

Learning Invariances in Dynamical System

Cheng-Cheng Lao

CID: 01353756

Supervised by Dr. Andrew Duncan and Dr. Mark van der Wilk

August 2, 2022

Submitted in partial fulfilment of the requirements for the MSc in Statistics of
Imperial College London

The work contained in this thesis is my own work unless otherwise stated.

Signed: Cheng-Cheng Lao

Date: August 2, 2022

Abstract

ABSTRACT GOES HERE

Acknowledgements

ANY ACKNOWLEDGEMENTS GO HERE

1 Introduction

Dynamical system is one of the richest theory that models the real world with wide application across all science and engineering. Of course, it is an interesting area of mathematics to study on its own. In practice, it will be very valuable to be able to reliably predict the evolution of a system for many purpose; both from a theoretical and practical point of view. Recently, data intensive machine learning methods, such as deep learning, has shown some very powerful results. However, the amount of data required are often enormous, and that will involve lots of manual labour to label and process the data. Sometimes, the data could be very difficult to come by. Therefore, data efficiency is the key in the future if we wish to study more complicated systems. Inductive bias is one way to achieve the goal. What was usually done is to embed the inductive bias into the model in some way so that the model "understands" this knowledge before seeing the data. Assuming the inductive bias is correct, we would expect the model to perform better, either in a data efficiency way or prediction performance way. One very powerful inductive bias is the use of symmetry and invariances, that puts strong constraints on the model to obey.

2 Background

Here we will cover the background knowledge required to understand the remaining thesis, including Gaussian Process, dynamical systems as well as symmetry and invariances. This section will also cover related work in the field.

2.1 Gaussian Process

Gaussian Process (GP) can be thought of a distribution over functions. Rasmussen and Williams (2006) More formally,

GP is a collection of random variables, any finite number of which have a joint Gaussian distribution

We will be able to specify a GP on $f(\mathbf{x})$ by mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$ as well as the kernel function $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[f(\mathbf{x}) - m(\mathbf{x})(f(\mathbf{x}') - m(\mathbf{x}'))]$. We then write

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')).$$

If we assume the signal is not noise free such that we assume a Gaussian noise $y = f + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma_n^2)$. Then we will have An important identity we use throughout the thesis is the Gaussian conditional formula; if

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix} \right)$$

, then

$$x|y \sim \mathcal{N}(\mu_x + CB^{-1}(y - \mu_y), A - CB^{-1}C^T) \quad (2.1)$$

2.1.1 GPflow

GPflow is used throughout the project (Matthews et al. (2017)). It is a Python package built on TensorFlow to allow efficient and fast computation on GPUs to do GP inferences. A particular powerful feature is that it allows the choosing of hyperparameter such as lengthscales and variance of a kernel automatic by calculating the gradient (using automatic differentiation feature of Tensorflow) and perform optimisation. In later sections, we will need to optimise more parameters than aforementioned ones with respect to maginal likelihood objective.

2.2 Dynamical Systems

Dynamical system, which describes the evolution of a system in time, can be broadly classified to two categories: differential equations and difference equations, with the former being continuous time and the latter being discrete time. (Strogatz (2019)). We will formalise our notion and notation using continuous time (differential equation) description. There are many real world systems that can be modelled by a dynamical system, for example a huge class of problems in engineering and physics. A simple example, which will also be explored later, is a damped pendulum, with equation of motion

$$m \frac{d^2 x}{dt^2} + b \frac{dx}{dt} + kx = 0.$$

In general an n^{th} degree differential equations can be written as

$$\frac{d^n x}{dt^n} = F\left(t, x, \frac{dx}{dt}, \dots, \frac{d^{n-1}x}{dt^{n-1}}\right),$$

(Glendinning (1994)) and we will need n initial conditions to specify a solution for the differential equations. We can also write n^{th} degree ordinary differential equations (ODEs) as n degree one coupled ODE. For example, in the above case, we can write it as

$$\begin{cases} \frac{dx}{dt} = v \\ \frac{dv}{dt} = \frac{-bv - kx}{m} \end{cases}$$

More generally, we can write

$$\begin{cases} \dot{x}_1 = f_1(x_1, \dots, x_n) \\ \vdots \\ \dot{x}_n = f_n(x_1, \dots, x_n) \end{cases}$$

We can then visualise the state of the system described by the n variables in the "phase space", and this system will evolve or move in the phase space as time evolves. For instance, the damped pendulum is of degree 2 so there will be two variables describing its evolution, which in this case will be x and $\frac{dx}{dt} = \dot{x}$, which is the angular displacement as well as angular velocity. Another example would be a naive model describing exponential growth of a organism, such as some sort of bacteria population in a petri dish

$$\dot{x} = rx,$$

and this time it is one-dimensional with the phase variable being the population number. Another important concept is linearity and nonlinearity.

2.3 Symmetry and Invariances

Invariances are functions of the states that describe a dynamical system, and is unchanged throughout the evolution of the system over time. An example would be from the field of physics, the conservation of energy, which will be unchanged throughout the trajectories of the system. There are other types of symmetry, such as translational symmetry, rotational symmetry, permutation symmetry, which are widely explored in geometric deep learning literature. More formally, we can describe symmetry using group theory. In physics, there is a deeper connection between symmetry and invariance, namely the Noether's Theorem. In plain words, every conservation law is associated with an underlying symmetry. For instance, conservation of momentum is obeyed when there is a translation symmetry of the system. Similarly, rotational symmetry is associated with the conservation of angular momentum. There are more abstract conservation laws, such as conservation of electrical charge is related to Gauge Symmetry, an advanced physics idea that we will not go into too much detail. Here the more relevant one is the conservation of energy, which is related to the symmetry of the system under temporal translation.

2.4 Related Work

There are many existing literatures in the field of learning invariance as well as physics informed machine learning.

2.4.1 Physics Informed Machine Learning

Hamiltonian Neural Network

Symbolic approach

2.4.2 ODE approach

2.4.3 Symmetry and Invariance

Learning Invariance using Marginal Likelihood

Geometric Deep Learning

2.4.4 GP in dynamical systems

3 Invariance Kernel

In this chapter we will start with general theoretical construction of invariance kernel given the knowledge of the invariance of the system. We will then apply the construction to various systems, namely linear and nonlinear system in both one and two dimensions.

3.1 General Construction

If we are given a dynamical system of dimension d with variables \mathbf{q}, \mathbf{p} , where $\mathbf{q} = (q_1, q_2, \dots, q_d)$ is the vector of positional coordinates, and $\mathbf{p} = (p_1, p_2, \dots, p_d)$ is the vector of velocity coordinates of the states of the dynamical system. We are interested to predict the future trajectories of the state of the system. Therefore, we would like to know the time derivative of these coordinates so we can update them according to Euler's integrator. These time derivatives are referred to as the dynamics of a dynamical system, which governs the evolution of the dynamical system, and is a function of the coordinates \mathbf{q} and \mathbf{p} . For our systems, we will denote the dynamics of \mathbf{p} as $\frac{d\mathbf{p}}{dt} = \mathbf{a}(\mathbf{q}, \mathbf{p})$, and that of \mathbf{q} as $\frac{d\mathbf{q}}{dt} = \mathbf{v}(\mathbf{q}, \mathbf{p})$. Once we have the dynamics, we can integrate up to obtain the future trajectories. Notation wise, we will collect the two dynamics term and call them

$$\mathbf{f}(\mathbf{q}, \mathbf{p}) = \begin{pmatrix} \mathbf{a}(\mathbf{q}, \mathbf{p}) \\ \mathbf{v}(\mathbf{q}, \mathbf{p}) \end{pmatrix}$$

For simplicity of notation, the dependence on \mathbf{q}, \mathbf{p} is now implicit. We will put independent GP prior on \mathbf{a} and \mathbf{v} since there are no reason we should assume they are correlated. For the choice of kernel, we chose the standard smooth kernel, the squared exponential kernel, or RBF, and we denote this kernel to be K_{RBF} . We will also denote any set of coordinates of length n and dimension d as

$$X \equiv \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} = \begin{pmatrix} q_{11} & q_{21} & \dots & q_{d1} & p_{11} & p_{21} & \dots & p_{d1} \\ q_{12} & q_{22} & \dots & q_{d2} & p_{12} & p_{22} & \dots & p_{d2} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ q_{1n} & q_{2n} & \dots & q_{dn} & p_{1n} & p_{2n} & \dots & p_{dn} \end{pmatrix}.$$

Without loss of generality, we will choose to stack the dynamics vertically in \mathbf{f} ; for example in a d dimensional system, we have

$$\mathbf{f}(X) = \begin{pmatrix} a_1(\mathbf{x}_1) \\ \vdots \\ a_1(\mathbf{x}_n) \\ \vdots \\ a_d(\mathbf{x}_1) \\ \vdots \\ a_d(\mathbf{x}_n) \\ v_1(\mathbf{x}_1) \\ \vdots \\ v_1(\mathbf{x}_n) \\ \vdots \\ v_d(\mathbf{x}_1) \\ \vdots \\ v_d(\mathbf{x}_n) \end{pmatrix}$$

For our naive independent RBF GP prior, we have block diagonal

$$K = \text{Cov}(\mathbf{f}(X), \mathbf{f}(X')) = \begin{pmatrix} K_{RBF, a_1}(X, X') & \dots & \dots & \dots & 0 \\ \vdots & \ddots & \dots & \ddots & \vdots \\ 0 & \dots & K_{RBF, a_d}(X, X') & \dots & 0 \\ \vdots & \ddots & \vdots & K_{RBF, v_1}(X, X') & \vdots \\ 0 & \dots & 0 & \dots & K_{RBF, v_d}(X, X') \end{pmatrix},$$

with all the off diagonal terms being zero block matrix because of the independence prior assumption. This naive kernel will be our baseline to be compared to throughout the whole project. Now we can start considering the invariance. If we assume the invariance is true throughout the input space \mathbb{R}^{2d} , then we can condition the GP on the a grid of points, which we referred to invariance points, X_L , and assume there are ℓ of them, and we will call the invariance constraints L .

The form of L will depend on the system as well as X_L , and examples will be given in the following sections to make it clear. Since the invariance function will always be a linear function on the dynamics by the fact that it is the time derivative

$$\text{energy} = E(\theta, \mathbf{p}, \mathbf{q}); \frac{dE}{dt} = \sum_{i=1}^d \frac{\partial E}{\partial p_i} a_i + \sum_{i=1}^d \frac{\partial E}{\partial q_i} v_i$$

We can see it is always a linear combination of something independent of the dynamics. If we apply this linear transformation L on $\mathbf{f}(X_L)$, $L(\mathbf{f}(X_L))$ will again be a GP with a

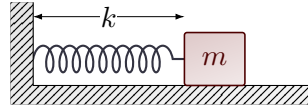
transformed kernel. Therefore, we have

$$\begin{pmatrix} \mathbf{f}(X) \\ L[\mathbf{f}(X_L)] \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0_{2nd} \\ 0_\ell \end{pmatrix}, \begin{pmatrix} K & LK \\ KL^T & LKL^T \end{pmatrix} \right)$$

3.2 1D System

3.2.1 Linear

We will first examine one of the most simple dynamical system, an 1D simple harmonic motion (SHM). An example would be a mass spring system as shown in figure with mass m and spring constant k , an example trajectory is shown in figure 3.1.



The equation of motion of the system is

$$m \frac{d^2 q}{dt^2} = -kq,$$

where q is the displacement. And the analytical solution would be of the form $q = A \sin(\omega_0 t + \phi)$, where $\omega = \frac{k}{m}$ and A, ϕ depends on the initial condition, which dictates the amplitude and phase of the motion. For this case, we have

$$\mathbf{f}(X) = \begin{pmatrix} \mathbf{a}(X) \\ \mathbf{v}(X) \end{pmatrix}$$

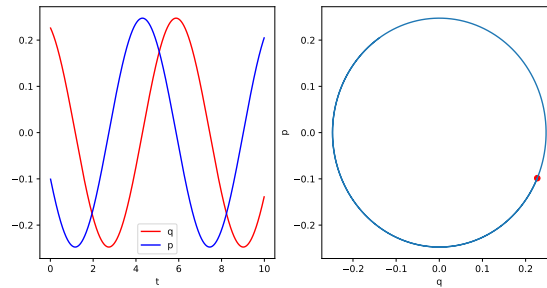


Figure 3.1: Example trajectory of SHM

We also have the energy, $E = \frac{mp^2}{2} + \frac{kq^2}{2}$, Therefore, to obtain our invariance L , we use the conservation of energy $\frac{dE}{dt} = 0$ so we finally have

$$L(a, v) = mpa + kqv = 0$$

While L is not able to be into a matrix form, it is a linear operator as shown below. If we have

$$X_L \equiv \begin{pmatrix} \mathbf{x}_{L,1} \\ \vdots \\ \mathbf{x}_{L,\ell} \end{pmatrix} = \begin{pmatrix} q_{L,1} & p_{L,1} \\ \vdots & \vdots \\ q_{L,\ell} & p_{L,\ell} \end{pmatrix} \equiv \begin{pmatrix} \vdots & \vdots \\ q_L & p_L \\ \vdots & \vdots \end{pmatrix},$$

then we have

$$L([\mathbf{f}(X_L)]) = \begin{pmatrix} mp_{L,1}a(q_{L,1}, p_{L,1}) + kq_{L,1}v(q_{L,1}, p_{L,1}) \\ \vdots \\ mp_{L,\ell}a(q_{L,\ell}, p_{L,\ell}) + kq_{L,\ell}v(q_{L,\ell}, p_{L,\ell}) \end{pmatrix},$$

which can be readily checked to be linear. Combine with original GP prior assumption, we have

$$\begin{pmatrix} \mathbf{f}(X) \\ L([\mathbf{f}(X_L)]) \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0_{2n} \\ 0_\ell \end{pmatrix}, \begin{pmatrix} A & B \\ C & D \end{pmatrix} \right),$$

where

$$A = K(X, X), B = \begin{pmatrix} K_{RBF,a} \\ K_{RBF,v} \end{pmatrix} \odot \begin{pmatrix} mP_L \\ kQ_L \end{pmatrix}, C = B^T, D = K_{RBF,a} \odot m^2(p_L \otimes p_L) + K_{RBF,v} \odot k^2(q_L \otimes q_L),$$

where \odot is the elementwise product and \otimes is the Kronecker product so that

$$P_L = \begin{pmatrix} p_{L,1} & \cdots & p_{L,\ell} \\ \vdots & \text{repeats n rows} & \vdots \\ p_{L,1} & \cdots & p_{L,\ell} \end{pmatrix}, Q_L = \begin{pmatrix} q_{L,1} & \cdots & q_{L,\ell} \\ \vdots & \text{repeats n rows} & \vdots \\ q_{L,1} & \cdots & q_{L,\ell} \end{pmatrix}$$

and we have

$$p_L \otimes p_L = \begin{pmatrix} p_{L,1}^2 & p_{L,1}p_{L,2} & \cdots & p_{L,1}p_{L,\ell} \\ \vdots & \vdots & \vdots & \vdots \\ p_{L,\ell}p_{L,1} & p_{L,\ell}p_{L,2} & \cdots & p_{L,\ell}^2 \end{pmatrix}, q_L \otimes q_L = \begin{pmatrix} q_{L,1}^2 & q_{L,1}q_{L,2} & \cdots & q_{L,1}q_{L,\ell} \\ \vdots & \vdots & \vdots & \vdots \\ q_{L,\ell}q_{L,1} & q_{L,\ell}q_{L,2} & \cdots & q_{L,\ell}^2 \end{pmatrix},$$

These matrix expressions are derived as follows by computing the covariance manually. For B , we wish to calculate

$$\begin{aligned} B_{ij} &= \text{Cov}(\mathbf{f}(X), L[\mathbf{f}(X_L)])_{ij} \\ &= \text{Cov}(\mathbf{f}(X)_i, L\mathbf{f}(X_L)_j) \\ &= \begin{cases} \text{Cov}(a(q_i, p_i), mp_{L,j}a(q_{L,j}, p_{L,j}) + kq_{L,j}v(q_{L,j}, p_{L,j})) & i \leq n \\ \text{Cov}(v(q_i, p_i), mp_{L,j}a(q_{L,j}, p_{L,j}) + kq_{L,j}v(q_{L,j}, p_{L,j})) & i > n \end{cases} \\ &= \begin{cases} K_{RBF,a}(\mathbf{x}_i, \mathbf{x}_{L,j})mp_{L,j} & i \leq n \\ K_{RBF,v}(\mathbf{x}_i, \mathbf{x}_{L,j})kq_{L,j} & i > n \end{cases}, \end{aligned}$$

and hence we have the form above. For D , we have

$$\begin{aligned} D_{ij} &= \text{Cov}(L[\mathbf{f}(X_L)], L[\mathbf{f}(X_L)])_{ij} \\ &= \text{Cov}(mp_{L,i}a(q_{L,i}, p_{L,i}) + kq_{L,i}v(q_{L,i}, p_{L,i}), mp_{L,i}a(q_{L,i}, p_{L,i}) + kq_{L,i}v(q_{L,i}, p_{L,i})) \\ &= m^2 p_{L,i} p_{L,j} K_{RBF,a}(\mathbf{x}_{L,i}, \mathbf{x}_{L,j}) + k^2 q_{L,i} q_{L,j} K_{RBF,v}(\mathbf{x}_{L,i}, \mathbf{x}_{L,j}) \end{aligned}$$

using the bilinear property of the covariance operator and the fact that v and a are independent. Since we assume invariance on these invariance points, we will condition on $L([\mathbf{f}(X_L)]) = 0$. Now we can simply use the Gaussian conditional formula to obtain the Schur Complement using equation 2.1.

$$\mathbf{f}(X)|L[\mathbf{f}(X_L)] = 0 \sim \mathcal{N}(0_{2n}, A - BD^{-1}C),$$

we will then call the resulting covariance our Invariance Kernel for 1D SHM, K_L .

3.2.2 Non Linear

The idea is pretty much the same for nonlinear system compared to linear; however, the fitting is expected to be more difficult. A simple nonlinear system in every day life is a simple pendulum as shown in figure below. The governing equation is

$$\frac{d^2 q}{dt^2} = -\frac{g}{\ell} \sin q,$$

where q is the angle of displacement this time and we have example trajectory in 3.2, where we can see the effect of nonlinearity on the shape. The nonlinear dynamics arises from the sine term above, which will complicate the derivation for invariance kernel slightly. However, since $\sin x \approx x$ at small angle, this system is approximately linear under small displacement. There is no analytical solution to this nonlinear problem in general, besides the small displacement case, which will be 1D SHM.

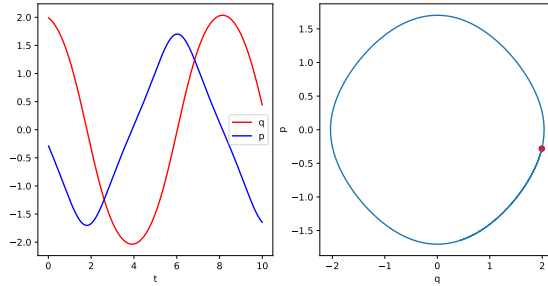


Figure 3.2: Example trajectory of pendulum

This time we have energy, $E = \frac{m\ell^2 p^2}{2} + mg\ell(1 - \cos q)$, and by setting the time

derivative to 0, we have

$$L(a, v) = \frac{dE}{dt} = m\ell^2 pa + mg\ell(\sin q)v = 0.$$

If we cancel out the common term $m\ell$ since their product cannot be zero, we have

$$L(a, v) = \ell pa + g(\sin q)v = 0$$

Most of the terms are unchanged from the linear case. However, this time,

$$B = \begin{pmatrix} K_{RBF,a} \\ K_{RBF,v} \end{pmatrix} \odot \begin{pmatrix} \ell P_L \\ g \sin(Q_L) \end{pmatrix}, D = K_{RBF,a} \odot \ell^2 (p_L \otimes p_L) + K_{RBF,v} \odot g^2 (\sin(q_L) \otimes \sin(q_L)),$$

where

$$\sin(Q_L) = g \begin{pmatrix} \sin(q_{L,1}) & \dots & \sin(q_{L,\ell}) \\ \vdots & \text{repeats n rows} & \vdots \\ \sin(q_{L,1}) & \dots & \sin(q_{L,\ell}) \end{pmatrix},$$

$$\sin(q_L) \otimes \sin(q_L) = \begin{pmatrix} \sin(q_{L,1})^2 & \sin(q_{L,1}) \sin(q_{L,2}) & \dots & \sin(q_{L,1}) \sin(q_{L,\ell}) \\ \vdots & \vdots & \vdots & \vdots \\ \sin(q_{L,\ell}) \sin(q_{L,1}) & \sin(q_{L,\ell}) \sin(q_{L,2}) & \dots & \sin(q_{L,\ell})^2 \end{pmatrix},$$

derived in exactly the same way as the linear case.

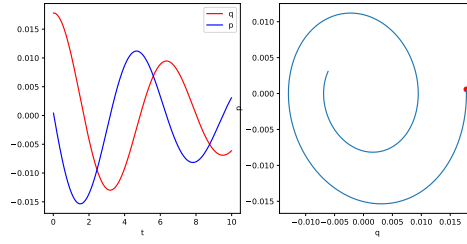
3.3 Damped System

In previous sections, we have considered a perfect system, i.e. ideal world without frictions and the conservation of energy is completely obeyed such that $L = 0$ exactly. However, in a real system, there will be dissipation where energy is lost in the form of heat etc. Therefore, to model that, we need to allow "approximate" invariance, such that the invariance L is noisy and not always equal to zero, i.e. soft invariance. To account for that, in a similar spirit to when we add noise to the underlying latent GP to represent noise signal in equation, where we have $K + \sigma_n^2 \mathbb{I}$; Here, since the white noise is uncorrelated, the only place the noise will enter is D , so we simply have to add a noise term to D and replace all D with $\tilde{D} = D + \epsilon \mathbb{I}$, where ϵ is a parameter to be learnt. For the data, we will model the system as a damped SHM or pendulum with linear velocity dependent force. If we denote $\omega_0^2 = k/m$ or $\omega_0^2 = g/\ell$ We will have the equation of motion for SHM and pendulum respectively

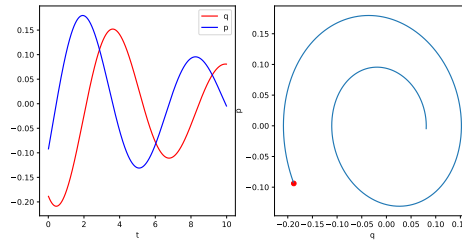
$$\frac{d^2 q}{dt^2} + 2\gamma \frac{dq}{dt} + \omega_0^2 q = 0; \quad \frac{d^2 q}{dt^2} + 2\gamma \frac{dq}{dt} + \omega_0^2 \sin q = 0,$$

and the γ is the damping factor that controls how much frictional force there is. Note that we will only focus on underdamping case ($\gamma < \omega_0^2$) so that the system still oscillate but with gradually reduced amplitude instead of critical damped ($\gamma = \omega_0^2$) or overdamped

$(\gamma > \omega_0^2)$ case where the system simply slowly decays to still. An example trajectory for damped SHM and damped pendulum is shown below in figure 3.3.



(a) Example trajectory of damped SHM



(b) Example trajectory of damped pendulum

Figure 3.3: Example trajectories of damped systemes

3.4 2D System

A 2D system is not that different from an 1D system. The major difference being there will be two more variables. We will again look at two simple examples, a linear 2D SHM, and a nonlinear double pendulum.

3.4.1 Linear

A simple extension to 1D SHM to 2D is just allowing the spring to be at an angle so that the system is now in a 2D space. Again, the equations are simple that we have two equations for the two coordinates.

$$\begin{cases} \frac{d^2 q_1}{dt^2} = -\frac{k}{m} q_1 \\ \frac{d^2 q_2}{dt^2} = -\frac{k}{m} q_2 \end{cases}$$

The analytical solution is exactly the same, but now there are two of them with differing amplitudes and phases depending on the initial condition. And now we have

$$\mathbf{f}(X) = \begin{pmatrix} \mathbf{a}_1(X) \\ \mathbf{a}_2(X) \\ \mathbf{v}_1(X) \\ \mathbf{v}_2(X) \end{pmatrix},$$

where \mathbf{a}_1 and \mathbf{a}_2 are the dynamics or time derivative for \mathbf{p}_1 and \mathbf{p}_2 respectively; similarly \mathbf{v}_1 and \mathbf{v}_2 are the time derivative of \mathbf{q}_1 and \mathbf{q}_2 . In this system, the energy is the sum of energy in the two directions so $E = \frac{m(p_1^2 + p_2^2)}{2} + \frac{k(q_1^2 + q_2^2)}{2}$ and so the invariance $L(a_1, a_2, v_1, v_2) = \frac{dE}{dt} = mp_1 a_1 + mp_2 a_2 + kq_1 v_1 + kq_2 v_2 = 0$. Now our naive baseline GP has the kernel

$$K(X, X') = \begin{pmatrix} K_{RBF, a_1}(X, X') & 0 & 0 & 0 \\ 0 & K_{RBF, a_2}(X, X') & 0 & 0 \\ 0 & 0 & K_{RBF, v_1}(X, X') & 0 \\ 0 & 0 & 0 & K_{RBF, v_2}(X, X') \end{pmatrix}.$$

We will then have the joint distribution of

$$\begin{pmatrix} \mathbf{f}(X) \\ L[\mathbf{f}(X_L)] \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0_{4n} \\ 0_\ell \end{pmatrix}, \begin{pmatrix} A & B \\ C & D \end{pmatrix} \right),$$

with

$$A = K(X, X'), B = \begin{pmatrix} K_{RBF, a_1} \\ K_{RBF, a_2} \\ K_{RBF, v_1} \\ K_{RBF, v_2} \end{pmatrix} \odot \begin{pmatrix} mP_{1L} \\ mP_{2L} \\ kQ_{1L} \\ kQ_{2L} \end{pmatrix}, C = B^T$$

$$D = K_{RBF, a_1} m^2 \odot (p_{1L} \otimes p_{1L}) + K_{RBF, a_2} m^2 \odot (p_{2L} \otimes p_{2L}) \\ + K_{RBF, v_1} k^2 \odot (q_{1L} \otimes q_{1L}) + K_{RBF, v_2} k^2 \odot (q_{2L} \otimes p_{2L})$$

where the terms are defined the exact same way as before, but now with respect to different coordinates, and the derivation are also the same. We will then again take the Schur Complement.

3.4.2 Non Linear

Now we will to introduce double pendulum, which is a fairly nonlinear system and quite complicated. It has two mass blobs m_1 and m_2 as well as two lengths for the pendulum

stem ℓ_1 and ℓ_2 . The defining equations are as follows:

$$\begin{cases} \frac{d^2 q_1}{dt^2} = \frac{-g(2m_1+m_2)\sin q_1 - m_2 g \sin(q_1-2q_2) - 2\sin(q_1-q_2)m_2(p_2^2 l_2 + p_1^2 l_1 \cos(q_1-q_2))}{l_1(2m_1+m_2-m_2 \cos(2q_1-2q_2))} \\ \frac{d^2 q_2}{dt^2} = \frac{2\sin(q_1-q_2)(p_1^2 l_1(m_1+m_2) + g(m_1+m_2)\cos q_1 + p_2^2 l_2 m_2 \cos(q_1-q_2))}{l_2(2m_1+m_2-m_2 \cos(2q_1-2q_2))} \end{cases}$$

As we can see, it is very complicated in term of the form. We also have form of energy

$$E = -(m_1+m_2)gl_1 \cos q_1 - m_2 gl_2 \cos q_2 + \frac{m_1 l_1^2 p_1^2}{2} + \frac{m_2}{2}(l_1^2 p_1^2 + l_2^2 p_2^2 + 2l_1 l_2 p_1 p_2 \cos(q_1 - q_2))$$

While the form is more complicated, the underlying principle to construct the invariance kernel. We again differentiate with respect to time and set it to zero to obtain the invariance equation to obtain

$$L(a_1, a_2, v_1, v_2) = \frac{dE}{dt} = (m_1 + m_2)gl_1 \sin \theta_1 v_1 + m_2 gl_2 \sin \theta_2 v_2 + m_1 l_1^2 \dot{\theta}_1 a_1 + m_2(l_1^2 \dot{\theta}_1 a_1 + l_2^2 \dot{\theta}_2 a_2 + l_1 l_2(\dot{\theta}_2 \cos(\theta_1 - \theta_2)a_1 + \dot{\theta}_1 \cos(\theta_1 - \theta_2)a_2 - \dot{\theta}_1 \dot{\theta}_2 \sin(\theta_1 - \theta_2)(v_1 - v_2))) = 0$$

The form of the invariance matrix is too cubersome to write down, but the derivation is as straightforward as before.

4 Learning Invariance

Previously, we have assumed the knowledge of the dynamics equation from Newton's Law, which allowed us to derive the energy and thereafter the invariance equation. However, for the kernel to be useful in real life, we cannot assume the equations of the system beforehand, since if we have the knowledge, then the system is pretty much solved that we can just use an ODE integrator to obtain the trajectories. Therefore, we should find a way for the system to learn the form of invariance from data directly. One way to do so is to parameterise our invariance function and allow the system to learn the coefficients by maximising the log marginal likelihood.

4.1 1D system

A simple way to parameterise an unknown function, when the problem is relatively bounded is to use polynomial basis to parameterise the unknown function $L(a, v)$. In 1D cases, for our examples, they are simple enough so that $L(a, v)$ is a sum of a function of a and p only, $f_p(a)$ and a different function of v and q only, $g_q(v)$ so that $L(a, v) = f_p(a) + g_q(v)$. In the 1D SHM case, $f_p(a) = mpa$ and $g_q(v) = kqv$. For simple pendulum, we have $f_p(a) = \ell pa$ and $g_q(v) = \sin qv$. We will therefore set

$$f_p(a) = \sum_{i=0}^{n_f} c_{f,i} p^i, g_q(a) = \sum_{i=0}^{n_g} c_{g,i} q^i,$$

where $c_{f,i}, c_{g,i}$ are the parameter we will optimise based on marginal likelihood and we will evaluate different degree of polynomial n_f, n_g and see which gives the best results.

4.2 2D system

In 2D, things are a bit more complicated. For example, it is very clear that in double pendulum, the invariance function is not separable. Instead, we need to consider all combinations of polynomials; i.e. multivariate polynomial. Therefore, for 2D we have $L(a, v) = f_1 a_1 + f_2 a_2 + g_1 v_1 + g_2 v_2$. We will parameterise each of f_1, f_2, g_1, g_2 with

$$\sum_{i,j,k,l=0}^n c_{ijkl} p_1^i p_2^j q_1^k q_2^l$$

Therefore, there will be many more coefficients.

5 Experiments

5.1 Data Generation

To generate the data, we use RK4 integrator to generate data using the ODEs for each problem. We will also add some small Gaussian noise 10^{-8} to the dynamics to reflect the real world better. Since we aim to predict the value of dynamics at any point in the input space, we will need the dynamics as targets. This is done by combining forward and backward difference to compute the finite gradient. We also choose the time step to be 0.01 to ensure sufficient accuracy for the integrator. For training data, we will choose a few random starting point within a range(exact details depends on the task as will be explained later) in the input space and generate data from the point for a number of timesteps. For testing data, we will do the same both in the same range and outside the range, we will pick a few starting points and then average the result. Without loss of generality, we choose $m = k = 1$ and $g = \ell = 1$ for 1D and 2D SHM and simple pendulum for simplicity. For the damping, we will choose $\gamma = 0.1$. For double pendulum, we also let $m_1 = m_2 = \ell_1 = \ell_2 = g = 1$.

5.2 Evaluation Method

We will use a few different evaluation metrics.

1. Generate future trajectories from a random starting point and compare to true trajectories
2. Evaluate the difference between the true analytical dynamics and the predicted dynamics on a grid
3. Log marginal likelihood

If we are learning the invariance, we can also further compare the true invariance and the learnt invariance. Or we can assess the conservation of energy over time.

5.3 Implementation Technicalities

We need to add jitter to the D matrix in computation due to the fact that we need the invert it and because of numerical issues, it will be more stable if we add a little bit of jitter in the range of 10^{-6} and its effect will be explained in the sections below. Similarly, to stabilise the algorithm, when backpropagating the hyperparameters, we need to narrow range of search so it's not too crazy. For example, we wouldn't expect

a lengthscale of 100 in our example with maximum range of 5; therefore, we will set for example, the maximum range of search is 10^{-3} to 5 for the length scale and similarly for the variance term. Another changes we made in the 2D cases is to scale the data using MinMax scaler to squash them to be between ± 1 . This is done since the polynomial behaves better around this region.

Local Invariance

Conditioning on a grid of invariance points is simply not scalable in high dimensions. If we have an invariance density of 1 in a ± 5 range for 2D SHM, then we will have 10000 points. We will easily run out of memory and also taking the matrix inverse are very difficult and time consuming. As a result, to make the approach scalable to 2D space, we will use the idea of local invariance. Which is to take samples in the local region around training and testing points and condition on these points, where we sample uniformly around the spaces of these points.

5.4 1D

We will illustrate a lot of examples in 1D cases since the results are easy to visualise, but the principle readily extend to high dimensional spaces.

5.4.1 SHM

We start with the 1D SHM case, and for the data we will have only two trajectory starting randomly for $-3 \leq q \leq 3$ and p to be between ± 0.3 so that it does not overshoot the boundary, and we allow it to run for 0.1 seconds so there will be in total 10 training points. We will then draw three test point from the same range as training, to test the performance of the future predictive power for 3 seconds (300 points) and average over the three trajectories. In figure 5.1, we will choose one at random and draw the future trajectories. The performance is summarised in table 5.1 with jitter 1×10^{-5} with invariance density of 20 within ± 3 . We have also plotted the posterior distribution of using invariance kernel and compare that to using naive RBF kernel in figure 5.2. We can see the invariance has much better predictive power and generalisability.

We will test all combination of polynomial of degree 0, 1, 2, 3 so there will be 16 of them, up to the cubic terms, and we will choose the ones with highest marginal likelihood. We initialise the coefficients with normal distribution with variance of 10^{-3} , we also restrict the search between ± 1 . Data is again between ± 3 so is the invariance conditioning. The best polynomial pair is

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	74.66	86.28	83.32
MSE	3.5563	0.1057	0.0927

Table 5.1: SHM Invariance Kernel Performance

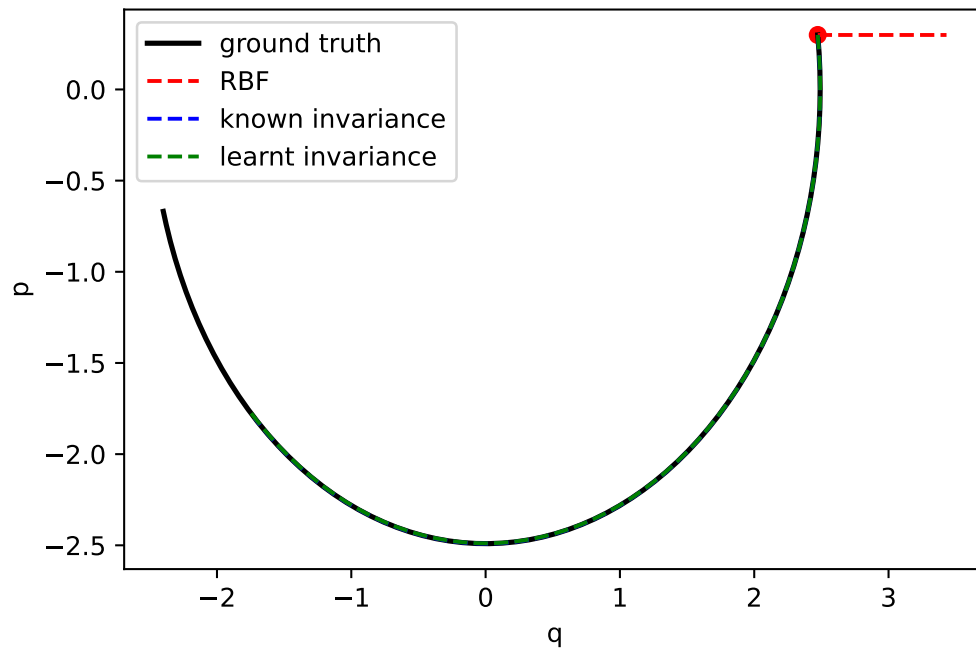
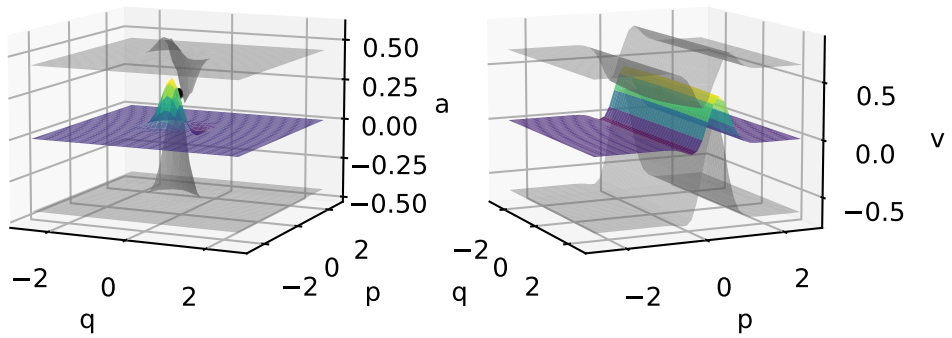


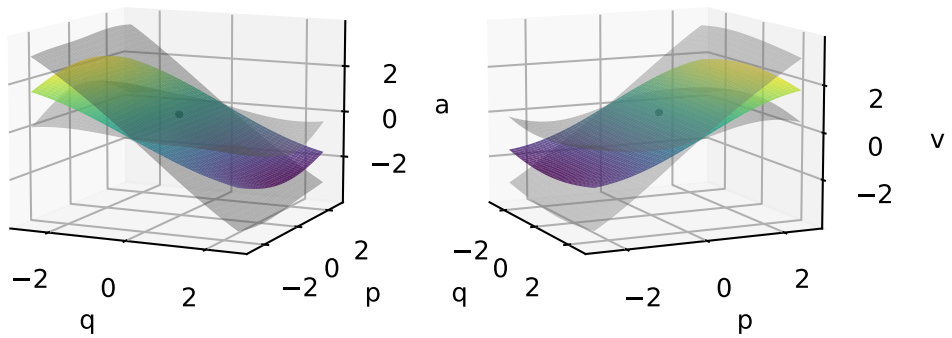
Figure 5.1: SHM predicted trajectory

Invariance GP Posterior of a Invariance GP Posterior of v



(a) SHM RBF posterior

Invariance GP Posterior of a Invariance GP Posterior of v



(b) SHM invariance posterior

Figure 5.2: SHM posteriors

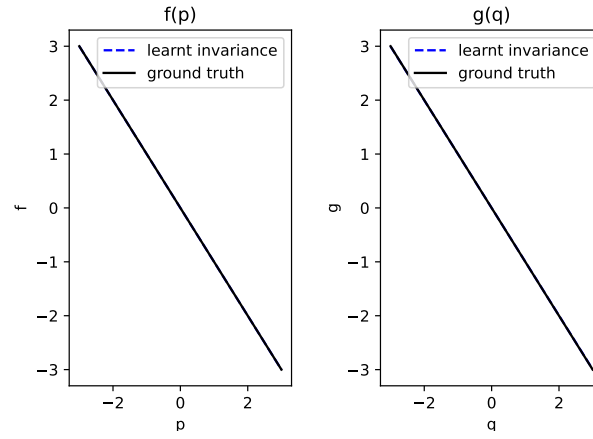


Figure 5.3: Learnt invariance for SHM

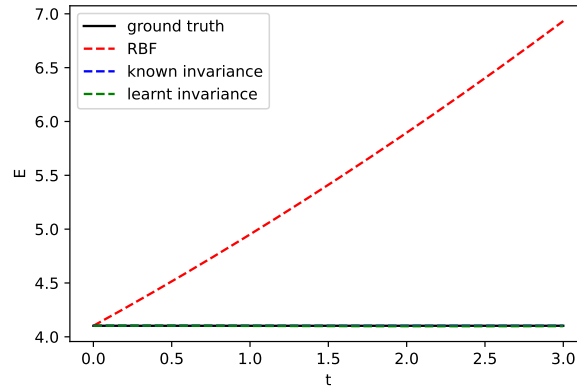


Figure 5.4: Energy conservation for SHM

5.4.2 Pendulum

The procedure is pretty much the same here that we will draw 3 single training example between $q = \pm 150^\circ$ and $\pm 10^\circ$ for p . The rest of the procedure is the same as the SHM case. Again, the invariance is conditioned between $\pm 150^\circ$ with invariance density of 30 with jitter 1×10^{-5} . We again have the prediction performance summarised in table 5.2 as well as figures 5.5, 5.6.

The best polynomial degree is

We again see a much better improvement in performance even in this nonlinear case.

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	233.89	262.59	255.76
MSE	0.05	0.0025	0.0020

Table 5.2: Pendulum

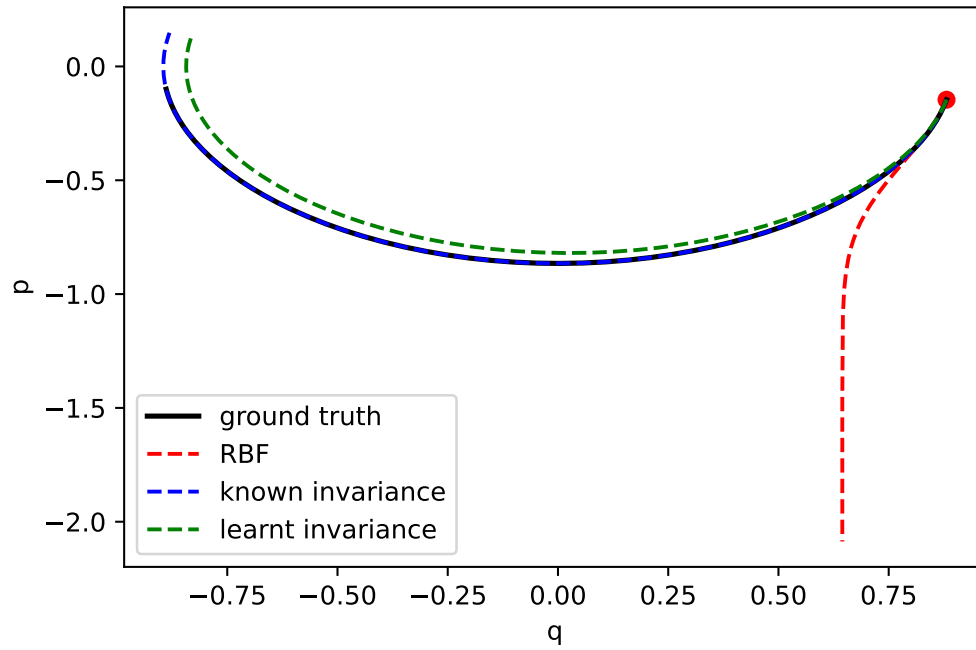
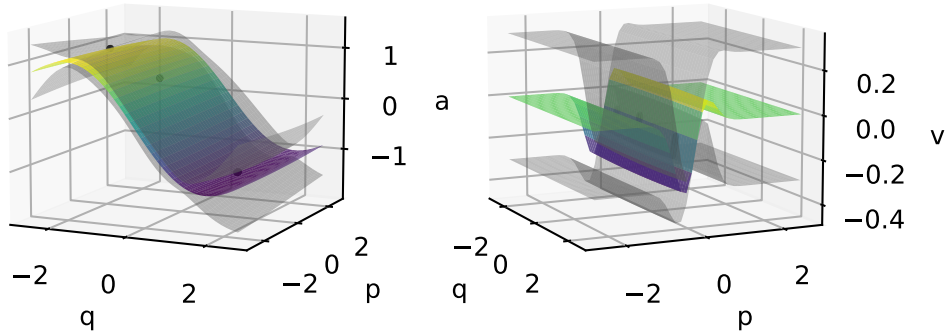


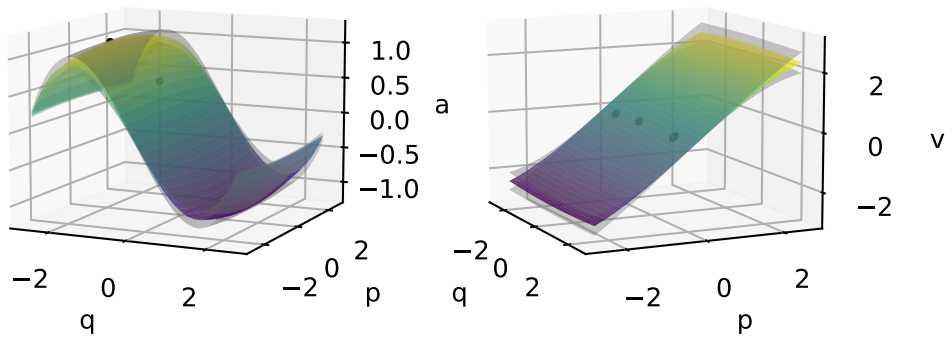
Figure 5.5: pendulum predicted trajectory

Invariance GP Posterior of a Invariance GP Posterior of v



(a) pendulum RBF posterior

Invariance GP Posterior of a Invariance GP Posterior of v



(b) pendulum invariance posterior

Figure 5.6: pendulum posteriors

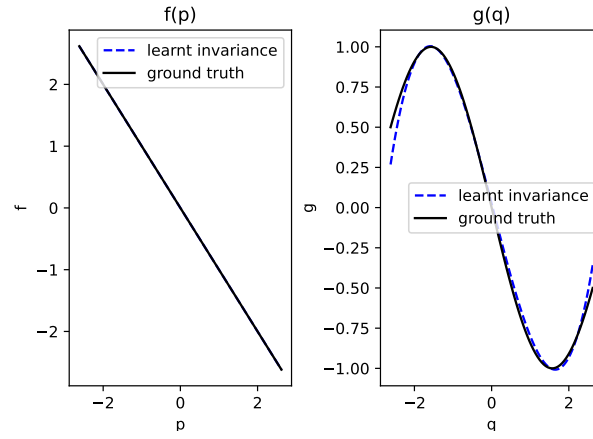


Figure 5.7: Learnt invariance for pendulum

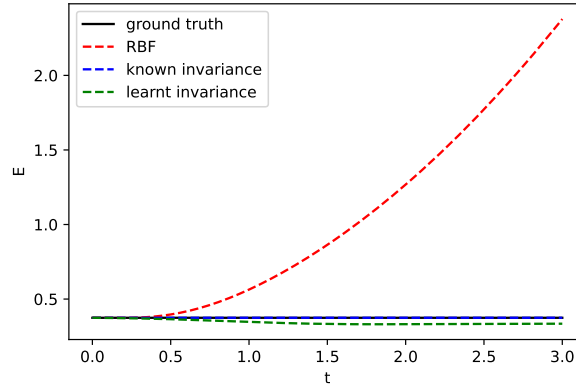
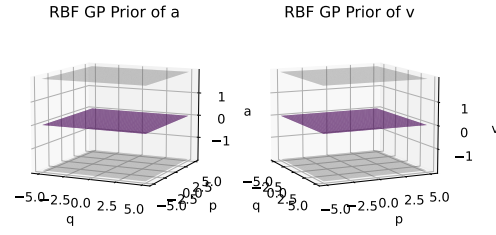
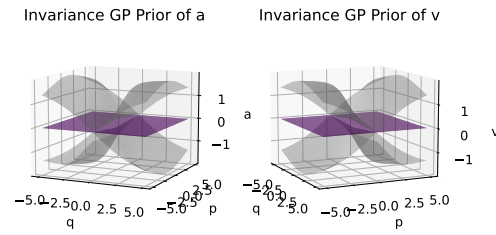


Figure 5.8: Energy conservation for pendulum

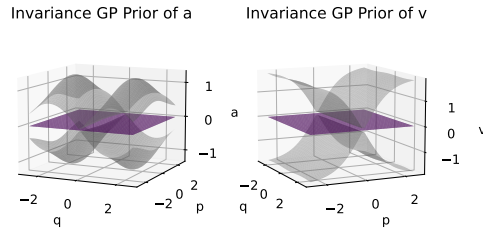
Here we also include the prior distribution of RBF kernel in figure 5.9a, SHM invariance kernel in figure 5.9b and pendulum invariance kernel in figure 5.9c. We can see the prior knowledge of the dynamics, such as the linearity in SHM as well as the sinusoidal nature of pendulum dynamics.



(a) RBF prior



(b) SHM invariance prior



(c) Pendulum invariance prior

Figure 5.9: Priors

5.4.3 Invariance Density

Here we will to explore how the change of invariance density could affect the performance of the kernel as well as degree of freedom. Since 1D SHM is rather trivial, I will test it on the nonlinear pendulum case. We will assess its prediction power using the same method above and we will also calculate its degree of freedom using equation. We can

see that beyond a certain point, the increase in performance is not that significant but the computation cost in terms of time and space all increases dramatically. Therefore, from here, we will only use the minimum number of invariance points that ensure the performance of the invariance kernel.

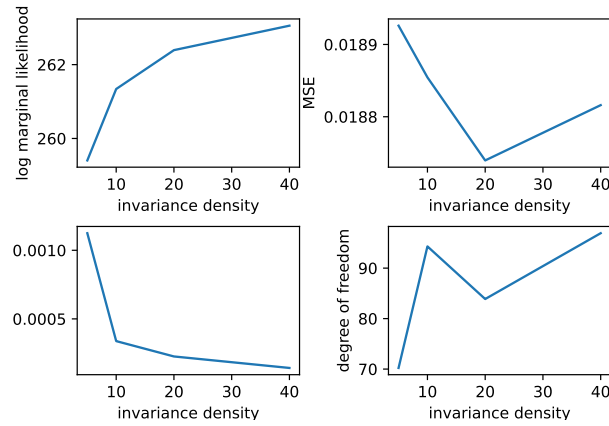


Figure 5.10: Effect of varying invariance density

5.4.4 Effect of Jitter

Since we need a little bit of jitter to stabilise the computation of our invariance kernel as explained above, and that adding jitter is the same as adding noise to the invariance (since we are adding it on the diagonal of D submatrix), we need to make sure it is not too noisy. So we will again assess the performance using different amount of jitter. We can again see beyond a certain point, the performance does not change. While the lower the jitter, the higher the performance should be in theory; the instability of matrix inversion and numerical errors could cause a drop in performance. As a result, now we will also choose the smallest amount of jitter that remains stable.

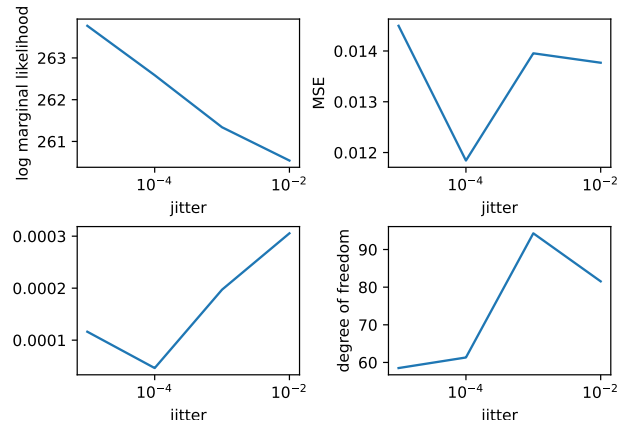


Figure 5.11: Effect of varying jitter

5.4.5 Damped SHM

In damped system, we introduce a damping of $\gamma = 0.1$, which is very high already. The rest of the setting is essentially the same as undamped SHM case. However, this time we use double the amount of data (20 points) since it is a more difficult problem. The results is again summarised in table 5.3 and figure 5.12, and the learnt ϵ is 0.0297, which is not too high a noise. We see that even in damped cases, the knowledge of conservation of energy is helpful.

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	67.69	74.84	77.43
MSE	0.9226	0.4772	0.2405

Table 5.3: Damped SHM performance

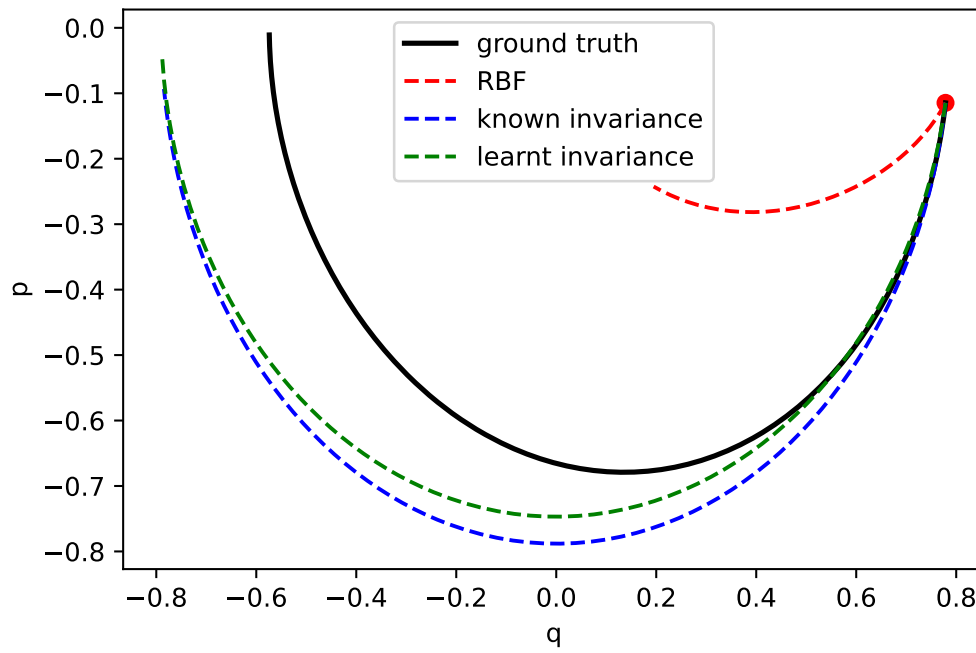


Figure 5.12: damped shm predicted trajectory

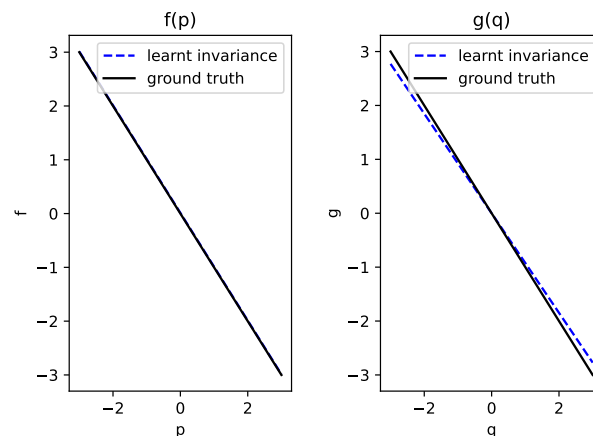


Figure 5.13: Learnt invariance for damped shm

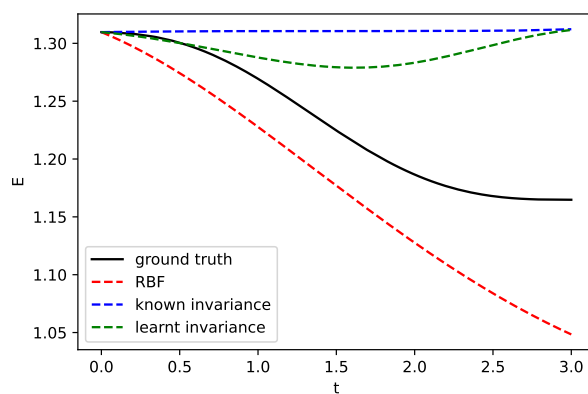


Figure 5.14: Energy conservation for damped shm

5.4.6 Damped Pendulum

Again, the setting is almost the same as the undamped pendulum, and $\gamma = 0.1$. Also, this time we kept at 10 data points since the phase space is smaller than that of SHM. The results is summarised in table 5.4 and figure 5.15, and the learnt ϵ is 0.0014, again not too high. We again see an improvement using invariance kernel in this damped nonlinear case.

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	229.96	243.41	243.48
MSE	0.5983	0.0323	0.0877

Table 5.4: Damped Pendulum Performance

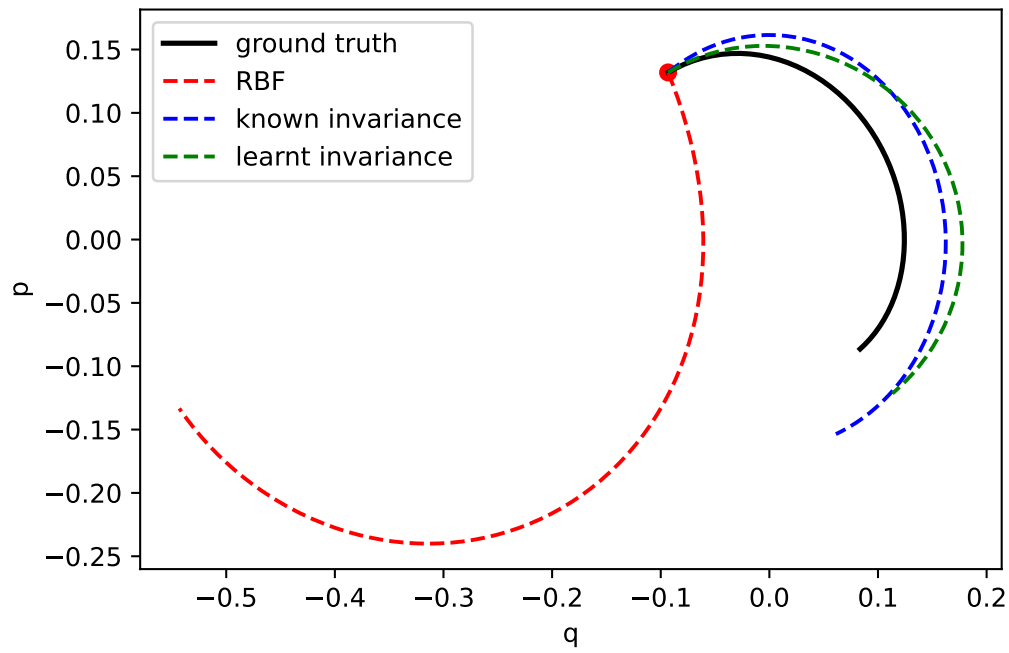


Figure 5.15: damped pendulum predicted trajectory

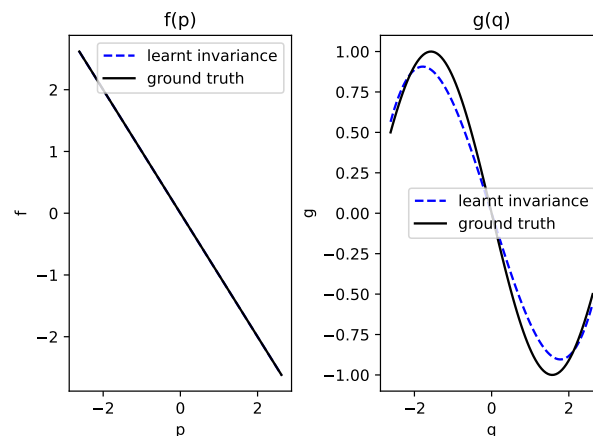


Figure 5.16: Learnt invariance for damped pendulum

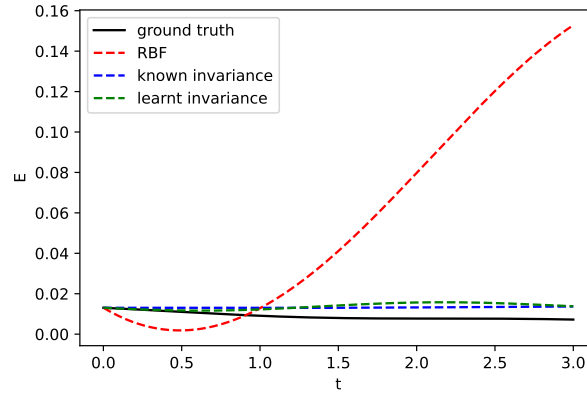


Figure 5.17: Energy conservation for damped pendulum

5.5 2D

5.5.1 SHM

Here we again use the same set of input space as the 1D space, but double the dimensions. The results are summarised in table 5.5 and figure 5.18. Since 4 dimensional space is difficult to visualise, only individual prediction of the trajectory is provided.

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	382.50	439.74	430.80
MSE	8.6260	2.8882	4.1416

Table 5.5: 2D SHM Invariance performance

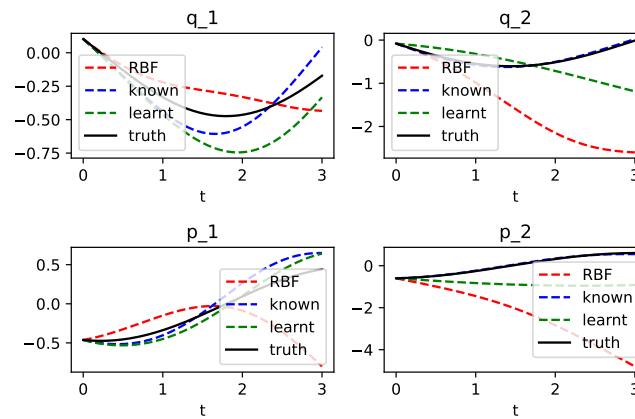


Figure 5.18: 2D SHM prediction

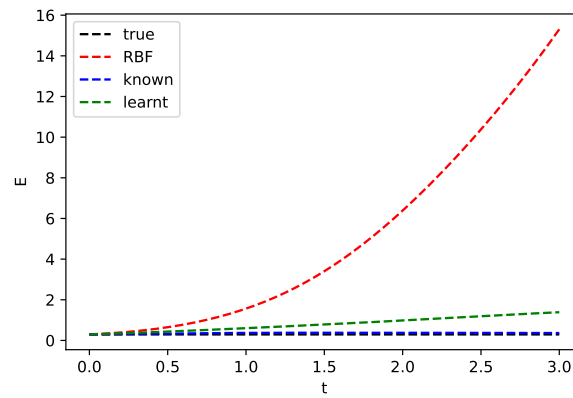


Figure 5.19: 2D SHM energy

5.5.2 Double Pendulum

With double pendulum, we have angles being between $\pm 60^\circ$ and angular velocity between $\pm 10^\circ$. The results are summarised in table 5.6 and figure 5.20.

Method	RBF	Known Invariance	Learnt Invariance
Log Marginal Likelihood	452.34	469.75	465.85
MSE	0.5592	0.2030	0.1995

Table 5.6: Double pendulum Invariance performance

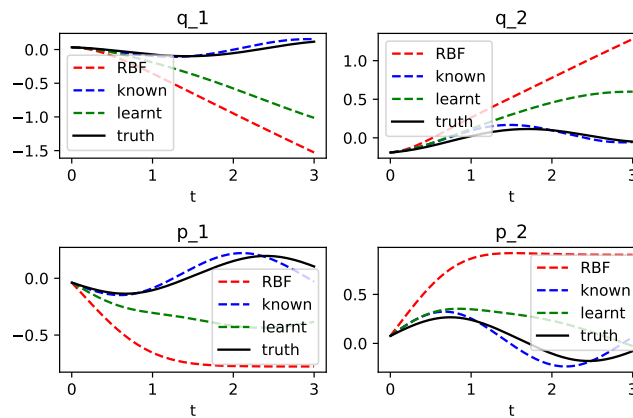


Figure 5.20: 2D SHM prediction

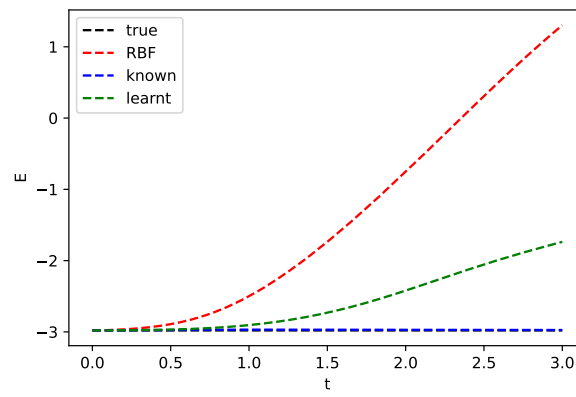


Figure 5.21: 2D SHM energy

6 Conclusion

Conclusion goes here.

References

- Paul Glendinning. *Stability, instability and Chaos: An introduction to the theory of nonlinear differential equations*. Cambridge University Press, 1994.
- Alexander G. de G. Matthews, Mark van der Wilk, Tom Nickson, Keisuke. Fujii, Alexis Boukouvalas, Pablo León-Villagr , Zoubin Ghahramani, and James Hensman. GPflow: A Gaussian process library using TensorFlow. *Journal of Machine Learning Research*, 18(40):1–6, apr 2017. URL <http://jmlr.org/papers/v18/16-537.html>.
- Carl Edward Rasmussen and Christopher K I Williams. *Gaussian process for machine learning*. The MIT Press, 2006.
- Steven Strogatz. *Nonlinear Dynamics and Chaos: With applications to physics, Biology, Chemistry and Engineering*. CRC Press, 2019.