

Diffeomorphic Learning Based on ODE Flows

Scientific work within the Ingenieurspraxis
from the Department of Electrical and Computer Engineering at the
Technical University of Munich.

Supervised by Univ.-Prof. Dr.-Ing. Sandra Hirche
M. Sc. Petar Bevanda
Chair of Information-Oriented Control

Submitted by cand. Ing. Adrian Brünger
Willi-Graf Str. 17/0704
80805 Munich
0049 1631925390

Submitted on Munich, 11.04.2021

April 7, 2022

INGENIEURPRAXIS
for

Adrian Brünger
Student ID 03715546, Degree EI

Diffeomorphic learning based on ODE flows

Problem description:

When considering stable systems, stability guarantees are of great importance to allow for safety and physical consistency even for unseen states. Conventional learning-based approaches via function approximators such as neural networks or kernel methods for modeling dynamical systems need not reflect the system's physical properties such as stability [2]. Often, however stability conditions are enforced only locally on data points and verifying stability in the whole state space leads to the curse of dimensionality. However, simplifying the conditions by structure and equivalences has shown to be a promising direction in recent years. For that, we use a smooth invertible transform (diffeomorphism) to fit the data to a latent linear system whose stability is trivially ensured. To learn such transforms, a method based on ordinary differential equations (ODE) solution curves is considered as it ensures invertibility. Together with a sufficiently smooth ODE learner, it forms a diffeomorphism. A seminal method employing such an approach is Neural ODE (NODE) [1], learning a neural network that delivers an invertible map when integrated over time. To simplify the training process, however, employing a kernel-based linear regressor in an inducing-state kernel approach is studied along with an approach based on random Fourier features approximating the kernel of interest. The performance of the aforementioned methods is evaluated on a handwriting dataset.

Work schedule:

- Literature research
- Designing a kernel-based learner
- Verification on a handwriting dataset
- Comparison to state-of-the-art

- [1] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [2] Muhammad Asif Rana, Anqi Li, Dieter Fox, Byron Boots, Fabio Ramos, and Nathan D. Ratliff. Euclideanizing flows: Diffeomorphic reduction for learning stable dynamical systems. In *L4DC*, volume 120 of *Proceedings of Machine Learning Research*, pages 630–639. PMLR, 2020.

Supervisor: M. Sc. Petar Bevanda
Start: 22.11.2021
Delivery: 11.04.2022

(S. Hirche)
Univ.-Professor

Abstract

Considering an imitation learning problem, smooth, bijective maps (diffeomorphisms) can be utilized to guarantee stability for the learned dynamic system, by linking it to a stable *latent* system. In this work, an *inducing state* kernel method and a *random Fourier feature* kernel approximation method, constructing such mappings, are presented. The kernel-based approaches are compared to the state-of-the-art *neural ordinary differential equations* method. Evaluating the approaches on a handwriting dataset, results show that the methods are capable of learning dynamic systems reproducing serverly curved trajectories, closely rembling the demonstraions. Further, due to its simplicity, the *random Fourier feature* method proves to be computationally faster then the state-of-the-art approach.

Contents

1	Introduction	5
1.1	Problem Setup	6
1.1.1	Learning Dynamical Systems from Demonstrations	6
1.1.2	Stability	6
2	Technical Approach	7
2.1	Theoretical Foundations	7
2.1.1	Diffeomorphisms and Riemannian Geometry	7
2.1.2	Natural Gradient Descent	9
2.1.3	Stability	10
2.1.4	Flow-based Construction of Diffeomorphisms	11
2.2	Infinitesimal Generators	11
2.2.1	Kernel Methods	11
2.2.2	Neural Ordinary Differential Equations (NODE)	13
2.3	Implementation	14
3	Evaluation and Conclusion	15
3.1	Experimental Results and Performance	15
3.2	Conclusion and Future Work	16
	Bibliography	19

Chapter 1

Introduction

Diffeomorphic learning is an imitation learning approach, in which a stable dynamic system is learned in order to reproduce human expert demonstrations. In many robot applications, imitation learning is a natural approach for teaching a robot motions, without having to model explicit functional relationships or cost functions. Hereby, it is of major importance to guarantee stability for the learned system, to avoid uncontrolled behavior, which is especially crucial in environments involving human-robot interaction.

By introducing a smooth bijective map, a diffeomorphism, linking the target system to a stable *latent* system, one can ensure stability of the target system. When defining such a map via flows generated by ordinary differential equations (ODEs), computations of the inverse mapping are trivial. In particular when dealing with highly dynamic surroundings, fast evaluations of learned dynamics are necessary and adaptability - in terms of e.g. obstacle avoidance - is requested.

Previous approaches, such as [RLF⁺20, UGTP20], construct a diffeomorphism a layered fashion by composing invertible coupling neural networks. These approaches succeed in learning stable dynamic systems, however, the layered structure gives little insight in how the networks construct the diffeomorphism, which increases the difficulty of adapting the target system to changes in the environment. Another downside is computational complexity. A different approach is presented in [PSC16], which realizes a fast evaluation of the learned dynamics, while only being able to learn from one example trajectory.

This work introduces a diffeomorphic learning framework and focuses on explaining and evaluating kernel-approximation-based methods, to simplify training and minimize prediction times. Utilizing the presented framework, adaptability can be achieved by simply composing learned and customizable diffeomorphic mappings to incorporate desired behavior while maintaining stability.

1.1 Problem Setup

1.1.1 Learning Dynamical Systems from Demonstrations

Let \mathbb{R}^d be a d dimensional state space. A nonlinear time-invariant first order dynamical system of the form

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t)), \mathbf{x}(t_0) = \mathbf{x}_0 \quad (1.1)$$

can be interpreted as a mapping of a system state \mathbf{x} , e.g. cartesian position or angular coordinates of a robot manipulator, to a velocity $\dot{\mathbf{x}}$ through the non-linear map $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^d$. The initial state of the system is \mathbf{x}_0 . Without loss of generality t_0 is set to $t_0 = 0$ in the following.

Given a set of N discretized trajectories (demonstrations) $\Xi = \{\mathbf{x}_{i,1}, \mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,l_i}\}_{i=1}^N$ of length l_i , with $\mathbf{x}_{i,j} \in \mathbb{R}^d \forall i, j$, the goal is to find a stable dynamic system (1.1), that reproduces the demonstrations best. \mathbf{f} is constructed to be globally Lipschitz continuous to guarantee existence and uniqueness of solutions.

Assuming the velocities $\dot{\mathbf{x}}_{i,j}$ and \mathbf{W} parametrizes \mathbf{f} , the least squares loss

$$MSE(\mathbf{W}) = \sum_{i=1}^N \sum_{j=1}^{l_i} \|\dot{\mathbf{x}}_{i,j} - \hat{\dot{\mathbf{x}}}_{i,j}\|^2 = \sum_{i=1}^N \sum_{j=1}^{l_i} \|\dot{\mathbf{x}}_{i,j} - \mathbf{f}_{\mathbf{W}}(\mathbf{x}_{i,j})\|_2^2 \quad (1.2)$$

is minimized with respect to \mathbf{W} , under the constraint that

$$\dot{\hat{\mathbf{x}}} = \mathbf{f}_{\mathbf{W}}(\mathbf{x}) \quad (1.3)$$

is stable.

Trajectory predictions are then generated by integrating the learned dynamic system (1.3).

$$\hat{\boldsymbol{\xi}}(t, \mathbf{x}_0) = \mathbf{x}_0 + \int_0^t \hat{\dot{\mathbf{x}}}(u) du. \quad (1.4)$$

1.1.2 Stability

An equilibrium \mathbf{x}^* is a state in which the velocity is zero $\dot{\mathbf{x}}^* = \mathbf{f}(\mathbf{x}^*) = \mathbf{0}$. It further is a Lyapunov stable equilibrium if there exists a $\delta > 0$ for every $\epsilon > 0$, so that for any initial state \mathbf{x}_0 , that satisfies $\|\mathbf{x}_0 - \mathbf{x}^*\| < \delta$, the state satisfies $\|\mathbf{x}(t) - \mathbf{x}^*\| < \epsilon$ for all $t > 0$. In other words, an equilibrium is stable if for every neighborhood U of \mathbf{x}^* all trajectories starting in a neighborhood $V \subseteq U$ of \mathbf{x}^* remain in U . An equilibrium is asymptotically stable if it is Lyapunov stable and all trajectories starting in V converge to the equilibrium $\lim_{t \rightarrow \infty} \|\mathbf{x}(t) - \mathbf{x}^*\| \rightarrow 0$. A dynamic system is asymptotically stable if the equilibrium is unique and $U = \mathbb{R}^d$.

Linear dynamic system of the form $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$, $\mathbf{x}(0) = \mathbf{x}_0$, with \mathbf{A} being non-singular, is stable for negative definite \mathbf{A} .

Chapter 2

Technical Approach

The overall goal is to learn a stable dynamic system that imitates certain motions by minimizing a distance metric between demonstrated and predicted trajectories. For $\mathbf{z} \in \mathbb{R}^d$ consider a gradient system

$$\dot{\mathbf{z}} = \mathbf{g}(\mathbf{z}) = -\nabla_{\mathbf{z}}\Phi(\mathbf{z}) \quad (2.1)$$

with convex potential

$$\Phi = \|\mathbf{z} - \mathbf{z}_{goal}\|_2^2. \quad (2.2)$$

The goal of the following section is now to briefly explain the elegant approach of assuming that there exists a bijective mapping \mathbf{d} between d -dimensional Euclidean space \mathbb{R}^d (*latent* space) and a d -dimensional non-Euclidean topological space \mathcal{M} , so that trajectories generated by a simple, *latent* gradient descent system (2.1) in \mathbb{R}^d and mapped to \mathcal{M} closely match the demonstrated trajectories. Hereby it is presumed, that the demonstrations are generated by a globally asymptotically stable first order system with one unique equilibrium, so that is exactly linearizable by such a coordinate change \mathbf{d} .

2.1 Theoretical Foundations

2.1.1 Diffeomorphisms and Riemannian Geometry

In this section, based on [Lee06], basic concepts of Riemannian geometry and properties of diffeomorphisms are discussed to gain some intuition on the methodology used in this work. A d -dimensional manifold \mathcal{M} in general is a topological space that locally resembles Euclidean space \mathbb{R}^d . Riemannian geometry studies Riemannian manifolds, which are smooth manifolds with a Riemannian metric. A Riemannian metric g is a positive definite 2-tensor field that determines an inner product on each tangent space of $T_p\mathcal{M}$. This is often written as $g(\mathbf{X}, \mathbf{Y}) = \langle \mathbf{X}, \mathbf{Y} \rangle$ for $\mathbf{X}, \mathbf{Y} \in T_p\mathcal{M}$. The Riemannian metric allows for definitions of geometric concepts,

such as lengths of curves and consequently distances between points on \mathcal{M} . A diffeomorphism $\mathbf{d} : \mathcal{M}_1 \rightarrow \mathcal{M}_2$ is a smoothly differentiable, bijective mapping between two smooth d -dimensional manifolds \mathcal{M}_1 and \mathcal{M}_2 with a smoothly differentiable inverse mapping $\mathbf{d}^{-1} : \mathcal{M}_2 \rightarrow \mathcal{M}_1$. To avoid confusion, these notations hold for the remainder of this report.

$\mathbf{z} \in \mathbb{R}^d$	Point/State in d -dimensional <i>latent</i> Euclidean space
$\mathbf{x} \in \mathcal{M}$	Point/State on d -dimensional Riemannian manifold
$\mathbf{d} : \mathcal{M} \rightarrow \mathbb{R}^d : \mathbf{d}(\mathbf{x}) = \mathbf{z}$	Diffeomorphism
$\mathbf{d}^{-1} : \mathbb{R}^d \rightarrow \mathcal{M} : \mathbf{d}^{-1}(\mathbf{z}) = \mathbf{x}$	Inverse diffeomorphism
$g : T_p\mathcal{M} \times T_p\mathcal{M} \rightarrow \mathbb{R}$	Positive definite Riemannian metric
$\mathbf{G} : \mathcal{M} \rightarrow \mathbb{R}^{d \times d}$	Positive definite Riemannian metric matrix
\mathbf{W}	Learning parameters

Following the approach sketched above, it can be concretized, that the diffeomorphism maps from $\mathcal{M}_1 = \mathbb{R}^d$ to $\mathcal{M}_2 = \mathcal{M}$. A Riemannian metric in euclidean space evaluated in $\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^d$ can be represented as an inner product between the Euclidean tangent space vectors in $\mathbf{z}_1, \mathbf{z}_2$, which are the velocities $\dot{\mathbf{z}}_1, \dot{\mathbf{z}}_2$.

$$g^{\mathbb{R}}(\dot{\mathbf{z}}_1, \dot{\mathbf{z}}_2) := \dot{\mathbf{z}}_1^T \dot{\mathbf{z}}_2 := \dot{\mathbf{z}}_1^T \mathbf{I} \dot{\mathbf{z}}_2 \quad (2.3)$$

$\mathbf{I} \in \mathbb{R}^{d \times d}$ is the identity matrix. Since the diffeomorphism $\mathbf{d}(\mathbf{x}) = \mathbf{z}$ is differentiable, a tangent space vector/velocity vector on the manifold can be written as

$$\frac{d}{dt} \mathbf{d}(\mathbf{x}(t)) = \mathbf{J}_d(\mathbf{x}(t)) \dot{\mathbf{x}}(t) = \dot{\mathbf{z}}(t) . \quad (2.4)$$

With $\mathbf{J}_d(\mathbf{x}(t))$ being the Jacobian of the diffeomorphism evaluated in $\mathbf{x}(t)$. In the following time dependencies and state dependencies of Jacobians are hidden to simplify notation. Utilizing (2.4) a Riemannian metric $g^{\mathcal{M}}$ of the Riemannian manifold \mathcal{M} can be defined as

$$\begin{aligned} g^{\mathbb{R}}(\dot{\mathbf{z}}_1, \dot{\mathbf{z}}_2) &= g^{\mathbb{R}} \left(\frac{d}{dt} \mathbf{d}(\mathbf{x})|_{\mathbf{x}_1}, \frac{d}{dt} \mathbf{d}(\mathbf{x})|_{\mathbf{x}_2} \right) = \left(\frac{d}{dt} \mathbf{d}(\mathbf{x})|_{\mathbf{x}_1} \right)^T \frac{d}{dt} \mathbf{d}(\mathbf{x})|_{\mathbf{x}_2} \\ &= (\mathbf{J}_d \dot{\mathbf{x}}_1)^T \mathbf{J}_d \dot{\mathbf{x}}_2 = \dot{\mathbf{x}}_1^T \mathbf{J}_d^T \mathbf{J}_d \dot{\mathbf{x}}_2 = \dot{\mathbf{x}}_1^T \mathbf{G} \dot{\mathbf{x}}_2 \\ &:= g^{\mathcal{M}}(\dot{\mathbf{x}}_1, \dot{\mathbf{x}}_2) . \end{aligned} \quad (2.5)$$

Here, \mathbf{G} is the positive definite Riemannian metric matrix associated with $g^{\mathcal{M}}$, which - in this case ($g^{\mathbb{R}}, \mathbf{d} \in \mathbb{C}^1$) - also is the *pullback* of $g^{\mathbb{R}}$ by \mathbf{d} . In differential geometry there are *pushforwards* and *pullbacks*. Using a map, in this case the diffeomorphism $\mathbf{d} : \mathcal{M} \rightarrow \mathbb{R}^d$, with $\mathcal{M} \rightarrow \mathbb{R}^d$ being the "forward direction", it can be observed, that points are "pushed forward" by \mathbf{d} (as well as velocities by $\dot{\mathbf{d}}$) while functions are "pulled back", $g^{\mathbb{R}} : \mathbb{R} \rightarrow \mathbb{R}$ and $g^{\mathbb{R}} \circ \mathbf{d} : \mathcal{M} \rightarrow \mathbb{R}$. In other words $g^{\mathcal{M}}$ is the metric on the manifold \mathcal{M} induced by the map \mathbf{d} , by "pulling" the metric $g^{\mathbb{R}}$ from \mathbb{R} onto

\mathcal{M} through \mathbf{d} .

An important observation is that \mathbf{d} is an isometry due to $g^{\mathbb{R}} \circ \mathbf{d} = g^{\mathcal{M}}$. This distance preserving feature of \mathbf{d} can be sketched as follows. On a Riemannian manifold \mathcal{M} the *distance* between two points $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{M}$ can be defined as the infimum length of all piecewise continuous paths $\gamma^{\mathcal{M}} : [t_1, t_2] \rightarrow \mathcal{M}$ connecting $\mathbf{x}_1, \mathbf{x}_2$:

$$d(\mathbf{x}_1, \mathbf{x}_2) = \inf\{L(\gamma^{\mathcal{M}}) \mid \gamma^{\mathcal{M}} : [t_1, t_2] \rightarrow \mathcal{M}, \gamma^{\mathcal{M}}(t_1) = \mathbf{x}_1, \gamma^{\mathcal{M}}(t_2) = \mathbf{x}_2\} \quad (2.6)$$

L is a cost functional assigning a length scalar to each input path:

$$\begin{aligned} L(\gamma^{\mathcal{M}}) &= \int_{t_1}^{t_2} \sqrt{g^{\mathcal{M}}(\dot{\gamma}^{\mathcal{M}}(t), \dot{\gamma}^{\mathcal{M}}(t))} dt \\ &= \int_{t_1}^{t_2} \sqrt{g^{\mathbb{R}}(\dot{\gamma}^{\mathbb{R}}(t), \dot{\gamma}^{\mathbb{R}}(t))} dt \\ &= L(\gamma^{\mathbb{R}}) \\ &= L(\mathbf{d}(\gamma^{\mathcal{M}})) \end{aligned} \quad (2.7)$$

So the length of each path $\gamma^{\mathcal{M}}$ on \mathcal{M} between $\mathbf{x}_1, \mathbf{x}_2$ equals the length of its mapping to Euclidean space $\gamma^{\mathbb{R}}$ between $\mathbf{z}_1, \mathbf{z}_2$ via \mathbf{d}^{-1} . The equation (2.7) holds since (2.5) holds and \mathbf{d}^{-1} is continuously differentiable, which allows to invert (2.4):

$$\frac{d}{dt} \mathbf{d}^{-1}(\mathbf{z}) = \mathbf{J}_{\mathbf{d}}^{-1} \dot{\mathbf{z}} = \dot{\mathbf{x}}. \quad (2.8)$$

It leads to the conclusion that \mathbf{d} maps L -minimizing paths (shortest paths/geodesics) $\gamma^{\mathcal{M}}$ on between $\mathbf{x}_1, \mathbf{x}_2$ to L -minimizing paths (shortest paths/geodesics) $\gamma^{\mathbb{R}}$ between $\mathbf{z}_1, \mathbf{z}_2$. This intuition is further exploited in [2.1.2]. Furthermore, mapped points have the same distance in both spaces

$$\begin{aligned} l(\mathbf{x}_1, \mathbf{x}_2) &= \inf\{L(\gamma^{\mathbb{R}}) \mid \gamma^{\mathbb{R}} : [t_1, t_2] \rightarrow \mathbb{R}^d, \gamma^{\mathbb{R}}(t_1) = \mathbf{z}_1, \gamma^{\mathbb{R}}(t_2) = \mathbf{z}_2\} \\ &= l(\mathbf{z}_1, \mathbf{z}_2). \end{aligned} \quad (2.9)$$

Another extremely useful feature of diffeomorphisms is that compositions of diffeomorphisms are diffeomorphisms, which allows for a layered construction of diffeomorphisms as in [RLF⁺20].

2.1.2 Natural Gradient Descent

A simple approach to generating *paths* of a desired shape would be to assume that the paths of the given trajectories are geodesics on some manifold \mathcal{M} and to learn a diffeomorphism $\mathbf{d}_{\mathbf{W}}$ that maps from the manifold to Euclidean space. By mapping the start and goal point to Euclidean space one could then generate the desired path on the manifold by mapping a straight line (geodesic in \mathbb{R}^d), connecting the mapped start and goal point, back onto the manifold via $\mathbf{d}_{\mathbf{W}}^{-1}$ as shown in Figure

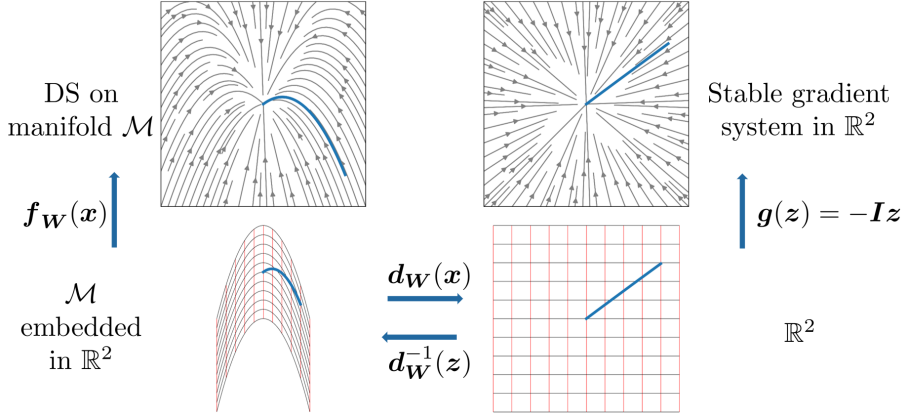


Figure 2.1: Sketch of generating a desired trajectory by a learned dynamic system (DS) on a manifold \mathcal{M} (*upper left*), by mapping a trajectory, generated by a stable, *latent* gradient system in \mathbb{R}^d (*upper right*) onto \mathcal{M} via a learned diffeomorphism d_W^{-1} (*lower right to lower left*).

2.1. Note that a straight line in \mathbb{R}^d can be interpreted as a trajectory generated by a (stable) gradient system (2.1).

Now, this approach is extended in order to be able to learn a dynamic system that generates trajectories, whose paths are geodesics on \mathcal{M} that match the example trajectories and traces them in similar time. In other words, the velocities along the given and generated trajectories should match. To achieve this, one can further note, that in a gradient system $\dot{z} = -\nabla_z \Phi(z)$ the velocity in a point z equals the steepest descent of the potential Φ in z . To mimic this behavior on \mathcal{M} , the goal is to find a similar gradient system $\dot{x} = f_W(x)$ on \mathcal{M} , which modifies the steepest gradient on the potential Φ of a mapped point $d_W(x)$ in \mathbb{R}^d so that it matches some given velocity \dot{x} , which then is the steepest gradient on \mathcal{M} , respective the potential $\Phi \circ d_W$. Such a system is called *Natural Gradient System* and can be formally written as:

$$\dot{x} = f_W(x) = -G^{-1}(x) \nabla_x \Phi(d_W(x)) = -\nabla_G \Phi(d_W(x)) \quad (2.10)$$

One can show [Ama98] that the negative *Natural Gradient* $-\nabla_G \Phi(d_W(x))$ is the steepest descent of $\Phi \circ d_W$ on \mathcal{M} .

With the potential $\Phi = \|z - z_{goal}\|_A^2$, weighted with a positive negative A , (2.10) can be written as

$$f_W(x) = -J_{d_W} A d_W(x) \stackrel{A=I}{=} J_{d_W} I d_W(x), \quad (2.11)$$

due to the potential choice (2.2) in this work.

2.1.3 Stability

An intuition on why the *Natural Gradient System* is stable if the *latent* gradient system is stable can be gained as follows. Since the *latent* system is stable, all

paths of trajectories in *latent* space are straight lines starting in a mapped state $\mathbf{z}_0 = \mathbf{d}_{\mathbf{W}}(\mathbf{x}_0)$ and ending in an equilibrium $\mathbf{z}^* = \mathbf{z}_{goal}$. Since \mathbf{d} is bijective, every mapped-back path on \mathcal{M} ends in $\mathbf{x}^* = \mathbf{d}_{\mathbf{W}}^{-1}(\mathbf{z}^*)$. Furthermore, since every trajectory in *latent* space is generated by moving along the steepest gradient and converges towards \mathbf{z}^* , every trajectory on M , generated by moving along the steepest gradient, will converge towards \mathbf{x}^* . This feature is proven in [RLF⁺20] and guarantees stability of the learned system by structure.

2.1.4 Flow-based Construction of Diffeomorphisms

To guarantee diffeomorphic properties and a straight forward computation of the inverse mapping $\mathbf{d}_{\mathbf{W}}^{-1}$, the diffeomorphism is constructed in a flow-based fashion. Given the ordinary differential equation (ODE)

$$\dot{\gamma}(\mathbf{s}, t) = \mathbf{V}_{\mathbf{W}}(\gamma(\mathbf{s}, t), t) \quad (2.12)$$

a diffeomorphism can be defined as

$$\mathbf{d}_{\mathbf{W}}(\mathbf{x}) = \gamma(\mathbf{x}, T) = \mathbf{x} + \int_0^T \mathbf{V}_{\mathbf{W}}(\gamma(\mathbf{x}, t)) dt = \mathbf{z} \quad (2.13)$$

if $\mathbf{V}_{\mathbf{W}}$ is a smoothly differentiable vector field from some *separable Hilberspace* \mathcal{H} [DGM98]. In other words, the diffeomorphism can be defined as a smooth flow generated by integrating the vector field $\mathbf{V}_{\mathbf{W}}$ for time T . As in 1.1.1, due to smooth differentiability, every partial derivative of $\mathbf{V}_{\mathbf{W}}$ is bounded on every bounded neighborhood. Therefore $\mathbf{V}_{\mathbf{W}}$ is Lipschitz continuous. The existence and uniqueness of a solution to the ODE and therefore the diffeomorphism follows by applying Picard-Lindelöf. By noting that for any initial point $\mathbf{s}_0 \in \mathcal{M}$ and time t an inverse of $\mathbf{s}_t = \mathbf{d}_{\mathbf{W}}(\mathbf{s}_0) = \gamma(\mathbf{s}_0, t)$ can be defined as $\mathbf{s}_0 = \mathbf{d}_{\mathbf{W}}^{-1}(\mathbf{s}_t) = \gamma(\mathbf{s}_t, -t)$. That leads to the following definition of the inverse diffeomorphism:

$$\mathbf{d}_{\mathbf{W}}^{-1}(\mathbf{z}) = \gamma(\mathbf{z}, -T) = \mathbf{z} + \int_{-T}^0 \mathbf{V}_{\mathbf{W}}(\gamma(\mathbf{x}), t) dt = \mathbf{x} \quad (2.14)$$

2.2 Infinitesimal Generators

In the following, three approaches are shown to define $\mathbf{V}_{\mathbf{W}}$, utilized to implicitly learn diffeomorphisms $\mathbf{d}_{\mathbf{W}}$.

2.2.1 Kernel Methods

To learn a scalar nonlinear function, in this case $y = v_{\mathbf{w}}(\mathbf{s})$, *kernel methods* [HSS08] can be employed to implicitly use a very high dimensional feature space \mathcal{H} , in which the targets y may be a linear function of the mapped input $\varphi(\mathbf{s})$, $\varphi : \mathcal{M} \rightarrow \mathcal{H}$. That

allows for utilizing simple methods such as *linear regression* in the *lifted* domain \mathcal{H} . Considering a *linear regression* approach of the form

$$\hat{y} = v_{\hat{\mathbf{w}}}(\mathbf{s}) = \langle \hat{\mathbf{w}}, \mathbf{s} \rangle . \quad (2.15)$$

Given a data matrix $\mathbf{S} \in \mathbb{R}^{N \times d}$ with full column rank d (linearly independent data points $\{\mathbf{s}\}_{i=1}^N \in \mathcal{M}$ as rows), corresponding to a target vector \mathbf{y} , $\hat{\mathbf{w}}$ can be obtained by minimizing the mean square error \hat{R} with respect to \mathbf{w}

$$\hat{R}(\mathbf{w}) = \frac{1}{N} \|\mathbf{S}^T \mathbf{w} - \mathbf{y}\|_2^2 \quad (2.16)$$

as follows.

$$\min_{\mathbf{w}}(\hat{R}(\mathbf{w})) \implies \hat{\mathbf{w}} = (\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T \mathbf{y} . \quad (2.17)$$

Which allows to make a prediction \hat{y} by

$$\begin{aligned} \hat{y} &= \hat{\mathbf{w}}^T \mathbf{s} \\ &= [\mathbf{S}^T \mathbf{S}]^{-1} \mathbf{S}^T \mathbf{y}]^T \mathbf{s} \\ &= \mathbf{s}^T \mathbf{S}^T (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{y} \end{aligned} \quad (2.18)$$

Applying the feature map φ and noting that the prediction \hat{y} only depends on pairs of mapped input points, the *Representer Theorem* allows for an implicit formulation of the form

$$\mathbf{v}_{\mathbf{w}}(\mathbf{s}) = \langle \mathbf{w}, \varphi(\mathbf{s}) \rangle_{\mathcal{H}} = \sum_{i=1}^N \alpha_i \mathbf{k}(\mathbf{s}, \mathbf{s}_i) \quad (2.19)$$

with $\mathbf{k} : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$, $\mathbf{k}(\mathbf{s}, \mathbf{s}') = \langle \varphi(\mathbf{s}), \varphi(\mathbf{s}') \rangle_{\mathcal{H}}$ being the positive definite *kernel* corresponding to the feature map φ . Note, that if φ is infinite dimensional, \mathbf{w} would have to be infinite dimensional as well, which is infeasible. Utilizing the *Representer Theorem* reduces the number of learnable parameters to N , depending only on the number of input data-points. This can be shown for the sketched linear regression approach by rewriting (2.18) as

$$\hat{y} = \boldsymbol{\kappa}(\mathbf{s}) \mathbf{K}(\mathbf{S}, \mathbf{S})^{-1} \mathbf{y} \quad (2.20)$$

with $\boldsymbol{\kappa}_i = \mathbf{k}(\mathbf{s}, \mathbf{s}_i)$ and the kernel matrix $\mathbf{K} : K_{i,j} = \mathbf{k}(\mathbf{s}_i, \mathbf{s}_j)$. Here, the expression $\mathbf{K}(\mathbf{S}, \mathbf{S})^{-1} \mathbf{y} \in \mathbb{R}^N$ represents the learned implicit N weights α_i from (2.19).

In 2 dimensions $\mathbf{y} = \mathbf{V}_{\mathbf{w}}$ can be defined element-wise as

$$\mathbf{V}_{\mathbf{w}}(\mathbf{s}) = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} \langle \mathbf{w}_1, \varphi(\mathbf{s}) \rangle_{\mathcal{H}} \\ \langle \mathbf{w}_2, \varphi(\mathbf{s}) \rangle_{\mathcal{H}} \end{bmatrix} = \langle \mathbf{W}, \varphi(\mathbf{s}) \rangle_{\mathcal{H}} . \quad (2.21)$$

In the following approaches the Gaussian radial-basis function (RBF) kernel

$$\mathbf{k}(\mathbf{s}, \mathbf{s}') = \exp(-l \|\mathbf{s} - \mathbf{s}'\|_2^2) \quad (2.22)$$

is used due to its smoothness properties [JHS⁺15]. The tunable parameter $l > 0$ can be interpreted as the *kernel width*.

To avoid having to compute the inverse of the kernel matrix \mathbf{K} as in (2.20), which has a computational complexity of $\mathcal{O}(N^3)$, two kernel-approximation approaches are presented.

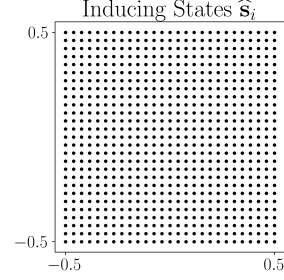


Figure 2.2: Equally spaced grid of $m = 729$ inducing states on $[-0.5, 0.5] \times [-0.5, 0.5]$.

Induced Gaussian Kernel (IGK)

A first, rather primitive approach to approximate the Gaussian kernel is to define the vector field as an inner product between a weight matrix $\mathbf{W} \in \mathbb{R}^{m \times d}$ and a *kernel vector*

$$\mathbf{V}_{\mathbf{W}}(\mathbf{s}) = \mathbf{W}^T \tilde{\mathbf{k}}(\mathbf{s}, \hat{\mathbf{s}}), \quad \tilde{\mathbf{k}}(\mathbf{s}, \hat{\mathbf{s}}) = \begin{bmatrix} \exp(l \|\mathbf{s} - \hat{\mathbf{s}}_1\|_2^2) \\ \exp(l \|\mathbf{s} - \hat{\mathbf{s}}_2\|_2^2) \\ \dots \\ \exp(l \|\mathbf{s} - \hat{\mathbf{s}}_m\|_2^2) \end{bmatrix}. \quad (2.23)$$

The states $\hat{\mathbf{s}}$, called *inducing states*, are predefined states on an equally spaced grid, covering the range of demonstrations. An example is shown in Figure 2.2.

Random Fourier Features (RFF)

To circumvent specifying *inducing states* as in the previous section, one can use *random Fourier features* (RFF) to approximate the Gaussian kernel as shown in [RR07]

$$\mathbf{V}_{\mathbf{W}}(\mathbf{s}) = \mathbf{W}^T \boldsymbol{\tau}(\mathbf{s}), \quad \boldsymbol{\tau}(\mathbf{s}) = \sqrt{\frac{2}{m}} \begin{bmatrix} \cos(\boldsymbol{\alpha}_1^T \mathbf{s} + \beta_1) \\ \cos(\boldsymbol{\alpha}_2^T \mathbf{s} + \beta_2) \\ \dots \\ \cos(\boldsymbol{\alpha}_m^T \mathbf{s} + \beta_m) \end{bmatrix}, \quad \begin{matrix} \boldsymbol{\alpha}_i \sim \mathcal{N}(\mathbf{0}, l\mathbf{I}) \\ \beta_i \sim \mathcal{U}(0, 2\pi) \end{matrix}. \quad (2.24)$$

The approximation $\langle \boldsymbol{\varphi}(\mathbf{s}), \boldsymbol{\varphi}(\mathbf{s}') \rangle_{\mathcal{H}} \approx \boldsymbol{\tau}(\mathbf{s})^T \boldsymbol{\tau}(\mathbf{s}')$, $\boldsymbol{\tau}$ being an m dimensional vector as described in (2.24), avoids computing the inverse kernel matrix by instead computing an inverse RFF matrix $(\boldsymbol{\tau}(\mathbf{s})^T \boldsymbol{\tau}(\mathbf{s}'))^{-1}$, which has a computational complexity of $\mathcal{O}(m^3)$. With $N \gg m$ a significant increase in performance can be achieved.

2.2.2 Neural Ordinary Differential Equations (NODE)

Neural Ordinary Differential Equations [CRBD18] can be seen as the state-of-the-art approach to ODE-based diffeomorphic learning. Utilizing NODE, $\mathbf{V}_{\mathbf{W}}$ is defined as a Neural Network with smooth activation functions, to guarantee that the vector field is smooth as desired.

$$\mathbf{V}_{\mathbf{W}}(\mathbf{s}) = \mathbf{N} \mathbf{N}_{\mathbf{W}}^{ELU}(\mathbf{s}) \quad (2.25)$$

2.3 Implementation

For implementing the three methods, the Python [PyTorch](#) library is used. The code is CUDA optimized to run on a graphics card for accelerating matrix operations. ADAM [KB14] is chosen for optimzing and gradient computations are performed using the PyTorch [autograd](#) package. ODE solvers are chosen from the [torchdiffeq](#) library.

Chapter 3

Evaluation and Conclusion

The presented algorithms are evaluated on the [LASA handwriting dataset](#), consisting of 30 human handwriting motions. For each motion/shape, 7 trajectories of 1000 data points are given.

Firstly, determined heuristically, the five most complex motions are chosen in order to test and compare the three methods. The selection can be seen in Figure ?? and focuses on either motions with heavily curved trajectories, such as the second, third and last shape, motions represented by frequently intersecting trajectories (fourth shape), or shapes with closely distanced starting and goal points.

Secondly, the data is preprocessed to remove noise and to smooth the trajectories as follows. The trajectories are downsampled and filtered. Further, the data points at the beginning and end of each trajectory are dropped to avoid random deflections when placing the pen on -or lifting the pen from the drawing surface. Lastly the velocities along the trajectories are recomputed via infinitesimal differences to circumvent noisy velocity sensor data. After preprocessing each trajectory consists of $l_i = 235$ points, resulting in a total of $N = 1645$ training points.

3.1 Experimental Results and Performance

To account for a *fair* training set-up, every method is trained on the $N = 1645$ data points for 5000 epochs. The kernel methods are constructed with $m = 900$ inducing states/RFFs respectively, to approximate a Gaussian RBF kernel with a *kernel width* $l = 60$. The NODE model has a layer structure of $[2, 50, 50, 2]$. For generating the diffeomorphism as in 2.1.4, 4 steps of a fourth order Runge Kutta (RK4) method are used for approximating a solution to the ODE in the time-frame $t \in [0, 2s]$.

Analyzing the test results, one can see that every approach is able to learn stable nonlinear systems, that are able to reproduce severely curved trajectories. One can further notice that the learned first order dynamical systems are incapable of modeling intersections in the data. To measure the reproduction quality, the root mean square error (RMSE) between generated and example trajectories is computed.

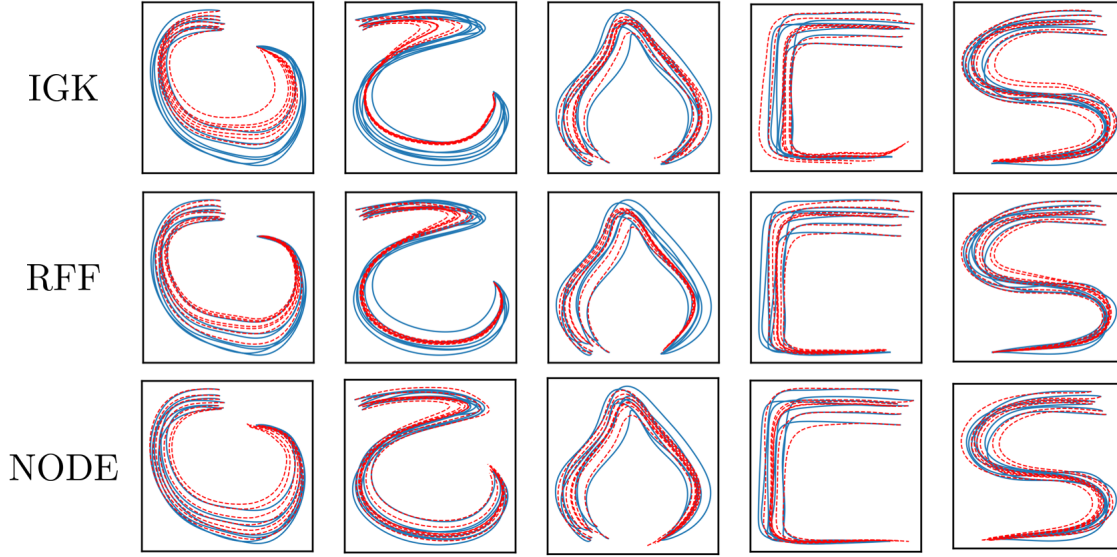


Figure 3.1: Five LASA handwriting shapes, human expert demonstrations (blue) and reproduced trajectories by learned dynamic systems (red). Induced gaussian kernel (IGK) approach is evaluated in *first row*, random Fourier feature (RFF) approach in *second row* and neural ordinary differential equation (NODE) approach in *third row*.

By examining the mean RMSE per shape in Figure 3.2, one can see that the RFF approach performs similarly well compared to the state-of-the-art NODE approach. The IGK approach has a slightly worse reproduction quality.

To compare computational performance, the respective algorithms are run on a *NVIDIA TITAN V* graphics card for 1000 epochs. The training time per epoch includes a prediction/evaluation of the learned dynamics and the computation of a gradient descent step to update the weights \mathbf{W} . In Figure 3.2 the averaged training and prediction times per epoch are shown respective the fastest approach, which is the RFF method.

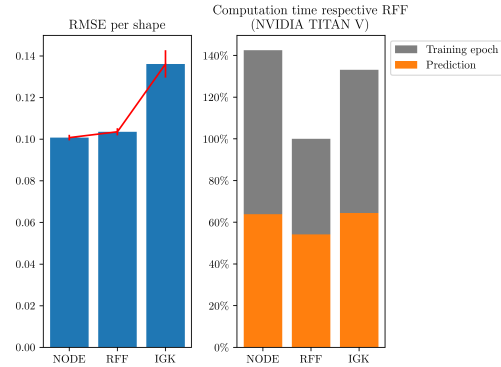


Figure 3.2: Mean RMSE per shape (*left*) and computation times respective RFF approach (*right*).

3.2 Conclusion and Future Work

Overall, regarding reproduction quality, the RFF and the state-of-the-art NODE approach perform similarly well and with some further tuning of the *kernel width* l a slight performance increase of the RFF approach is most likely possible. Respective computation time, a large margin is caused by the computation of the batch-wise inverse Jacobian $\mathbf{J}_{d\mathbf{w}}^{-1}$, which is furthermore computed multiple times in each ODE-

solving step. Therefore, when reducing the number of steps, a major decrease in computation time could be achieved. In 3.1, using a simple explicit Euler step approximating the ODE solution, the number of computations of \mathbf{J}_{dw}^{-1} can be reduced from 16 (4 RK4 steps) to 1. With some further testing one may find a better trade-off between approximation quality of the ODE (in terms of integration methods and number of steps) and computation time. Furthermore the computation time of the RFF approach may be improved without losing reproduction quality by reducing the number of features (probably also by a large margin).

Structurally, the methods may be improved by further learning a positive definite \mathbf{A} matrix (2.11), corresponding to the *latent* system $\mathbf{g}(\mathbf{z}) = -\mathbf{A}\mathbf{z}$, to reduce the *deformation* the diffeomorphism has to cope for. To be able to learn systems, that can model intersections in the data, second order dynamical systems may be considered.

Concluding, the RFF and NODE approach both perform well, with the RFF approach offering a lot of fine-tuning options and future work, which might make it the method of choice for learning stable dynamic systems, especially if fast evaluations of the dynamics are desired.

Bibliography

- [Ama98] Shun-ichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276, 02 1998. URL: <https://doi.org/10.1162/089976698300017746>, arXiv:<https://direct.mit.edu/neco/article-pdf/10/2/251/813415/089976698300017746.pdf>, doi:10.1162/089976698300017746.
- [CRBD18] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. *CoRR*, abs/1806.07366, 2018. URL: <http://arxiv.org/abs/1806.07366>, arXiv:1806.07366.
- [DGM98] PAUL DUPUIS, ULF GRENANDER, and MICHAEL I. MILLER. Variational problems on flows of diffeomorphisms for image matching. *Quarterly of Applied Mathematics*, 56(3):587–600, 1998. URL: <http://www.jstor.org/stable/43638248>.
- [HSS08] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine learning. *The annals of statistics*, 36(3):1171–1220, 2008.
- [JHS⁺15] Sadeep Jayasumana, Richard Hartley, Mathieu Salzmann, Hongdong Li, and Mehrtaash Harandi. Kernel methods on riemannian manifolds with gaussian rbf kernels. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(12):2464–2477, 2015. doi:10.1109/TPAMI.2015.2414422.
- [KB14] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2014. URL: <https://arxiv.org/abs/1412.6980>, doi:10.48550/ARXIV.1412.6980.
- [Lee06] John M Lee. *Riemannian manifolds: an introduction to curvature*, volume 176. Springer Science & Business Media, 2006.
- [PSC16] Nicolas Perrin and Philipp Schlehüser-Caissier. Fast diffeomorphic matching to learn globally asymptotically stable nonlinear dynamical systems. *Systems Control Letters*, 96:51–59, 2016. URL: <https://www.sciencedirect.com/science/article/pii/S0167691116300846>, doi:<https://doi.org/10.1016/j.sysconle.2016.06.018>.

-
- [RLF⁺20] Muhammad Asif Rana, Anqi Li, Dieter Fox, Byron Boots, Fabio Ramos, and Nathan D. Ratliff. Euclideanizing flows: Diffeomorphic reduction for learning stable dynamical systems. *CoRR*, abs/2005.13143, 2020. URL: <https://arxiv.org/abs/2005.13143>, arXiv:2005.13143.
- [RR07] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL: <https://proceedings.neurips.cc/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf>.
- [UGTP20] Julen Uraïn, Michele Ginesi, Davide Tateo, and Jan Peters. Imitation-flow: Learning deep stable stochastic dynamic systems by normalizing flows. *CoRR*, abs/2010.13129, 2020. URL: <https://arxiv.org/abs/2010.13129>, arXiv:2010.13129.

License

This work is licensed under the Creative Commons Attribution 3.0 Germany License. To view a copy of this license, visit <http://creativecommons.org> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California 94105, USA.