MATHEMATISCHES FORSCHUNGSINSTITUT OBERWOLFACH

# Constrained Dynamics, Stochastic Numerical Methods and the Modeling of Complex Systems

Organized by
Benedict Leimkuhler, Edinburgh
Richard Tsai, Austin
Gilles Vilmart, Geneva
Rachel Ward, Austin

26 May – 31 May 2024

ABSTRACT. The workshop aimed to unite researchers from diverse fields of mathematics and statistics to explore the foundations of high-dimensional modeling and computational studies. It addressed recent advancements in numerical analysis, dynamical systems, and stochastic differential equations that support model reduction for large-scale complex systems.

Incorporating targeted geometric structures, such as Riemannian manifolds, into large-scale statistical models is known to enhance the stability, reliability, and efficiency of numerical methods. However, algorithms are often presented in application contexts without adequate attention to their fundamental properties, limiting the adoption of these advanced modeling methods.

The workshop emphasized understanding the fundamental properties of these structures, their impact on dynamics and stochastic dynamics, and the need to redesign algorithms to capture essential properties, aiming for robustness and suitability for high-performance computation.

By bringing together numerical analysts, statisticians, and modelers, the workshop sought to improve the quality of methods and identify new model frameworks to guide future development.

# Introduction by the Organizers

The aim of this workshop was to convene researchers from various disciplines within mathematics and statistics to discuss recent advancements and innovative ideas in the modeling of complex systems. These discussions focused on the integration of model reduction, machine learning, statistical inference, and numerical analysis.

In recent years, these fields have been creatively combined to tackle complex problems, such as characterizing large datasets through machine learning. As mathematical modelers evolve their frameworks to better incorporate real-world observations and improve uncertainty quantification, statistical algorithms have become crucial components within larger schemes. Ensuring the convergence and stability of these methods is essential, particularly in addressing the scale couplings and stochastic perturbations inherent in large-scale systems.

Our workshop provided a platform to explore the incorporation of geometric structures, such as Riemannian manifolds, into statistical models to enhance their stability, reliability, and efficiency. We emphasized understanding the fundamental properties of these structures, their impact on dynamics and stochastic processes, and the redesign of algorithms to capture essential properties for robust and efficient computation.

Bringing together numerical analysts, statisticians, and modelers, the workshop aimed to improve existing methods and identify new frameworks to guide future developments. Topics discussed included the role of Bayesian paradigms in data approximation, the influence of numerical errors on sampling processes, and the performance of state-of-the-art algorithms.

The compilation of extended abstracts in this report captures the diverse insights and visions shared during the workshop, offering a roadmap for future research in this critical and evolving field.

## Workshop: Constrained Dynamics, Stochastic Numerical Methods and the Modeling of Complex Systems

## Table of Contents

# Abstracts

## Stochastic partial differential equations on surfaces and evolving random surfaces: a computational approach

Annika Lang

(joint work with David Cohen, Erik Jansson, Mike Pereira and Christoph Schwab)

Looking at ice crystals in clouds or moving cells as shown in data in [7] and [4], we observe that the shapes of different individuals are similar but not equal. What are possible ways to model this difference and to efficiently generate many such similar shapes that are possibly changing over time?

The main tool of the talk is the use of stochastic models and more specifically random fields and solutions to stochastic partial differential equations on spheres and other surfaces. These are used to transform a sphere to a random surface. With this approach, shapes are similar in the sense of being samples of the same distribution that are described by the mean and the covariance in the considered Gaussian examples.

An important question when modeling with uncertainty is how to add the randomness such that surface structures are not destroyed. We start by taking the spherical harmonic functions $(Y_{\ell,m}, \ell \in \mathbb{N}, m = -\ell, \ldots, \ell)$ on the unit sphere $\mathbb{S}^2$ as basis of $L^2(\mathbb{S}^2)$ and looking at the basis expansion

$$(1) \qquad U = \sum_{\ell=0}^{\infty} \sum_{m=-\ell}^{\ell} a_{\ell,m} Y_{\ell,m},$$

where we perturb each mode with a standard Gaussian random number.

In a simulation of the truncated series at some mode $L$, we see that the samples get more and more spiky the more modes we include as shown in Figure 1. For



(a) $L = 1$    (b) $L = 2$    (c) $L = 4$    (d) $L = 8$    (e) $L = 16$    (f) $L = 32$    (g) $L = 64$

FIGURE 1. Random field generated by the series expansion (1) with different truncation levels $L$.

a better intuition for the spikes and therefore roughness of the field, we introduce next a transformation of the random field to a random surface as was used for the modeling of ice crystals in [7]. We compute first the lognormal random field

$\exp(U)$ of (1) and shift every point $y$ on the sphere in normal direction to the point $x$ given by

(2) $$x = \exp(U(y))\, y, \qquad y \in \mathbb{S}^2.$$

This transformation generates out of the random samples in Figure 1 the random surfaces shown in Figure 2.



(a) $L = 1$     (b) $L = 2$     (c) $L = 4$     (d) $L = 8$     (e) $L = 16$     (f) $L = 32$     (g) $L = 64$

FIGURE 2. Random surfaces generated by the transformation (2) with different truncation levels $L$.

A natural question is if the obtained surfaces make sense when taking $L$ to infinity. The generation of a surface without cracks and holes requires the continuity of the random field $U$. To approach this question, we first extend the class of random fields by adding weights corresponding to the inverse of Bessel potentials

(3) $$U = \sum_{\ell=0}^{\infty} \sum_{m=-\ell}^{\ell} a_{\ell,m} (1 + \ell(\ell+1))^{-\alpha/2}\, Y_{\ell,m}$$

with a smoothness parameter $\alpha$. The variant of the Kolmogorov–Chentsov theorem developed in [6] tells us that the random field $U$ has a continues modification for $\alpha > 1$ and that it is differentiable for $\alpha > 2$. This means that looking at the sample in Figure 2(g) with different parameters $\alpha$ in Figure 3, we see the smoothing and that for $\alpha = 1/2$ the continues surface just exists due to the truncation of the series expansion (1).



(a) $\alpha = 1/2$     (b) $\alpha = 1$     (c) $\alpha = 3/2$     (d) $\alpha = 2$     (e) $\alpha = 5/2$     (f) $\alpha = 3$     (g) $\alpha = 7/2$

FIGURE 3. Random surfaces generated by the transformation (2) with different smoothness index $\alpha$ in (3).

Next, we simulate isotropic Wiener processes $W$ on the sphere by adding up the Gaussian random fields (3) scaled by the square root of the time step size similarly to the simulation of sample paths of one-dimensional Brownian motions, i.e.,

$$W(t_{n+1}) - W(t_n) \sim \sqrt{t_{n+1} - t_n}\, U.$$

With this construction at hand, we are able to solve the linear stochastic wave equation on the sphere

$$\partial_{tt} u(t) - \Delta_{\mathbb{S}^2} u(t) = \dot{W}(t)$$

based on an expansion with respect to the spherical harmonic functions and its truncation. We discuss strong and weak convergence that depend on the smoothness of $W$ and of the initial condition as shown in [1].

In the last part of the presentation, we look at alternatives to generate Gaussian random fields on surfaces when the eigenexpansion of the covariance operator is unknown. In [3], we extend our earlier results from [2, 5] and consider random fields of the form

$$u = \gamma(\mathcal{L})\mathcal{W}$$

on a hypersurface $\mathcal{M}$ of dimension $d = 1, 2$ with the differential operator $\mathcal{L} = -\nabla_{\mathcal{M}} D \nabla_{\mathcal{M}} + V$ and a power spectral density $\gamma$ that decays asymptotically with rate $\alpha$. Here, $\mathcal{W}$ denotes white noise on $\mathcal{M}$. We approximate the solution using a surface finite element method with linear elements and a Chebyshev approximation as in [5]. We give strong convergence results of essentially order $\min(\alpha - d/4, 2)$ in the mesh width of the surface approximation and exponentially in the degree of the Chebyshev polynomials.

## References

[1] D. Cohen and A. Lang, *Numerical approximation and simulation of the stochastic wave equation on the sphere*, Calcolo **59**(3):32, 2022.

[2] E. Jansson, M. Kovács, and A. Lang, *Surface finite element approximation of spherical Whittle–Matérn Gaussian random fields*, SIAM J. on Sci. Comp. **44**(2):A825–A842, 2022.

[3] E. Jansson, A. Lang, and M. Pereira, *Non-stationary Gaussian random fields on hypersurfaces: Sampling and strong error analysis*, arXiv:2406.08185 [math.NA], 2024.

[4] Johns Hopkins Medicin, *Ever wonder how cells move?* YouTube, https://www.youtube.com/watch?v=HPdl-tBYXHY, 2017.

[5] A. Lang and M. Pereira, *Galerkin–Chebyshev approximation of Gaussian random fields on compact Riemannian manifolds*, BIT Num. Math. **63**(4):51, 2023.

[6] A. Lang and Ch. Schwab, *Isotropic Gaussian random fields on the sphere: regularity, fast simulation and stochastic partial differential equations*, Ann. Appl. Probab. **25**(6):3047–3094, 2015.

[7] T. Nousiainen and G. M. McFarquhar, *Light scattering by quasi-spherical ice crystals*, J. Atmos. Sci. **61**(18):2229–2248, 2004.

# Convergence of kinetic Langevin samplers for non-convex potentials

Katharina Schuh

(joint work with Peter A. Whalley)

We study three kinetic Langevin samplers to sample a given target measure

$$\pi(\mathrm{d}x) \propto e^{-U(x)}\mathrm{d}x \quad \text{in } \mathbb{R}^d.$$

The samplers include the Euler discretization, the BU and the UBU splitting scheme and form approximations of the kinetic Langevin dynamics given by

$$\begin{cases} \mathrm{d}X_t = V_t\mathrm{d}t \\ \mathrm{d}V_t = -\nabla U(X_t)\mathrm{d}t - \gamma V_t\mathrm{d}t + \sqrt{2\gamma}\mathrm{d}B_t, \end{cases}$$

where $\gamma > 0$ is the friction parameter and $(B_t)_{t \geqslant 0}$ is a $d$-dimensional standard Brownian motion. The unique invariant measure of the continuous dynamics is the Boltzmann-Gibbs measure $\tilde{\pi}(\mathrm{d}x\mathrm{d}v) \propto e^{-U(x)-1/2|v|^2}\mathrm{d}x\mathrm{d}v$ where the marginal measure in the position component is the target measure.

We are interested in how well the numerical schemes sample $\pi$ if $U$ is a non-convex potential which includes double-well potentials. In particular, we assume that the potential is only strongly convex outside a Euclidean ball and has a Lipschitz continuous gradient.

We show contraction results in $L^1$-Wasserstein distance for these schemes. The results are based on a carefully tailored distance function taken from the continuous framework [1] and an appropriate coupling construction. If the two copies of the coupling are far apart a synchronous coupling is considered. Then due to the friction term and the convexity local contraction is obtained. If the two copies are close to each other a coupling containing partially a reflection coupling is constructed. Together with the concavity of the distance function this results in local contraction on average. By combining the two local contraction results we obtain global contraction in $L^1$-Wasserstein distance.

Additionally, we analyse the error in the $L^1$-Wasserstein distance between the target measure and the invariant measure of the discretization schemes. To get an $\varepsilon$-accuracy in $L^1$-Wasserstein distance, we show complexity guarantees of order $\mathcal{O}(\sqrt{d}/\varepsilon)$ for the Euler scheme and $\mathcal{O}(d^{1/4}/\sqrt{\varepsilon})$ for the UBU scheme under appropriate regularity assumptions on the target measure. Here, a global error result is obtained by first showing local error bounds and combining them with the contraction result.

The results can also be applied to interacting particle systems and provide bounds for sampling probability measures of mean-field type.

## References

[1] K. Schuh, *Global contractivity for Langevin dynamics with distribution-dependent forces and uniform in time propagation of chaos*, Ann. Inst. H. Poincaré Probab. Statist. 60(2): 753-789 (May 2024). DOI: 10.1214/22-AIHP1337.

[2] K. Schuh and P. A. Whalley, *Convergence of kinetic Langevin samplers for non-convex potentials*, arXiv preprint arXiv:2405.09992 (2024).

## Stable generative modelling using diffusion maps

SEBASTIAN REICH

(joint work with Georg Gottwald, Fengyi Li and Youssef Marzouk)

Generative modelling is the process of learning a mechanism for synthesizing new samples that resemble those of the original data-generating distribution, given only a finite set of samples. In my blackboard talk, I sketched out a new nonparametric approach to generative modelling that combines ideas from diffusion maps and Schrödinger bridges [2] with discretised reversible Langevin dynamics.

More specifically, suppose we are given $M$ training samples $x^{(i)} \sim \pi$, $i = 1, \ldots, M$, from an unknown distribution $\pi$ on $\mathbb{R}^d$. Then the key ideas put forward in [1] are

(i) to approximate the conditional expectation value $\mathbb{E}[X_\epsilon | X_0 = x]$, where $X_t$ satisfies the reversible Langevin process with invariant measure $\pi$ (see (1) below), by

$$\mathbb{E}[X_\epsilon | X_0 = x] \approx \sum_{i=1}^{M} x^{(i)} p_{i,\epsilon}(x),$$

where $p_\epsilon(x) \in \mathbb{R}^M$ is a state-dependent probability vector which can be obtained from a Schrödinger bridge problem of the data $\{x^{(i)}\}_{i=1}^{M}$ with itself [2, 1], and

(ii) to note that, for $\epsilon > 0$ sufficiently small and $M$ sufficiently large, we can approximate the score $s(x) = \nabla_x \log \pi(x)$ by

$$s(x) \approx \frac{\sum_{i=1}^{M} x^{(i)} p_{i,\epsilon}(x) - x}{\epsilon}.$$

These approximations suggest to consider the following split-step time-stepping scheme with step-size $\Delta t = \epsilon$:

$$\text{noising:} \quad X_{n+1/2} = X_n + \sqrt{2\epsilon}\Xi_n,$$

$$\text{denoising:} \quad X_{n+1} = \sum_{i=1}^{M} x^{(i)} p_{i,\epsilon}(X_{n+1/2})$$

for the reversible Langevin process

$$\text{(1)} \qquad \mathrm{d}X_t = \nabla \log \pi(X_t)\mathrm{d}t + \sqrt{2}\mathrm{d}B_t$$

with invariant measure $\pi$, where $\Xi_n \sim \mathrm{N}(0, I)$ and $B_t$ denotes standard multi-dimensional Brownian motion.

Provided the support of $\pi$ is compact, it can be shown that the proposed time-stepping scheme is stable and geometric ergodic. The accuracy of the scheme will depend on the chosen time-step, $\epsilon$, and the number, $M$, of data points. See the theoretical analysis in [2]. It is also clear by construction that the generated samples $X_n$ will be contained in the convex hull of the samples $\{x^{(i)}\}_{i=1}^{M}$.

Numerical results can be found in [1] including implementations with variable bandwidth. The proposed methodology naturally extends to Bayesian inference with negative log-likelihood $l(x, y)$, i.e.,

$$(2a) \qquad\qquad X_{n+1/2} = X_n - \epsilon \nabla l(X_n, y) + \sqrt{2\epsilon}\Xi_n,$$

$$(2b) \qquad\qquad X_{n+1} = \sum_{i=1}^{M} x^{(i)} p_{i,\epsilon}(X_{n+1/2}),$$

conditional sampling, and score generative modelling. The time-stepping (2) is particularly attractive in the context of sequential data assimilation where the prior $\pi$ at any data assimilation cycle is provided by a forecast ensemble only [3]. Future challenges include extensions to high-dimensional problems building upon the concept of localisation [3]. In this context, we note the resemblance of (2b) with the ensemble transform particle filter [3].

### References

[1] G. Gottwald, F. Li, Y. Marzouk, and S. Reich, *Stable generative modeling using diffusion maps*, arXiv:2401.04372 (2024).

[2] C.L. Wormell and S. Reich, *Spectral convergence of diffusion maps: Improved error bounds and an alternative normalization*, SIAM J. Numer. Anal. **59** (2021), 1687–1734.

[3] S. Reich and C. Cotter, *Probabilistic forecasting and Bayesian data assimilation*, Cambridge University Press, Cambridge, 2015

## Learning of neural networks with low-dimensional and multiscale structures in data

### Juncai He

#### (joint work with Lewis Liu, Richard Tsai and Rachel Ward)

The low-dimensional manifold hypothesis posits that data found in many applications, such as those involving natural images, lie (approximately) on low-dimensional manifolds embedded in a high-dimensional Euclidean space. In this setting, a typical neural network defines a function that takes a finite number of vectors in the embedding space as input. However, one often needs to consider evaluating the optimized network at points outside the training distribution. In the work of [1], we consider the case in which the training data are distributed in $\mathcal{M}$, a linear subspace of $\mathbb{R}^d$. We derive estimates on the variation of the learning function $f_\theta : \mathbb{R}^d \mapsto \mathbb{R}$, defined by a parameterized family of both linear and ReLU neural networks, in the direction transversal to the subspace. We study the potential regularization effects associated with the network's depth and noise in the codimension of the data manifold. Furthermore, we demonstrate the multiscale structure in the training dynamics of learning $f_\theta$ when the noise has a small positive variance in the orthogonal complement of $\mathcal{M}$ for deep linear neural networks based on the result in [2].

In the next work of [3], we present that the data manifold's extrinsic geometry can lead to a multiscale structure in linear regressions. Specifically, we analyze

the impact of the manifold's curvatures (or higher-order nonlinearity in the parameterization when the curvatures are locally zero) on the uniqueness of the regression solution. Our findings suggest that the corresponding linear regression does not have a unique solution when the embedded submanifold is flat in some dimensions. Otherwise, the manifold's curvature (or higher-order nonlinearity in the embedding) may contribute significantly, particularly in the solution associated with the normal directions of the manifold. Our findings thus reveal the data manifold's geometry in ensuring the stability of regression models for out-of-distribution inferences. All these results are essentially established based on the multiscale structure of loss functions in terms of the manifold's curvatures.

In our most recent work [4], we investigate the impact of multiscale structure in data on machine learning algorithms, particularly in the context of deep learning. A dataset is multiscale if its distribution shows large variations in scale across different directions. This work reveals multiscale structures in the loss landscape, including its gradients and Hessians inherited from the data for both linear and logistic regressions and deep neural networks. More precisely, let us assume the data have the structure

$$x_i = \left(x_i^0, \epsilon_1 x_i^1, \cdots, \epsilon_m x_i^m\right) \sim (\mathcal{O}(1), \mathcal{O}(\epsilon_1), \cdots, \mathcal{O}(\epsilon_m))$$

with $1 \gg \epsilon_1 \gg \cdots \gg \epsilon_m > 0$. Then, for both linear and logistic regressions, one can have the following multiscale expansion of the gradients of loss functions

$$\frac{\partial \mathcal{L}}{\partial \theta} = (\mathcal{O}(1), \mathcal{O}(\epsilon_1), \cdots, \mathcal{O}(\epsilon_m)).$$

However, for deep neural networks, the explicit multiscale expansion only exists for parameters in the first hidden layer. For general parameters in the $\ell$-th layer, we have

$$\frac{\partial \mathcal{L}}{\partial W^\ell} = \frac{1}{N} \sum_{k=0}^{m} \epsilon^k A_k^\ell \left(x^0, \cdots, x^k\right),$$

where $\epsilon_k = \epsilon^k$, $A_k^\ell$ is of $\mathcal{O}(1)$ and more details can be found in [4]. Given this structure, correspondingly, we introduce a novel gradient descent approach called Multirate Gradient Descent (MrGD), drawing inspiration from multiscale algorithms used in scientific computing [5]. This approach seeks to transcend empirical learning rate selection, offering a more systematic, data-informed strategy to enhance training efficiency, especially in the later stages. The key to the success of this method is to choose the number of iterations for different learning rates with different scales. Theoretically, we establish a comprehensive and rigorous theory demonstrating that the MrGD scheme achieves a quasi-optimal convergence rate for linear problems and can be extended to convex functions. Numerical examples are also provided to demonstrate the efficiency of MrGD.

References

[1] J. He, R. Tsai, and R. Ward, *Side effects of learning from low-dimensional data embedded in a Euclidean space*, Research in the Mathematical Sciences **10** (2023).
[2] S. Arora, N. Cohen, and E. Hazan, *On the optimization of deep networks: Implicit acceleration by overparameterization*, Proceedings of the 35th International Conference on Machine Learning, PMLR **80** (2018), 244–253.
[3] L. Liu, J. He, and R. Tsai, *Linear regression on manifold structured data: the impact of extrinsic geometry on solutions*, Proceedings of 2nd Annual Workshop on Topology, Algebra, and Geometry in Machine Learning (TAG-ML), PMLR **221** (2023), 557–576.
[4] J. He, L. Liu, and R. Tsai, *Data-induced multiscale losses and efficient multirate gradient descent schemes*, ArXiv:2402.03021 (2023).
[5] B. Engquist and R. Tsai, *Heterogeneous multiscale methods for stiff ordinary differential equations*, Mathematics of Computation **74** (2005), 1707–1742.

# A structure-preserving kernel method for learning Hamiltonian systems on symplectic and Poisson manifolds

Juan-Pablo Ortega

(joint work with Jianyu Hu and Daiying Yin)

Hamiltonian systems are essential tools to model physical systems [1, 6, 2]. In the simplest case in which the phase space is Euclidean and is endowed with a constant symplectic form, Hamiltonian systems are determined by a scalar-valued Hamiltonian function $H : \mathbb{R}^{2d} \longrightarrow \mathbb{R}$, $d \in \mathbb{N}$, and when using the so-called canonical Darboux coordinates, the corresponding dynamics is governed by the well-known **Hamilton's equations**

$$(1) \qquad\qquad \dot{\mathbf{z}}(t) = J\nabla H(\mathbf{z}(t)),$$

where $\mathbf{z} = (\mathbf{q}^{\top}, \mathbf{p}^{\top})^{\top} \in \mathbb{R}^{2d}$ is the phase space vector comprising the positions and the momenta of the system, and $J$ is the canonical symplectic matrix. Modern technology has made collecting trajectory data directly from physical systems increasingly feasible. This motivates us to address the fundamental inverse problem: *determining the underlying Hamiltonian function and the governing Hamilton's equations from trajectory data.*

## 1. Kernel ridge regression setup

Let $K \in C_b^3(\mathbb{R}^{2d} \times \mathbb{R}^{2d})$ be a Mercer kernel and let $\mathcal{H}_K$ be the reproducing kernel Hilbert space (RKHS) associated to it. the symbol $C_b^s(\mathbb{R}^d)$ denotes the set of bounded $s$-continuously differentiable functions with bounded derivatives. The main purpose of this work is to learn in a structure-preserving fashion the unknown Hamiltonian function $H : \mathbb{R}^{2d} \longrightarrow \mathbb{R}$ of the system (1) out of realizations of random samples containing $N$ noisy observations of the Hamiltonian vector field. More explicitly, the observed data consists of $N$ independent random samples of states in the phase space and noisy observations of the Hamiltonian vector fields

at the corresponding $N$ states. We shall write the associated random samples as:

$$\mathbf{Z}_N := \mathrm{Vec}\left(\mathbf{Z}^{(1)} | \cdots | \mathbf{Z}^{(N)}\right) \in \mathbb{R}^{2dN},$$

(2)

$$\mathbf{X}_{\sigma^2,N} := \mathrm{Vec}\left(\mathbf{X}_{\sigma^2}^{(1)} | \cdots | \mathbf{X}_{\sigma^2}^{(N)}\right) \in \mathbb{R}^{2dN},$$

where $\mathbf{Z}^{(n)} \in \mathbb{R}^{2d}$ is the phase space vector containing the position and the conjugate momenta of the system, and $\{\mathbf{Z}^{(1)}, \cdots, \mathbf{Z}^{(N)}\}$ are IID random variables with the same distribution $\mu_{\mathbf{Z}}$. The symbol 'Vec' stands for the vectorization of the corresponding matrices and $\mathbf{X}_{\sigma^2}^{(n)} \in \mathbb{R}^{2d}$ denotes a noisy vector field value at $\mathbf{Z}^{(n)}$, that is, $\mathbf{X}_{\sigma^2}^{(n)} = J\nabla H(\mathbf{Z}^{(n)}) + \varepsilon^{(n)}$, where $\varepsilon^{(n)}$ are IID $\mathbb{R}^{2d}$-valued random variables with mean zero and variance $\sigma^2$. In the sequel, if $f : \mathbb{R}^{2d} \to \mathbb{R}^s$ is a function, we then shall denote the value $\mathrm{Vec}\left(f(\mathbf{Z}^{(1)}) | \cdots | f(\mathbf{Z}^{(N)})\right) \in \mathbb{R}^{sN}$ by $f(\mathbf{Z}_N)$.

To address the above-mentioned learning problem, we propose a structure-preserving kernel ridge regression method. In contrast to traditional kernel ridge regressions, our approach guarantees that the learned vector field is indeed Hamiltonian. Structure-preservation is achieved by searching for vector fields $\mathbf{f} : \mathbb{R}^{2d} \to \mathbb{R}^{2d}$ with Hamiltonian form, that is, $\mathbf{f}_h := X_h = J\nabla h$, where $h : \mathbb{R}^{2d} \longrightarrow \mathbb{R}$ is an element of $\mathcal{H}_K$. More precisely, we will be studying the following optimization problem

(3)
$$\widehat{h}_{\lambda,N} := \underset{h \in \mathcal{H}_K}{\arg\min} \frac{1}{N} \sum_{n=1}^{N} \left\| X_h(\mathbf{Z}^{(n)}) - \mathbf{X}_{\sigma^2}^{(n)} \right\|^2 + \lambda \|h\|_{\mathcal{H}_K}^2,$$

where $X_h = J\nabla h$ and $\lambda \geqslant 0$ is a Tikhonov regularization parameter. We call the solution $\widehat{h}_{\lambda,N}$ of the optimization problem (3) the **structure-preserving kernel estimator of the Hamiltonian function**.

## 2. Some Results

The following result shows that the optimization problem (3) can be cast as the solution of a convex Gramian regression. This convexity feature is a comparative advantage with the (potentially non-convex) maximum likelihood problem introduced that we would face when using Gaussian processes or neural networks.

**Theorem 2.1 (Differential Representer Theorem for Symplectic Vector Spaces).** *For every* $\lambda > 0$, *the optimization problem* (3) *has a unique solution* $\widehat{h}_{\lambda,N} \in \mathcal{H}_K$ *that can be represented as*

(4)
$$\widehat{h}_{\lambda,N} = \sum_{i=1}^{N} \langle \widehat{\mathbf{c}}_i, \nabla_1 K(\mathbf{Z}^{(i)}, \cdot) \rangle,$$

*with* $\widehat{\mathbf{c}}_1, \ldots, \widehat{\mathbf{c}}_N \in \mathbb{R}^{2d}$, $\langle \cdot, \cdot \rangle$ *the Euclidean inner product in* $\mathbb{R}^{2d}$, *and where* $\nabla_1 K(\mathbf{z}, \cdot) \in \mathbb{R}^{2d}$ *denotes the gradient of* $K$ *with respect to the* $\mathbf{z}$ *variable. Moreover, if we denote by* $\widehat{\mathbf{c}} \in \mathbb{R}^{2dN}$ *the vectorization of* $(\widehat{\mathbf{c}}_1 | \cdots | \widehat{\mathbf{c}}_N)$, *then we have*

$$\widehat{\mathbf{c}} = (\nabla_{1,2} K(\mathbf{Z}_N, \mathbf{Z}_N) + \lambda NI)^{-1} \mathbb{J}^\top \mathbf{X}_{\sigma^2, N}.$$

*The matrix* $\nabla_{1,2}K(\mathbf{Z}_N, \mathbf{Z}_N)$ *(the symbol* $\nabla_{1,2}$ *denotes partial derivatives with respect to all the entries in* $K$*) is the differential Gram matrix that can be shown to be positive semidefinite.*

PAC bounds and convergence upper rates for the total reconstruction error can be formulated using additional conditions. The first one is a customary restriction on the target Hamiltonian function introduced in [3] under the name of **source condition**, an alternative way to handle the approximation error using universal kernels. Let $\gamma \in (0,1)$, $S > 0$, and $B = A^*A$ with $A : \mathcal{H}_K \longrightarrow \mathcal{H}_K^{2d}$ given by $Ah = J\nabla h$ with $h \in \mathcal{H}_K$. The source condition consists of assuming that

$$(5) \qquad H \in \Omega_S^\gamma := \{h \in \mathcal{H}_K \mid h = B^\gamma \psi, \psi \in \mathcal{H}_K, \|\psi\|_{\mathcal{H}_K} < S\}.$$

**Theorem 2.2** (**PAC bounds of the total reconstruction error**). *Let* $\widehat{h}_{\lambda,N}$ *be the unique minimizer of the optimization problem (3). Suppose that* $H \in \Omega_S^\gamma$ *as defined in (5). Then, for any* $\varepsilon, \delta > 0$*, there exist* $\lambda > 0$ *and* $n \in \mathbb{N}_+$ *such that for all* $N > n$*, it holds that*

$$\mathbb{P}\left(\left\|\widehat{h}_{\lambda,N} - H\right\|_{\mathcal{H}_K} > \varepsilon\right) < \delta.$$

In order to get a convergence upper rate of $\|\widehat{h}_{\lambda,N} - H\|_{\mathcal{H}_K}$ as $N \to \infty$, we shall work not with a fixed, but with a dynamical $\lambda$ that is adapted with respect to the sample size $N$. More specifically, we shall assume that $\lambda \propto N^{-\alpha}$, $\alpha > 0$.

**Theorem 2.3** (**Convergence upper rate of the total reconstruction error**). *Let* $\widehat{h}_{\lambda,N}$ *be the unique minimizer of the optimization problem (3). Suppose that* $H$ *satisfies the source condition (5), that is,* $H \in \Omega_S^\gamma$*. Then for all* $\alpha \in (0, \frac{1}{3})$*, and for any* $0 < \delta < 1$*, with probability as least* $1 - \delta$*, it holds that*

$$\left\|\widehat{h}_{\lambda,N} - H\right\|_{\mathcal{H}_K} \leqslant C(\gamma, \delta, \kappa) \ N^{-\min\{\alpha\gamma, \frac{1}{2}(1-3\alpha)\}},$$

*where*

$$C(\gamma, \delta, \kappa) = \max\left\{\|B^{-\gamma}H\|_{\mathcal{H}_K}, 8\sqrt{4\log(8/\delta)}d^{\frac{3}{2}}\kappa^3\|H\|_{\mathcal{H}_K}\right\}.$$

## 3. A Numerical Illustration

It is generally a challenging task to learn a Hamiltonian function that has a highly non-convex potential function. We showcase our algorithm by learning the following Hamiltonian function

$$H(q_1, q_2, p_1, p_2) = \frac{1}{2}(p_1^2 + p_2^2) + \sin\left(\frac{2\pi}{3} \cdot q_1\right)\cos\left(\frac{2\pi}{3} \cdot q_2\right) + \frac{\sin(\sqrt{q_1^2 + q_2^2})}{\sqrt{q_1^2 + q_2^2}},$$

whose potential function is visualized below in Figure 1 together with the solution given in Theorem 2.1 with $N = 1500$.

FIGURE 1. Learning with $N = 1500$ (a) Ground truth potential (b) Potential of the learned Hamiltonian (c) Mismatch error after vertical shift

## 4. ADDITIONAL RESULTS

Improved convergence rates can be formulated by invoking additional hypotheses like the so-called coercivity hypothesis. See [5] for a detailed presentation.

The results presented in this abstract for symplectic vector spaces admit a generalization to symplectic and Poisson manifolds (see [4]). In that case, the learning problem is far more degenerate since any modification of the original Hamiltonian function using a Casimir of the Poisson algebra yields the same Hamiltonian vector field and, hence, the same dynamics. It can nevertheless be proved that the ridge regularization term using the RKHS norm also guarantees the uniqueness of the solution of the learning problem.

## REFERENCES

[1] R. Abraham and J. E. Marsden, *Foundations of Mechanics*. Addison-Wesley, Reading, MA, 2nd Ed., 1978

[2] V. I. Arnol'd, *Mathematical Methods of Classical Mechanics*. Vol. 60, Springer Science & Business Media, 2013.

[3] J. Feng, C. Kulick, Y. Ren and S. Tang, *Learning particle swarming models from data with Gaussian processes*, Mathematics of Computation, 2023.

[4] J. Hu, J.-P. Ortega, and D. Yin. *A global structure-preserving kernel method for the learning of Poisson systems*, In preparation, 2024.

[5] J. Hu, J.-P. Ortega and D. Yin. *A structure-preserving kernel method for learning Hamiltonian systems*, ArXiv Preprint ArXiv:2403.10070, 2024.

[6] J. E. Marsden and Ratiu, T. *Introduction to mechanics and symmetry*. Springer-Verlag, New York, second Ed., 1999.

# Using exotic aromatic forests to construct order two scheme for the invariant measure sampling of Langevin dynamics with variable diffusion

EUGEN BRONASCO

(joint work with Benedict Leimkuhler, Dominic Phillips and Gilles Vilmart)

Exotic aromatic forests [6] are an extension of aromatic forests into the stochastic context and serve pivotal roles in generating order conditions for invariant measure sampling and studying the algebraic properties [2] of stochastic integrators.

Let $\phi$ denote a test function $\mathbb{R}^d \to \mathbb{R}$. Consider a system of stochastic differential equations with multiplicative noise with smooth vector field $F : \mathbb{R}^d \to \mathbb{R}^d$ and smooth diffusion $D : \mathbb{R}^d \to \mathbb{R}^{d \times d}$:

$$dX = F(X)dt + \sigma D(X)dW, \quad X(t) \in \mathbb{R}^d,$$

where $W(t) \in \mathbb{R}^d$ is a standard Wiener process. The weak Taylor expansion of the solution $X(t)$ is given by

$$\mathbb{E}[\phi(X(h))] = \phi(X_0) + h\mathcal{L}\phi(X_0) + \cdots + \frac{h^k}{k!}\mathcal{L}^{\circ k}\phi(X_0) + \cdots,$$

with generator, using Hessian matrix $\nabla^2\phi$, given by

$$\mathcal{L}\phi = \phi'F + \frac{\sigma^2}{2}\sum_{a=1}^{d}\phi''(D_a, D_a) = F \cdot \nabla\phi + \frac{\sigma^2}{2}Tr((\nabla^2\phi)DD^T).$$

An integrator $X_1 = \Phi_h(X_0)$ with the weak Taylor expansion

(1) $$\mathbb{E}[\phi(X_1)] = \phi(X_0) + h\mathcal{A}_1\phi(X_0) + \cdots + h^k\mathcal{A}_k\phi(X_0) + \cdots,$$

has weak order $p$ if $\mathcal{A}_k = \frac{1}{k!}\mathcal{L}^{\circ k}$ for $k = 1, \ldots, p$. [9]

For an *ergodic model* (e.g. overdamped Langevin dynamics where $F = -\nabla V$ and mild assumptions) with invariant measure $\mu$, the solution $X(t)$ satisfies

$$\lim_{T \to \infty}\frac{1}{T}\int_0^T \phi(X(t))dt = \int_{\mathbb{R}^d}\phi(x)d\mu(x), \quad \text{a.s.}$$

An *ergodic integrator* $X_n \mapsto X_{n+1}$ has order $q$ with respect to invariant measure sampling if

(2) $$\left|\lim_{N \to \infty}\frac{1}{N+1}\sum_{k=0}^{N}\phi(X_k) - \int_{\mathbb{R}^d}\phi(x)d\mu(x)\right| \leqslant Ch^q,$$

Given the differential operators $\mathcal{A}_k$ from the weak Taylor expansion (1) of $\mathbb{E}[\phi(X_1)]$, the condition (2) is satisfied if,

(3) $$\int_{\mathbb{R}^d}\mathcal{A}_k\phi(x)d\mu(x) = 0, \quad k = 1, \ldots, q.$$

For details see [1, 9].

## 1. Exotic aromatic forests

An *exotic aromatic forest* is a forest with edges oriented from top to bottom. This forest can contain cycles with edges oriented counterclockwise, and some of its vertices may be paired. For example:



In these forests, vertices represent vector fields, and edges represent directional derivatives. Cycles allow us to represent divergences, while paired vertices correspond to Laplacians. [2, 6]

**Using forests to check inv. measure sampling order.** We can use *integration by parts* denoted by $\sim$ to modify $\mathcal{A}_k$ without changing the value of the integral in (3)[2, 5]. The order conditions become

$$(a \circ A)(\tau) = 0, \quad \text{for all } \tau \in EAT, |\tau| \leqslant q,$$

where $A$ is an adjoint operation of the integration by parts. For example, we obtain among the order two conditions:



$$\ldots$$

In our recent work, we use integration by parts to develop a new order 2 method for the sampling of the invariant measure of Langevin dynamics with variable diffusion.

## 2. New scheme

Let $D(x) \in \mathbb{R}^{d \times d}$ be a symmetric matrix with $D = (D_1, \ldots, D_d)$ being smooth with respect to $x \in \mathbb{R}^d$ with columns $D_i = (D_i^1, \ldots, D_i^d)^T$. Then, we consider an SDE model with variable diffusion in $\mathbb{R}^d$ of the following form:

$$(4) \qquad dX = D^2(X)f(X)dt + \frac{\sigma^2}{2}\mathrm{div}(D^2)(X)dt + \sigma D(X)dW,$$

where $X(t) \in \mathbb{R}^d$ with $X(0) = X_0$ being deterministic, $f = -\nabla V$ with $V : \mathbb{R}^d \to \mathbb{R}$ being a smooth and globally Lipschitz potential, $\sigma > 0$ is a constant, and $W(t)$ is $d$-dimensional Wiener process fulfilling the usual assumptions. Divergence of the symmetric matrix $D$ is defined as

$$\mathrm{div}(D) = \Big(\sum_{j=1}^{d} \partial_j D_i^j\Big)_{i=1}^{d} = \begin{pmatrix} \mathrm{div}D_1 \\ \vdots \\ \mathrm{div}D_d \end{pmatrix}$$

In [7, 8], Leimkuhler-Matthews scheme is introduced for constant $D$ case. It has order 2 with respect to the invariant measure, requires only one evaluation of $f$, and has the following form:

$$(5) \qquad X_{n+1} = X_n + hD^2 f(X_n) + \sqrt{h}\sigma D \frac{R_n + R_{n+1}}{2}.$$

We generalize the Leimkuhler-Matthews scheme to the Langevin equation with variable diffusion. The new method has order 2 w.r.t. the invariant measure sampling, requires only one evaluation of the drift, and has the form:

$$X_{n+1} = X_n + hF(\overline{X}_n) + \hat{\Phi}_h^D(X_n + \frac{1}{4}hF(\overline{X}_{n-1})),$$

$$(6) \qquad \overline{X}_n = X_n + \frac{1}{2}\sqrt{h}\sigma D(X_n)R_n, \quad \text{with } \overline{X}_{-1} = X_0,$$

where $\Phi_h^D(X_n) = X_n + \hat{\Phi}_h^D(X_n)$ is an integrator of weak order 2 applied to the problem,

$$dX = \sigma D(X)dW,$$

where $\Phi^D(X_0) = X_0 + \sqrt{h}\sigma D(X_0)R_n + O(h)$. We study its stability properties and use the exotic aromatic forests framework to prove its convergence in the invariant measure. This work will be published in [5].

**Related ongoing work.**

(1) study of the algebraic properties of exotic aromatic forests and description of the backward error analysis and modified equation of ergodic stochastic differential equations in collaboration with Adrien Laurent, [4]

(2) implementation of a software package to automate the computations involving exotic aromatic forests in collaboration with Jean-Luc Falcone and Gilles Vilmart. [3]

REFERENCES

[1] A. Abdulle, and G. Vilmart, and K. Zygalakis, *High order numerical approximation of the invariant measure of ergodic SDEs*, SIAM Journal on Numerical Analysis **52**, 4 (2014), 1600–1622.

[2] E. Bronasco, *Exotic B-Series and S-Series: Algebraic Structures and Order Conditions for Invariant Measure Sampling*, Foundations of Computational Mathematics (2024).

[3] E. Bronasco, J.-L. Falcone, and G. Vilmart, *GraphAlgebra.hs: A Haskell library for the algebraic manipulation of graphs*, in preparation.

[4] E. Bronasco, A. Laurent, *Hopf algebra structures for the backward error analysis of ergodic stochastic differential equations*, in preparation.

[5] E. Bronasco, B. Leimkuhler, D. Phillips, and G. Vilmart, *Order 2 scheme for the invariant measure sampling of Langevin dynamics with variable diffusion*, in preparation.

[6] A. Laurent, and G. Vilmart, *Exotic aromatic B-series for the study of long time integrators for a class of ergodic SDEs*, Mathematics of Computation **89**, 321 (2019), 169–202.

[7] B. Leimkuhler, and C. Matthews, *Rational construction of stochastic numerical methods for molecular sampling*, Applied Mathematics Research Express. AMRX **1** (2013), 34–56.

[8] B. Leimkuhler, and C. Matthews, and M. V. Tretyakov, *On the long-time integration of stochastic gradient systems*, Proceedings of The Royal Society of London. Series A. Mathematical, Physical and Engineering Sciences **470**, 2170 (2014), 16.

[9] D. Talay, and L. Tubaro, *Expansion of the global error for numerical schemes solving stochastic differential equations*, Stochastic Analysis and Applications **8**, 4 (1990), 483–509.

# Wasserstein convergence and bias estimates for discretized kinetic Langevin dynamics

PETER A. WHALLEY

(joint work with Benedict Leimkuhler and Daniel Paulin)

We consider kinetic Langevin dynamics on $\mathbb{R}^d$ within the context of sampling. The dynamics are given by

$$\mathrm{d}X_t = V_t\mathrm{d}t$$
$$\mathrm{d}V_t = -\nabla U(X_t)\mathrm{d}t - \gamma V_t\mathrm{d}t + \sqrt{2\gamma}\mathrm{d}W_t,$$

where $\gamma > 0$, $(W_t)_{t \geqslant 0}$ is a $d$-dimensional standard Brownian motion and $U : \mathbb{R}^d \to \mathbb{R}$ is the potential energy function. The dynamics has invariant measure $\pi$ with density proportional to $\exp\left\{-U(x) + \frac{1}{2}\|v\|^2\right\}$. Ultimately, one is interested in sampling from $\pi$, which is typically done by discretizing the dynamics. We consider the setting where $U$ is $M$-$\nabla$Lipschitz and $m$-convex.

For some initial measure $\pi_0$ and a discretization with transition kernel $P_h$, stepsize $h > 0$ and invariant measure $\pi_h$, one is often interested in the distance to the target measure $\pi$ after $n \in \mathbb{N}$ steps

$$\mathcal{W}_2(\pi_0 P_h^n, \pi) \leqslant \underbrace{\mathcal{W}_2(\pi_0 P_h^n, \pi_h)}_{\text{Convergence Rate}} + \underbrace{\mathcal{W}_2(\pi_h, \pi)}_{\text{Bias}},$$

which can be split up in terms of the convergence rate of the discretization and the bias of the invariant measure.

We introduce methods to study the convergence rate and bias in the invariant measure separately. We provide convergence rates of $\mathcal{O}(m/M)$, with explicit stepsize restrictions, which are of the same order as the stability threshold for Gaussian targets and are valid for a large interval of the friction parameter. We apply this methodology to various integration schemes which are popular in the molecular dynamics and machine learning communities. Further, we introduce the property "$\gamma$-limit convergent" (GLC) to characterize underdamped Langevin schemes that converge to the overdamped dynamics in the high-friction limit and which have stepsize restrictions that are independent of the friction parameter; we show that this property is not generic by exhibiting methods from both the class and its complement.

We next consider the invariant measure bias of the BAOAB scheme [3], typical approaches for quantifying the asymptotic bias in Wasserstein distance rely on strong order estimates. However, BAOAB is only strong order one in stepsize with respect to the continuous dynamics, but it has weak order two in stepsize in the

FIGURE 1. Contour plots of convergence rate for various schemes in the case of an anisotropic Gaussian with parameters $m = 1$ and $M = 10$. Regions of white indicate instability.

invariant measure. Our approach to achieve second-order estimates is to strongly approximate the BAOAB scheme by a modified stochastic dynamics which preserves the invariant measure. In particular, we introduce the HOH scheme, where H corresponds to exact Hamiltonian steps and this approximates the BAOAB scheme up to second order in stepsize uniformly in time, whilst preserving the invariant measure.

We refer to the works [1] and [2] for more details.

## REFERENCES

[1] B. J. Leimkuhler, D. Paulin, and P. A. Whalley, *Contraction and convergence rates for discretized kinetic Langevin dynamics*, SIAM Journal on Numerical Analysis **62**, 3 (2024), 1226–1258.
[2] B. Leimkuhler, D. Paulin, and P. A. Whalley, *Contraction rate estimates of stochastic gradient kinetic Langevin integrators*, to appear in ESAIM: Mathematical Modelling and Numerical Analysis (2024).
[3] B. Leimkuhler, and C. Matthews, *Rational construction of stochastic numerical methods for molecular sampling*, Applied Mathematics Research Express. AMRX **1** (2013), 34–56.

## Constrained and partitioned training of neural networks

### TIFFANY VLAAR

(joint work with Jonathan Frankle, Matthias Hein, Benedict Leimkuhler, Maximilian Müller, Timothée Pouchon, David Rolnick and Amos Storkey)

In the first half of the talk I discussed the use of constrained stochastic differential equations to train deep neural networks. Common techniques used to improve the generalization performance of deep neural networks (such as e.g. L2 regularization [1, 2] and batch normalization [3]) are tantamount to imposing a parameter constraint, but despite their widespread use are often not well understood [4, 5]. I

described an approach for efficiently incorporating hard constraints into a stochastic gradient Langevin dynamics framework [6]. Our constraints offer direct control of the parameter space, which allows us to study their effect on generalization. In the second half of the talk, I focused on the role played by individual layers and substructures of neural networks: layer-wise sensitivity to the choice of initialization and optimizer hyperparameter settings varies [7] and training different neural network layers differently may lead to enhanced generalization and reduced computational cost [8]. In particular, I showed that a multirate approach can be used to train deep neural networks for transfer learning applications in half the time, without reducing the generalization performance of the model [9].

## REFERENCES

[1] A. Hoerl and R. Kennard, *Ridge Regression: Biased Estimation for Nonorthogonal Problems*, Technometrics **12** (1970), 55–67.

[2] A. Krogh and J. Hertz, *A Simple Weight Decay can Improve Generalization*, Advances in Neural Information Processing Systems (1991).

[3] S. Ioffe and C. Szegedy, *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*, International Conference on Machine Learning, PMLR **37** (2015).

[4] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, *How Does Batch Normalization Help Optimization?*, Advances in Neural Information Processing Systems (2018).

[5] M. Andriushchenko, F. D'Angelo, A. Varre, and N. Flammarion, *Why Do We Need Weight Decay in Modern Deep Learning?*, preprint arXiv:2310.04415 (2023).

[6] B. Leimkuhler, T. Vlaar, T. Pouchon, and A. Storkey, *Better Training using Weight-Constrained Stochastic Dynamics*, International Conference on Machine Learning, PMLR **139** (2021).

[7] T. Vlaar and J. Frankle, *What Can Linear Interpolation of Neural Network Loss Landscapes Tell Us?*, International Conference on Machine Learning, PMLR **162** (2022).

[8] M. Müller, T. Vlaar, D. Rolnick, and M. Hein, *Normalization Layers Are All That Sharpness-Aware Minimization Needs*, Advances in Neural Information Processing Systems **36** (2023).

[9] T. Vlaar and B. Leimkuhler, *Multirate Training of Neural Networks*, International Conference on Machine Learning, PMLR **162** (2022).

# High order integration of stochastic dynamics in $\mathbb{R}^d$, on manifolds, and in the neighbourhood of manifolds

ADRIEN LAURENT

(joint work with Eugen Bronasco, Hans Z. Munthe-Kaas and Gilles Vilmart)

On a Riemannian manifold $(\mathcal{M}, g)$, the overdamped Langevin dynamics write as the following equation defined using the Eells-Elworthy-Malliavin construction (see [9]):

$$(1) \qquad dX(t) = f(X(t))dt + dW^{\mathcal{M}}(t), \quad f = -\nabla V.$$

We are interested in particular in the Euclidean case $\mathcal{M} = \mathbb{R}^d$ with the standard Euclidean overdamped Langevin dynamics,

$$(2) \qquad dX(t) = f(X(t))dt + dW(t),$$

and in the case $\mathcal{M} = \{x \in \mathbb{R}^d | \zeta(x) = 0\}$, where the manifold is defined as the zero-level set of a smooth constraint $\zeta \colon \mathbb{R}^d \to \mathbb{R}^q$,

$$(3) \qquad dX(t) = \Pi_{\mathcal{M}}(X(t))f(X(t))dt + \Pi_{\mathcal{M}}(X(t)) \circ dW(t),$$

with $\Pi_{\mathcal{M}}$ the orthogonal projection on $T\mathcal{M}$. These models appear naturally when modelling the motion of a set of particles subject to a potential $V$ in a high friction regime. The constraints can for instance represent fixed distances or angles in molecules. The constraints encoded by $\zeta$ are often satisfied only up to a small parameter $\varepsilon$, yielding dynamics that evolve in the neighbourhood of $\mathcal{M}$,

$$dX^\varepsilon(t) = f(X^\varepsilon(t))dt + dW(t) + \frac{1}{4}\nabla \ln(\det(G))(X^\varepsilon(t))dt$$

$$(4) \qquad\qquad - \frac{1}{\varepsilon}(gG^{-1}\zeta)(X^\varepsilon(t))dt, \quad g = \nabla\zeta, \quad G = g^T g.$$

Our focus is the simulation of the law of overdamped Langevin dynamics in finite time (weak approximation) and in long time (approximation of the invariant measure $d\mu_\infty \propto e^{-2V} d\mathrm{vol}_{\mathcal{M}}$). Following [20], a method of high weak order also samples the invariant measure with high order. It is, however, well known that there exists method of low weak order, typically one, that sample the invariant measure with high order. For instance, the Leimkuhler-Matthews method [16] has order two for sampling the invariant measure of (2) for a similar cost as the Euler Maruyama method. The methodology for the creation of such high order methods for sampling the invariant measure of (2) was studied in [6, 1] and in [7, 11] in the case of embedded manifolds (3). We introduce the algebraic formalism of exotic aromatic Butcher series in [14, 15] to deal with the tedious calculations. This generalisation of the standard Butcher series [8] was later studied for its far-reaching algebraic and geometric properties [4, 13, 5].

The contributions presented in the talk are the following. We uncover in [14] a class of stochastic Runge-Kutta integrators with the order conditions for reaching up to any order for sampling the invariant measure of (2). We extend these results for sampling (3) on embedded manifolds in [15], where we introduce the first stochastic projection methods of order two (see Figure 1). We emphasize that one cannot use standard Euclidean methods in this context as the measure sampled by the numerical scheme would then be absolutely continuous w.r.t. the Lesbesgue measure on $\mathbb{R}^d$ and not w.r.t. the measure on the manifold. It is thus crucial that the chosen integrator lies on $\mathcal{M}$.

We then present a new approach for the creation of projection methods of high weak order for the dynamics (4) in the neighbourhood of manifolds. In this context, a constrained method works only for small $\varepsilon$, while a Euclidean method will face severe timestep restriction for small $\varepsilon$. Taking advantage of the geometry of the problem, we propose a new projection method in [12] with accuracy and cost independent of $\varepsilon$. In the case of a manifold of codimension one, and under

technical assumptions (see [12]), the method is the following:

$$X_{n+1}^\varepsilon = X_n^\varepsilon + \sqrt{h}\sigma\xi_n + hf(X_n^\varepsilon) + \frac{(1-e^{-h/\varepsilon})^2}{2}(\zeta^2 G^{-2}g'(g))(X_n^\varepsilon)$$

$$+ \frac{\sigma^2\varepsilon}{4}(1-e^{-2h/\varepsilon})(G^{-1}g'(g))(X_n^\varepsilon) + g(X_n^\varepsilon)\lambda_{n+1}^\varepsilon,$$

$$\zeta(X_{n+1}^\varepsilon) = e^{-h/\varepsilon}\zeta(X_n^\varepsilon) + \sigma\sqrt{\frac{\varepsilon}{2}(1-e^{-2h/\varepsilon})}g^T(X_n^\varepsilon)\xi_n$$

$$+ \varepsilon(1-e^{-h/\varepsilon})(g^T f + \frac{\sigma^2}{2}G^{-1}g^T g'(g) + \frac{\sigma^2}{2}\mathrm{div}(g))(X_n^\varepsilon).$$



FIGURE 1. A trajectory (left) and the convergence curves for the invariant measure (right) (from [15]).

While the approach using projection methods is computationally efficient, easy to implement, and widely used in applications, projection methods rely on an embedding of $\mathcal{M}$ into a possibly high-dimensional vector space and on the use of non-intrinsic quantities. In the spirit of [3], we shall propose in future works new stochastic Lie-group methods [10, 17] of high order in the weak sense and for sampling the invariant measure of (1). The extension of Lie group methods for the approximation of stochastic dynamics on general Riemannian manifolds is a challenging problem as its deterministic counterpart is already an active field of research [2, 19, 18]. This is exciting matter for future work.

## REFERENCES

[1] Assyr Abdulle, Gilles Vilmart, and Konstantinos C. Zygalakis. High order numerical approximation of the invariant measure of ergodic SDEs. *SIAM J. Numer. Anal.*, 52(4):1600–1622, 2014.

[2] M. J. H. Al-Kaabi, K. Ebrahimi-Fard, D. Manchon, and H. Z. Munthe-Kaas. Algebraic aspects of connections: from torsion, curvature, and post-lie algebras to Gavrilov's double exponential and special polynomials. *Submitted*, 2022.

[3] Karthik Bharath, Alexander Lewis, Akash Sharma, and Michael V Tretyakov. Sampling and estimation on manifolds using the Langevin diffusion. *arXiv preprint arXiv:2312.14882*, 2023.

[4] Eugen Bronasco. Clumped forests for the description of the Hopf algebra structures of exotic aromatic series and their applications in stochastic numerical analysis. *Submitted*, 2022.

[5] Eugen Bronasco and Adrien Laurent. Hopf algebra structures for the backward error analysis of ergodic stochastic differential equations. *In preparation*, 2023.

[6] Arnaud Debussche and Erwan Faou. Weak backward error analysis for SDEs. *SIAM J. Numer. Anal.*, 50(3):1735–1752, 2012.

[7] Erwan Faou and Tony Lelièvre. Conservative stochastic differential equations: mathematical and numerical analysis. *Math. Comp.*, 78(268):2047–2074, 2009.

[8] Ernst Hairer, Christian Lubich, and Gerhard Wanner. Geometric numerical integration, volume 31 of *Springer Series in Computational Mathematics.* Springer-Verlag, Berlin, second edition, 2006. Structure-preserving algorithms for ordinary differential equations.

[9] Elton P. Hsu. *Stochastic analysis on manifolds*, volume 38 of *Graduate Studies in Mathematics.* American Mathematical Society, Providence, RI, 2002.

[10] Arieh Iserles, Hans Z. Munthe-Kaas, Syvert P. Nørsett, and Antonella Zanna. Lie-group methods. In *Acta numerica*, 2000, volume 9 of *Acta Numer.*, pages 215–365. Cambridge Univ. Press, Cambridge, 2000.

[11] Adrien Laurent. *Algebraic Tools and Multiscale Methods for the Numerical Integration of Stochastic Evolutionary Problems.* PhD thesis, University of Geneva, 2021.

[12] Adrien Laurent. A uniformly accurate scheme for the numerical integration of penalized Langevin dynamics. *SIAM J. Sci. Comput.*, 44(5):A3217–A3243, 2022.

[13] Adrien Laurent and Hans Z. Munthe-Kaas. The universal equivariance properties of exotic aromatic B-series. *Submitted*, arXiv:2305.10993, 2023.

[14] Adrien Laurent and Gilles Vilmart. Exotic aromatic B-series for the study of long time integrators for a class of ergodic SDEs. *Math. Comp.*, 89(321):169–202, 2020.

[15] Adrien Laurent and Gilles Vilmart. Order conditions for sampling the invariant measure of ergodic stochastic differential equations on manifolds. *Found. Comput. Math.*, 22(3):649–695, 2022.

[16] Benedict Leimkuhler and Charles Matthews. Rational construction of stochastic numerical methods for molecular sampling. *Appl. Math. Res. Express. AMRX*, 2013(1):34–56, 2013.

[17] Simon J. A. Malham and Anke Wiese. Stochastic Lie group integrators. *SIAM J. Sci. Comput.*, 30(2):597–617, 2008.

[18] Hans Munthe-Kaas. Geometric integration on symmetric spaces. *J. Comput. Dyn.*, 11(1):43–58, 2024.

[19] Hans Munthe-Kaas and Jonatan Stava. Lie admissible triple algebras: The connection algebra of symmetric spaces. *Submitted*, 2023.

[20] Denis Talay and Luciano Tubaro. Expansion of the global error for numerical schemes solving stochastic differential equations. *Stochastic Anal. Appl.*, 8(4):483–509 (1991), 1990.

# Retraction-based simulations of Brownian motion on Riemannian and sub-Riemannian manifolds

SIMON SCHWARZ

(joint work with Michael Herrmann, Anja Sturm and Max Wardetzky)

Probabilistic models in continuous time with geometric constraints lead to Brownian motion and stochastic differential equations (SDEs) on Riemannian manifolds. A naive approach to simulate such a process is to use geodesic random walks on a Riemannian manifold $(M, g)$: In each step one can

(1) sample a tangent vector at the current position and
(2) follow a geodesic for some fixed stepsize in the sampled direction.

Jørgensen proved in [2] that under some mild assumptions geodesic random walks converge to diffusion processes on manifolds (depending on the mean and variance of the probability measure on the tangent bundle) as the stepsize decreases. Computing geodesics, however, is computationally expensive – we therefore propose to use *retractions* in Step (2) of the algorithm. Retractions are computationally efficient approximations of the exponential map that were originally introduced for numerical optimization on manifolds, see [1].

**Definition 1.** *Let* $\mathrm{Ret} : TM \to M$ *be a smooth map, and denote the restriction to the tangent space* $T_xM$ *by* $\mathrm{Ret}_x$ *for any* $x \in M$. $\mathrm{Ret}$ *is a* retraction *if the following two conditions are satisfied for all* $x \in M$ *and all* $v \in T_xM$:

(1) $\mathrm{Ret}_x(0) = x$, *where* $0$ *is the zero element in* $T_xM$ *and*
(2) $\frac{d}{d\tau} \mathrm{Ret}_x(\tau v)\big|_{\tau=0} = v$ *(where we identify* $T_0 T_xM \simeq T_xM$*).*

*A retraction is a* second-order retraction *if it additionally satisfies that for all* $x \in M$ *and for all* $v \in T_xM$ *one has that*

$$\frac{D}{d\tau}\left( \frac{d}{d\tau} \mathrm{Ret}_x(\tau v) \right)\bigg|_{\tau=0} = \frac{D}{d\tau}\left( \frac{d}{d\tau} \mathrm{Exp}_x(\tau v) \right)\bigg|_{\tau=0} = 0 \ ,$$

*where* $\frac{D}{d\tau}(\frac{d}{d\tau}\gamma(\tau))$ *denotes covariant differentiation of the tangent vector field* $\dot{\gamma}(\tau) = \frac{d}{d\tau}\gamma(\tau)$ *along the curve* $\gamma$.

We prove in [4] that retraction-based random walks converge to the correct limiting process if and only if the respective retraction is a second-order retraction. Moreover, we give several examples of second-order retractions and therefore provide an *efficient* and *convergent* way of simulating diffusion processes on Riemannian manifolds. Our retraction-based algorithms can also be generalized to sub-Riemannian manifolds, see [3].

## References

[1] P.-A. Absil, R. Mahony, R. Sepulchre *Optimization algorithms on matrix manifolds*, Princeton University Press, 2008.
[2] E. Jørgensen *The Central Limit Problem for Geodesic Random Walks*, Z. Wahrscheinlichkeitstheorie verw. Gebiete **32** (1975), 1–64.
[3] M. Herrmann, P. Neumann, S. Schwarz, A. Sturm, M. Wardetzky, *Sub-Riemannian Random Walks: From Connections to Retractions* (2023), arXiv:2311.17289.
[4] S. Schwarz, M. Herrmann, A. Sturm, M. Wardetzky, *Efficient Random Walks on Riemannian Manifolds*, Foundations of Computational Mathematics (2023), 1–17.

## XAI meets quantum chemistry

Klaus-Robert Müller

(joint work with ML group, QCML, BIFOLD and others)

In this talk I describe the recently developed field of explainable AI (XAI) [5, 6, 7], which has now many applications in the sciences and industry [5, 4]. XAI allows to study how machine learning methods (ML) such as deep learning, LSTMs and

kernel methods come to their decision (in terms of input variables) on a single sample basis despite of their non-linearities. The concept of Clever Hans is introduced for learning models [8]. Using this, we analyse ML for quantum chemistry and demonstrate that novel scientific insights emerge from trained ML models [2, 1, 10, 11, 12, 3, 9]. Finally ML for PDEs are discussed using a particularly challenging application, namely, the control of the curling robot Curly [13].

### REFERENCES

[1] K. Schütt, F. Arbabzadah, S. Chmiela, K.-R. Müller, A. Tkatchenko. *Quantum-chemical insights from deep tensor neural networks*, Nature Communication. **8** (2017), 13890.

[2] M. Rupp, A. Tkatchenko, K.-R. Müller, O.A. Von Lilienfeld, *Fast and accurate modeling of molecular atomization energies with machine learning*, Physical Review Letters **108(5)** (2012), 058301.

[3] O.T. Unke, S. Chmiela, H.E. Sauceda, M. Gastegger, I. Poltavsky, K.T. Schütt, A. Tkatchenko and K.-R. Müller, *Machine learning force fields*, Chemical Reviews **121(16)** (2021), 10142-10186.

[4] F. Klauschen, J. Dippel, P. Keyl, P. Jurmeister, M. Bockmayr, A. Mock, O. Buchstab, M. Alber,L. Ruff, G. Montavon, K.-R. Müller, *Toward explainable artificial intelligence for precision pathology*, Annual Review of Pathology: Mechanisms of Disease **19** (2024), 541-70.

[5] W. Samek, G. Montavon, S. Lapuschkin, C.J. Anders and K.-R. Müller, *Explaining deep neural networks and beyond: A review of methods and applications*, Proceedings of the IEEE **109(3)** (2021), 247-278.

[6] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller and W. Samek, *On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation*, PloS one **10(7)** (2015), e0130140.

[7] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen and K.-R. Müller, *How to explain individual classification decisions*, The Journal of Machine Learning Research **11** (2010), 1803-1831.

[8] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek and K.-R. Müller, *Unmasking Clever Hans predictors and assessing what machines really learn*, Nature communications **10(1)** (2019), 1096.

[9] O.T. Unke, M. Stöhr, S. Ganscha, T. Unterthiner, H. Maennel, S. Kashubin, D. Ahlin, M. Gastegger, L. Medrano Sandonas, J.T. Berryman, A. Tkatchenko and K.-R. Müller, *Biomolecular dynamics with machine-learned quantum-mechanical force fields trained on diverse chemical fragments*, Science Advances **10(14)** (2024), eadn4397.

[10] S. Chmiela, A. Tkatchenko,H.E. Sauceda, I. Poltavsky, K.T. Schütt and K.-R. Müller, *Machine learning of accurate energy-conserving molecular force fields*, Science advances **3(5)** (2017), e1603015.

[11] S. Chmiela, H.E. Sauceda, K.-R. Müller and A. Tkatchenko, *Towards exact molecular dynamics simulations with machine-learned force fields*, Nature communications **9(1)** (2018), 3887.

[12] S. Chmiela, V. Vassilev-Galindo, O.T. Unke, A. Kabylda, H.E. Sauceda, A. Tkatchenko and K.-R. Müller, *Accurate global machine learning force fields for molecules with hundreds of atoms*, Science Advances **9(2)** (2023), eadf0873.

[13] D.O. Won, K.-R. Müller and S.W. Lee, *An adaptive deep reinforcement learning framework enables curling robots with human-like performance in real-world conditions*, Science Robotics **5(46)** (2020), eabb9764.

## Sampling and generative modeling on Lie group

Molei Tao

(joint work with Lingkai Kong, Yuchen Zhu and Tianrong Chen)

My talk and this report are based on two articles, [1] and [2]. Interested readers are strongly encouraged to read the full articles.

### 1. Sampling of Probability Distributions on Lie Groups

Explicit, momentum-based dynamics for optimizing functions defined on Lie groups was recently constructed, based on techniques such as variational optimization and left trivialization. We appropriately add tractable noise to the optimization dynamics to turn it into a sampling dynamics, leveraging the advantageous feature that the trivialized momentum variable is Euclidean despite that the potential function lives on a manifold. We then propose a Lie-group MCMC sampler, by delicately discretizing the resulting kinetic-Langevin-type sampling dynamics. The Lie group structure is exactly preserved by this discretization. Exponential convergence with explicit convergence rate for both the continuous dynamics and the discrete sampler are then proved under $W_2$ distance. Only compactness of the Lie group and geodesically $L$-smoothness of the potential function are needed. To the best of our knowledge, this is the first convergence result for kinetic Langevin on curved spaces, and also the first quantitative result that requires no convexity or, at least not explicitly, any common relaxation such as isoperimetry.

### 2. Generative Modeling of Data on Lie Groups

The generative modeling of data on manifold is an important task, for which diffusion models in flat spaces typically need nontrivial adaptations. We demonstrate how the trivialization technique used in Part I can transfer the effectiveness of diffusion models in Euclidean spaces to Lie groups. In particular, an auxiliary momentum variable was algorithmically introduced to help transport the position variable between data distribution and a fixed, easy-to-sample distribution. Normally, this would incur further difficulty for manifold data because momentum lives in a space that changes with the position. However, our trivialization technique creates to a new momentum variable that stays in a simple **fixed vector space**. This design, together with a manifold preserving integrator, simplifies implementation and avoids inaccuracies created by approximations such as projections to tangent space and manifold, which were typically used in prior work, hence facilitating generation with high-fidelity and efficiency. The resulting method achieves state-of-the-art performance on protein and RNA torsion angle generation and sophisticated torus datasets. We also, arguably for the first time, tackle the generation of data on high-dimensional Special Orthogonal and Unitary groups, the latter essential for quantum problems.

REFERENCES

[1] Kong, L. & Tao, M., *Convergence of Kinetic Langevin Monte Carlo on Lie groups*, COLT (2024).
[2] Zhu, Y., Chen, T., Kong, L., Theodorou, E. & Tao, M., *Trivialized Momentum Facilitates Diffusion Generative Modeling on Lie Groups*, ArXiv:2405.16381 (2024).

## Computational methods for Bayesian inverse problems

KONSTANTINOS ZYGALAKIS

(joint work with Yoann Altmann, Paul Dobson, Teresa Klatzer, Marcelo Pereyra and Jesus Maria Sanz-Serna)

Bayesian (imaging) inverse problems provide a coherent mathematical and algorithmic framework that enables researchers to combine mathematical models with data. A problem of typical interest in this setting is the recovery of an unknown image $x \in \mathbb{R}^d$ through a measurement $y \in \mathbb{R}^m$ (typically $m \leqslant d$) which is assumed to be related with $x$ by a statistical model $p(y|x)$. The simplest example of such relation between the image $x$ and the measurement $y$ is the following one

$$y = Ax + w$$

where $A \in \mathbb{R}^{m \times d}$ (rank-deficient) and $w$ is additive (Gaussian noise).

The fact that $A$ is rank deficient implies that the recovery of $x$ from $y$ is ill-posed resulting in significant uncertainty about $x$. In order to reduce the overall uncertainty and deliver accurate results one needs a prior distribution $p_r(x)$ that is meant to reflecting the properties of true image $x$. Given the prior distribution $p_r(x)$ and the statistical model $p(y|x)$ we can write down the posterior distribution using Bayes rule:

$$\pi(x|y) \propto p(y|x)p_r(x),$$

which models our knowledge of $x$ after observing $y$. The ability thus to solve such (imaging) inverse problems depends crucially on the efficient calculation of quantities relating to the posterior distribution. In particular, one might be interested in calculating

$$\hat{x}_{\mathrm{MAP}} := \operatorname*{argmax}_{x \in \mathbb{R}^d} \pi(x|y),$$

the point estimate that maximises the posterior distribution or some other statistic of the form

$$(1) \qquad \mathbb{E}_\pi(f) := \int_{x \in \mathbb{R}^d} f(x)\pi(x|y)dx.$$

Calculations of terms of the form (1) can be done using the discretizations of the Langevin equation

$$(2) \qquad dX_t = \nabla \log \pi(X_t|y)dt + \sqrt{2}dW_t,$$

where $W_t$ is the standard $d$-dimensional Brownian motion. One of the main computational challenges in the case of imaging inverse problems is the non-smoothness

of the prior $p_r(x)$. To deal with this, one replaces $\pi(x|y)$ with a smooth approximation $\pi^\lambda(x|y)$ for which $\lambda \log \pi^\lambda(x|y)$ is gradient Lipschitz with constant that behaves like $\lambda^{-1}$ which results in time-step restrictions for standard stochastic integrators. To address this issue one needs to use tailored stochastic numerical integrators that are based either on explicit stabilised solvers [1] or they are implicit [2]. In the case of the implicit methods each step of the stochastic integrator corresponds to solving an optimization step. In particular, we have

$$X_{n+1} = \operatorname{argmin}_{x \in \mathbb{R}^d} F(x; X_n; \xi_{n+1}), \qquad \xi_{n+1} \sim \mathcal{N}(0, I_d)$$

$$F(x; u, z) := -\theta^{-1} \log \pi^\lambda(\theta x + (1-\theta)u) + \frac{1}{2\delta}\|x - u - \sqrt{2\delta}z\|^2.$$

which for the case $\theta = 1/2$ we call implicit midpoint Langevin algorithm. One can show that this integrator is exact for Gaussian targets while in the case of strongly log-concave potentials converges towards the *biased* invariant measure with similar speed to accelerated optimization methods. An interesting open research question is to study this algorithm in the case where an accept-reject correction is introduced to remove the bias from the discretization.

<div align="center">REFERENCES</div>

[1] L. Vargas, M. Pereyra, K.C. Zygalakis, *Accelerating proximal Markov Chain Monte Carlo by using explicit stabilised methods.*, SIAM J. Imaging Sci. **13**(2) (2020), 905–935.

[2] T. Klatzer, P. Dobson, Y. Altmann, M. Pereyra, J. M. Sanz-Serna, and K. C. Zygalakis, *Accelerated Bayesian imaging by relaxed proximal-point Langevin sampling*, SIAM J. Imaging Sci. **17**(2) (2024), 1078–1117.

<div align="center">

## Random batch methods for interacting particle systems and molecular dynamics
### SHI JIN

</div>

We first develop random batch methods for classical interacting particle systems with large number of particles [1, 5]. These methods use small but random batches for particle interactions, thus the computational cost is reduced from $O(N^2)$ per time step to $O(N)$, for a system with $N$ particles with binary interactions. For one of the methods, we give a particle number independent error estimate under some special interactions [1, 2, 7].

This method is also extended to molecular dynamics with Coulomb interactions, in the framework of Ewald summation [3]. We will show its superior performance compared to the current state-of-the-art methods (for example PPPM) for the corresponding problems, in the computational efficiency and parallelizability [4, 6].

<div align="center">REFERENCES</div>

[1] Shi Jin, Lei Li and Jian-Guo Liu, *Random Batch Methods (RBM) for interacting particle systems*, J. Comp. Phys. **400** (2020), 108877.

[2] Shi Jin, Lei Li and Jian-Guo Liu, *Convergence of Random Batch Method for interacting particles with disparate species and weights*, SIAM J. Numer. Anal. **59** (2021), 746–768.

[3] Shi Jin, Lei Li, Zhenli Xu, and Yue Zhao, *A random batch Ewald method for particle systems with Coulomb interactions*, SIAM J. Sci. Comp., **43** (2021), B937–B960.

[4] Jiuyang Liang, Pan Tan, Yue Zhao, Lei Li, Shi Jin, Liang Hong, Zhenli Xu, *Super-Scalable Molecular Dynamics Algorithm*, J. Chem. Phys. **156** (2022), 014114.

[5] Shi Jin and Lei Li, *Random Batch Methods for classical and quantum interacting particle systems and statistical sampling*, Active Particles, III, Birkhäuser, Cham, (2022), 153-200, 2022 (ed. N. Bellomo, J. Carrillo, and E. Tadmor).

[6] Jiuyang Liang, Pan Tan, Liang Hong, Shi Jin, Zhenli Xu, Lei Li, *A random batch Ewald method for charged particles in the isothermal-isobaric ensemble,* J. Chem. Phys. **156** (2022), 014114.

[7] Shi Jin, Lei Li, Xuda Ye, Zhennan Zhou, *Ergodicity and long-time behavior of the Random Batch Method for interacting particle systems*, Math. Model. Math. Appl. Sci. **33** (2023), 67–102.

# Entropy of convex hulls revisited

SARA VAN DE GEER

Let $Q$ be a probability measure on a measurable space $\mathbf{X}$. We consider the convex hull $\mathbf{F}$ of a set $X \in L_2(Q)$ which is defined as

$$\mathbf{F} := \{f = \sum_{i=1}^{p} \beta_j x_j : \ x_j \in X, \ \beta_j \geqslant 0, \ j = 1, \ldots, p, \ \sum_{j=1}^{p} \beta_j = 1, \ p \in \mathbb{N}\}.$$

We study the (metric) entropy of $\mathbf{F}$. Metric entropy is a key concept in approximation theory, with numerous applications in various areas such as statistics, signal processing, probability and information theory. For example, the entropy of the parameter space in a statistical model typically (in an asymptotic sense where the number of samples goes to infinity) determines the estimation error.

We assume that $X$ is "small" in the sense that it has polynomial covering numbers.

**Definition 1.** *For $\epsilon > 0$, the $\epsilon$-covering number $N(\epsilon, T)$ of a set $T \subset L_2(Q)$ is defined as the minimum number of balls with radius $\epsilon$, necessary to cover $T$. The entropy of $T$ is $\mathbf{H}(\cdot, T) := \log N(\cdot, T)$.*

Let $\|\cdot\|$ be the $L_2(Q)$-norm. We assume that $\sup_{x \in X} \|x\| \leqslant 1$ and that, for some constant $V > 0$,

$$N(\epsilon, X) \lesssim \epsilon^{-V}.$$

We cite the following result of [1]:

**Theorem 1.**

$$\mathbf{H}(\epsilon, \mathbf{F}) \lesssim \epsilon^{-\frac{2V}{2+V}}.$$

The above bound is tight in certain cases, but it is not tight in general. We therefore introduce another concept, namely the approximation number.

**Definition 2.** *For $\epsilon > 0$ the $\epsilon$-approximation number $M(\epsilon, T)$ of a set $T \subset L_2(Q)$ is defined as the smallest dimension $M$ such that there exists a linear space $\mathbf{V} \subset$*

$L_2(Q)$ with dimension $M$ that has $\delta(T, \mathbf{V}) \leqslant \epsilon$, where $\delta(T, \mathbf{V})$ is the distance of $T$ to $\mathbf{V}$:

$$\delta(T, \mathbf{V}) := \max_{t \in T} \min_{v \in \mathbf{V}} \|t - v\|.$$

Clearly, $M(\epsilon, X) \leqslant N(\epsilon, X)$ for all $\epsilon > 0$.

**Theorem 2.** *If, for some constants $W > 0$ and $w \geqslant 0$,*

$$M(\epsilon, X) \lesssim \epsilon^{-W} \log^w(1/\epsilon).$$

*Then*

$$\mathbf{H}(\epsilon, \mathbf{F}) \quad \lesssim \quad \epsilon^{-\frac{2W}{2+W}} \log^{\frac{2w}{2+W}}(1/\epsilon) \log^{\frac{W}{2+W}}(1/\epsilon).$$

Here are two examples.

**Example 1.** *Let $\mu$ be Lebesgue measure on $[0, 1]$ and $\Psi = \{\psi_v : v \in [0, 1]\}$ be the collection of heaviside functions, that is the collection of indicators of halfintervals*

$$\psi_v(\cdot) := 1\{\cdot \geqslant v\}, \ v \in [0, 1].$$

*Then*

$$X := \underbrace{\Psi \otimes \cdots \otimes \Psi}_{d \text{ times}}$$

*is the collection of indicators of halfintervals in $[0, 1]^d$, and $\mathbf{F}$ is the collection of all d-dimensional distribution functions. Let $Q = \mu \times \cdots \times \mu$ be Lebesgue measure on $[0, 1]^d$. Then it can be shown that*

$$M(\epsilon, X) \lesssim \epsilon^{-2} \log^{2(d-1)}(1/\epsilon).$$

*Hence*

$$\mathbf{H}(\epsilon, \mathbf{F}) \lesssim \epsilon^{-1} \log^{d-\frac{1}{2}}(1/\epsilon).$$

*Thus we recover the bound of [2].*

Note that in the above example, the dimension $d$ occurs only in the logarithmic term. The same is true for the second example.

**Example 2.** *Consider the convolution $U + V$, where $U$ and $V$ are independent random vectors in $[0, 1]^d$ and the components of $U$ are independent and all have a Beta$(q_1, q_2)$-distribution with $(q_1, q_2) \in \mathbb{N}^2$. Let $\mathbf{F}$ be the class of all possible densities of $U + V$, with respect to Lebesgue measure on $[0, 2]^d$. Then for $q := \min\{q_1, q_2\}$ one can show that*

$$M(\epsilon, X) \lesssim \epsilon^{-\frac{2}{2q-1}} \log^{\frac{2q(d-1)}{2q-1}}(1/\epsilon).$$

*We conclude that*

$$\mathbf{H}(\epsilon, \mathbf{F}) \lesssim \epsilon^{-\frac{1}{q}} \log^{d-\frac{1}{2q}}(1/\epsilon).$$

References

[1] K. Ball, and A. Pajor, *The entropy of convex bodies with "few" extreme points*, London Math. Soc. Lecture Note Series **158** (1990), 25–32.
[2] R. Blei and F. Gao and W. Li, *Metric entropy of high dimensional distributions*, Proceedings of the American Mathematical Society **135** (2007), 4009–4018.

# Infinite-dimensional Wishart Processes

SONJA COX

(joint work with Christa Cuchiero and Asma Khedher)

The goal is to introduce and analyse infinite-dimensional Wishart processes. An infinite-dimensional Wishart processes is a stochastic process $X = (X_t)_{t \geqslant 0}$ taking values in $S_1^+(H)$, the cone of positive self-adjoint trace class operators on a separable real Hilbert space $H$, and satisfying (in some sense) the following stochastic differential equation:

$$(1) \quad dX_t = (\alpha Q + X_t A + A^* X_t) \, dt + \sqrt{X_t} \, dW_t \sqrt{Q} + \sqrt{Q} \, dW_t^* \sqrt{X_t}, \, t \geqslant 0, \, X_0 = x_0.$$

Here $\alpha \in \mathbb{R}$, $A \colon D(A) \subset H \to H$ is the generator of a $C_0$-semigroup, $x_0$ and $Q$ are positive self-adjoint bounded operators, and $(W_t)_{t \geqslant 0}$ is an $L_2(H)$-cylindrical Brownian motion (where $L_2(H)$ is the space of Hilbert Schmidt operators on $H$).

*Finite-dimensional Wishart processes*, i.e., processes taking values in $S^+(\mathbb{R}^n)$, the cone of positive semidefinite $n \times n$ matrices, have been thoroughly studied: in [1, 2] the existence of finite-dimensional Wishart processes was established under certain conditions on the parameters, and stochastic differential equations were derived for the eigenvalues and eigenvectors. It was soon recognised that these finite-dimensional Wishart processes are *affine*, i.e., Markov processes whose Laplace transform depends in an exponentially affine way on the initial value.

Wishart processes are popular because they provide tractable stochastic covariance models. Indeed, one important application of finite-dimensional Wishart processes is multivariate asset price modelling with stochastic covariances. The fact that certain models for bond and commodity markets call for *infinite-dimensional stochastic covariance models* inspired us to consider the infinite-dimensional analogue. Moreover, so-called matrix-valued Volterra-Wishart processes can be interpreted as infinite-dimensional Wishart processes by an appropriate lift.

The difficulty we face when studying infinite-dimensional Wishart processes is that there are strong indications from the finite-dimensional theory that in the presence of a non-degenerate diffusion part, i.e., *when $Q$ in (1) is of infinite rank, then an infinite-dimensional Wishart process is necessarily of finite rank almost everywhere*. To explain this statement, let us return for a moment to the finite-dimensional setting: a Wishart process $X$ taking values in $S^+(\mathbb{R}^n)$ is a process satisfying

$$(2) \quad dX_t = (\alpha Q + A X_t + X_t A^*) \, dt + \sqrt{X_t} dW_t \sqrt{Q} + \sqrt{Q} \, dW_t^* \sqrt{X_t}, \quad t \geqslant 0,$$

where $A \in \mathbb{R}^{n \times n}$, $Q \in S^+(\mathbb{R}^n)$, and $W$ is a standard $\mathbb{R}^{n \times n}$-valued Brownian motion. It is well-known (see, e.g., [2, 4, 5, 7, 8]) that if $Q$ is injective, then such a finite-dimensional Wishart process exists if and only if either $\alpha \in [n - 1, \infty)$, or $\alpha \in \{0, \ldots, n - 2\}$ and $\mathrm{rank}(x_0) \leqslant \alpha$. In case of the latter one has $\mathrm{rank}(X_t) \leqslant \alpha$ a.s. for all $t \geqslant 0$. When translated to the infinite-dimensional setting this suggests that Wishart processes of infinite rank are hard to come by. Indeed, we prove the following:

**Theorem.** Assume that $H$ is infinite-dimensional. If $Q$ is of trace class and injective and $A$ is bounded then an analytically and probabilistically weak solution to (1) exists if and only if $\alpha \in \mathbb{N}$ and $\mathrm{rank}(x_0) \leqslant \alpha$. In this case, $\mathrm{rank}(X_t) = \alpha$ a.s. for almost all $t > 0$.

In fact, our results go beyond the realm of the Theorem above. In general, we only assume that $A \colon D(A) \subset H \to H$ is the generator of a $C_0$-semigroup $(e^{tA})_{t \geqslant 0}$ and $Q$ a bounded positive self-adjoint operator on $H$ satisfying $\int_0^t \| e^{sA} \sqrt{Q} \|_{L_2(H)} \, ds < \infty$ for all $t > 0$. For this setting we have the following results:

(1) If $\alpha \in \mathbb{N}$ and $\mathrm{rank}(x_0) \leqslant \alpha$, then there exists a probabilistically and analytically weak solution $X$ to (1). By construction, this solution is necessarily of rank at most $\alpha$.

(2) Formulas for the Fourier-Laplace transform (below, tr denotes the trace)

$$\mathbb{E}\left[\exp(-\operatorname{tr}((u - iv)X_t)) \mid x_0\right]$$

of a Wishart process $X$ for
   (a) $u \in S^+(H)$ (the positive self-adjoint operators) and $v = 0$;
   (b) $v \in S^+(H)$ or $-v \in S^+(H)$ and $u = 0$;
   (c) $u \in S^+(H)$, and $v \in S(H)$ (the self-adjoint operators) and $u$, $v$, $Q$, $A$, and $x_0$ all jointly diagonalizable;
   (d) $\alpha \in \mathbb{N}$, $u \in S(H)$, $v \in S(H)$, and $t$ sufficiently small.
   In all cases the Fourier-Laplace transform is of exponential affine form, i.e.,

$$\mathbb{E}\left[\exp(-\operatorname{tr}((u - iv)X_t)) \mid x_0\right] = \exp(-\operatorname{tr}(\psi(t, u - iv)x_0) - \phi(t, u - iv)), \qquad t \geqslant 0,$$

where $\psi$ and $\phi$ are solutions of operator valued Riccati equations that can be solved explicitly in all the cases listed above. As a consequence we obtain that an infinite-dimensional Wishart process is an affine process satisfying the Markov property and is thus unique in law.

(3) If $Q$ is injective and if there exists a $t > 0$ such that $e^{tA}$ is injective, then the existence of a probabilistically and analytically weak solution $X$ to (1) *implies* that $\alpha \in \mathbb{N}$, $\mathrm{rank}(x_0) \leqslant \alpha$ and $\mathrm{rank}(X_t) = \alpha$ a.s. for almost all $t > 0$.

(4) If there exists a probabilistically and analytically weak solution $X$ to (1), then either $\mathrm{rank}(X_t) \geqslant \mathrm{rank}(Q)$ a.s. for almost all $t > 0$, or $\alpha \in \mathbb{N}$ and $\mathrm{rank}(X_t) = \alpha$ a.s. for almost all $t > 0$. This provides new insights even in the finite-dimensional setting, since a characterisation of Wishart processes

in $S(\mathbb{R}^n)$ that are of rank at most $k$ is only known when $A \equiv 0$ and $Q = \mathrm{id}_{\mathbb{R}^n}$ see [5, Theorem 3.10]. Moreover, in the infinite-dimensional setting under the condition that $\mathrm{rank}(Q) = \infty$, this result implies that finite-rank Wishart processes exist if and only if $\alpha \in \mathbb{N}$.

(5) If $e^{tA}$ is injective for all $t \geqslant 0$, then a probabilistically and analytically weak solution to (1) is Feller with respect to (a minor refinement of) the weak-$*$-topology on the space of self-adjoint trace class operators. We make use of the Feller property to prove the existence of unique limit distribution of the constructed Wishart process.

Our results also give rise to various intriguing questions. Firstly, we were not able to rule out the existence of a Wishart process when $Q$ is injective, $\alpha \notin \mathbb{N}$, and $e^{tA}$ is not injective for all $t > 0$:

**Open probem 1.** Let $Q$ in (1) be injective. Does there exist an unbounded operator $A$ such that (1) allows for a solution for some $\alpha \notin \mathbb{N}$?

Secondly, we have little insight (even for $H = \mathbb{R}^n$) of existence of Wishart processes when $Q$ is not injective and $\alpha \notin \mathbb{N}$:

**Open probem 2.** If $Q$ is not injective, for what $\alpha \in \mathbb{R} \setminus \mathbb{N}$ and what $x_0 \in S_1^+(H)$ does a solution to (1) exits? Which role does the operator $A$ play?

## References

[1] M.-F. Bru, *Diffusions of perturbed principal component analysis*, J. Multivariate Anal. **29**(1):127–136, 1989.

[2] M.-F. Bru, *Wishart processes*, Journal of Theoretical Probability **4**(4):725–751, 1991.

[3] S. Cox, Ch. Cuchiero, A. Khedher, *Infinite-dimensional Wishart processes*, arXiv:2304.03490.

[4] C. Cuchiero, D. Filipović, E. Mayerhofer, and J. Teichmann, *Affine processes on positive semidefinite matrices*, The Annals of Applied Probability **21**(2):397–463, 2011.

[5] P. Graczyk, J. Małecki, and E. Mayerhofer, *A characterization of Wishart processes and Wishart distributions*, Stochastic Processes and their Applications **128**(4):1386–1404, 2018.

[6] G. Letac and H. Massam, *The noncentral Wishart as an exponential family, and its moments*, Journal of Multivariate Analysis **99**(7):1393–1417, 2008.

[7] G. Letac and H. Massam, *The Laplace transform* $(\det s)^{-p} \exp \mathrm{tr}(s^{-1}w)$ *and the existence of non-central Wishart distributions*, J. Multivariate Anal. **163**:96–110, 2018.

[8] E. Mayerhofer, *On Wishart and noncentral Wishart distributions on symmetric cones*, Transactions of the American Mathematical Society **371**(10):7093–7109, 2019.

## HJ-sampler: a Bayesian sampler for inverse problems of a stochastic process by leveraging Hamilton–Jacobi PDEs and score-based generative models

TINGWEI MENG[1]

(joint work with Zongren Zou[1], Jérôme Darbon and George Em Karniadakis)

The interplay between stochastic processes and optimal control has been extensively explored in the literature. With the recent surge in the use of diffusion models, stochastic processes have increasingly been applied to sample generation. This paper builds on the log transform, known as the Cole-Hopf transform in Brownian motion contexts, and extends it within a more abstract framework that includes a linear operator. Within this framework, we found that the well-known relationship between the Cole-Hopf transform and optimal transport is a particular instance where the linear operator acts as the infinitesimal generator of a stochastic process. We also introduce a novel scenario where the linear operator is the adjoint of the generator, linking to Bayesian inference under specific initial and terminal conditions. Leveraging this theoretical foundation, we develop a new algorithm, named the HJ-sampler, for Bayesian inference for the inverse problem of a stochastic differential equation with given terminal observations. The HJ-sampler involves two stages: solving viscous Hamilton-Jacobi (HJ) partial differential equations (PDEs) and sampling from the associated stochastic optimal control problem. Our proposed algorithm naturally allows for flexibility in selecting the numerical solver for viscous HJ PDEs. We introduce two variants of the solver: the Riccati-HJ-sampler, based on the Riccati method, and the SGM-HJ-sampler, which utilizes diffusion models. Numerical examples demonstrate the effectiveness of our proposed methods.

## B-stability of geodesic integrators on Riemannian manifolds

BRYNJULF OWREN

(joint work with Elena Celledoni and Ergys Çokaj)

Nonlinear stability of numerical methods for differential equations in Euclidean spaces is often realised in the form of non-expansiveness. The idea is to consider the behaviour of a numerical approximation method when it is applied to a vector field $X$ which satisfies a monotonicity condition

$$(1) \qquad \langle X(y) - X(z), y - z \rangle \leqslant \nu |y - z|^2$$

for any pair $y, z$ belonging to the space $V$ where $\nu \leqslant 0$ is a monotonicity constant, and where $\langle \cdot, \cdot \rangle$ is an inner product and $|\cdot|$ the induced norm. Flows, $\varphi_{t,X}$ of vector fields $X$ satisfying (1) yield the bound

$$|\varphi_{t,X}(y) - \varphi_{t,X}(z)| \leqslant e^{\nu t}|y - z|$$

---

[1]These authors contributed equally to this work.

and are thus non-expansive when $\nu \leqslant 0$. A one-step integrator $\phi_{h,X}$ is called B-stable if

$$(2) \qquad |\phi_{h,X}(y) - \phi_{h,X}(z)| \leqslant |y - z|, \quad \forall y, z, \; h > 0.$$

whenever (1) holds with $\nu \leqslant 0$. In the literature, see e.g. [1, 2], conditions for B-stability of Runge–Kutta methods can be found.

These definitions can be extended to Riemannian manifolds. A manifold $M$ with a metric $g$, induces a Riemannian distance $d(y, z)$ between points on $M$. A straightforward generalisation of B-stability is to demand that, given any pair of points $y, z$ in some appropriately chosen subset of $M$, whenever the exact flow, $\varphi_{t,X}$ of a vector field $X$ satisfies

$$d(\varphi_{t,X}(y), \varphi_{t,X}(z)) \leqslant d(y, z), \quad \forall t > 0,$$

a B-stable numerical method $\phi_{h,X}$ should verify

$$d(\phi_{h,X}(y), \phi_{h,X}(z)) \leqslant d(y, z), \quad \forall h > 0.$$

We need a monotonicity condition which generalises (1) to $M$. In order to do this we use the Levi-Civita connection $\nabla_X Y$ associated to the metric $g$, and writing $\langle u, v \rangle$ for $g(u, v)$ we introduce the condition

$$(3) \qquad \langle \nabla_Y X, Y \rangle \leqslant \nu |Y|^2, \quad \forall Y \in \mathcal{X}(\mathcal{U}), \quad \mathcal{U} \subseteq M,$$

where $\mathcal{U}$ is a suitably chosen open subset of $M$ to be made precise later. See also [4]. We recall from [3] that the flow of a vector field $X$ satisfying (3) has the property that

$$d(\varphi_{t,X}(y), \varphi_{t,X}(z)) \leqslant e^{\nu t} d(y, z)$$

Geodesic integrators. On a Riemannian manifold, we can use geodesics to define numerical integrators. Examples of these are

**GIE.:** The geodesic implicit Euler method, defined implicitly as $\phi_{IE} : y \mapsto y_1$ where

$$\exp_{y_1}(-hX(y_1)) = y$$

**GIMP.:** The Geodesic implicit midpoint rule defined in terms of a mid point $\bar{y}$ as $y \mapsto y_1$ through

$$y = \exp_{\bar{y}}\left(-\frac{h}{2}X(\bar{y})\right), \qquad y_1 = \exp_{\bar{y}}\left(\frac{h}{2}X(\bar{y})\right)$$

Stability and uniqueness results. The following result was proved in [3]: *The GIE method is B-stable if the manifold $(M, g)$ has non-negative sectional curvature.* It is of interest to understand the behaviour of such geodesic integrators also on positively curved spaces. For this, a number of numerical experiments were conducted on the 2-sphere $S^2$. One can then find instances of non-expansive vector fields ($\nu \leqslant 0$), such as Killing fields, and initial points $y, z \in S^2$ which cause the GIE method to exhibit an expansive behaviour. The following plot from [3] shows this, since for small stepsizes, $h$, clearly $d(\phi_{h,X}(y), \phi_{h,X}(z)) > d(y, z)$. For the exact flow, the distance between $y(t)$ and $z(t)$ remains constant.

FIGURE 1. The Riemannian distance between to initial points $y, z \in S^2$ for increasing values of the step size $h$ in the GIE method applied to a Killing vector field.

Another issue with positively curved manifolds, such as $S^2$, is the non-uniqueness of the solution to the equations which arise due to the implicitness of the scheme. For the same Killing vector field used above, the GIE method exhibits multiple solutions. It is interesting to observe that in the Euclidean case, non-expansivity ensures unique solutions for the implicit Euler method, see [2, Theorem 14.4]. The solutions can be depicted in a bifurcation plot that is borrowed from [3].



FIGURE 2. The solutions to the third cartesian coordinate for a range of step sizes $h$ when the GIE method is applied to a Killing vector field.

Some technical assumptions hitherto omitted. As indicated above, one needs to constrain the B-stability to some subset $\mathcal{U}$ of the manifold $M$, defined in terms of the vector field $X$. We here list some conditions that must be imposed on such a set:

- $\mathcal{U}$ is geodesically convex
- The vector field $X$ is *forward complete* on $\mathcal{U}$, meaning that $\varphi_{t,X}(y)$ exists for every $y \in \mathcal{U}$ and $t > 0$
- The vector field $X$ is *forward invariant* on $\mathcal{U}$, meaning that $\varphi_{t,X}(y) \in \mathcal{U}$ for every $y \in \mathcal{U}$ and $t > 0$
- The numerical solution is also defined and invariant on $\mathcal{U}$ for all $h > 0$

Conclusion and further work. We have introduced the notion of B-stability on Riemannian manifolds. This is an unconditional stability definition, meaning that no upper bound on the stepsize $h > 0$ is imposed. The main result is that the Geodesic Implicit Euler method is B-stable on manifolds $(M, g)$ of non-positive sectional curvature. We also show through numerical experiments that B-stability does not seem to hold in general for positively curved spaces, here demonstrated on the compact manifold $S^2$. Future work will aim at considering explicit integrators and conditional stability, and we intend to apply the theory in the design of neural networks with manifold valued features and parameters.

#### References

[1] J.C. Butcher, *A stability property of implicit Runge-Kutta methods*, BIT **15** (1975), 358–361.
[2] E. Hairer and G. Wanner, *Solving Ordinary Differential Equations II, Stiff and Differential-Algebraic Problems*, Second Revised Edition, Springer-Verlag, 1996.
[3] M. Arnold and E. Celledoni and E. Çokaj and B. Owren and D. Tumiotto, *B-stability of numerical integrators on Riemannian manifolds*, J. Comput. Dyn. **11** (2024) 92–107
[4] J. W. Simpson-Porco and F. Bullo, *Contraction theory on Riemannian manifolds*, Systems & Control Letters **65**, (2014), 74–80

## A dynamical systems view to deep learning: contractivity and structure preservation

ELENA CELLEDONI

(joint work with Davide Murari and Brynjulf Owren)

### 1. INTRODUCTION

Deep learning neural networks have recently been interpreted as discretisations of an optimal control problem subject to an ordinary differential equation constraint. There is a growing effort to mathematically understand the structure in existing deep learning methods and to design new approaches preserving (geometric) structure in neural networks. The (discrete) optimal control point of view to neural networks offers an interpretation of deep learning from a numerical analysis perspective and opens the way to mathematical insight [9, 8, 1]. We discuss a number of interesting directions of current and future research in structure preserving deep learning [2].

## 2. STABLE NETWORKS

Some deep neural networks can be designed to have desirable properties such as invertibility and group equivariance or can be adapted to problems of manifold value data. Equivariant neural networks are effective in reducing the amount of data for solving certain imaging problems [3].

We show how classical results of stability of ODEs are useful to construct contractive neural networks architectures. Thus, neural networks can be designed with guaranteed stability properties. This can be used to ensure robustness against adversarial attacks and to obtain converging "Plug-and-Play" algorithms for inverse problems in imaging [2, 7, 12].

We consider vector fields of the type

$$X(t, y(t)) = -A(t)^T \sigma(A(t)y(t) + b(t)),$$

with $\sigma$ increasing activation function[1], $A \in \mathbb{R}^{n \times k}$, $b \in \mathbb{R}^n$, and use the forward Euler discretizations of the corresponding differential equations to build the neural network. These vector fields are contractive in the $\ell^2$ norm [2], which means that there is $\nu < 0$ such that for all $y_1, y_2 \in \mathbb{R}^n$ and $t \in [0, T]$:

$$\langle X(t, y_2) - X(t, y_1), y_2 - y_1 \rangle \leqslant \nu \|y_2 - y_1\|^2.$$

As a consequence for any two integral curves solutions of $\dot{y} = X(t, y)$, $y(0) = y_0$ and $\dot{z} = X(t, z)$, $z(0) = z_0$ we have that

$$\|y(t) - z(t)\| \leqslant e^{t\nu} \|y(0) - z(0)\|.$$

For the vector field $X(t, y(t)) = -A(t)^T \sigma(A(t)y(t) + b(t))$, one can also prove that the following, stronger monotonicity condition

$$\langle X(t, y_2) - X(t, y_1), y_2 - y_1 \rangle \leqslant \bar{\nu} \|X(t, y_2) - X(t, y_1)\|^2, \quad \bar{\nu} < 0$$

holds with $\bar{\nu} = -\frac{1}{\|A\|^2 L}$ where $L$ is the Lipschitz constant of the activation function $\sigma$ [12]. The latter condition is essential to guarantee also contractivity of the forward Euler method for small enough step-sizes.

Stable networks can be constructed composing layers of contractive and expansive vector fields and using the above ideas to restrict the Lipschitz constant of the overall network.



Accuracy versus size of the perturbation for plain ResNet and stable neural networks, [7].

---

[1]For $x \in \mathbb{R}^n$, $\sigma(x)$ is a scalar function applied to each component of the vector $x$.

## 3. Optimal shape parametrisation and learning Hamiltonian systems from data

Shape analysis is a framework for treating complex data and obtaining metrics on spaces of data. Examples are spaces of unparametrized curves, time-signals, surfaces and images [10, 4].

A computationally demanding task for estimating distances between shapes, e.g. in object recognition, is the computation of optimal reparametrizations. This is an optimisation problem on the infinite dimensional group of orientation preserving diffeomorphisms $\mathrm{Diff}^+(\Omega)$, for some domain $\Omega$. The optimisation problem takes the form

$$\inf_{\varphi \in \mathrm{Diff}^+(\mathrm{I})} E(\varphi).$$

We approximate the diffeomorphisms with neural network parametrizations (with a finite number of parameters) and where each layer is a diffeomorphism, and we solve a finite dimensional optimisation problem for the parameters on a finite dimensional space [5]. It is useful to consider geometric properties in this context e.g. reparametrization invariance of the distance function, invertibility and contractivity of the neural networks.

We can show that finite compositions of diffeomorphisms of the type

$$\varphi_\ell = \mathrm{id} + f_\ell, \qquad \ell = 1, \dots, L$$

with $f_\ell$ a 1-Lipschitz vector field can be used to describe the whole group of diffeomorphisms fixing the boundary of $\Omega$ a compact, convex subset of $\mathbb{R}^d$; and of diffeomorphisms on a cube $\Omega = [0,1]^d$.

We also consider applications of deep learning to mechanical systems, for learning Hamiltonians on manifolds and from noisy data [6, 11].

## References

[1] M. Benning, E. Celledoni, M. J. Ehrhardt, B. Owren, and C. B. Schönlieb, *Deep learning as optimal control problems: models and numerical methods.* Journal of Computational Dynamics **6**(2):171–198, 2019.

[2] E. Celledoni, M. J. Ehrhardt, C. Etmann, R. I. McLachlan, B. Owren, C. B. Schönlieb, F. Sherry, *Structure preserving deep learning*, European Journal of Applied Mathematics (2021).

[3] E. Celledoni, M. J. Ehrhardt, C. Etmann, B. Owren, C. B. Schönlieb, F. Sherry, *Equivariant neural networks for inverse problems*, Inverse Problems (2021).

[4] E. Celledoni, M. Eslitzbichler, and A. Schmeding, *Shape analysis on Lie groups with applications in computer animation*, J. Geom. Mech. **8**(3):273–304, 2016.

[5] E. Celledoni, H. Glöckner, J. Riseth, A. Schmeding, *Deep neural networks on diffeomorphism groups for optimal shape reparameterization*, BIT Numerical Mathematics **63** (4), 1-38.

[6] E. Celledoni, A. Leone, D. Murari, and B. Owren, *Learning Hamiltonians of constrained mechanical systems*, J. Comput. Appl. Math. **417**, Paper No. 114608, 12 pp.

[7] E. Celledoni, D. Murari, B. Owren, C.-B. Schönlieb, and F. Sherry, *Dynamical systems' based neural networks*, SISC, 2023.

[8] W. E, *A Proposal on Machine Learning via Dynamical Systems*, Commun. Math. Stat. **5**, 1–11, 2017.

[9] E. Haber and L. Ruthotto, *Stable architectures for Deep Neural Networks*, Inverse Problems **34** (1), 2017.

[10] S. Kurtek, E. Klassen, J. C. Gore, Z. Ding, A. Srivastava, *Elastic geodesic paths in shape space of parameterized surfaces*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2012.

[11] E. Celledoni, S. Eidnes, H. Noren Myhr, *Learning Dynamical Systems from Noisy Data with Inverse-Explicit Integrators*, arXiv:2306.03548, 2023.

[12] F. Sherry, E. Celledoni, M. J. Ehrhardt, D. Murari, B. Owren, and C.-B. Schönlieb, *Designing Stable Neural Networks using Convex Analysis and ODEs*, Physica D: Nonlinear Phenomena **463**, 2024.

## Reservoir kernels and Volterra series

LYUDMILA GRIGORYEVA

(joint work with Lukas Gonon and Juan-Pablo Ortega)

**Reservoir computing.** In this talk we consider input/output systems determined by SSSs. The symbols $\mathcal{Z}$ and $\mathcal{Y}$ denote the *input* and the *output spaces*, respectively, and $\mathcal{X}$ is the *state space* (typically finite or infinite-dimensional manifolds). A *discrete-time SSS* is determined by:

$$(1) \qquad \mathbf{x}_t = F(\mathbf{x}_{t-1}, \mathbf{z}_t), \quad \mathbf{y}_t = h(\mathbf{x}_t), \quad t \in \mathbb{Z},$$

with $\mathbf{z} \in \mathcal{Z}^{\mathbb{Z}}$, $\mathbf{y} \in \mathcal{Y}^{\mathbb{Z}}$, $\mathbf{x} \in \mathcal{X}^{\mathbb{Z}}$. The map $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$ is called the *state map* and $h : \mathcal{X} \to \mathcal{Y}$ the *readout/observation* map. Cases in which $F$ is randomly generated and $h$ is functionally simple are known as *reservoir computing (RC)* ([6, 7, 4, 5]; a popular family is *echo state networks*).

We focus on (1) that determine an *input/output* system, which happens in the presence of the *echo state property (ESP)*: when for any $\mathbf{z} \in \mathcal{Z}^{\mathbb{Z}}$ there exists a unique $\mathbf{y} \in \mathcal{Y}^{\mathbb{Z}}$ s.t. (1) holds. In that case, the *state-space (SS) filter* $U_h^F : \mathcal{Z}^{\mathbb{Z}} \to \mathcal{Y}^{\mathbb{Z}}$ is defined by $U_h^F(\mathbf{z}) := \mathbf{y}$, with $\mathbf{z} \in \mathcal{Z}^{\mathbb{Z}}$, $\mathbf{y} \in \mathcal{Y}^{\mathbb{Z}}$ linked by (1) via the ESP. If the ESP holds at the level of (1), we can define a *state* or *reservoir filter* $U^F : \mathcal{Z}^{\mathbb{Z}} \to \mathcal{X}^{\mathbb{Z}}$ and, in that case, we have that $U_h^F := h \circ U^F$. The SS filters are causal (C) and time-invariant (TI) [1] and it suffices to work with their restriction $U_h^F : \mathcal{Z}^{\mathbb{Z}_-} \to \mathcal{Y}^{\mathbb{Z}_-}$. Moreover, $U_h^F$ determines an SS (reservoir) *functional* $H_h^F : \mathcal{Z}^{\mathbb{Z}_-} \to \mathcal{Y}$ as $H_h^F(\mathbf{z}) := U_h^F(\mathbf{z})_0$, $\forall \mathbf{z} \in \mathcal{Z}^{\mathbb{Z}_-}$. The same holds for $U^F$ and $H^F$ with the ESP at the level of the state equation. If $\mathcal{Z}$ is a compact metric space and $\mathcal{Z}^{\mathbb{Z}_-}$ is endowed with the product topology we say that $H_h^F$ or $H^F$ have the *fading memory property (FMP)* when they are continuous [2].

**Reservoirs with a linear readout and the RKHS associated to a state system.** This talk is based on [3] in which RKHS was associated to any ESP SSS $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$, using $H^F : \mathcal{Z}^{\mathbb{Z}_-} \to \mathcal{X}$ as a feature map. This allows, using the Represeter Theorem, to reduce the search for a linear $h$ which is optimal with respect to the regularized empirical risk minimization (ERM) associated to any loss to the search for $h$ defined on the linear span of reachable states $\overline{\mathcal{X}_R} := \mathrm{span}\{\mathcal{X}_R\} = \mathrm{span}\{H^F(\mathbf{z}) \mid \mathbf{z} \in \mathcal{Z}^{\mathbb{Z}_-}\}$. Let $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$ be a state map s.t. the

pair $(\mathcal{X}, \langle \cdot, \cdot \rangle_{\mathcal{X}})$ is a Hilbert space and $F$ has the ESP. Let $H^F : \mathcal{Z}^{\mathbb{Z}_-} \to \mathcal{X}$ be the corresponding state functional. Define the *reservoir kernel map*

$$(2) \qquad \begin{array}{rccc} K : & \mathcal{Z}^{\mathbb{Z}_-} \times \mathcal{Z}^{\mathbb{Z}_-} & \to & \mathbb{R} \\ & (\mathbf{z}, \mathbf{z}') & \longmapsto & \langle H^F(\mathbf{z}), H^F(\mathbf{z}') \rangle_{\mathcal{X}}. \end{array}$$

The reservoir kernel $K$ is symmetric and positive semidefinite. Let $(\mathbb{H}, \langle \cdot, \cdot \rangle_{\mathbb{H}})$ be the corresponding RKHS given by

$$(3) \qquad \mathbb{H} := \overline{\mathrm{span}\{K_{\mathbf{z}} := K(\mathbf{z}, \cdot) : \mathcal{Z}^{\mathbb{Z}_-} \to \mathbb{R} \mid \mathbf{z} \in \mathcal{Z}^{\mathbb{Z}_-}\}}.$$

**Proposition 1.** [3] *Let* $(\mathcal{X}, \langle \cdot, \cdot \rangle_{\mathcal{X}})$ *be a finite-dimensional Hilbert space and let* $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$ *be a state map that satisfies the ESP. Let* $(\mathbb{H}, \langle \cdot, \cdot \rangle_{\mathbb{H}})$ *be the associated RKHS in* (3)*. Then*

$$(4) \qquad \mathbb{H} = \{ \langle \mathbf{W}, H^F(\cdot) \rangle_{\mathcal{X}} \mid \mathbf{W} \in \overline{\mathcal{X}_R} \}.$$

*Moreover, for any* $\mathbf{W}_1, \mathbf{W}_2 \in \overline{\mathcal{X}_R}$, $\langle \langle \mathbf{W}_1, H^F(\cdot) \rangle_{\mathcal{X}}, \langle \mathbf{W}_2, H^F(\cdot) \rangle_{\mathcal{X}} \rangle_{\mathbb{H}} = \langle \mathbf{W}_1, \mathbf{W}_2 \rangle_{\mathcal{X}}$, *and the map* $\Psi : (\overline{\mathcal{X}_R}, \langle \cdot, \cdot \rangle_{\mathcal{X}}) \to (\mathbb{H}, \langle \cdot, \cdot \rangle_{\mathbb{H}})$, $\mathbf{W} \mapsto \langle \mathbf{W}, H^F(\cdot) \rangle_{\mathcal{X}} =: H^F_{\mathbf{W}}(\cdot)$ *is an isometric isomorphism.*

For infinite-dimensional state-space representations, the RKHS $\mathbb{H}$ is infinite-dimensional and (4) $\{ \langle \mathbf{W}, H^F(\cdot) \rangle_{\mathcal{X}} \mid \mathbf{W} \in \overline{\mathcal{X}_R} \} \subset \mathbb{H}$. Moreover, if $\mathcal{X}$ is infinite-dimensional, $\Psi$ is an injective isometry but ceases to be surjective in general.

**Estimation of the empirical risk minimizing readout.** Consider the ESP map $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$ and a finite sample $\{(\mathbf{Z}_{-i}, Y_{-i})\}_{i \in \{0, \dots, n-1\}}$ of input/output observations. For $i \in \{0, \dots, n-1\}$ define the truncated input samples $\mathbf{Z}_{-i}^{-n+1} := (\dots, \mathbf{0}, \mathbf{0}, \mathbf{Z}_{-n+1}, \dots, \mathbf{Z}_{-i-1}, \mathbf{Z}_{-i})$, and the *empirical risk* $\widehat{R}_n(H^F_{\mathbf{W}})$ associated to the loss $L : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ for the system $H^F_{\mathbf{W}}(\cdot) = \langle \mathbf{W}, H^F(\cdot) \rangle_{\mathcal{X}}$ with readout $\mathbf{W} \in \mathcal{X}$ as $\widehat{R}_n(H^F_{\mathbf{W}}) = \frac{1}{n} \sum_{i=0}^{n-1} L(\langle \mathbf{W}, H^F(\mathbf{Z}_{-i}^{-n+1}) \rangle_{\mathcal{X}}, Y_{-i})$.

**Proposition 2** (Implicit reduction and kernelization). *Let* $F : \mathcal{X} \times \mathcal{Z} \to \mathcal{X}$ *be an ESP state system and let* $L : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ *be a loss function. Let* $\{(\mathbf{Z}_{-i}, Y_{-i})\}_{i \in \{0, \dots, n-1\}}$ *be a sample and let* $\Omega : (0, \infty) \to \mathbb{R}$ *be a strictly increasing function. Then, the regularized ERM problem admits the following reformulations:*

- **Implicit reduction:** *Let* $\mathbb{H}$ *be the RKHS in* (3)*. Then,*

$$\min_{\mathbf{W} \in \mathcal{X}} \{ \widehat{R}_n(H^F_{\mathbf{W}}) + \Omega(\|\mathbf{W}\|_{\mathcal{X}}^2) \} = \min_{\mathbf{W} \in \overline{\mathcal{X}_R}} \{ \widehat{R}_n(H^F_{\mathbf{W}}) + \Omega(\|\mathbf{W}\|_{\mathcal{X}}^2) \}$$

$$(5) \qquad \qquad \qquad = \min_{H^F_{\mathbf{W}} \in \mathbb{H}} \{ \widehat{R}_n(H^F_{\mathbf{W}}) + \Omega(\|H^{\mathbf{W}}_F\|_{\mathbb{H}}^2) \}.$$

*If* $\mathcal{X}$ *is finite-dimensional, these minima coincide with* $\min_{f \in \mathbb{H}} \{ \widehat{R}_n(f) + \Omega(\|f\|_{\mathbb{H}}^2) \}$.

- **Kernelization:** *The problem* (5) *in* $\mathcal{X}$ *can be written in terms of the Gramian* $\mathcal{K} \in \mathbb{M}_n$, *defined as* $\mathcal{K}_{i,j} = K(\mathbf{Z}_{-(n-i)}^{-n+1}, \mathbf{Z}_{-(n-j)}^{-n+1})$, $i, j \in \{1, \dots, n\}$, *that is*

$$\min_{f \in \mathbb{H}} \{ \widehat{R}_n(f) + \Omega(\|f\|_{\mathbb{H}}^2) \} = \min_{\boldsymbol{\alpha} \in \mathbb{R}^n} \{ \frac{1}{n} \sum_{i=1}^{n} L((\mathcal{K}\boldsymbol{\alpha})_i, Y_{-(n-i)}) + \Omega(\boldsymbol{\alpha}^{\top} \mathcal{K} \boldsymbol{\alpha}) \}.$$

# 1. Main results

Consider $\mathcal{Z} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}^m$, and denote by $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ the Euclidean inner product and norm, as well as certain inner products and norms induced by it. For any $\mathbf{z} \in (\mathbb{R}^d)^{\mathbb{Z}_-}$ define the $p$-norms as $\|\mathbf{z}\|_p := (\sum_{t \in \mathbb{Z}_-} \|\mathbf{z}_t\|^p)^{1/p}$, for $1 \leqslant p < \infty$, $\|\mathbf{z}\|_\infty := \sup_{t \in \mathbb{Z}_-} \{\|\mathbf{z}_t\|\}$ for $p = \infty$. Given $M > 0$, $K_M := \{\mathbf{z} \in (\mathbb{R}^d)^{\mathbb{Z}_-} \mid \|\mathbf{z}_t\| \leqslant M$ for all $t \in \mathbb{Z}_-\}$. It is easy to see that $K_M = \overline{B_M} \subset \ell_-^\infty(\mathbb{R}^d)$, with $B_M := B_{\|\cdot\|_\infty}(\mathbf{0}, M)$ and $\ell_-^\infty(\mathbb{R}^d) := \{\mathbf{z} \in (\mathbb{R}^d)^{\mathbb{Z}_-} \mid \|\mathbf{z}\|_\infty < \infty\}$. We define $\widetilde{B}_M := B_M \cap \ell_-^1(\mathbb{R}^d)$, $\ell_-^1(\mathbb{R}^d) := \{\mathbf{z} \in (\mathbb{R}^d)^{\mathbb{Z}_-} \mid \|\mathbf{z}\|_1 < \infty\}$ and write $L(V, W)$ for the space of linear maps between the real vector spaces $V$ and $W$.

**Theorem 1** (The Volterra reservoir). *Let $M > 0$ and let $\tau > 0$ be s.t. $\tau^2 M^2 < 1$. Let $0 < \lambda < \sqrt{1 - \tau^2 M^2}$. Consider the state system with inputs in $K_M \subset \ell_-^\infty(\mathbb{R}^d)$ and states in $T_2(\mathbb{R}^d)$ given by the recursion*

$$(6) \qquad \mathbf{x}_t = \lambda \mathbf{x}_{t-1} \otimes \widetilde{\mathbf{z}}_t + 1, \quad t \in \mathbb{Z}_-,$$

*where $\widetilde{\mathbf{z}}_t \in \overline{T(\mathbb{R}^d)}$ is the $\tau$-tensorization of $\mathbf{z}_t$. Then,*

- *(i) This system has the ESP and defines a unique CTI and FMP filter $U_\lambda^{\mathrm{Volt}} : K_M \subset \ell_-^\infty(\mathbb{R}^d) \to \{\mathbf{x} \in (T_2(\mathbb{R}^d))^{\mathbb{Z}_-} \mid \|\mathbf{x}_t\| \leqslant L$ for all $t \in \mathbb{Z}_-\}$ for $t \in \mathbb{Z}_-$ is given by $U_\lambda^{\mathrm{Vlt}}(\mathbf{z})_t = 1 + \sum_{j=0}^\infty \lambda^{j+1} \widetilde{\mathbf{z}}_{t-j} \otimes \widetilde{\mathbf{z}}_{t-(j-1)} \otimes \cdots \otimes \widetilde{\mathbf{z}}_t$, with $L = 1/(1 - \lambda/\sqrt{1 - \tau^2 M^2})$.*
- *(ii) Let $U : K_M \subset \ell_-^\infty(\mathbb{R}^d) \to K_L \subset \ell_-^\infty(\mathbb{R}^m)$ be a CTI, FMP filter whose restriction $U|_{B_M}$ is analytic as a map between open sets in the Banach spaces $\ell_-^\infty(\mathbb{R}^d)$ and $\ell_-^\infty(\mathbb{R}^m)$, and $U(\mathbf{0}) = \mathbf{0}$. Then, there exists a unique map $\mathbf{W} \in L(T_2(\mathbb{R}^d), \mathbb{R}^m)$, s.t. $U(\mathbf{z})_t = \mathbf{W} U_\lambda^{\mathrm{Vlt}}(\mathbf{z})_t$, for any $\mathbf{z} \in \widetilde{B}_M$, $t \in \mathbb{Z}_-$.*

*We refer to $F_\lambda^{\mathrm{Vlt}} : T_2(\mathbb{R}^d) \times \{\mathbf{z} \in \mathbb{R}^d \mid \|\mathbf{z}\| \leqslant M\} \to T_2(\mathbb{R}^d)$ given by $F_\lambda^{\mathrm{Vlt}}(\mathbf{x}, \mathbf{z}) := \lambda \mathbf{x} \otimes \widetilde{\mathbf{z}} + 1$ as the Volterra reservoir map, to the filter $U_\lambda^{\mathrm{Vlt}}$ as the Volterra filter, and to the equality in (ii) as the Volterra filter representation of the FMP filter $U$. We call the corresponding functional $H_\lambda^{\mathrm{Vlt}}(\mathbf{z}) := U_\lambda^{\mathrm{Vlt}}(\mathbf{z})_0$ the Volterra functional.*

Let $\tau \in \mathbb{Z}_-$ and define the $\tau$-*time delay operator* $T_{-\tau} : (\mathbb{R}^d)^{\mathbb{Z}_-} \to (\mathbb{R}^d)^{\mathbb{Z}_-}$ by $T_{-\tau}(\mathbf{z})_t := \mathbf{z}_{t+\tau}$, for any $\mathbf{z} \in (\mathbb{R}^d)^{\mathbb{Z}_-}$, $t \in \mathbb{Z}_-$. For any $t \in \mathbb{Z}_-$, denote by $p_t : (\mathbb{R}^d)^{\mathbb{Z}_-} \to \mathbb{R}^d$ the $t$-*projection* given by $p_t(\mathbf{z}) = \mathbf{z}_t$. Finally, recall that a filter $U : (\mathbb{R}^d)^{\mathbb{Z}_-} \to (\mathbb{R}^m)^{\mathbb{Z}_-}$ is TI when $U \circ T_{-\tau} = T_{-\tau} \circ U$, for any $\tau \in \mathbb{Z}_-$.

**Proposition 3** (The Volterra reservoir kernel). *Let $M > 0$, $\tau > 0$ be s.t. $\tau^2 M^2 < 1$. Let $0 < \lambda < \sqrt{1 - \tau^2 M^2}$. The reservoir kernel $K^{\mathrm{Vlt}} : K_M \times K_M \to \mathbb{R}$ of the Volterra reservoir in (6) defined using $\lambda$, $\tau$-tensorization, and with inputs in $K_M \subset \ell_-^\infty(\mathbb{R}^d)$ is well-defined and can be calculated for $\mathbf{z}, \mathbf{z}' \in K_M$ using*

$$(7) \qquad K^{\mathrm{Vlt}}(\mathbf{z}, \mathbf{z}') = 1 + \frac{\lambda^2 K^{\mathrm{Vlt}}(T_1(\mathbf{z}), T_1(\mathbf{z}'))}{1 - \tau^2 \langle p_0(\mathbf{z}), p_0(\mathbf{z}') \rangle}.$$

**The Volterra kernel recursion.** Expression (7) can be used to compute the Gram matrix of $K^{\mathrm{Vlt}} \in \mathbb{M}_n$ for a truncated input $\{\mathbf{Z}_{-i}\}_{i \in \{0, \ldots, n-1\}}$, defined as

$\mathcal{K}_{i,j}^{\mathrm{Vlt}} = K^{\mathrm{Vlt}}(\mathbf{Z}_{-(n-i)}^{-n+1}, \mathbf{Z}_{-(n-j)}^{-n+1})$, $i, j \in \{1, \ldots, n\}$. By setting $\mathcal{K}_{0,0}^{\mathrm{Vlt}} = \mathcal{K}_{i,0}^{\mathrm{Vlt}} = 1/(1 - \lambda^2)$, $\forall i \in \{1, \ldots, n\}$, for $j \in \{1, \ldots, i\}$ it holds

$$\mathcal{K}_{i,j}^{\mathrm{Vlt}} = 1 + \lambda^2 \frac{\mathcal{K}_{i-1,j-1}^{\mathrm{Vlt}}}{1 - \tau^2 \langle \mathbf{Z}_{-(n-i)}, \mathbf{Z}_{-(n-j)} \rangle}.$$

The Gram matrix can be completed by using its symmetry.

**Theorem 2** (Universality of the Volterra reservoir kernel). *Let $K^{\mathrm{Vlt}} : K_M \times K_M \to \mathbb{R}$ be the reservoir kernel map in Proposition 3 and let $K^{\mathrm{Vlt}}(K_M)$ be the associated space of kernel sections. Then*

$$K^{\mathrm{Vlt}}(K_M) = C^0(K_M),$$

*that is, the Volterra reservoir kernel is universal.*

REFERENCES

[1] L. Grigoryeva and J.-P. Ortega, *Echo state networks are universal*, Neural Networks **108**, pp. 495–508, 2018.
[2] L. Grigoryeva and J.-P. Ortega, *Differentiable reservoir computing*, Journal Of Machine Learning Research **20**, No. 179, pp. 1–62, 2019.
[3] L. Grigoryeva and J.-P. Ortega, *Dimension reduction in recurrent networks by canonicalization* Journal of Geometric Mechanics **13**, No. 4, pp. 647–677, 2021.
[4] H. Jaeger, *Short term memory in echo state networks*, Fraunhofer Institute for Autonomous Intelligent Systems, Technical Report, **152**, 2002.
[5] H. Jaeger and H. Haas, *Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication*, Science **304**, No. 5667, pp. 78–80, 2004.
[6] H. Jaeger, *The 'echo state' approach to analysing and training recurrent neural networks with an erratum note*, Tech. Rep., German National Research Center for Information Technology, 2010.
[7] W. Maass, T. Natschläger, and H. Markram, *Real-time computing without stable states: a new framework for neural computation based on perturbations*, Neural Computation **14**, pp. 2531–2560, 2002.

# Mean field limits for interacting particle systems: phase transitions, inference and control

GRIGORIOS A. PAVLIOTIS

We report on a recent research program on the study of interacting particle systems, described by systems of weakly interacting diffusions, and of their mean field limit. We consider $\{X_t^i\}_{i=1,\ldots,N} \subset \mathbb{R}^d$, the positions of $N$ indistinguishable interacting particles at time $t \geqslant 0$, satisfying the following system of SDEs:

$$(1) \quad \begin{cases} dX_t^i = -\nabla V(X_t^i) \, dxt - \dfrac{1}{N} \sum_{j=1}^{N} \nabla_1 W(X_t^i, X_t^j) \, dxt + \sqrt{2\beta^{-1}} \, dB_t^i, \\ \mathrm{Law}(X_0^1, \ldots, X_0^N) = \rho_{\mathrm{in}}^{\otimes N} \in \mathcal{P}_{2,\mathrm{sym}}((\mathbb{R}^d)^N), \end{cases}$$

where $V : \mathbb{R}^d \to \mathbb{R}$, $W : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, $\beta^{-1} > 0$ is the inverse temperature, $B_t^i, i = 1, \ldots, N$ are independent $d$-dimensional Brownian motions, and the initial position of the particles is i.i.d with law $\rho_{\mathrm{in}}$. We consider the dynamics both in the

whole space, as well as on the torus $\mathbb{T}^d$ with periodic boundary conditions. Under appropriate assumptions on the confining and interaction potentials, we can pass to the mean field limit to show that the empirical measure $\mu_t^N := \frac{1}{N} \sum_{j=1}^{N} \delta_{X_t^j}$ converges, in an appropriate sense, to $\rho_t(x)\,dx$ the solution of the McKean-Vlasov PDE

$$(2a) \qquad \frac{\partial \rho_t}{\partial t} = \nabla \cdot (\nabla V \rho_t) + \nabla \cdot ((\nabla W \star \rho_t)\rho_t) + \beta^{-1}\Delta\rho_t,$$

$$(2b) \qquad \rho_t(x,0) = \rho_{\text{in}}.$$

The mean field McKean SDE is

$$(3) \qquad dX_t = -\nabla V(X_t)\,dt - (\nabla W \star \rho_t)(X_t)\,dt + \sqrt{2\beta^{-1}}\,dW_t, \quad X_0 \sim \rho_{\text{in}}.$$

In a series of recent papers we studied phase transitions for the McKean-Vlasov PDE on the torus [4], the stability of multipeak steady states [1], $\Gamma-$convergence approaches to the study of propagation of chaos [3], multiscale problems and the noncommutativity between the diffusive and mean field limits [5, 9], mean field limits or interacting generalized Langevin particles [7], we developed spectral numerical methods for solving mean field PDEs with colored noise [8], explored the link between (the absence of) uniform propagation of chaos and (the absence of) uniform logarithmic Sobolev inequalities [6]. In addition, we have developed inference methodologies for the McKean SDE given observations of trajectories of the interacting particle system [11, 12, 13] and we have also studied the fully nonparametric problem of inferring the interaction potential from noisy measurements of the solution to the McKean-Vlasov PDE in a purely nonparametric framework [10].

Finally, in recent work we developed techniques for identifying all, stable and unstable, steady states of the McKean-Vlasov PDE and for stabilizing unstable steady states [2]. In particular, we developed an efficient numerical scheme for identifying all steady states (both stable and unstable) of the mean field McKean-Vlasov PDE, based on a spectral Galerkin approximation combined with a deflated Newton's method to handle the multiplicity of solutions of the Kirkwood-Monroe integral equation

$$(4) \qquad \rho_\infty = \frac{1}{Z}e^{-\beta(V+W*\rho_\infty)}, \quad Z = \int_{\mathbb{T}^d} e^{-\beta(V+W*\rho_\infty)}\,dx.$$

Having found all possible stationary states, we then formulate an optimal control strategy for steering the dynamics towards a chosen unstable steady state. The control is computed using iterated open-loop solvers in a receding horizon fashion. Our proposed methodology was applied to several examples, including the noisy Hegselmann-Krause model for opinion dynamics and the Haken-Kelso-Bunz model from biophysics. For the Heglselmann-Krause model, it is well known that it exhibits a discontinuous phase transition: at sufficiently high temperatures, the uniform distribution, describing non-consensus, is the unique stationary state of the mean field systems. A discontinuous phase transition occurs at at critical temperature, below which a localized state, describing consensus formation, becomes a stable steady state.

The proposed inference and computational framework opens up new possibilities for the calibration of mean field models and for understanding and controlling the collective behavior of noise-driven interacting particle systems, with potential applications in various fields such as social dynamics, biological synchronization, and collective behavior in physical and social systems.

## References

[1] B. Bertoli, B. D. Goddard, and G. A. Pavliotis, *Stability of stationary states for mean field models with multichromatic interaction potentials*, 2024.

[2] S. Bicego, D. Kalise, and G. A. Pavliotis, *Computation and control of unstable steady states for mean field multiagent systems*, 2024.

[3] J. A. Carrillo, M. G. Delgadino, and G. A. Pavliotis, *A λ-convexity based proof for the propagation of chaos for weakly interacting stochastic particles*, J. Funct. Anal. **279**(10):108734, 2020.

[4] J. A. Carrillo, R. S. Gvalani, G. A. Pavliotis, and A. Schlichting, *Long-time behaviour and phase transitions for the McKean-Vlasov equation on the torus*, Arch. Ration. Mech. Anal. **235**(1):635–690, 2020.

[5] M. G. Delgadino, R. S. Gvalani, and G. A. Pavliotis, *On the diffusive-mean field limit for weakly interacting diffusions exhibiting phase transitions*, Arch. Ration. Mech. Anal. **241**(1):91–148, 2021.

[6] M. G. Delgadino, R. S. Gvalani, G. A. Pavliotis, and S. A. Smith, *Phase transitions, logarithmic Sobolev inequalities, and uniform-in-time propagation of chaos for weakly interacting diffusions*, Comm. Math. Phys. **401**(1):275–323, 2023.

[7] MH Duong and GA Pavliotis, *Mean field limits for non-Markovian interacting particles: Convergence to equilibrium, generic formalism, asymptotic limits and phase transitions*, Communications in Mathematical Sciences **16**:2199–2230, 2018.

[8] S. N. Gomes, G. A. Pavliotis, and U. Vaes *Mean field limits for interacting diffusions with colored noise: phase transitions and spectral numerical methods*, Multiscale Model. Simul. **18**(3):1343–1370, 2020.

[9] S.N. Gomes and G.A. Pavliotis, *Mean field limits for interacting diffusions in a two-scale potential*, J. Nonlin. Sci. **28**(3):905–941, 2018.

[10] R. Nickl, G. A. Pavliotis, and K. Ray, *Bayesian nonparametric inference in McKean-Vlasov models*, 2024.

[11] G. A. Pavliotis and A. Zanoni, *Eigenfunction martingale estimators for interacting particle systems and their mean field limit*, SIAM J. Appl. Dyn. Syst. **21**(4):2338–2370, 2022.

[12] G. A. Pavliotis and A. Zanoni, *A method of moments estimator for interacting particle systems and their mean field limit*, SIAM/ASA J. Uncertain. Quantif. **12**(2):262–288, 2024.

[13] L. Sharrock, N. Kantas, P. Parpas, and G. A. Pavliotis, *Online parameter estimation for the McKean-Vlasov SDE*, Stoch. Proc. Appl. **162**:481–546, 2023.

## Unbiased kinetic Langevin Monte Carlo with inexact gradients

Daniel Paulin

(joint work with Neil Chada, Benedict Leimkuhler, and Peter A. Whalley)

We present the unbiased UBU (UBUBU) method for Bayesian posterior means based on kinetic Langevin dynamics that combines advanced splitting methods with enhanced gradient approximations. Our approach avoids Metropolis correction by coupling Markov chains at different discretization levels in a multilevel

Monte Carlo approach. Theoretical analysis demonstrates that our proposed estimator is unbiased, attains finite variance, and satisfies a central limit theorem. It can achieve accuracy $\epsilon > 0$ for estimating expectations of Lipschitz functions in $d$ dimensions with $O(d^{1/4}\epsilon^{-2})$ expected gradient evaluations, without assuming warm start. We exhibit similar bounds using both approximate and stochastic gradients, and our method's computational cost is shown to scale independently of the size of the dataset. The proposed method is tested using a multinomial regression problem on the MNIST dataset and a Poisson regression model for soccer scores. Experiments indicate that the number of gradient evaluations per effective sample is independent of dimension, even when using inexact gradients. For product distributions, we give dimension-independent variance bounds. Our results demonstrate that the unbiased algorithm we present can be much more efficient than the "gold-standard" randomized Hamiltonian Monte Carlo.
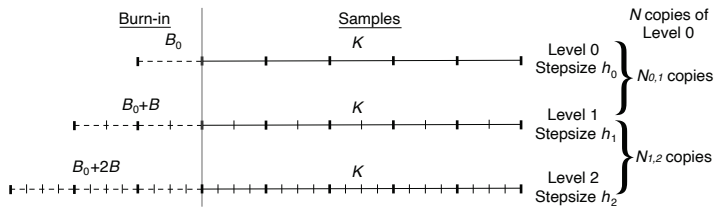


FIGURE 1. Elimination of bias by increasing burn-in lengths at higher discretization levels.
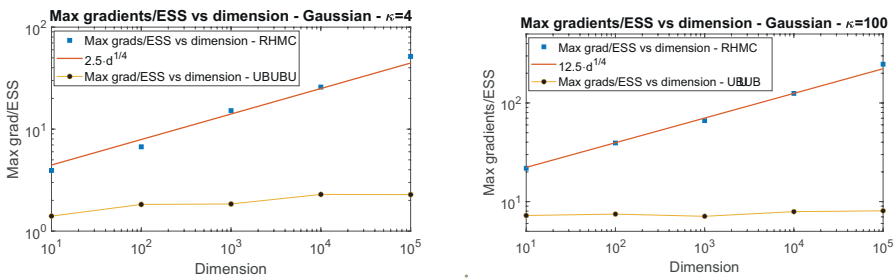


FIGURE 2. Dimensional dependence of gradients/ESS over all components for Gaussian targets.

## Unbiased sampling using reversibility checks

TONY LELIÈVRE

(joint work with Mathias Rousset, Régis Santet, Gabriel Stoltz and Wei Zhang)

Let us consider a probability measure $\pi(dx)$, on a measurable space $(X, \mathcal{F})$ and a Markov kernel $T(x, dx')$ on this same space. The standard Metropolis Hastings algorithm writes: for a given sample $X^n$ at iteration $n$,

$$
\begin{cases}
\text{Sample } \tilde{X}^{n+1} \sim T(X^n, dx') \\
\text{If } U^n \leqslant \min(1, r(X^n, \tilde{X}^{n+1})) \\
\qquad \text{accept the proposal: } X^{n+1} = \tilde{X}^{n+1} \\
\qquad \text{else reject the proposal: } X^{n+1} = X^n
\end{cases}
$$

where $U^n$ is a random variable uniformly distributed, and $r(x, x') = \frac{\pi(dx')T(x',dx)}{\pi(dx)T(x,dx')}$ is the so-called Metropolis-Hastings ratio. It is standard to check that the Markov chain $(X^n)_{n \geqslant 0}$ is invariant with respect to $\pi$. Our objective is to generalize this algorithm in the case when the sampling $\tilde{X}^{n+1} \sim T(X^n, dx')$ requires to solve an implicit problem. In particular, $\tilde{X}^{n+1}$ may be ill-defined (no solution to the implicit problem, or multiple solutions to the implicit problem).

We have more precisely two applications in mind, which are related to the (Generalized) Hamiltonian Monte Carlo algorithm:

- (Generalized) Hamiltonian Monte Carlo algorithm with non-separable hamiltonians;
- Sampling measures on submanifolds with the projected (Generalized) Hamiltonian Monte Carlo algorithm.

In the former case, the generalized Verlet algorithm, which is used to build the proposed move, is implicit because of the non-separability of the hamiltonian. In the latter case, the projection steps used in the Rattle discretization of the constrained Hamiltonian dynamics also require to solve an implicit problem, in order to compute the Lagrange multipliers associated with the projection on the submanifold of interest. We propose a generic method to build a reversible Markov chain in both cases, by modifying the proposal by a so-called reversibility check, very much inspired by [4]. In the specific case when the implicit problems can be solved exactly, for example by root finding softwares for polynomial equations such as Bertini and HomotopyContinuation.jl, we show how this reversibility check can be bypassed.

We refer to the works [1, 2, 3] for more details and numerical illustrations.

REFERENCES

[1] T. Lelièvre, M. Rousset, and G. Stoltz, *Hybrid Monte Carlo methods for sampling probability measures on submanifolds*, Numerische Mathematik **143**(2):379–421, 2019.
[2] T Lelièvre, R. Santet, and G. Stoltz, *Unbiasing Hamiltonian Monte Carlo algorithms for a general Hamiltonian function*, https://arxiv.org/abs/2303.15918.

[3] T. Lelièvre, G. Stoltz, and W. Zhang, *Multiple projection Markov Chain Monte Carlo algorithms on submanifolds*, IMA Journal of Numerical Analysis **43**(2):737–788, 2023.

[4] E. Zappa, M. Holmes-Cerfon, and J. Goodman, *Monte Carlo on manifolds: Sampling densities and integrating functions*, Commun. Pure Appl. Math. **71**(12):2609–2647, 2018.

# Analysis of a positivity-preserving splitting scheme for some semilinear stochastic heat equations

### DAVID COHEN

(joint work with Charles-Edouard Bréhier and Johan Ulander)

**Introduction and motivation.** Let $\Omega$ be a nice domain in $\mathbb{R}^d$. Consider the deterministic linear heat equation with homogeneous Dirichlet boundary conditions

$$\begin{cases} u_t(t,x) = \Delta u(t,x) \ , & t > 0, x \in \Omega, \\ u(0,x) = u_0(x) \ , & x \in \Omega, \end{cases}$$

where $\Delta = \partial^2_{x_1 x_1} + \ldots + \partial^2_{x_d x_d}$ and where the given initial value $u_0$ is non-negative.

Using the semi-group notation, the solution to the above partial differential equation (PDE) may be written

$$u(t) = E(t)u_0,$$

with the solution operator $E(t) = e^{\Delta t}$ which is non-negative by the maximum principle, see for instance [7]. This implies that the solution to the above PDE remains non-negative: $u(t) = E(t)u_0 \geqslant 0$ in $\Omega$ for all time $t \geqslant 0$.

In this extended abstract, we are interested in investigating if the above property of $u(t,x)$ remains valid when the PDE is driven by a random noise, that is $u_t(t,x) = \Delta u(t,x) + NOISE$. Furthermore, we derive and analyse a positivity-preserving numerical integrator for stochastic partial differential equations (SPDEs).

In order to explain the main ideas and obtained theoretical results, we will first consider a linear heat equation driven by a standard Brownian motion (see [3] for details) and then present the results in the case of a linear heat equation driven by a space-time white noise in dimension one (see [2] for details).

**Heat equations driven by a standard Brownian motion.** Let $\big(\beta(t)\big)_{t \geqslant 0}$ be a standard real-valued Brownian motion defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ satisfying the usual conditions. Let an integer $d \geqslant 1$. Define the spatial domain $\mathcal{D} = (0,1)^d$. Let $T \in (0, \infty)$ denote the finite time horizon.

In this part, we consider the time discretisation of the semilinear stochastic heat equation driven by a purely time-dependent Brownian motion (Itô sense)

$$\text{(1)} \quad \begin{cases} \mathrm{d}u(t,x) = \Delta u(t,x)\,\mathrm{d}t + g(u(t,x))\,\mathrm{d}\beta(t) \ , & t > 0, x \in \mathcal{D}, \\ u(t,x) = 0 \ , & t \geqslant 0, \ x \in \partial\mathcal{D}, \\ u(0,x) = u_0(x) \ , & x \in \mathcal{D}, \end{cases}$$

where the (non-random) initial value $u_0 : \overline{\mathcal{D}} \to \mathbb{R}$ is non-negative and the nonlinearity satisfies $g(0) = 0$. Under some technical assumptions on the nonlinearity

$g : \mathbb{R} \to \mathbb{R}$ and the initial value $u_0$, it is known that solutions to the SPDE (1) are Hölder continuous with exponent $1/2-$ in time and $1-$ in space, furthermore these solutions remain non-negative almost surely: $u(t, x) \geqslant 0$ for all $(t, x) \in [0, T] \times \overline{\mathcal{D}}$, see for instance [4].

We now present the efficient time integrator for the SPDE (1) studied in [3]. Let us define the time-step size $\tau = T/M$ where $M \in \mathbb{N}$ is an integer. Set $t_m = m\tau$ for all $m \in \{0, \ldots, M\}$. Define the increments $\delta\beta_m = \beta(t_{m+1}) - \beta(t_m)$ for all $m \in \{0, \ldots, M-1\}$. Next, introduce the bounded and continuous function $f : \mathbb{R} \to \mathbb{R}$ defined by

$$f(v) = \frac{g(v)}{v} \mathbb{1}_{v \neq 0} + g'(0) \mathbb{1}_{v=0}.$$

The positivity-preserving Lie–Trotter splitting scheme $u_m^{\mathrm{LT}}(\cdot) \approx u(t_m, \cdot)$ is based on a Lie–Trotter splitting strategy: given $u_m^{\mathrm{LT}}$ for some $m \in \{0, \ldots, M-1\}$, the numerical solution $u_{m+1}^{\mathrm{LT}}$ is obtained by solving successively two subsystems on the time interval $[t_m, t_{m+1}]$:

- first, a family of linear Itô stochastic differential equations (corresponding to geometric Brownian motions)
- second, a deterministic PDE (the linear heat equation).

The Lie–Trotter splitting scheme can be written as

$$(2) \qquad u_{m+1}^{\mathrm{LT}}(x) = \int_{\mathcal{D}} G(\tau, x, y) \left( \exp\left( f(u_m^{\mathrm{LT}}(y)) \delta\beta_m - \frac{f(u_m^{\mathrm{LT}}(y))^2 \tau}{2} \right) \right) \, \mathrm{d}y,$$

where $G(t, x, y)$ is the fundamental solution to the deterministic heat equation.

The time integrator (2) satisfies the following, see [3] for precise statements:

- Almost surely, $u_m^{\mathrm{LT}}(x) \geqslant 0$ for all $x \in \mathcal{D}$ and $m \in \{0, \ldots, M\}$,
- It is exact when applied to the SPDE (1) when $g(v) = v$, i.e. $f(v) = 1$,
- It has mean-square order of convergence $1/2$, that is:

$$\sup_{0 \leqslant m \leqslant M} \sup_{x \in \overline{\mathcal{D}}} \left( \mathbb{E}[|u_m^{\mathrm{LT}}(x) - u(t_m, x)|^2] \right)^{1/2} \leqslant C \, \tau^{1/2}.$$

These properties have been illustrated numerically in [3].

We conclude the first part of the extended abstract by mentioning some possible future research questions: It has been noted that the LT splitting scheme remains positive and has a rate of mean-square convergence $1/2$ even in the case of the non-globally Lipschitz nonlinearity $g(v) = v^{1.25}$. To prove this remains an open question. In addition, the rate of weak convergence of the LT splitting scheme has been investigated in [1] for finite-dimensional SDEs. To prove a rate of weak convergence of the LT splitting scheme when applied to the SPDE (1) is still open.

**Heat equations driven by a space-time white noise.** In this second part, we consider the discretisation in space and time of the semilinear stochastic heat

equation (in the Itô sense)

(3)
$$\begin{cases} \partial_t u(t,x) = \partial_{xx}^2 u(t,x) + g(u(t,x))\,\dot{W}(t,x), \\ u(t,0) = u(t,1) = 0, \\ u(0,x) = u_0(x) \end{cases}$$

for $t \in [0,T]$ and $x \in [0,1]$, where $\dot{W}(t,x)$ is a space-time white noise, see for instance [8].

The spatial discretisation of the SPDE (3) is analysed in [5]. For any integer $N \in \mathbb{N}$, let $h = 1/N$ be the space mesh size, and let $x_n = nh$ for $0 \leqslant n \leqslant N$ be the grid points. A standard finite difference discretisation gives the $N-1$-dimensional stochastic differential equation (SDE)

$$\mathrm{d}u^N(t) = N^2 D^N u^N(t)\ \mathrm{d}t + \sqrt{N} g(u^N(t))\ \mathrm{d}W^N(t)$$

with initial value $u^N(0) = \big(u(0,x_n)\big)_{1\leqslant n\leqslant N-1}$, $D^N = \mathrm{diag}(1,-2,1)$ is a tridiagonal matrix and $W_n^N(t) = \sqrt{N}\big(W(t,x_{n+1}) - W(t,x_n)\big)$ for $1 \leqslant n \leqslant N-1$. We use the notation $\big(g(u^N(t))\ \mathrm{d}W^N(t)\big)_n = g(u_n^N(t))\ \mathrm{d}W_n^N(t)$.

This system of SDEs is then discretised in time by the Lie–Trotter splitting scheme presented in the first part of this extended abstract. The fully-discrete numerical approximation of the SPDE (3) then reads:

(4)
$$u_{m+1}^{\mathrm{LT}} = e^{\tau N^2 D^N} \left( \exp\!\Big(\sqrt{N} f(u_{m,n}^{\mathrm{LT}})\Delta_{m,n}W - \frac{N f(u_{m,n}^{\mathrm{LT}})^2 \tau}{2}\Big) u_{m,n}^{\mathrm{LT}} \right)_{1\leqslant n\leqslant N-1},$$

where $\Delta_{m,n}W = W_n^N(t_{m+1}) - W_n^N(t_m)$.

Under some technical assumptions, the fully-discrete numerical scheme (4) satisfies the following, see [2] for precise statements:

- Almost surely, $u_{m,n}^{\mathrm{LT}} \geqslant 0$ for all $m \in \{0,\dots,M\}$ and $n \in \{1,\dots,N-1\}$,
- It has mean-square order of convergence $1/4$, that is: For all $\gamma \in (0,\infty)$ and $T \in (0,\infty)$, there exists $C = C_{\gamma,T}(u_0) \in (0,\infty)$ such that for all $\tau = T/M$ and $h = 1/N$ satisfying the condition $\tau \leqslant \gamma h^2$,

$$\sup_{0\leqslant m\leqslant M}\ \sup_{0\leqslant n\leqslant N}\ \big(\mathbb{E}[|u_{m,n}^{\mathrm{LT}} - u(t_m,x_n))|^2]\big)^{\frac{1}{2}} \leqslant C\tau^{\frac{1}{4}}.$$

It is possible to get an upper bound, not uniformly with respect to $h$, for the mean-square error under the condition $\tau \leqslant \gamma h$, see [2] for details.

These properties have been illustrated numerically in [2].

We conclude this extended abstract by mentioning some possible future research questions: It has been observed numerically in [2] that the LT splitting scheme remains positive and has a rate of mean-square convergence $1/4$ even in the case of the non-globally Lipschitz nonlinearity $g(v) = v^{1.25}$. To prove these properties for the SPDE (3) with $g$ less regular than in [2] remains an open question. In addition, the rate of weak convergence of the LT splitting scheme when applied to the SPDE (3) has not yet been investigated.

REFERENCES

[1] M. Bossy, J.-F. Jabir, K. Martínez, *On the weak convergence rate of an exponential Euler scheme for SDEs governed by coefficients with superlinear growth*, Bernoulli **27**, 312–347 (2021).
[2] C.-E. Bréhier, D. Cohen, J. Ulander, *Analysis of a positivity-preserving splitting scheme for some semilinear stochastic heat equations*, to appear in ESAIM:M2AN (2024).
[3] C.-E. Bréhier, D. Cohen, J. Ulander, *Positivity-preserving schemes for some nonlinear stochastic PDEs*, the proceedings of the Sixteenth International Conference Zaragoza-Pau on Mathematics and its Applications 2022 (2024).
[4] J. Cresson, M. Efendiev, S. Sonner, *On the positivity of solutions of systems of stochastic PDEs*, ZAMM Z. Angew. Math. Mech. **93**, 414–422 (2013).
[5] I. Gyöngy, *Lattice approximations for stochastic quasi-linear parabolic partial differential equations driven by space-time white noise I*, Potential Anal. **9**, 1–25 (1998).
[6] G.J. Lord, C.E. Powell, T. Shardlow *An Introduction to Computational Stochastic PDEs*, Cambridge University Press (2014).
[7] V. Thomée, *On positivity preservation in some finite element methods for the heat equation*, Lecture Notes in Comput. Sci. **8962**, 13–24, Springer (2015).
[8] J.B. Walsh *An Introduction to Stochastic Partial Differential Equations*, Springer (1986).

## Extended dynamic mode decomposition and error rates

### CAROLINE WORMELL

Computation of nonlinear dynamical systems is a complicated area, and there are few universal tools at our disposal. One of these tools is the Koopman operator, which allows us to study the action of a dynamical system as a linear object. Given a map $f : M \to M$ (which could be stochastic and Markov), the Koopman operator acts on functions $\varphi : M \to \mathbb{R}$ as

$$\mathcal{K}\varphi = \mathbb{E}[\varphi \circ f].$$

Various ergodic and geometric properties of the map can be revealed by studying the Koopman operator on appropriate spaces of functions: for instance, autocorrelation functions, Birkhoff variances, and almost-invariant subsets of phase space. This is particularly useful for quantifying e.g. outputs of the dynamical system that drive other systems.

Given its utility, we might wonder how we can compute with the Koopman operator, especially when our information on the dynamical system is relatively limited. The most natural idea is to project $\mathcal{K}$ onto a subspace of functions, for example those spanned by a *dictionary of observables* $\{\psi_n\}_{|n|<N}$ (or some other set of in). Computationally, we can achieve this most easily by least squares approximation, and that ideally with respect to a finite set of data points $\{(x_m, f(x_m))\}_{m=1,...,M}$. In practice, these data points usually are taken from some sampling measure $\mu$ (If

the points form an orbit, so $f(x_m) = x_{m,+1}$, this measure is an ergodic measure of $f$).

To do this, one constructs a design matrix and fitting matrix

$$\Psi_X = (\psi_n(x_m))_{m=1,\ldots,M,|n|<N}, \Psi_Y = (\psi_n(f(x_m)))_{m=1,\ldots,M,|n|<N}$$

and the Koopman matrix is given as $K_{M,N} = \Psi_X^+ \Psi_Y$, where $A^+$ denotes the pseudoinverse of $A$. Such a method of approximating the Koopman operator is known as Extended Dynamical Mode Decomposition (EDMD). Thousands of papers are published about it, but little is known about the interpretation or quality output, particularly in chaotic dynamics.

The goal of this talk is to understand what the output of EDMD converges to, and some principles governing its corresponding error rates. As a simple metric for "error", we'll talk about the error in eigenvalues of the Koopman operator.

As a starting point, we choose our phase space to be the periodic interval $M = \mathbb{R}/2\pi\mathbb{Z}$, and on this space compare stochastic dynamics (with, say, smooth kernel) with uniformly expanding dynamics (i.e. smooth maps with $|f'| > \gamma > 1$). We will choose our observables to be the Fourier basis $\{e^{inx}\}_{|n|<N}$, but the sampling measure may not be Lebesgue—although for both kinds of dynamical system, the physically observed ergodic measure has a smooth Lebesgue density.

For these kinds of dynamics, it is possible to show that, as the amount of data converges to infinity the Koopman matrix $K_{N,M}$ converges to a limiting matrix $K_N$, with an $\mathcal{O}(1/\sqrt{M})$ error in each matrix entry [2]. For these observables, it is also possible to see that the errors in each entries are sufficiently uncorrelated that the $\ell^2$ norm of this $2N-1 \times 2N-1$ matrix is of order $\mathcal{O}(\sqrt{N}/\sqrt{M})$: by Fourier duality this is the corresponding $L^2(\mathrm{d}x)$ distance in function space between the operators $\mathcal{K}_{N,M}$ and $\mathcal{K}_N$.

This limiting operator $\mathcal{K}_N$ is at first sight fairly easy to characterise: if we define $\mathcal{P}_N$ to be the orthogonal projection onto the subspace spanned by the observables in $L^2(\mu)$, then we have $\mathcal{K}_N = \mathcal{P}_N \mathcal{K}|_{\mathrm{im}\,\mathcal{P}_N}$. For our choice of observables, it turns out the effectiveness of this projection is quite similar to (but cannot be reduced to) truncation of Fourier modes (which is what you get when $\mu$ is equal to Lebesgue measure). We can quantify this by studying the operator error in Sobolev spaces $H^s$, consisting of functions whose $s$th derivative is in $L^2$.

**Theorem 1** (Theorem 1.1,[3])**.** *There exists $C$ such that for all $s > r \geqslant 0$,*

$$\|I - \mathcal{P}_N\|_{H^s \to H^r} \leqslant CN^{-(s-r)}.$$

In the case of our stochastic dynamics, the smoothness of the kernel implies that for all $s$ $\mathcal{K}$ is bounded from $L^2 \to H^s$, and so we have that $\|\mathcal{K} - \mathcal{K}_N\|_{L^2} \leqslant C'N^{-s}$. Since $\mathcal{K}$ is compact on $L^2$ for stochastic maps, we can put this together with the noise bound to get the following:

**Corollary 1.** *For every simple eigenvalue $\lambda \in \sigma(\mathcal{K}, L^2)$, there exists a sequence $\{\lambda_{N,M} \in \sigma(\mathcal{K}_{N,M})\}_{N,M}$ so that*

$$|\lambda_{N,M} - \lambda| \leqslant \mathcal{O}(N^{-s} + \sqrt{N}/\sqrt{M})$$

*for each $N, M$ with high probability, for every $s > 0$.*

With some more work, we could probably remove the $\sqrt{N}$ factor in the last term.

We would hope that the same is possible for chaotic systems. However, in this case, there are some obstacles. The first is that $\mathcal{K}$ is not bounded from a stronger Sobolev space to a weaker Sobolev space. However, there are ways to get around this, and we can show the following resolvent bound

**Theorem 2.** *For every $\lambda \in \mathbb{C}\backslash\sigma(\mathcal{K}, H^s)$ and every $s, t > 0$,*

$$\|R_\lambda(\mathcal{K}_N) - R_\lambda(\mathcal{K})|_{\operatorname{im}\mathcal{P}_N}\|_{H^s} = \mathcal{O}(N^{-t}).$$

This implies that

**Corollary 2.** *For every $\lambda \in \mathbb{C}\backslash\sigma(\mathcal{K}, H^s)$ and every $s, t > 0$,*

$$\|R_\lambda(\mathcal{K}_{N,M}) - R_\lambda(\mathcal{K})|_{\operatorname{im}\mathcal{P}_N}\|_{H^s} = \mathcal{O}(N^{-t} + \|\mathcal{K}_{N,M} - \mathcal{K}_M\|_{H^s})$$
$$\leqslant \mathcal{O}(N^{-t} + N^{1/2+s}/\sqrt{M}),$$

*with high probability for each $N, M$.*

with the last inequality obtained by comparing the Sobolev norms in Fourier space. (It seems for these maps that there is no special covariance structure in the "random matrix" that allows a better bound.) Setting $H^0 = 0$ we recover the same kind of decay we had for stochastic systems.

However, the other problem we have is that the spectrum of $\mathcal{K}$ in $L^2$ is very boring (essential spectrum filling the unit disk), and isn't even the limit of $\sigma(\mathcal{K}_N)$ as $N \to \infty$. However, if we study the spectrum of $H^s$ as $s > 0$, we get something more interesting. There are a discrete set of values (with multiplicity) that are known as the Ruelle–Pollicott resonances $\sigma(f)$, which can be defined by means independent of function spaces. It is known that for $s > 0$,

$$\sigma(\mathcal{K}, H^s) = \sigma(f) \cup \sigma_{\operatorname{ess}}(\mathcal{K}_N, H^s).$$

The essential spectrum $\sigma_{\operatorname{ess}}(\mathcal{K}_N, H^s)$ is contained in a ball of radius $e^{-P(s)}$ about zero, where $P(s)$ can be quantified in terms of expansion rates of $f$, and is $\sim L_{\operatorname{exp}}s + o(s)$ for $s$ small [1]. Furthermore, increasing $s$ shrinks the essential spectrum.

We can use Corollary 2 to bound the eigenvalue error between some $\lambda_{N,M} \in \sigma(\mathcal{K}_{N,M})$ and $\lambda \in \sigma(f)$ as

$$\inf_{s:\lambda\in\sigma_{\operatorname{pt}}(\mathcal{K}, H^s)} \mathcal{O}(N^{-t} + N^{1/2+s}/\sqrt{M})$$

This is obviously (asymptotically) increasing in $s$, so we want to minimise $s$. By the bounds on the essential spectrum, we know we can take $s = P^{-1}(\log|\lambda|) + \epsilon$, so, for any $t$,

$$|\lambda - \lambda_{N,M}| \leqslant \mathcal{O}\left(N^{-t} + \frac{N^{1/2 + P^{-1}(\log|\lambda|)}}{\sqrt{M}}\right).$$

The upshot of this is that the dependence on the size of $\lambda$ here suggests that smaller eigenvalues require much larger amounts of data to be resolvable. The variable essential spectral radius phenomenon is a standard property of chaotic systems, so we can expect to see this phenomenon to more realistic dynamics. Fortunately, smaller eigenvalues tend to be less physically meaningful than larger ones, and the extra penalty on eigenvalues close to the unit circle is relatively small.

The results in this talk suggest that when it comes to estimating Koopman operator spectra, there is a disjunction between stochastic and chaotic dynamics: all eigenfunctions of stochastic systems are $L^2$ objects, whereas the internal spectrum of chaotic systems are creatures of higher regularity, even though can have set-based interpretations. This means it is harder to study the internal spectra of chaotic systems.

<div align="center">REFERENCES</div>

[1] V. Baladi. *Dynamical zeta functions and dynamical determinants for hyperbolic maps*. Berlin: Springer International Publishing (2018).

[2] S. Klus and C. Schütte. *Towards tensor-based methods for the numerical approximation of the Perron–Frobenius and Koopman operator*, Journal of Computational Dynamics **3(2)** (2016) 139–161.

[3] C. Wormell, *Orthogonal polynomial approximation and Extended Dynamical Mode Decomposition in chaos*, preprint (2024).

<div align="center">

## Sampling on Riemannian manifolds via stochastic differential equations

ALEXANDER LEWIS

(joint work with Karthik Bharath, Akash Sharma, Michael V. Tretyakov)

</div>

The general objective we wish to pursue is the sampling from a general probability measure $d\mu_\phi \propto e^{-\phi}\text{dvol}$ on a Riemannian manifold $M$ of dimension $q$, where dvol is the canonical volume form. When $M$ is Euclidean, it is well known that sampling may be carried out by discretizing an overdamped Langevin SDE.

Our analysis centres around the intrinsic (overdamped) Langevin equation

$$(1) \qquad dX_t = dB_t^M - \frac{1}{2}\nabla\phi(X_t)dt, \quad X_0 = x,$$

where $B_t^M$ is an $M$-valued Brownian motion constructed by the Eells–Elworthy–Malliavin approach (e.g. [3]). When $M$ is compact, the geometry guarantees that

the distribution of the solution $X_t$ converge exponentially to the stationary distribution $\mu_\phi$. So by discretizing (1), it should be possible to approximately sample from the target measure $\mu_\phi$.

Previous methods (e.g. [4, 5]) for numerical schemes on manifolds have relied upon embedding the manifold in some higher dimensional Euclidean space (typically with $\mathrm{codim}(M) = 1$). The constrained diffusion then lies on the manifold. This technique may be inefficient since the discretized process lives in a higher-dimensional space, and needs to be projected onto $M$. Moreover numerically, as time increases, the repeated application of the projection will accumulate error that can interfere with the order of convergence. Instead we shall rely on an intrinsic based approach, where we do not rely on any particular choice of embedding and the number of random variables driving the discretization at each step is equal to the intrinsic dimension of the manifold. An intrinsic approach which provided orders of convergence of an Euler discretization was first developed in [2], however their method of proof relied heavily upon choice of coordinate charts. Calculus was performed on the entirety of charts, and the exit time of a chart had to be tracked. For this reason, the random variables driving the Euler discretization had to be restricted. Our approach on the other hand enables coordinate-free analysis of the algorithm.

We define the geodesic Langevin algorithm as

$$(2) \qquad X_{n+1}^h = \exp_{X_n^h}\left( -\frac{h}{2}\nabla\phi(X_n^h) + \sqrt{h}g^{-1/2}(X_n^h)\xi_{n+1} \right),$$

where $\xi$ is i.i.d. $\mathbb{R}^q$-valued random vector that satisfies the following moment matching conditions:

$$\mathbb{E}[\xi^i] = 0, \ \mathbb{E}[\xi^i\xi^j] = \delta_{ij}, \ \mathbb{E}[\xi^i\xi^j\xi^k] = 0, \ \mathbb{E}[(\xi^i)^2(\xi^j)^2] < \infty.$$

Examples of such random variables that are usable are standard Gaussian and the discrete random variable $\xi_i = \pm 1$ w.p. $1/2$.

**Theorem 1** (Global convergence theorem). *Let $M$ be a compact Riemannian manifold and assume $g^{ij} \in C^3(M)$, $\phi \in C^3(M)$ and $\varphi \in C^4(M)$. Define the estimator $\hat{\mu}_{\phi,N}(\varphi) = \frac{1}{N}\sum_{n=1}^N \varphi(X_N^{(n),h})$. Then,*

$$(3) \qquad |\mathbb{E}[\hat{\mu}_{\phi,N}(\varphi)] - \mathbb{E}_{\mu_\phi}[\varphi]| \leqslant C\left( h + e^{-\lambda T} \right).$$

For a global order 1 convergence, we aim for a local, one step approximation, in order 2.

**Lemma 1** (One-step approximation lemma). *Let $M$ be a compact Riemannian manifold and assume $g^{ij} \in C^3(M)$, $\phi \in C^3(M)$ and $\varphi \in C^4(M)$. Given $x \in M$, let $X_1^h$ be computed according to the following formula:*

$$X_1^h = \exp_x\left( -\frac{h}{2}\nabla\phi(x) + \sqrt{h}g^{-1/2}(x)\xi \right).$$

*Let $u(t,x)$ be the solution to the backward Kolmogorov equation of the diffusion (1). Then*

$$\mathbb{E}[u(t+h, X_1^h) - u(t,x)] \leqslant Ch^2 e^{-\lambda(T-t)},$$

*where $C > 0$ is independent of $T$ and $h$.*

The proof of Lemma 1 (see e.g. [1]) exploits the Taylor expansion of $u(t+h, X_1^h)$ along a geodesic.

Hence, a key ingredient in our coordinate-free approach to the proofs is to write the Taylor expansion of a function $f \in C^4(M)$ along geodesics. Let $\gamma : [0, \sqrt{h}] \to M$ be a geodesic with $\gamma(0) = x$ and $\dot\gamma(0) = V$, then we write the Taylor series $f(\gamma(s))$ in the parameter $s$ as

$$
\begin{aligned}
f(\gamma(s)) = f(\gamma(0)) + s\frac{\mathrm{d}}{\mathrm{d}s}f(\gamma(s))\Big|_{s=0} + \frac{s^2}{2}\frac{\mathrm{d}^2}{\mathrm{d}s^2}f(\gamma(s))\Big|_{s=0} \\
+ \frac{s^3}{6}\frac{\mathrm{d}^3}{\mathrm{d}s^3}f(\gamma(s))\Big|_{s=0} + \frac{s^4}{24}\frac{\mathrm{d}^4}{\mathrm{d}s^4}f(\gamma(s))\Big|_{s=\alpha}
\end{aligned}
$$

for some $\alpha \in [0, \sqrt{h}]$.

Then by recognising that $\frac{\mathrm{d}^k}{\mathrm{d}s^k}f(\gamma(s))|_{s=0} = D^k f(V, \ldots, V)$, we see that the only random part in the derivative terms are the vector fields $V$. Hence, the problem of evaluating expectations on the manifold is changed into evaluating expectations on the tangent space, and because $T_x M$ is isomorphic to $\mathbb{R}^q$, makes calculating the expectations much simpler.

We obtain the same order (both local and global) of convergence when compared to the Euclidean case [6].

For non-compact manifolds, the issue of convergence of the diffusion (1) is a subtle one, tied intimately to the geometry of $M$. The Bakry–Èmery criterion is one such condition that is sufficient for exponential erogidicty. However, we must impose very stringent conditions on both the geometry and density to ensure it is satisfied; when setting $M = \mathbb{R}^q$, this is akin to demanding that the target density is strongly log-concave. To broaden the class of allowed densities we can theoretically sample from, we are able to split the Bakry–Èmery criterion into two separate assumptions on the Ricci curvature and Hessian: For constants $b, c > 0$,

$$(4) \qquad\qquad \mathrm{Ric} \geqslant (-c - b^2\rho_o^2)g,$$

where $\rho_o = \rho(o, x)$, $o, x \in M$. And

$$(5) \qquad\qquad \mathrm{Hess}^\phi \geqslant \delta g$$

outside of a compact set in $M$, where the constant $\delta$ relates to the constants from (4) as $\delta > (1 + \sqrt{2})b\sqrt{q-1} > 0$. If both (4) and (5) are satisfied, then the log-Sobolev inequality on $M$ holds [8].

Experiments were performed on $\mathbb{S}^2$ and $\mathcal{P}_3$ (the space of $3\times3$ symmetric positive definite matrices) to numerically illustrate the convergence rate of the algorithm (2).

Theorem 1 and Lemma 1 can also be adapted when the algorithm (2) moves along retractions [7]. Define the first order retraction $F_x : T_x M \to M$ approximating the exponential map in (2), and further define $c(s)$ a corresponding curve satisfying $c(0) = x$ and $\dot{c}(0) = -\frac{\sqrt{h}}{2}\nabla\phi(x) + g^{-1/2}(x)\xi$. Then the convergence rates in Theorem 1 and Lemma 1 are the same when $\exp_x$ is replace with $F_x$ in (2) if $|\mathbb{E}[D_{\dot{c}(s)}\dot{c}(s)|_{s=0}]| \leqslant Ch$ and $|\mathbb{E}[D_{\dot{c}(s)}D_{\dot{c}(s)}\dot{c}(s)|_{s=0}]| \leqslant Ch^{1/2}$.

Further avenues of research include; e.g.,

(1) Higher order intrinsic weak order algorithms;
(2) Alternative sampling methods such as splitting in the underdamped Langevin diffusion;
(3) Modifying the algorithm to cater to manifolds with a boundary by introducing a reflection term.

Further details, proofs for Theorem 1 and Lemma 1, and numerical experiments can be found in [1].

REFERENCES

[1] K. Bharath, A. Lewis, A. Sharma, M.V. Tretyakov, *Sampling and estimation on manifolds using the Langevin diffusion*, Preprint (2023), arXiv:2312.14882.
[2] A. Grorud, D. Talay, *Approximation of Lyapunov exponents of nonlinear stochastic differential equations*, SIAM Journal on Applied Mathematics, **56** (1996), no. 2, 627-650.
[3] E.P. Hsu, *Stochastic analysis on manifolds*, American Mathematical Society, **37** (2002).
[4] A. Laurent, G. Vilmart, *Order conditions for sampling the invariant measure of ergodic stochastic differential equations on manifolds*, Foundations of Computational Mathematics **22** (2022), no. 3, 649–695.
[5] T. Lelièvre, M. Rousset, G. Stoltz, *Hybrid Monte Carlo methods for sampling probability measures on submanifolds*, Numerische Mathematik **143** (2019), no. 2, 379–421
[6] G.N. Milstein, M.V. Tretyakov, *Stochastic numerics for mathematical physics*, Springer, **39** (2004).
[7] D. Talay, L. Tubaro, *Expansion of the global error for numerical schemes solving stochastic differential equations*, Stochastic analysis and applications, **8** (1990), no. 4, 483–509.
[8] F-Y. Wang, *Log-Sobolev inequalities: Different roles of Ric and Hess*, Annals of Probability, **37** (2009), no. 4, 1587-1604.

## Deterministic simulation of SDEs driven by $\alpha$-stable processes

GEORG A. GOTTWALD

(joint work with Ian Melbourne)

We present a novel framework to numerically integrate stochastic differential equations (SDEs) which are driven by Lévy noise [1]. Whereas classical methods discretize continuous-time SDEs using Taylor expansions, we view an SDE as a limiting slow equation of a discrete slow-fast map using rigorous theory for homogenization of deterministic weakly chaotic systems [2, 3]. Our method naturally deals with the notorious Marcus-integral appearing in SDEs driven by multiplicative Lévy-noise. As a by-product this implies an entirely deterministic procedure

to generate $\alpha$-stable noise. In particular, we devise an explicit method to integrate $\alpha$-stable stochastic differential equations (SDEs) with nonglobally Lipschitz coefficients. To mitigate against numerical instabilities caused by unbounded increments of the Lévy noise, we use a deterministic map which has the desired SDE as its homogenised limit. We present an example of an SDE with a natural boundary showing that our method naturally respects the boundary whereas Euler-Maruyama discretisation fails to do so and shows leakage into forbidden regions.

We further report on recent results for multi-dimensional decorated $\alpha$-stable processes which were recently discovered [4].

## References

[1] G.A. Gottwald and I. Melbourne, *Simulation of non-Lipschitz stochastic differential equations driven by $\alpha$-stable noise: a method based on deterministic homogenisation*, Multiscale Modeling & Simulation **19**(2) (2021), 665–687.

[2] G.A. Gottwald and I. Melbourne, *Homogenization for deterministic maps and multiplicative noise*, Proceedings of the Royal Society A **469** (2013), 20130201.

[3] I. Chevyrev, P.K. Friz, A. Korepanov, and I. Melbourne, *Superdiffusive limits for deterministic fast–slow dynamical systems*, Probability Theory and Related Fields **178**(3) (2020), 735–770.

[4] I. Chevyrev, P.K. Friz, A. Korepanov, and I. Melbourne, *Superdiffusive limits beyond the Marcus regime for deterministic fast-slow systems*, ArXiv arXiv:2312.15734 [math.DS] (2024).

## A mathematical analysis of autoencoders for free energy calculations

Gabriel Stoltz

(joint work with Tony Lelièvre, Thomas Pigeon and Wei Zhang)

Autoencoders fall into the class of unsupervised machine learning methods. They can be used a dimension reduction tool, to extract salient features of the data at hand. They can be seen as a nonlinear extension to principal component analysis. They are used in particular in molecular dynamics to find so-called collective variables, in the context of free energy calculations [11].

In essence, autoencoders aim at representing the identity function with a model of limited capacity. More precisely, for a given input data point $x \in \mathcal{X} \subset \mathbb{R}^D$, we denote by $f_\theta(x)$ the prediction of the neural network. The parameters $\theta \in \Theta$ (weights and biases for each layer) are chosen to minimize the loss function

$$(1) \qquad \mathscr{L}(\theta) = \mathbb{E}\left[\|X - f_\theta(X)\|^2\right],$$

where the expectation is over the realizations of the input data $X$ distributed according to some probability measure denoted by $\mu$. In practice, the population loss $\mathscr{L}$ is replaced by the empirical loss over a training set of $N_{\text{data}}$ given input

data points $\{x^1, \ldots, x^{N_{\text{data}}}\}$:

$$\widehat{\mathscr{L}}(\theta) = \frac{1}{N_{\text{data}}} \sum_{n=1}^{N_{\text{data}}} \|x^n - f_\theta(x^n)\|^2 .$$

There are various classes of autoencoders, as reviewed in [10, Section 2.1]. It is useful to distinguish between undercomplete and overcomplete models. Undercomplete models have a limited capacity that prevents them from achieving zero training loss. The most prominent example is provided by (feedforward) bottleneck autoencoders for which

$$(2) \qquad\qquad f_\theta = f_{\text{dec},\theta_2} \circ f_{\text{enc},\theta_1},$$

where the parameters $\theta = (\theta_1, \theta_2)$ have been decomposed into parameters used in the encoder and decoder parts, respectively. The limitation in the capacity of the autoencoder arises from the fact that the encoding function $f_{\text{enc},\theta_1}$ has values in a latent space $\mathcal{Z} \subset \mathbb{R}^d$ of dimension $d$ strictly smaller than the dimension $D$ of the input/output space $\mathcal{X}$, usually much smaller in fact. Autoencoders are often symmetric in their structures. In some cases, tied weights are being used, *i.e.* the weights $\theta_2$ are the transpose of the weights $\theta_1$ when writing the prediction function as (2). However, there is no particular motivation to use symmetric architectures, and the results in Section 1.2 below in fact suggest to resort to very expressive decoders.

## 1. INTERPRETATIONS OF THE LOSS FUNCTION

We discuss in this section various reformulations and reinterpretations of the loss function (1) for bottleneck autoencoders (2) when the loss function is the square loss, and discuss in particular the relationship with principal curves/manifolds [9, 14], and conditional expectations.

1.1. **Three viewpoints on the loss function.** We consider an ideal setting where we minimize upon all measurable functions $f_{\text{enc}} : \mathcal{X} \to \mathcal{Z}$ and $f_{\text{dec}} : \mathcal{Z} \to \mathcal{X}$. We denote by $\mathcal{F}_{\text{enc}}$ and $\mathcal{F}_{\text{dec}}$ the sets of measurable functions from $\mathcal{X}$ to $\mathcal{Z}$ and from $\mathcal{Z}$ to $\mathcal{X}$, respectively; and by $\mathcal{F}$ the set of measurable functions obtained by composing functions of $\mathcal{F}_{\text{enc}}$ with functions of $\mathcal{F}_{\text{dec}}$:

$$\mathcal{F} = \{f = f_{\text{dec}} \circ f_{\text{enc}}, \ f_{\text{enc}} \in \mathcal{F}_{\text{enc}}, \ f_{\text{dec}} \in \mathcal{F}_{\text{dec}}\} .$$

Minimizing the reconstruction error over the set of functions in $\mathcal{F}$ can be then rewritten as

$$(3) \qquad\qquad \inf_{f \in \mathcal{F}} \mathbb{E}\left[\|X - f(X)\|^2\right] .$$

Note that we do not consider a regularization term here, so that overfitting may occur in practice (for instance, even with $\mathcal{Z}$ of dimension 1, $f_{\text{dec}}$ can parametrize a space-filling curve).

As discussed in [6] (which complements [8] which was already hinting at autoencoders), the unsupervised least-square problem (3) can be thought of in various

ways. In particular, there is some duality in the way the minimization over $f \in \mathcal{F}$ is performed, as one can decide to either

   (i) simultaneously minimize over $f_{\mathrm{enc}}$ and $f_{\mathrm{dec}}$, which is the standard way to proceed when training neural networks;

  (ii) minimize first over the encoder part, which allows to reformulate the minimization as the well-known problem of finding principal manifolds;

 (iii) minimize first over the decoder part, which is natural when thinking of the reconstruction error as some total variance to be decomposed using a conditioning on the values of the encoder. This approach is discussed more precisely in Section 1.2.

The chosen numerical approach has a natural impact on the topology of the networks which are considered: in situation (i), encoders and decoders are treated on an equal footing, and it is therefore natural to consider them to be of a similar complexity; whereas options (ii) and (iii) suggest to consider asymmetric autoencoders. For instance, in option (iii), the minimization over the decoder part, which is performed first, could be done more carefully, with more expressive networks in order to better approximate the optimal decoder for a given encoder.

1.2. **Reformulating autoencoders with conditional expectations.** We discuss here how to reformulate the training of autoencoders with conditional expectations, and provide alternative interpretations to the reconstruction error. We minimize the reconstruction error (3) by first minimizing over the decoder part for a given encoding function, as already considered in [7]. This approach is natural in molecular dynamics, as it is reminiscent of free energy computations [3, 11] where average quantities are computed for a fixed value of the collective variable $f_{\mathrm{enc}}$. From a mathematical viewpoint, it corresponds to introducing conditional averages associated with fixed values of the encoder.

    The loss function for unsupervised least-squares can be rewritten as

$$\inf_{f \in \mathcal{F}} \mathbb{E}\left[\|X - f(X)\|^2\right] = \inf_{f_{\mathrm{enc}} \in \mathcal{F}_{\mathrm{enc}}} \left\{ \inf_{f_{\mathrm{dec}} \in \mathcal{F}_{\mathrm{dec}}} \mathbb{E}\left[\|X - f_{\mathrm{dec}} \circ f_{\mathrm{enc}}(X)\|^2\right] \right\}$$

$$(4) \qquad\qquad\qquad = \inf_{f_{\mathrm{enc}} \in \mathcal{F}_{\mathrm{enc}}} \mathbb{E}\left[\|X - g^{\star}_{f_{\mathrm{enc}}} \circ f_{\mathrm{enc}}(X)\|^2\right],$$

where the ideal decoder $g^{\star}_{f_{\mathrm{enc}}}$ for a given encoder $f_{\mathrm{enc}}$ is the Bayes predictor associated with the least square regression problem (see [2, Section 2.2.3]):

$$(5) \qquad\qquad\qquad g^{\star}_{f_{\mathrm{enc}}}(z) = \mathbb{E}[\, X \mid f_{\mathrm{enc}}(X) = z].$$

Let us recall that, in all these expressions, expectations are taken with respect to the probability distribution $\mu$ of the input data (which is not necessarily the Boltzmann–Gibbs distribution). Equations (4)-(5) show that the question of finding the best autoencoder can be reduced to finding the best encoding function, provided that one is able to compute good approximations of the conditional expectation.

The reconstruction error (4) can be reinterpreted in terms of variances. Indeed,

$$
(6) \qquad
\begin{aligned}
\mathbb{E}\left[\left\|X - g^\star_{f_{\mathrm{enc}}} \circ f_{\mathrm{enc}}(X)\right\|^2\right] &= \mathbb{E}\left[\mathrm{Var}(X \,|\, f_{\mathrm{enc}}(X))\right] \\
&= \mathrm{Var}(X) - \mathrm{Var}\left[\mathbb{E}(X \,|\, f_{\mathrm{enc}}(X))\right].
\end{aligned}
$$

A consequence of (6) is that the minimization problem (4) can be reformulated as the following equivalent maximization problem:

$$
(7) \qquad
\sup_{f_{\mathrm{enc}} \in \mathcal{F}_{\mathrm{enc}}} \mathrm{Var}\left[\mathbb{E}(X \,|\, f_{\mathrm{enc}}(X))\right] = \sup_{f_{\mathrm{enc}} \in \mathcal{F}_{\mathrm{enc}}} \mathbb{E}\left[\left\|g^\star_{f_{\mathrm{enc}}} \circ f_{\mathrm{enc}}(X)\right\|^2\right].
$$

In words, this reformulation translates the equivalence between (the "classes" referring here to the level sets of $f_{\mathrm{enc}}$)

- minimizing the intraclass dispersion (4): the distribution of configurations $x \in \mathcal{X}$ for a fixed value $z$ of $f_{\mathrm{enc}}$ should concentrate around the mean value $g^\star_{f_{\mathrm{enc}}}(z)$ by having a variance as small as possible;
- maximizing the interclass dispersion (7): the values of the conditional averages of $X$ for fixed values of $f_{\mathrm{enc}}$ should be as spread out as possible over the range of $f_{\mathrm{enc}}$.

1.3. **Formal characterization of the optimal encoding function.** The alternative optimization problem (6) allows to characterize the optimal encoding function $f_{\mathrm{enc}}$ by some orthogonality condition similar to the self-consistency condition of principal curves, see [8, Section 2]. In fact, it is formally shown in [10] that critical points of (6) satisfy

$$
(8)
$$
$$
\forall j \in \{1, \dots, d\}, \quad \forall x \in \mathrm{Supp}(\mu), \qquad \left[x - g^\star_{f_{\mathrm{enc}}}(f_{\mathrm{enc}}(x))\right]^\top \partial_{z_j} g^\star_{f_{\mathrm{enc}}}(f_{\mathrm{enc}}(x)) = 0,
$$

where $\mathrm{Supp}(\mu)$ is the support of the probability measure $\mu$. The derivation of this condition, performed using the co-area formula [5, 1] together with the use of weak derivatives, can be seen as a variation of derivations of optimality conditions for principal curves, as written already in [9]; see also [8] where (8) is used to construct a new objective function to minimize in order to find $f_{\mathrm{enc}}$.

An interesting implication of (8) is that the intersection of $\mathrm{Supp}(\mu)$ and the submanifold

$$
(9) \qquad \Sigma_z = f_{\mathrm{enc}}^{-1}\{z\} = \{x \in \mathcal{X} \,|\, f_{\mathrm{enc}}(x) = z\}
$$

is in fact included in the $(D - d)$-dimensional hyperplane containing the point $g^\star_{f_{\mathrm{enc}}}(z)$ and orthogonal to the vectors $\partial_{z_1} g^\star_{f_{\mathrm{enc}}}(z), \dots, \partial_{z_d} g^\star_{f_{\mathrm{enc}}}(z)$ (recalling that $\mathcal{X}$ and $\mathcal{Z}$ have dimensions $D$ and $d$, respectively). As these hyperplanes generally have a non-empty intersection, finding a regular function $f_{\mathrm{enc}}$ which satisfies (8) is only possible for distributions $\mu$ which have a support sufficiently concentrated around the principal manifold.

We finally discuss the limit $\beta \to +\infty$ in (8), when the probability measure $\mu$ under consideration is the Boltzmann–Gibbs measure with a density proportional to $\mathrm{e}^{-\beta V(x)}$, with $V$ the potential energy function; and the latent space is one-dimensional. We consider two local minima $x_A$ and $x_B$ of the potential energy

function $V$, located respectively on $\Sigma_{z_A}$ and $\Sigma_{z_B}$ with $z_A = f_{\mathrm{enc}}(x_A)$ and $z_B = f_{\mathrm{enc}}(x_B)$ (assuming $z_A \leqslant z_B$ without loss of generality). Similarly to the discussion in [15] for principal curves, it can formally be shown in the low temperature limit that the decoder path $\{g^{\star}_{f_{\mathrm{enc}}}(z)\}_{z \in [z_A, z_B]}$ converges to a minimum energy path (MEP).

## 2. Extensions and generalizations

To conclude this brief presentation of autoencoders and some of their properties, we list some generalizations, described more precisely in [10].

In a situation where multiple transition paths link two metastable states, the autoencoder may fail to properly represent the system in the transition region between local minima for a one-dimensional latent space $\mathcal{Z}$, as it constructs only a single curve for the conditional expectations. An idea to address this issue is to consider multiple decoders associated with a common encoder, and to choose for a given configuration the decoder which best reconstructs the state through some assignment function reminiscent of the one considered for clustering.

Another idea is to put more emphasis on transition states to better describe the transition from one metastable state to another. This can be done by changing the reference measure from the Boltzmann–Gibbs measure to a probability measure putting more mass on regions between metastable states, such as the reactive trajectory measure [4, 12] (which is the distribution of configurations sampled by portions of trajectories switching from one metastable state to another).

A last idea is to use extra physical information encoded via additional terms in the loss functions (as considered in [13] for instance); see [10, Sections 2.6 and 2.7].

## References

[1] L. Ambrosio, N. Fusco and D. Pallara, *Functions of Bounded Variation and Free Discontinuity Problems*, Oxford Science Publications, (2000)

[2] F. Bach, *Learning Theory from First Principles*, MIT Press (2023)

[3] C. Chipot and A. Pohorille (editors), *Free Energy Calculations*, Springer Series in Chemical Physics **86** (2007)

[4] W. E and E. Vanden–Eijnden, *Towards a theory of transition paths*, J. Stat. Phys. **123** (2006), 503–523

[5] L.C. Evans and R.F. Gariepy, *Measure Theory and Fine Properties of Functions*, Studies in Advanced Mathematics, CRC Press (1992)

[6] S. Gerber, *Saddlepoints in unsupervised least squares*, arXiv preprint **2104.05000** (2021)

[7] S. Gerber, T. Tasdizen and R. Whitaker, *Dimensionality reduction and principal surfaces via kernel map manifolds*, 2009 IEEE 12th International Conference on Computer Vision (2009), 529–536

[8] S. Gerber and R. Whitaker, *Regularization-free principal curve estimation*, J. Mach. Learn. Res. **14** (2012), 1285–1302

[9] T. Hastie and W. Stuetzle, *Principal curves*, J. Amer. Statist. Assoc. **84** (1989), 502–516

[10] T. Lelièvre, T. Pigeon, G. Stoltz and W. Zhang, *Analyzing multimodal probability measures with autoencoders*, J. Phys. Chem. B **128**(11) (2024), 2607–2631

[11] T. Lelièvre, M. Rousset and G. Stoltz, *Free-energy Computations: A Mathematical Perspective*, Imperial College Press (2010)

[12] J. Lu and J. Nolen, *Reactive trajectories and the transition path process*, Probab. Theory Related Fields **161** (2015), 195–244

[13] M. Ramil, C. Boudier, A. Goryaeva, M.C. Marinica and J.-B. Maillet, *On sampling minimum energy path*, J. Chem. Theory and Comput. **18** (2022), 5864–5875
[14] R. Tibshirani, *Principal curves revisited*, Stat. Comput. **2** (1992), 183–190
[15] E. Vanden-Eijnden and M. Venturoli, *Revisiting the finite temperature string method for the calculation of reaction tubes and free energies.* J. Chem. Phys. **130** (2009), 194103

*Reporters: Eugen Bronasco, Peter A. Whalley*

# Participants

**Eugen Bronasco**
Section de Mathématiques
Université de Geneve
UNI DUFOUR
24, rue du Général Dufour
P.O. Box 64
1211 Genève 4
SWITZERLAND

**Prof. Dr. Elena Celledoni**
Department of Mathematical Sciences
Norwegian University of Science and
Technology
A. Getz vei 1
7491 Trondheim
NORWAY

**Dr. Neil Chada**
Department of Mathematics
Heriot-Watt University
Riccarton
Edinburgh EH14 4AS
UNITED KINGDOM

**David Cohen**
Mathematical Sciences
Chalmers University of Technology and
University of Gothenburg
412 96 Göteborg
SWEDEN

**Dr. Sonja G. Cox**
Korteweg-de Vries Institute for
Mathematics
Amsterdam University
Science Park 105
Postbus 94248
1090 GE Amsterdam
NETHERLANDS

**Dr. Carlos Esteve-Yagüe**
Centre for Mathematical Sciences
University of Cambridge
Wilberforce Road
Cambridge CB3 0WB
UNITED KINGDOM

**Prof. Dr. Georg A. Gottwald**
School of Mathematics and Statistics
The University of Sydney
Sydney 2206
AUSTRALIA

**Dr. Ludovic Goudenège**
Fédération de Mathématiques de
CentraleSupélec
Bâtiment Bouygues
9, rue Joliot Curie
91190 Gif-sur-Yvette
FRANCE

**Prof. Dr. Lyudmilla Grigoryeva**
Fachgruppe für Mathematik
und Statistik
University of St. Gallen
Bodanstr. 6
9000 St. Gallen
SWITZERLAND

**Ramona Häberli**
Section de Mathématiques
Université de Geneve
UNI DUFOUR
24, rue du Général Dufour
P.O. Box 64
1211 Genève 4
SWITZERLAND

**Prof. Dr. Carsten Hartmann**
Institut für Mathematik
Brandenburgische Technische Universität
Cottbus-Senftenberg
Konrad-Wachsmann-Allee 1
03046 Cottbus
GERMANY


**Dr. Juncai He**
Computer, Electrical and Mathematical
Science & Engineering Division
King Abdullah University of Science and
Technology
Jeddah 23955-6900
SAUDI ARABIA


**Dr. Franca Hoffmann**
California Institute of Technology
1200 E California Blvd
Pasadena, CA 91125
UNITED STATES


**Prof. Dr. Shi Jin**
Institute of Natural Sciences
Shanghai Jiao Tong University
No. 800 Dongchuan Road
Minhang District
Shanghai Shi 200 240
CHINA


**Aikaterini Karoni**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Bldg.
King's Buildings, Mayfield Road
Edinburgh EH9 3FD
UNITED KINGDOM


**George Kevrekidis**
Department of Applied Mathematics
and Statistics,
Johns Hopkins University
3400 North Charles Street
Baltimore, MD 21218-2689
UNITED STATES


**Prof. Dr. Annika Lang**
Department of Mathematical Sciences
Chalmers University of Technology
and the University of Gothenburg
412 96 Göteborg
SWEDEN


**Dr. Adrien Laurent**
Department of Mathematics
University of Bergen
Allégaten 41
5007 Bergen
NORWAY


**Prof. Dr. Benedict Leimkuhler**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Bldg.
King's Buildings, Mayfield Road
Edinburgh EH9 3FD
UNITED KINGDOM


**Prof. Dr. Tony Lelièvre**
CERMICS – ENPC
Cité Descartes, Champs-sur-Marne
6 et 8 avenue Blaise Pascal
77455 Marne-la-Vallée Cedex 2
FRANCE


**Dr. Alexander Lewis**
Institut für Mathematische Stochastik
Georg-August-Universität Göttingen
Goldschmidtstr. 7
37077 Göttingen
GERMANY


**Nicolas Masson**
Section de mathématiques
Université de Genève
Rue du Conseil-Général 9/7
1205 Genève
SWITZERLAND

**Dr. Tingwei Meng**
Department of Mathematics
University of California at
Los Angeles
405 Hilgard Avenue
Los Angeles, CA 90095-1555
UNITED STATES


**Prof. Dr. Klaus-Robert Müller**
Computer Science Department
Technical University of Berlin
Sekr. MAR 4-1
Marchstraße 23
10587 Berlin
GERMANY


**Stefan Oberdörster**
Institute for Applied Mathematics
University of Bonn
Endenicher Allee 60
53115 Bonn
GERMANY


**Prof. Dr. Juan-Pablo Ortega**
Division of Mathematical Sciences
School of Physical & Mathematical
Sciences
Nanyang Technological University
21 Nanyang Link
Singapore 637 371
SINGAPORE


**Dr. Michela Ottobre**
Department of Mathematics
Heriot-Watt University
Riccarton
Edinburgh EH14 4AS
UNITED KINGDOM


**Prof. Dr. Brynjulf Owren**
Department of Mathematical Sciences
NTNU
7491 Trondheim
NORWAY


**Dr. Daniel Paulin**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Bldg.
King's Buildings, Mayfield Road
Edinburgh EH9 3JZ
UNITED KINGDOM


**Prof. Dr. Grigorios Pavliotis**
Department of Mathematics
Imperial College London
Huxley Building
180 Queen's Gate
London SW7 2AZ
UNITED KINGDOM


**Dr. Anna Persson**
Department of Information Technology
University of Uppsala
P.O. Box 120
Lägerhyddsvägen 1
75237 Uppsala
SWEDEN


**Prof. Dr. Holger Rauhut**
Mathematisches Institut
LMU München
Theresienstr. 39
80333 München
GERMANY


**Prof. Dr. Sebastian Reich**
Institut für Mathematik
Universität Potsdam
Karl-Liebknecht-Straße 24-25
14476 Potsdam
GERMANY


**Dr. Matthias Sachs**
School of Mathematics
The University of Birmingham
Edgbaston
Birmingham B15 2TT
UNITED KINGDOM

**Dr. Katharina Schuh**
Institut für Analysis und Scientific
Computing
Technische Universität Wien
Wiedner Hauptstrasse 8 - 10
1040 Wien
AUSTRIA


**Simon Schwarz**
Institut für Mathematische Stochastik
Georg-August-Universität Göttingen
Goldschmidtstraße 7
37077 Göttingen
GERMANY


**Dr. Gabriel Stoltz**
CERMICS – ENPC
Cité Descartes, Champs-sur-Marne
6/8 Avenue Blaise Pascal
77455 Marne-la-Vallée Cedex 2
FRANCE


**Prof. Dr. Anja Sturm**
Institute of Mathematical Stochastics
University of Göttingen
Goldschmidtstraße 7
37077 Göttingen
GERMANY


**Dr. Molei Tao**
School of Mathematics
Georgia Institute of Technology
686 Cherry Street
Atlanta, GA 30332
UNITED STATES


**Prof. Dr. Yen-Hsi Richard Tsai**
Department of Mathematics
The University of Texas at Austin
1 University Station C1200
Austin, TX 78712-1082
UNITED STATES

**Prof. Dr. em. Sara van de Geer**
Seminar für Statistik
ETH Zürich
Rämistrasse 101
8092 Zürich
SWITZERLAND


**Dr. Gilles Vilmart**
Section de Mathématiques
Université de Genève
UNI DUFOUR
24, rue du Général Dufour
P.O. Box 64
1211 Genève 4
SWITZERLAND


**Dr. Tiffany Vlaar**
Department of Mathematics
University of Glasgow
University Gardens
Glasgow G12 8QW
UNITED KINGDOM


**Prof. Dr. Sven Wang**
Institut für Reine Mathematik
Fachbereich Mathematik
Humboldt-Universität Berlin
10099 Berlin
GERMANY


**Prof. Dr. Xu Wang**
Chinese Academy of Science
LSEC, ICMSEC, Academy of
Mathematics and Systems Science
55 Zhongguancun East Road
Beijing 100190
CHINA


**Prof. Dr. Rachel Ward**
Department of Mathematics
University of Texas at Austin
2515 Speedway
Austin, TX 78712
UNITED STATES

**Peter Whalley**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Bldg.
King's Buildings, Mayfield Road
Edinburgh EH9 3JZ
UNITED KINGDOM

**Prof. Dr. Konstantinos Zygalakis**
School of Mathematics
University of Edinburgh
James Clerk Maxwell Bldg.
Edinburgh EH9 3FD
UNITED KINGDOM

**Dr. Caroline Wormell**
School of Mathematics & Statistics
The University of Sydney
Sydney NSW 2006
AUSTRALIA