value function $Q_\pi(s, a)$ describes the expected return if the agent takes action $a$ in state $s$ with a policy $\pi$

$$
\begin{aligned}
Q_\pi(s, a) &= \mathbb{E}\left(G_t \big| s_t = s, a_t = a, \pi\right) \\
&= \mathbb{E}\left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \big| s_t = s, a_t = a, \pi\right)
\end{aligned}
\tag{4.39}
$$

The state-value function $V_\pi(s)$ describes the expected return from a given state $s$ under a policy $\pi$

$$
\begin{aligned}
V_\pi(s) &= \mathbb{E}\left(G_t \big| s_t = s, \pi\right) \\
&= \mathbb{E}\left(\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \big| s_t = s, \pi\right) \\
&= \sum_{a \in A} \pi(a|s) Q_\pi(s, a)
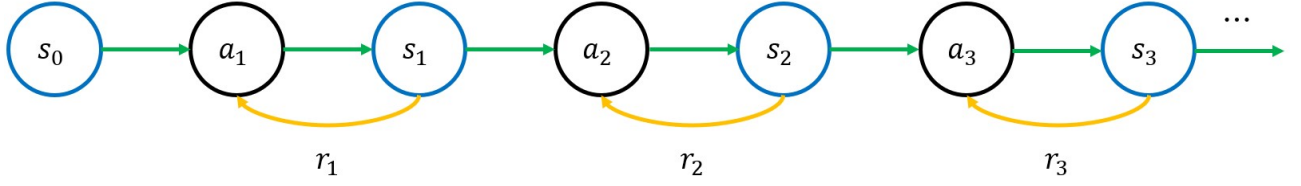\end{aligned}
\tag{4.40}
$$

Figure 4.4: Testing Actor with EMB environment after training

Almost all reinforcement learning algorithms based on the calculation of action-value function and state-value function. The relationship between those two function is build up by advantage function. In reinforcement learning, the advantage function $A$ is used to measure the advantage of the current action $a$ over the average behavior in the current state $s$ based on the current strategy $\pi$

$$A(s,a) = Q_\pi(s,a) - V_\pi(s) \tag{4.41}$$

Policy gradient methods is a kind of methods which directly optimize the policy $\pi$ [38]. The basic idea is to increase the probability of actions that can bring about positive advantages and decrease the probability of bringing about negative actions. Advance function is employed in this method to compute the gradient of the objective function

$$\nabla J(\pi) = \mathbb{E}\left[\nabla \log \pi(s,a) A(s,a)\right] \tag{4.42}$$

Traditional reinforcement learning methods, such as Q-learning and Temporal Difference (TD) methods, optimize an agent's policy by computing $Q$ and $V$ values. These methods have the advantage of being simple to construct and easy to understand, but they can only handle relatively simple problems. For example, in the case of Q-learning, maintaining a $Q$-table action space to be discrete, and in case of huge dimension of state space will directly lead to huge dimension of the $Q$-table, which limits its applicability to more complex scenarios.

With the development of deep learning and neural networks, reinforcement learning has gradually merged with deep learning. Since deep neural networks have the ability to fit arbitrary functions, Deep Reinforcement Learning (DRL) is theoretically almost unlimited in observation space and action space. This gives DRL the ability to solve complex problems in complex environments, such as robotic arm grasping and robot path planning problems.

So far, several DRL methods have been developed, such as Deep Q-Networks (DQN), Deep Deterministic Policy Gradient (DDPG), Twin Delayed DDPG (TD3), and Proximal Policy Optimization (PPO) [39]–[42]. Expect DQN and its variant, all otter DRL algorithm are based on Actor-Critic methods.

Compared to other DRL algorithms, PPO is known for its robustness. Since the application scenario in this thesis is a self-developed EMB model environment, unlike well-established reinforcement learning environments used wildly for benchmarking, the custom environment may be highly sensitive to the parameters of deep reinforcement learning algorithms, making tuning challenging. Compared to DDPG or A3C, PPO is relatively less sensitive to hyperparameter settings, making it a practical choice for this EMB environment that may be unstable or difficult to tune. Another advantage for PPO is, the clipping mechanism prevents large updates, making it more stable and robust compared to other actor-critic methods, which makes it suitable for the purposes of this thesis.

### 4.5.1 Actor-Critic methods

As showed in the name, two main components are included in this method: the Actor and the Critic. Actor is used to output the action, while Critic is responsible for give evaluation to those actions. In DRL, this two components are represented by two separate neural networks, policy network for Actor and value network for Critic. Policy network uses the policy gradient method combined advantage function $A(s,a)$ to update the network, as show in Equation (4.42).

For the critic network, the Temporal Difference (TD) error is used to minimize the difference between the predicted value ($V(s_t)$) and the actual return ($G_t$). The critic loss $L_V$ is typically computed as mean squared error (MSE) of the TD error

$$L_V = \frac{1}{2}\left(V(s_t) - G_t\right)^2 \tag{4.43}$$

### 4.5.2 Proximal Policy Optimization

Proximal Policy Optimization is a reinforcement learning algorithm base on Actor-Critic methods. It is an algorithm that achieves the data efficiency and reliable performance of Trust Region Policy Optimization (TRPO), while using only first-order optimization [42]. The object function of TRPO is introduced by [43]

$$J(\theta) = \hat{\mathbb{E}}_t\left[\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}\hat{A}_t\right] \tag{4.44}$$

TRPO introduces a Kullback-Leibler (KL) Divergence constraint to constrain the above objective function and uses fisher information matrix to compute the second-order approximation of the KL divergence. TRPO tries to convert this constrained optimization problem into an unconstrained optimization problem by introducing a hyperparameter $\beta$, but difficulty existing in selecting $\beta$.

The idea of PPO is to set a lower and upper bound for Equation (4.44), to represent the constraint and prevent the policy from making large updates. These bounds are converted into the following clip function with hyperparameter $\epsilon$

$$J^{CLIP}(\theta) = \hat{\mathbb{E}}_t\left[\min\left(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}\hat{A}_t, \operatorname{clip}(\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, 1-\epsilon, 1+\epsilon)\hat{A}_t\right)\right] \tag{4.45}$$

With this clipping mechanism, no KL divergence is needed for the calculation of the gradient of the object function, and the update of the police is strongly restricted. Making it have a more robust performance.

Different from other RL algorithms, the training of PPO doesn't based on episodes and stetting a maximum number of episode [44]. Instead, it use vectored environment and runs $N$ independent environments at the same time during the training. Each of these parallel environment collects $T$ timesteps of data in the rollout phase. If a done flag $d$ is triggered in any environment, that environment will be automatically reset. During the learning phase, the memory buffer with $NT$ steps data collected in the rollout phase will be used for the loss calculation and updating the networks.
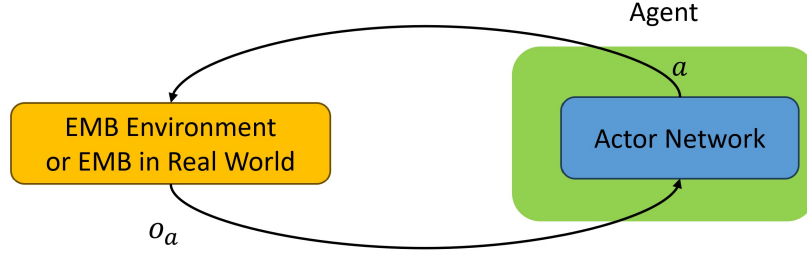
Figure 4.5: Testing Actor with EMB environment after training

### 4.5.3 Partially observable environments

In the reinforcement learning problem with EMB optimal input excitation, an important feature is partial observability. The FIM or sensitivity calculations described in the previous section cannot be accurately performed during actual parameter identification due to the lack of prior parameter values. An agent trained by reinforcement learning should only be able to acquire information from sensors to make decision and give the correct input trajectory. Thus, the gap between the theoretical simulation of optimal experimental design and its practical application needs to be taken into account in the framework design.

The goal of this reinforcement learning framework is to train the critic based on the full information from the EMB environment, which is used to accurately evaluate the behavior of the actor, thus guiding the actor to make the right decision. As showed in Figure 4.5, at the end of training, when the agent is actually used for generate the input excitation, only the actor network is employed to determine the corresponding system inputs based on the current system output measurements and time step, all the rest of the information in the environment stay unknown.

In fact, partial observability not unique to the EMB system, it also existed in many other area. For example in the robotics, robots are usually first trained in the simulator with all state information available. However, when the robot is deployed after training, it might only have images as input, which also leads to the partial observations problem. To address this problem, partially observable Markov decision processes (POMDPs) have been proposed [45]. A POMDPs is a MDP in which the agent is unable to observe the current state. Instead, it makes an observation based on the action and resulting state, and the goal of maximize expected discounted future reward remains the same. POMDPs are described by the tuple $\langle S, A, T, R, \Omega, O \rangle$ [45], while $S, A, T, R$ stay the same meaning with MDP, $\Omega$ is a set of possible observations that the agent can experience in its environment. $O : S \times A \to \Pi(\Omega)$ is the observation function, which provides a probability distribution over observations given a state and action.

POMDPs have been used in planning and control tasks [46], [47]. Most POMDPs applications rely on offline algorithms, the disadvantage of it is that the training time is too long, also it can not deal with small changes in the environment because the entire offline policy need to recalculated. To address these issues, several online POMDPs algorithms have been proposed and studied [48]–[50].

The Asymmetric Actor-Critic framework is variant of Actor-Critic framework, where full state in observation space is used to train the critic, while partial observation is used to trains the actor [51]. This framework has been applied to image-based robot learning and contextual learning tasks [52]. Although

| Network | Observation Space |
|---------|-------------------|
| Actor | $x_{1,m}, x_{2,m}, k$ |
| Critic | $x_1, x_2, k, \theta_{sample}, FIM$ |

Table 4.2: Asymmetric inputs for Actor and Critic

many asymmetric methods are effective, they are often evaluated in limited scenarios. Recent research has proved an asymmetric policy gradient theorem for partially observable control, an extension of the policy gradient Actor-Critic methods theorem[53].

Since PPO also based on Actor and Critic networks, it should be consider as a good way to introduce the symmetric Actor-Critic approach into PPO to deal with the existing partially observable space problem in EMB. The Detail of the asymmetric PPO framework will be conducted in the next section.

### 4.5.4 Asymmetric PPO framework

Inspired by the existing researches, in this thesis, an asymmetric PPO framework is introduced. As show in the Algorithm 1, compared to the classic PPO algorithm, this framework employs an asymmetric input strategy for the actor and critic networks. As showed in Figure 4.6, while the observation $o$ from the complete observation space is stored in memory after each step, the input to the actor network, $o_a$, contains only a partial observation. In contrast, the input to the critic network, $o_c$, contains more information and differs only slightly from the full observation $o$. Therefore, when designing and initializing the actor and critic networks, the different input structures must be taken into consideration.

In conjunction with the framework design specifically, the observation space of reinforcement learning contains all the information mentioned above. But the observation space of actor $o_a$ has only three dimensions, including the measured values of motor position and speed and the time step, while the observation space of critic $o_c$ contains the true values of motor position and speed, the a priori values of system parameters obtained from sampling the parameter distributions, as well as the elements of the time step and the elements of FIM. The difference of the observation space for those two network is listed in Table 4.2. During the development of the Asymmetric PPO framework, the PPO algorithm from $CleanRL$ is chosen as the fundamental reinforcement learning algorithm framework [54].

## 4.6 EMB environment setup for reinforcement learning

After solving the computational aspects of the optimal experimental design for the EMB system, one remaining challenge is build up an environment that can interact with reinforcement learning algorithms. In the field of reinforcement learning, a widely used standard interface is Gymnasium.

Gymnasium is an open-source library that provides a standard interface for reinforcement learning environments. It has features that allow different environments and training algorithms to work together easily, making it simpler to develop and test RL algorithms [55]. Gymnasium includes many classic reinforcement learning environments, such as $CartPole$, $Pendulum$, and $LunarLander$, etc. By using
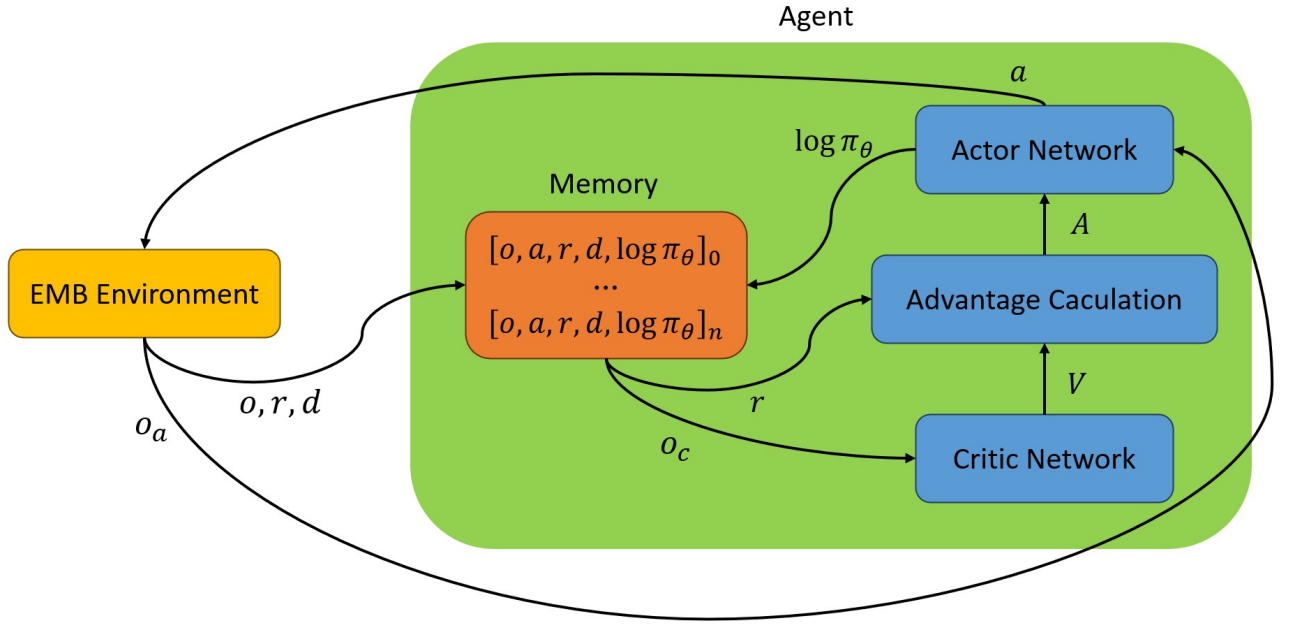
Figure 4.6: Training asymmetric PPO framework with EMB environment

these existing environments, interaction with RL algorithms can be easily achieved without the need to implement the environment, and most of the training results can be visualized.

In addition, Gymnasium provides standard tools and interfaces for reinforcement learning that allow for easy customization of environments. Since the EMB model in this thesis is a simplified simulation model based on a real test platform, building the corresponding reinforcement learning environment in Gymnasium is necessary for applying RL to the optimal experimental design.

The environment is the core part in Gymnasium. It is the instantiation of the RL environment and enables interaction with RL algorithms [55]. An environment consists of observation space, action space, $step$ method, and $reset$ method. The EMB model constraints introduced in Section 4.1 and the optimal experimental design requirements described in Section 4.3 will be implemented in these components. The following part describes how each section is built in a customized EMB environment.

The input to the EMB system is the motor voltage, which ranges from -6V to 6V, hence it corresponds to a one-dimensional action space. The design of the observation space determines the amount of information that the RL agent can obtain from the environment, and needs be considered in context of the task requirements and the structure of the RL algorithm. Parameters that can be contained in the observation space include the actual and measured values of motor position and speed, the prior values of parameters to be identified, the time step and elements in the FIM, etc. The $reset$ method is used to initialize the environment, which is mainly responsible for setting the parameters in the observation space to their initial values.

The $step$ function is the core of the reinforcement learning environment. In this part, the agent's action is input into the EMB system, and the motor position and speed in next step are updated. Additionally, the function used to calculate the FIM from the previous section is called and the states in the observation space are updated. The $step$ function also includes two conditions that automatically trigger the $reset$

**Algorithm 1** Asymmetric PPO Framework

---
1: Initialize Actor-Critic networks
2: Initialize memory buffer $\mathcal{B}$
3: Sample initial observation state $o_0$
4: **for** update $= 1, 2, \ldots$ **do**
5:     **for** actor $= 1, 2, \ldots, N$ **do**
6:         **for** $t = 0, T - 1$ **do**
7:             Get action $a_t$ and action distribution $\log \pi_\theta$ using actor policy $\pi_\theta(o_{a,t})$
8:             Execute action $a_t$, receive reward $r_t$, next observation $o_{t+1}$, and done flag $d_t$
9:             Get value $V_t$ from critic network
10:            Store $(o_t, a_t, r_t, d_t, V_t)$ in memory buffer $\mathcal{B}$
11:            **if** done flag $d_t == 1$ **then**
12:                Sample initial observation state $o_0$
13:            **end if**
14:         **end for**
15:     **end for**
16:     Compute returns $R$ and advantages $A$ with stored trajectories in memory buffer $\mathcal{B}$
17:     **for** epoch $= 1, 2, \ldots, K$ **do**
18:         Shuffle memory buffer $\mathcal{B}$ and divide into mini-batches
19:         **for** each minibatch **do**
20:             Update actor network using $(o_{a,t}, a_t, r_t, A_t)$
21:             Update critic network using $(o_{c,t}, a_t, r_t, V_t)$
22:         **end for**
23:     **end for**
24:     Clear memory buffer $\mathcal{B}$
25: **end for**

---

function. The first one is when the time step reaches 500, which means that the experiment has lasted for 500 ms and it is over. Furthermore, the agent must learn to control the motor so that it does not exceed speed and position limits. If this limit is broken in the step function, the $reset$ function is automatically triggered. This feature is implemented by the $terminated$ function within the environment. Another important part of the step function is to update the step reward based on the design of the reward function, which will directly affect the agent's training performance. The design of the reward function will be introduced in detail in the following sections.

## 4.7 Formulation of observation space for the optimization of input excitation

With the reinforcement learning algorithm chosen in Section 4.5 and the EMB environment for reinforcement learning implemented in Section 4.6, the details of the formulation of the observation space in reinforcement learning , which could assist in solving the optimal experiment design problem, will be introduced in this section.

Two structure of observation space design are considered in this thesis. In the first structure, the observation space contains all the information in the EMB environment, including the position of the motor, the real and measured values of the velocity, the time step, and the a priori values of the
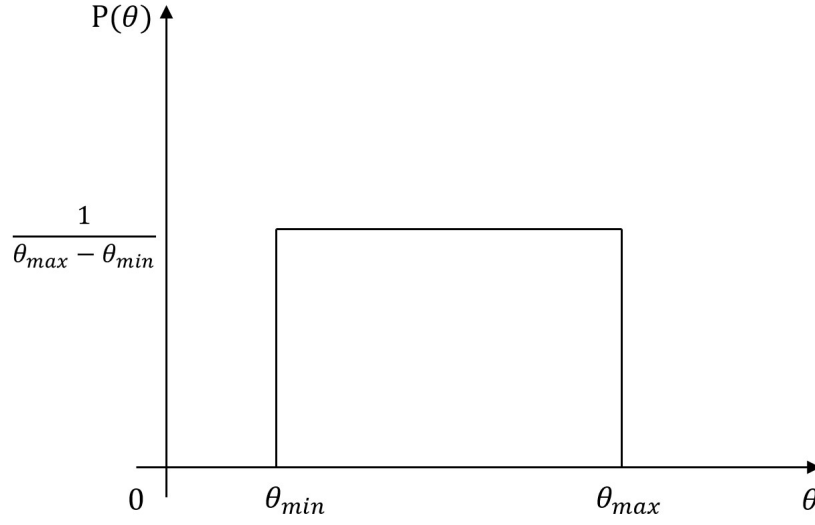
Figure 4.7: Structure of Asymmetric PPO framework with EMB environment

parameters and the FIM elements. The method employs fixed a priori values of the parameters, so it can be hypothesized that the trained agent will most likely not be able to do the right action when the actual parameter values in the EMB model is not equivalent to the priori values contained in the observation space during the training. However, this structure is a direct application of the conventional input excitation optimization problem to reinforcement learning, and can verify the correctness of the reinforcement learning framework as well as the learning potential of the agent, and is therefore essential.

The biggest difference between the second structure of observation space and the first one is that the prior values of the parameters to be identified are not included in the observation space. Instead, a distribution with the parameters to be identified is designed and used for training. Figure 4.7 illustrates the distribution of the one-dimensional parameters, the uniform distribution with a range from $\theta_{min}$ to $\theta_{max}$ is considered first. This distribution guarantees the same probability of sampling the parameter. Depending on the training requirements, other distributions such as the normal distribution can also be employed. Whenever the reset function is called in the EMB environment, a set of parameter values is sampled from the parameter distribution as a priori values and is saved in the observation space and used for training in this episode. This structure ensures that no assumptions about the a priori values of the parameters are needed for modeling, and the agent will no longer learn optimal the input base on specific parameter values during the training process, but will learn the optimal policy under the parameter distribution based on experience. The advantage of this structure is that it uses the random nature of reinforcement learning, helping the agent learn from different experiences.

# 5 Approaches of the single parameter Excitation

## 5.1 The Fisher information matrix of the vicious friction parameter

Firstly we choose the vicious friction $f_v$ as the parameter to be estimated. The calculation of the FIM is based on the formula given in equation (3.7):

$$\frac{\partial f_{d,k}}{\partial x_d^T} = \begin{bmatrix} 0 & 1 \\ -\frac{\gamma k_1}{J} & -\frac{f_v}{J} \end{bmatrix} \tag{5.1}$$

$$\frac{\partial f_{d,k}}{\partial \bar{f}_v} = \begin{bmatrix} 0 \\ -\frac{f_v x_2}{J} \end{bmatrix} \tag{5.2}$$

According to the real EMB model, the nominal range of position is $(-10, 100)$, the nominal range of velocity is $(-500, 500)$. Because of the calculation of the FIM is strongly relevant with the measurement value, to avoid heavy weighting on the velocity, the measured outputs need to be scaled with:

$$\bar{h}_k = \begin{bmatrix} \frac{1}{100} & 0 \\ 0 & \frac{1}{500} \end{bmatrix} h_k \tag{5.3}$$

Hence the Jacobian matrix of $\bar{h}_k$ is:

$$\frac{\partial \bar{h}_k}{\partial x_d^T} = \begin{bmatrix} \frac{1}{100} & 0 \\ 0 & \frac{1}{500} \end{bmatrix} \tag{5.4}$$

$$\frac{\partial \bar{h}_k}{\partial \bar{\theta}_i} = 0 \tag{5.5}$$

$$\chi_{\bar{\theta}_i,k+1} = \chi_{\bar{\theta}_i,k} + \Delta T \left( \begin{bmatrix} 0 & 1 \\ -\frac{\gamma k_1}{J} & -\frac{f_v}{J} \end{bmatrix} \chi_{\bar{\theta}_i,k} + \begin{bmatrix} 0 \\ -\frac{f_v x_2}{J} \end{bmatrix} \right) \tag{5.6}$$

So the calculation of $\chi_{\bar{\theta}_i,k}$ could be interpreted as a mass-damper system. When velocity $x_2$ stay zero after running a while, it is equal to no input added on the mass-damper system. The second term of $\chi_{\bar{\theta}_i,k}$ will oscillate with gradually decreasing amplitude until it stop at 0. But the first term of $\chi_{\bar{\theta}_i,k}$ will not to zero. This let the FIM always accumulating with time, even when the velocity already zero. The $\chi_{\bar{\theta}_i,k}$ was defined as:

$$\chi_{\bar{\theta}_i,k} = \begin{bmatrix} a \\ b \end{bmatrix} \tag{5.7}$$

We could have

$$\frac{d\bar{h}_k}{d\bar{f}_v} = \begin{bmatrix} \frac{1}{100} & 0 \\ 0 & \frac{1}{500} \end{bmatrix} \chi_{\bar{\theta}_i,k},$$
$$= \begin{bmatrix} \frac{a}{100} \\ \frac{b}{500} \end{bmatrix}$$

(5.8)

$$R^{-1} = \begin{bmatrix} 10^6 & 0 \\ 0 & 1 \end{bmatrix}$$

(5.9)

$$M_k = \left(\frac{d\bar{h}_k}{d\bar{f}_v}\right)^T R^{-1} \left(\frac{d\bar{h}_k}{d\bar{f}_v}\right),$$
$$= \begin{bmatrix} \frac{a}{100} & \frac{b}{500} \end{bmatrix} \begin{bmatrix} 10^6 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{a}{100} \\ \frac{b}{500} \end{bmatrix},$$
$$= 100a^2 + \frac{b^2}{250000}$$

(5.10)

## 5.2 Reward caculation

of the optimization in a time discrete system after $n$ steps, can be represented as the total reward $R_n$ in the reinforcement learning framework

$$R_n = \sum_{k=1}^{n} r_k = log|M_n|$$

(5.11)

From this equation, it is intuitive to define the step reward for the agent during training as

$$r_k = log|M_k| - log|M_{k-1}|$$

(5.12)

## 5.3 Integration of return to origin task

EMB needs to drive back to the origin when the experiment go to the end. This part could be done by a controller, but in this thesis it is integrated in the design of EMB environment for the training. Quintic Polynomial Interpolation

## 5.4 Comparison of performance between two observation space structures

Result from training with fixed parameter and parameter distribution. Training with parameter distribution is more robust
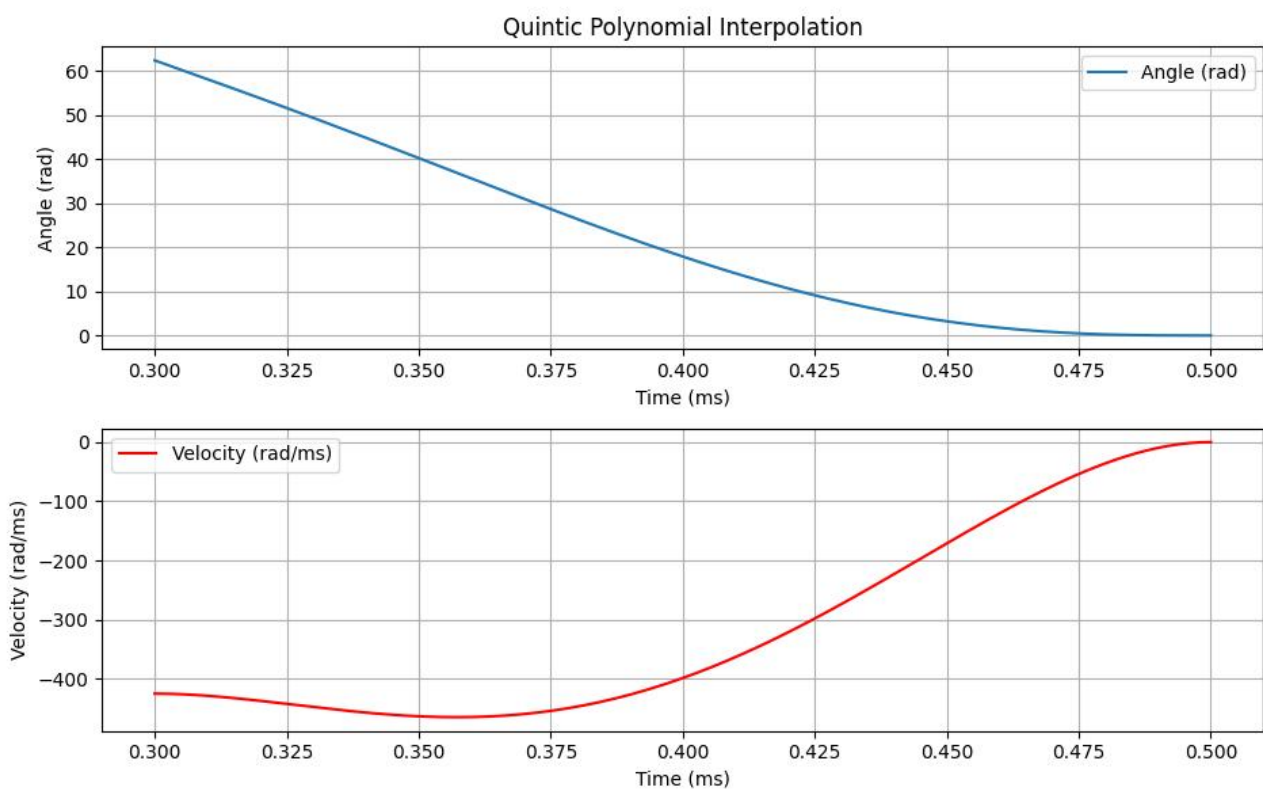
Figure 5.1: Qunitic Polynominal Interpolation for return to origin task

# 6 Approaches of the multi parameter Excitation

## 6.1 The Fisher information matrix of two parameter

After the succeed in the single parameter estimation, this RL structure was tested with two parameter. In the new turn of estimation, the vicious friction $f_v$ and the stiffness characteristics $k_1$ were chosen as the parameter to be estimated. Since significantly value difference could existed between different parameters, the parameters need to be scaled when the FIM was calculated. The normalised parameter vector could be written as:

$$\bar{\theta} = \left[ \frac{\theta_1}{\theta_{nom,1}}, \frac{\theta_2}{\theta_{nom,2}}, \ldots, \frac{\theta_n}{\theta_{nom,n}} \right]^T \tag{6.1}$$

The and the FIM could be calculated by the following equations:

$$\frac{\partial f_{d,k}}{\partial x_d^T} = \begin{bmatrix} 0 & 1 \\ -\frac{\gamma k_1}{J} & -\frac{f_v}{J} \end{bmatrix} \tag{6.2}$$

$$\frac{\partial f_{d,k}}{\partial \bar{\theta}} = \begin{bmatrix} 0 & 0 \\ -\frac{f_v x_2}{J} & -\frac{k_1 \gamma x_1}{J} \end{bmatrix} \tag{6.3}$$

$$\chi_{\bar{\theta},k+1} = \chi_{\bar{\theta},k} + \Delta T \left( \begin{bmatrix} 0 & 1 \\ -\frac{\gamma k_1}{J} & -\frac{f_v}{J} \end{bmatrix} \chi_{\bar{\theta}_i,k} + \begin{bmatrix} 0 & 0 \\ -\frac{f_v x_2}{J} & -\frac{k_1 \gamma x_1}{J} \end{bmatrix} \right) \tag{6.4}$$

Because there the parameter is not relevant with system output x1 and x2, so:

$$\frac{\partial \bar{h}_k}{\partial \bar{\theta}_i} = 0 \tag{6.5}$$

Then:

$$\begin{aligned}
\frac{d\bar{h}_k}{d\bar{\theta}_i} &= \frac{\partial \bar{h}_k}{\partial x_d^T} \chi_{\bar{\theta}_i,k} + \frac{\partial \bar{h}_k}{\partial \bar{\theta}_i}, \\
&= \begin{bmatrix} \frac{1}{100} & 0 \\ 0 & \frac{1}{500} \end{bmatrix} \chi_{\bar{\theta}_i,k}
\end{aligned} \tag{6.6}$$

The $\chi_{\bar{\theta},k}$ was defined as:

$$\chi_{\bar{\theta},k} = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \tag{6.7}$$

$$\frac{d\bar{h}_k}{d\bar{\theta}} = \left[ \frac{d\bar{h}_k}{d\bar{f}_v}, \frac{d\bar{h}_k}{dk_1} \right],$$

$$= \begin{bmatrix} \frac{1}{100} & 0 \\ 0 & \frac{1}{500} \end{bmatrix} \chi_{\bar{\theta},k}, \tag{6.8}$$

$$= \begin{bmatrix} \frac{a}{100} & \frac{c}{100} \\ \frac{b}{500} & \frac{d}{500} \end{bmatrix}$$

$$M_k = \left( \frac{d\bar{h}_k}{d\bar{\theta}} \right)^T R^{-1} \left( \frac{d\bar{h}_k}{d\bar{\theta}} \right),$$

$$= \begin{bmatrix} \frac{a}{100} & \frac{b}{500} \\ \frac{c}{100} & \frac{d}{500} \end{bmatrix} \begin{bmatrix} 10^6 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{a}{100} & \frac{c}{100} \\ \frac{b}{500} & \frac{d}{500} \end{bmatrix}, \tag{6.9}$$

$$= \begin{bmatrix} 100a^2 + \frac{b^2}{250000} & 100ac + \frac{bd}{250000} \\ 100ac + \frac{bd}{250000} & 100c^2 + \frac{d^2}{250000} \end{bmatrix}$$

## 6.2 Approaches for design of the reward function

## 6.3 Result of the agent

Compare the result between the different reward functions.

Also compare to the different input trajectories, such as step inputs, ramp inputs, and sinusoidal input etc.

Problem founded, the way try to fix them. Sth works, But still open problem

# 7 Conclusion and Outlook

Conclusion: compare to the related works, we are doing with mechanical area, we have some result

But an estimator is still needed.

# Bibliography

[1] C. L. J. Line, *Modelling and control of an automotive electromechanical brake*. University of Melbourne, Department of Mechanical and Manufacturing Engineering, 2007, pp. 63–79.

[2] C. F. Lee and C. Manzie, "Rapid parameter identification for an electromechanical brake", in *2013 Australian Control Conference*, IEEE, 2013, pp. 391–396.

[3] R. Huang, J. Fogelquist, and X. Lin, "Reinforcement learning of optimal input excitation for parameter estimation with application to li-ion battery", *IEEE Transactions on Industrial Informatics*, vol. 19, no. 11, pp. 11 160–11 170, 2023.

[4] A. Mesbah and S. Streif, "A probabilistic approach to robust optimal experiment design with chance constraints", *IFAC-PapersOnLine*, vol. 48, no. 8, pp. 100–105, 2015.

[5] K. Reif, *Automotive Mechatronics*. Springer, 2014, p. 13.

[6] H.-P. Schoner and P. Hille, "Automotive power electronics. new challenges for power electronics", in *2000 IEEE 31st Annual Power Electronics Specialists Conference. Conference Proceedings (Cat. No. 00CH37018)*, IEEE, vol. 1, 2000, pp. 6–11.

[7] L. Zhang, Z. Zhang, Z. Wang, J. Deng, and D. G. Dorrell, "Chassis coordinated control for full x-by-wire vehicles-a review", *Chinese Journal of Mechanical Engineering*, vol. 34, pp. 1–25, 2021.

[8] D. Li, C. Tan, W. Ge, J. Cui, C. Gu, and X. Chi, "Review of brake-by-wire system and control technology", in *Actuators*, MDPI, vol. 11, 2022, p. 80.

[9] C.-F. Zong, G. Li, H.-Y. Zheng, L. He, and Z.-X. Zhang, "Study progress and outlook of chassis control technology for x-by-wire automobile", *Zhongguo Gonglu Xuebao(China Journal of Highway and Transport)*, vol. 26, no. 2, pp. 160–176, 2013.

[10] C. Li, G. Zhuo, C. Tang, *et al.*, "A review of electro-mechanical brake (emb) system: Structure, control and application", *Sustainability*, vol. 15, no. 5, p. 4514, 2023.

[11] M. Zhang and J. Song, "A review of electromechanical brake (emb) system", *Mechanical Science and Technology*, vol. 24, no. 2, pp. 208–211, 2005.

[12] Y. Zhao, "Research and development of electro-mechanical brake system actuator", Ph.D. dissertation, Tsinghua University Beijing, China, 2010.

[13] S. Schrade, X. Nowak, A. Verhagen, and D. Schramm, "Short review of emb systems related to safety concepts", in *Actuators*, MDPI, vol. 11, 2022, p. 214.

[14] C. Line, C. Manzie, and M. C. Good, "Electromechanical brake modeling and control: From pi to mpc", *IEEE Transactions on Control Systems Technology*, vol. 16, no. 3, pp. 446–457, 2008. DOI: 10.1109/TCST.2007.908200.

[15] L. Chen, P. Wu, K. Chitta, B. Jaeger, A. Geiger, and H. Li, "End-to-end autonomous driving: Challenges and frontiers", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[16]   S. Kuutti, R. Bowden, Y. Jin, P. Barber, and S. Fallah, "A survey of deep learning applications to autonomous vehicle control", *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 2, pp. 712–733, 2020.

[17]   B. Desai and K. Patil, "Secure and scalable multi-modal vehicle systems: A cloud-based framework for real-time llm-driven interactions", *Innovative Computer Sciences Journal*, vol. 9, no. 1, pp. 1–11, 2023.

[18]   Y.-m. You, "Multi-objective optimal design of permanent magnet synchronous motor for electric vehicle based on deep learning", *Applied Sciences*, vol. 10, no. 2, p. 482, 2020.

[19]   R. R. Ardeshiri, B. Balagopal, A. Alsabbagh, C. Ma, and M.-Y. Chow, "Machine learning approaches in battery management systems: State of the art: Remaining useful life and fault detection", in *2020 2nd IEEE International conference on industrial electronics for sustainable energy systems (IESES)*, IEEE, vol. 1, 2020, pp. 61–66.

[20]   DeepL, *Deepl translator,* Accessed: 2024-11-18, 2024. [Online]. Available: `https://www.deepl.com/translator`.

[21]   OpenAI, *Chatgpt (gpt-4),* Accessed: 2024-11-18, 2024. [Online]. Available: `https://chat.openai.com/`.

[22]   T. Grigoratos, M. Mathissen, R. Vedula, *et al.*, "Interlaboratory study on brake particle emissions—part i: Particulate matter mass emissions", *Atmosphere*, vol. 14, no. 3, p. 498, 2023.

[23]   D. Lovrec and M. Kastrevc, "Modelling and simulating a controlled press-brake supply system", *International Journal of Simulation Modelling*, vol. 10, no. 3, pp. 133–144, 2011.

[24]   J. Kim, C. Jo, Y. Kwon, *et al.*, "Electro-mechanical brake for front wheel with back-up braking", *SAE International Journal of Passenger Cars-Mechanical Systems*, vol. 7, no. 2014-01-2538, pp. 1369–1373, 2014.

[25]   J. S. Cheon, J. Kim, and J. Jeon, "New brake by wire concept with mechanical backup", *SAE International Journal of Passenger Cars-Mechanical Systems*, vol. 5, no. 2012-01-1800, pp. 1194–1198, 2012.

[26]   C. Qi, "Kontaktereigniserkennung für ein elektromechanisches system", Master's Thesis, TU Berlin, 2024, p. 8.

[27]   G. Park, S. Choi, and D. Hyun, "Clamping force estimation based on hysteresis modeling for electro-mechanical brakes", *International Journal of Automotive Technology*, vol. 18, pp. 883–890, 2017.

[28]   R. Stribeck, "Die wesentlischen eigenschaften der gleit-und rollenlager. zeitschrift des vereines deutscher ingenieure", *vol*, vol. 46, pp. 1341–1348, 1902.

[29]   D. Karnopp, "Computer Simulation of Stick-Slip Friction in Mechanical Dynamic Systems", *Journal of Dynamic Systems, Measurement, and Control*, vol. 107, no. 1, pp. 100–103, Mar. 1985, ISSN: 0022-0434. DOI: `10.1115/1.3140698`. [Online]. Available: `https://doi.org/10.1115/1.3140698`.

[30]   H. Olsson, K. J. Åström, C. C. De Wit, M. Gäfvert, and P. Lischinsky, "Friction models and friction compensation", *Eur. J. Control*, vol. 4, no. 3, pp. 176–195, 1998.

[31]   C. C. De Wit, H. Olsson, K. J. Astrom, and P. Lischinsky, "A new model for control of systems with friction", *IEEE Transactions on automatic control*, vol. 40, no. 3, pp. 419–425, 1995.

[32]   Y. Bard, *Nonlinear parameter estimation*. Academic press New York, 1974, vol. 1209.

[33] R. Mehra, "Optimal input signals for parameter estimation in dynamic systems–survey and new results", *IEEE transactions on automatic control*, vol. 19, no. 6, pp. 753–768, 1974.

[34] S. Arimoto and H. Kimura, "Optimum input test signals for system identification—an information-theoretical approach", *International Journal of Systems Science*, vol. 1, no. 3, pp. 279–290, 1971.

[35] A. Paszke, S. Gross, S. Chintala, *et al.*, "Automatic differentiation in pytorch", 2017.

[36] D. A. Spielman, *Spectral and algebraic graph theory*, `http://cs-www.cs.yale.edu/homes/spielman/sagt/sagt.pdf`, Unpublished draft, 2019.

[37] R. Bellman, "A markovian decision process", *Journal of mathematics and mechanics*, pp. 679–684, 1957.

[38] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation", *Advances in neural information processing systems*, vol. 12, 1999.

[39] V. Mnih, "Playing atari with deep reinforcement learning", *arXiv preprint arXiv:1312.5602*, 2013.

[40] T. Lillicrap, "Continuous control with deep reinforcement learning", *arXiv preprint arXiv:1509.02971*, 2015.

[41] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods", in *International conference on machine learning*, PMLR, 2018, pp. 1587–1596.

[42] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms", *arXiv preprint arXiv:1707.06347*, 2017.

[43] J. Schulman, "Trust region policy optimization", *arXiv preprint arXiv:1502.05477*, 2015.

[44] S. Huang, R. F. J. Dossa, A. Raffin, A. Kanervisto, and W. Wang, "The 37 implementation details of proximal policy optimization", in *ICLR Blog Track*, https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/, 2022. [Online]. Available: `https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/`.

[45] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and acting in partially observable stochastic domains", *Artificial intelligence*, vol. 101, no. 1-2, pp. 99–134, 1998.

[46] M. Hauskrecht, "Value-function approximations for partially observable markov decision processes", *Journal of artificial intelligence research*, vol. 13, pp. 33–94, 2000.

[47] J. Pineau, G. Gordon, S. Thrun, *et al.*, "Point-based value iteration: An anytime algorithm for pomdps", in *Ijcai*, vol. 3, 2003, pp. 1025–1032.

[48] D. A. McAllester and S. Singh, "Approximate planning for factored pomdps using belief state simplification", *arXiv preprint arXiv:1301.6719*, 2013.

[49] G. Shani, R. Brafman, and S. Shimony, "Adaptation for changing stochastic environments through online pomdp policy learning", in *Proc. Eur. Conf. on Machine Learning*, 2005, pp. 61–70.

[50] S. Ross, J. Pineau, S. Paquet, and B. Chaib-Draa, "Online planning algorithms for pomdps", *Journal of Artificial Intelligence Research*, vol. 32, pp. 663–704, 2008.

[51] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel, "Asymmetric actor critic for image-based robot learning", *arXiv preprint arXiv:1710.06542*, 2017.

[52] W. Yue, Y. Zhou, X. Zhang, Y. Hua, Z. Wang, and G. Kou, "Aacc: Asymmetric actor-critic in contextual reinforcement learning", *arXiv preprint arXiv:2208.02376*, 2022.

[53] A. Baisero and C. Amato, "Unbiased asymmetric reinforcement learning under partial observability", *arXiv preprint arXiv:2105.11674*, 2021.

[54] S. Huang, R. F. J. Dossa, C. Ye, *et al.*, "Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms", *Journal of Machine Learning Research*, vol. 23, no. 274, pp. 1–18, 2022. [Online]. Available: `http://jmlr.org/papers/v23/21-1342.html`.

[55] M. Towers, A. Kwiatkowski, J. Terry, *et al.*, "Gymnasium: A standard interface for reinforcement learning environments", *arXiv preprint arXiv:2407.17032*, 2024.