

Proper Scoring Rules, Gradients, Divergences, and Entropies for Paths and Time Series

Patric Bonnier* and Harald Oberhauser†

Mathematical Institute, University of Oxford

Abstract

Many forecasts consist not of point predictions but concern the evolution of quantities. For example, a central bank might predict the interest rates during the next quarter, an epidemiologist might predict trajectories of infection rates, a clinician might predict the behaviour of medical markers over the next day, etc. The situation is further complicated since these forecasts sometimes only concern the approximate “shape of the future evolution” or “order of events”. Formally, such forecasts can be seen as probability measures on spaces of equivalence classes of paths modulo time-parametrization. We leverage the statistical framework of proper scoring rules with classical mathematical results to derive a principled approach to decision making with such forecasts. In particular, we introduce notions of gradients, entropy, and divergence that are tailor-made to respect the underlying non-Euclidean structure.

1 Introduction

Scoring rules provide a principled approach to form and evaluate probabilistic predictions. The earliest applications go at least back to the evaluation of weather forecasts [1], but have since then developed into a rich theoretical framework that plays a central part in modern statistical inference. We refer to [2] and [3] for a general background. The theoretical underpinnings of scoring rules are well-developed, but nearly all of the literature focuses on prediction of vector-valued or scalar-valued quantities. The aim of this article is develop a scoring rule framework for an important class of non-Euclidean data, namely sequential data – both in discrete and continuous time.

The Drawbacks of (Naïve) Vectorization. Given a dataset consisting of multivariate time series (TS), a common approach is to flatten each TS into a long vector and then use a standard pipeline for vector-valued data. However, this approach has several drawbacks. Firstly, there’s trouble whenever the different TS are irregularly sampled or are of different length since this embeds the different TS in Euclidean spaces of different

* bonnier@maths.ox.ac.uk

† oberhauser@maths.ox.ac.uk

dimension. Typically this addressed by ad-hoc approaches such as adding synthetic data by interpolation or dropping data points. Secondly, often the relevant information is independent of the time-parametrization (“*time-warping invariance*”), at least to a large degree; for instance, the meaning of a spoken word or an object being filmed are both independent of how fast or slow the audio signal or video is presented. Finally, many models are naturally formulated in continuous time rather than discrete time; for example, stochastic differential equations form a popular class of models in many applications and it is unclear how to evaluate such continuous time models in a scoring rule framework besides above naive vectorization on an arbitrary time grid.

A Non-Euclidean Data Domain. Key to our approach is that classic tools from pure mathematics faithfully capture the Non-Euclidean structure of the space of (unparametrized) paths. While there is no linear structure that allows for addition of paths of different length, any two paths can be concatenated into one path and any path can be run backwards. Both these operations – *concatenation and reversal* – are independent of the choice of parametrization, hence they also apply to equivalence classes of paths under reparametrization. We refer to an unparametrized path – that is an equivalence class of paths under reparametrization – as a *track*, as it is defined uniquely by the track it carves out in the space where it evolves. A classical result [4] is that there is a “*feature map*” from the set of tracks into a linear space that is *functorial* and *universal*; the former means that operations on tracks (concatenation and reversal) turn into algebraic operations in feature space, the latter means that any function of tracks can be approximated as a linear functional of this map. Moreover, this map is given as a series of iterated integrals which makes it amenable to computation and we refer to it as the *signature map*. In fact, the co-domain of this feature map (“*the feature space*”) is not only a linear space but forms a so-called Hopf algebra and it is the Hopf algebra structure that captures operations on tracks as algebraic operations. The third mathematical ingredient that we use are gradients of functions of tracks: the usual definition of linear (Fréchet) differentiability can sometimes be unsuitable for such functions due to the above lack of linear structure. However, Pansu generalized classical differentiation to a special class of groups and we leverage this to define *gradients of functions of (unparametrized) paths*. We show that the usual guarantees of gradient descent algorithms apply which allows us to compute quantities associated with our scoring rule framework (as even in the case of Euclidean data, many quantities are not given in closed form, but can be found by first order methods).

Outline. Section 2 recalls the basic definitions and general of scoring rules. Section 3 contains the theoretical background; it formalizes the structure of the spaces of (unparametrized paths), and introduces the signature feature map and the Hopf algebra structure of its co-domain (the feature space). Section 4 then shows how these quantities lead to natural scoring rules on the non-Euclidean space of tracks; in particular the so-called antipode of the Hopf algebra plays key role to relate the scoring rule framework to structural properties of tracks. From general results this then immediately leads to definitions of entropy, divergence, and mutual information that – unlike the (naïve) vectorization approach outlined above – are compatible with the structure of (probability measures) on

spaces of (unparametrized) paths. Section 5 then utilizes that one may identify a track as “group-like” element and shows that the concept of Pansu differentiability leads to a natural notion of gradient descent on the space of paths resp. tracks. Finally, Section 6 demonstrates that despite this approach being motivated by pure maths, the resulting quantities lead to efficiently computable quantities that have some advantageous properties compared to other methods with similar invariances.

1.1 Related Work

One of the earliest empirical insights for time series data was that time-parametrization (“time warping”) invariance is of key importance [5, 6]. Arguably the most popular way to address this invariance is via the classical dynamic time warping distance (DTW) and its many variations that introduce a distance between time series by searching over time changes. For example, [7, 8] introduce a regularised version of DTW, so-called soft DTW (s-DTW) which addresses the fact that the DTW distance is not differentiable, making it viable for use in deep learning pipelines. In the process of doing so it loses the invariance that DTW enjoys and introduces a trade-off of smoothness versus invariance. A more general point is that DTW and its variations do not aim to provide the full forecasting framework of scoring rules (divergence between measures, entropy, mutual information of TS) and although DTW approaches successfully deal with time-parametrization, they ignore other structural properties such as concatenation and reversal of TS. Another drawback is that while the focus of DTW is on discrete time it can be formulated in continuous time (so-called Fréchet distance) but the computation scales with quadratic complexity in the number of sequence entries which makes it too expensive for many sources of high-frequency data, whereas our distance can be computed in linear time for the price of higher complexity in the state space dimension. Ultimately the reason for this increase in efficiency is that the time-warpings are never explicitly computed or exhibited.

Another area that is directly related is kernel learning. Any kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ with reproducing kernel Hilbert space H_k induces a scoring rule by setting $s(x, \mu) := \|\delta_x - \mu\|_{H_k}^2$ where $\|\mu\|_{H_k}^2 = \int \int k(x, y) \mu(dx) \mu(dy)$, see [9, Section 5] and several kernels for sequences have been developed in the literature. Most relevant to our approach is the “signature kernel” introduced in [10]. However, for any scoring rule given by a kernel, the (generalized) divergence becomes simply the maximum mean discrepancy and the entropy simply the variance in the RKHS H_k . While kernels give rise to a powerful class of scoring rules, the success and popularity of non-kernel based scoring rules on $\mathcal{X} = \mathbb{R}^n$ is motivation enough to look for other interesting, non-linear scoring rules.

The technical key to our approach comes from mathematics where iterated integrals, so-called signatures, and non-commutative algebras are used to represent paths. This goes at least back to seminal work of Chen [4] in algebraic topology and subsequent applications in control theory [11, 12] and more recently rough path theory [13]. These results have been influential in stochastic analysis [14, 15] and only more recently have been started to be explored in a statistical and machine learning context. We mention pars-pro-toto [16, 17, 18] for inference about laws of stochastic processes; [10, 19, 20] for kernel learning; [21, 22, 23] for Bayesian approaches; [24, 25, 26] for generative modelling; [27, 28, 29] for applications in topology; [30, 31, 32] for algebraic perspectives.

Finally, we mention that the two topics that are central to us – invariances and non-Euclidean structure – have been considered in different contexts in scoring rule frameworks. For example, [33] studies equi- and in-variances for scoring rules for Euclidean data; non-vector valued data such as sets, contours, intervals, and quantiles have received attention [34, 35, 36].

2 Proper Scoring Rules, Entropies, and Divergences

We briefly recall general background on scoring rules following closely the notation in [3], see also [2, 37, 38, 39]. Let \mathcal{X} be a measurable space (*the outcome space*), \mathcal{A} be a set (*the action space*), and $\mathcal{L} : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ be a function (*the loss function*). Further, let X be a \mathcal{X} -valued random variable. We consider the following game between a Decision-maker and Nature: the task of the decision maker is to choose an action $a \in \mathcal{A}$, after which Nature reveals the outcome $x \in \mathcal{X}$ that is given by sampling X . The decision maker then suffers the loss $\mathcal{L}(x, a)$.

Given a set \mathcal{P} of probability measures on \mathcal{X} , a principled probabilistic (Bayesian) approach to this decision problem is to proceed as follows:

- (I) Associate with every $\mu \in \mathcal{P}$ its *Bayes act* $a_\mu \in \mathcal{A}$ defined by

$$a_\mu := \arg \min_{a \in \mathcal{A}} \mathbb{E}_{X \sim \mu} [\mathcal{L}(X, a)]$$

(assuming that a minimum exists; if it is not unique, choose a_μ arbitrary among the minimizers).

- (II) Use the Bayes act a_μ to define the *scoring rule* $s : \mathcal{X} \times \mathcal{P} \rightarrow \mathbb{R}$ on \mathcal{X} given by

$$s(x, \mu) := \mathcal{L}(x, a_\mu). \tag{1}$$

- (III) Use the scoring rule s to define the (generalised) *entropy* H , the (generalised) *divergence* $d : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$, and the (generalised) *mutual information* $I : \mathcal{P} \times \mathcal{Q} \rightarrow \mathbb{R}$ as

$$\begin{aligned} H : \mathcal{P} &\rightarrow \mathbb{R}, & \mu &\mapsto \mathbb{E}_{X \sim \mu} [s(X, \mu)], \\ d : \mathcal{P} \times \mathcal{P} &\rightarrow \mathbb{R}, & (\nu, \mu) &\mapsto \mathbb{E}_{X \sim \nu} [S(X, \mu)] - H(\nu), \\ I : \mathcal{P} \times \mathcal{Q} &\rightarrow \mathbb{R}, & (\mu, \nu) &\mapsto H(\mu) - \mathbb{E}_{U \sim \nu} [H(\mu|U)]. \end{aligned}$$

where \mathcal{Q} denotes a space of probability measures (not necessarily on \mathcal{X}) and $\mu|U$ denotes the law of μ conditioned on the random variable U .

The above definitions and nomenclature is justified as follows: firstly, it is an instructive exercise to check that for standard choices of state space \mathcal{X} , action space \mathcal{A} , and loss function \mathcal{L} , the above reduce to classical definitions of entropy, divergence and mutual information (e.g. if \mathcal{A} is the set of densities on \mathbb{R}^n then the *log score* $L(x, a) := -\log a(x)$ from [40] yields the the usual Shannon entropy, Kullback-Leibler divergence, and mutual information); see [3] for more examples. Secondly, Theorem 2.1 below shows that characteristic properties hold in the full generality of the above setup:

Theorem 2.1 ([3]). *Let \mathcal{X} , \mathcal{A} , and L be as above. Further, denote with H , I , and d the associated (generalized) scoring rule, entropy, and divergence. Then*

1. *the scoring rule (1) is proper, that is*

$$\mu \mapsto \mathbb{E}_{X \sim \nu}[s(X, \mu)]$$

is minimized at $\mu = \nu$.

2. *$\mu \mapsto H(\mu)$ is concave,*
3. *$(\nu, \mu) \mapsto \mathbb{E}_{X \sim \nu}[s(X, \mu)]$ is affine in ν for every μ ,*
4. *$d(\mu, \nu) \geq 0$ with equality for $\mu = \nu$,*
5. *$I(\mu, \nu) \geq 0$ with equality for $\mu \perp \nu$.*

Different applications areas demand different scoring rules. Classical choices for the Euclidean case $\mathcal{X} = \mathbb{R}^n$ are besides the already mentioned log score, the Brier score, the Tsallis score, the Bregman score, the Hyvärinen score, etc.; see [3] for details. The aim of the remainder of this article is to study the case when \mathcal{X} is a space of paths or a space of equivalence classes of paths (under reparametrisation).

A Toy Example: From Feature Maps to Bayes Actions. To motivate our scoring rule for paths let us first revisit the vector-valued case, $\mathcal{X} = \mathbb{R}^n$. To arrive at a proper scoring rule, the space of Bayes actions \mathcal{A} should be large enough to characterize any (sufficiently nice) probability measures on $\mathcal{X} = \mathbb{R}^n$. A classic way to characterize a probability measure, is to consider the sequence of moments,

$$(1, \mathbb{E}[X], \mathbb{E}[X^{\otimes 2}], \mathbb{E}[X^{\otimes 3}], \dots) \tag{2}$$

that is, $\mathbb{E}[X]$ is the mean vector, $\mathbb{E}[X^{\otimes 2}]$ is the covariance matrix, etc. (We tacitly assume that the sequence of moments is well-defined and decays quickly enough so that it characterizes the measure, see Remark 2.2). The sequence (2) is an element of the set

$$\mathcal{H} := \prod_{m \geq 0} (\mathbb{R}^n)^{\otimes m} \tag{3}$$

of sequences of tensor of increasing degree m . In fact, this set \mathcal{H} forms a vector space by element-wise addition of tensors of the same degree. If we define the “feature map”

$$\varphi : \mathbb{R}^n \rightarrow \mathcal{H}, \quad x \mapsto (1, x, x^{\otimes 2}, \dots).$$

With the above notation, the moment sequence (2) is simply the mean of $\varphi(X)$,

$$\mathbb{E}[\varphi(X)] = (1, \mathbb{E}[X], \mathbb{E}[X^{\otimes 2}], \mathbb{E}[X^{\otimes 3}], \dots) \in \mathcal{H}$$

Using a well-known characterization of the mean as minimizer of a quadratic we can introduce the loss function

$$\mathcal{L}(x, a) := \|\varphi(x) - a\|^2$$

which associates with a probability measure μ on \mathbb{R}^n the Bayes action

$$a_\mu \equiv \mathbb{E}[\varphi(X)] \equiv \operatorname{argmin}_{m \in \mathcal{H}} \mathbb{E}[\mathcal{L}(X, a)].$$

It follows from general principles that the resulting scoring rule on the state $\mathcal{X} = \mathbb{R}^n$ and action space $\mathcal{A} = \mathcal{H}$,

$$s(x, \mu) = \mathcal{L}(x, a_\mu)$$

is proper. Despite the elementary nature of this example it gives us a simple way to associate with any “feature map” a Bayes action and a scoring rule and already simple variations lead to interesting questions: for example, if the quadratic loss function is replaced by the absolute value, one ends with medians of moment as Bayes action and many other choices are possible. Such questions fall under the framework of “elicitation” of properties of probability measures with scoring rules which is an active research area, already in the classical vector-valued (even scalar) case; see [41] and [42, 43, 44] for some of the recent advances.

Remark 2.2. The question which probability measures on \mathbb{R}^n are characterized by the moment sequence (2) is classical but quite subtle in general. But for compactly supported measures this trivially holds. Our focus will soon shift to probability measures on pathspace but since spaces of paths are generically not even locally compact, compactness is a too strong assumption; in fact, important examples of measures on pathspace such as geometric Brownian motion are not characterized by their “signature moments” that we will use in Section 3 and Section 4. However, one can replace the moment sequence (2) by a normalized moment sequence that characterizes any probability measure on \mathbb{R}^n and this extends to path space and signature moments, see [17] for details. Hence, for simplicity, we assume throughout that the probability measures are characterized by their expected feature map (since this is possible by a slight modification of the feature map). \square

3 Structure of the Space of (Unparametrized) Paths

We review classic mathematical results about spaces of paths going back to seminal work of Chen [4]. The main result is the existence of a “feature map”

$$\Phi : \text{Tracks} \rightarrow \mathcal{H}, \quad \mathbf{x} \mapsto \Phi(\mathbf{x}) \tag{4}$$

that has as domain the set Tracks that consist of equivalence classes of paths that evolve in \mathbb{R}^n , and as co-domain the linear space \mathcal{H} .

We already encountered \mathcal{H} in Section 2 where it arose as the vector space of sequences of tensors $\mathbf{t}_m \in (\mathbb{R}^n)^{\otimes m}$ of increasing degree m , see (3). However, \mathcal{H} is not only a vector space but a so-called Hopf algebra: we can multiply elements of \mathcal{H} and take the “inverse” of elements of \mathcal{H} . One of the well-known and attractive properties of the map (4) is that these two algebraic operations (multiplication and inversion) capture the natural operations on Tracks (concatenation and reversal). Exploiting this correspondence will be essential for our main results in Section 4.

The Domain of Paths. A bounded variation path¹ \mathbf{x} in \mathbb{R}^n is a continuous map

$$\mathbf{x} : [0, T] \rightarrow \mathbb{R}^n \text{ such that } \|\mathbf{x}\| := \sup_{(t_1, \dots, t_L): 0 \leq t_1 < \dots < t_L < T} \sum \|\mathbf{x}(t_{i+1}) - \mathbf{x}(t_i)\| < \infty.$$

For $a, b \in \mathbb{R}^n$ we denote with $\text{Paths}(a, b)$ the set of all continuous bounded variation paths that start at a and end in b ,

$$\text{Paths}(a, b) := \{\mathbf{x} | \mathbf{x} : [0, T] \rightarrow \mathbb{R}^n \text{ is of bounded variation, } T > 0, \mathbf{x}(0) = a, \mathbf{x}(T) = b\}.$$

and by

$$\text{Paths} := \bigcup_{a, b \in \mathbb{R}^n} \text{Paths}(a, b)$$

the set of all bounded variation paths in \mathbb{R}^n . Although Paths is not a linear space, it has a rich structure given by concatenation and time reversal. Informally, this says that if one can go from a to b and from b to c then one can go from a to c and that if one can go from a to b then one can go from b to a . Formally, concatenation and reversal are defined as

1. For $\mathbf{x} \in \text{Paths}(a, b)$, $\mathbf{y} \in \text{Paths}(b, c)$, their *concatenation* $\mathbf{x} \star \mathbf{y} \in \text{Paths}(a, c)$ is defined by

$$(\mathbf{x} \star \mathbf{y})(t) := \begin{cases} \mathbf{x}(t), & \text{if } t \in [a, b] \\ \mathbf{x}(b) - \mathbf{y}(b) + \mathbf{y}(t), & \text{if } t \in [b, b + c] \end{cases}$$

2. for any $\mathbf{x} \in \text{Paths}(a, b)$ there exists an *inverse path* $\overleftarrow{\mathbf{x}} \in \text{Paths}(b, a)$ defined as

$$\overleftarrow{\mathbf{x}}(t) := \mathbf{x}(b - t).$$

The Domain of Tracks As discussed above, often we want to ignore the time parametrization, hence the fundamental object we care about is not the set of paths but equivalence classes of paths. It turns out that is useful to work with slightly more general equivalence relation, namely that of *tree-like equivalence* \sim . We define

$$\text{Tracks}(a, b) := \text{Paths}(a, b) / \sim, \text{ and } \text{Tracks} := \text{Paths} / \sim$$

With slight abuse of notation, we use the same notation \mathbf{x} for an element of Paths and an element of Tracks but emphasize that an element $\mathbf{x} \in \text{Tracks}$ is a whole equivalence class of paths. We give the precise definition of the equivalence relation \sim in Appendix A and only note here that if two paths $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^n$, $\mathbf{y} : [0, S] \rightarrow \mathbb{R}^n$ differ by time-parametrization, that is $\mathbf{x}(t) = \mathbf{y}(\varphi(t))$ for every t and an increasing function $\varphi : [0, T] \rightarrow [0, S]$, then $\mathbf{x} \sim \mathbf{y}$. However, in addition to time parametrization, tree-like equivalence also identifies

¹For simplicity we focus on (equivalence classes of) bounded variation paths but all the results immediately extend to paths with much less regularity such as trajectories of stochastic differential equations or (fractional) Brownian motion by replacing the iterated Riemann–Stieltjes integrals by stochastic integrals or rough path integrals [13]

paths that backtrack all their excursions, see Appendix A. We invite readers to think of elements of Tracks like animal tracks in nature: they provide shape and direction but not the speed at which the track was made. In particular, we note that the above operations of concatenation and reversal are well-defined for the elements of Tracks; after all, they do not depend on the time-parametrization. So again, with slight abuse of notation we have concatenation and reversal map,

$$\star : \text{Tracks}(a, b) \times \text{Tracks}(b, c) \rightarrow \text{Tracks}(a, c) \text{ and } \overleftarrow{\bullet} : \text{Tracks}(a, b) \rightarrow \text{Tracks}(b, a).$$

The co-domain \mathcal{H} . Our first encounter of \mathcal{H} was in Section 2 as the state space of the moment map (3). However, a more abstract way to introduce is by identifying it as the free algebra over \mathbb{R}^n . Informally, this means we want to keep the vector space structure of \mathbb{R}^n but we also would like to have a multiplication and do this in the most general way possible. Formally, this means \mathcal{H} is the free algebra over \mathbb{R}^n . Despite this abstract characterization as a free object, the space \mathcal{H}_n has a very concrete form which we will take as its definition,

$$\mathcal{H} := \prod_{m \geq 0} (\mathbb{R}^n)^{\otimes m} := \{\mathbf{t} = (\mathbf{t}_0, \mathbf{t}_1, \mathbf{t}_2, \mathbf{t}_3, \dots) : \mathbf{t}_m \in (\mathbb{R}^n)^{\otimes m}\}.$$

(one can then directly verify that this indeed is the free algebra, see [45]). That is, an element \mathbf{t} of \mathcal{H} is sequence of tensors $(\mathbf{t}_0, \mathbf{t}_1, \mathbf{t}_2, \dots)$ of increasing degree m where by convention $(\mathbb{R}^n)^{\otimes 0} = \mathbb{R}$. The vector space structure of \mathcal{H} is simply given as element-wise addition: addition of $\mathbf{s}, \mathbf{t} \in \mathcal{H}$ is defined as

$$\mathbf{s} + \mathbf{t} = (\mathbf{s}_0 + \mathbf{t}_0, \mathbf{s}_1 + \mathbf{t}_1, \dots)$$

and their multiplication is defined by $(\mathbf{s} \cdot \mathbf{t})_m = \sum_k \mathbf{s}_k \otimes \mathbf{t}_{m-k}$, i.e

$$\mathbf{s} \cdot \mathbf{t} := (1, \mathbf{s}_1 + \mathbf{t}_1, \mathbf{s}_2 + \mathbf{s}_1 \otimes \mathbf{t}_1 + \mathbf{t}_2, \dots)$$

where \otimes denotes the usual tensor (outer) product. Like matrix multiplication, this multiplication is associative but in general not commutative, $\mathbf{s} \cdot \mathbf{t} \neq \mathbf{t} \cdot \mathbf{s}$ and it has as multiplicative unit $(1, 0, \dots) \in \mathcal{H}$,

$$\mathbf{t} \cdot (1, 0, 0, \dots) = (1, 0, \dots) \cdot \mathbf{t} = \mathbf{t}.$$

The existence of a unit for multiplication naturally leads to the question of the existence of inverses, that is for $\mathbf{t} \in \mathcal{H}$ can one find another element in \mathcal{H} , denoted by $\mathbf{t}^{-1} \in \mathcal{H}$, such that

$$\mathbf{t} \cdot \mathbf{t}^{-1} = \mathbf{t}^{-1} \cdot \mathbf{t} = (1, 0, 0, \dots).$$

This is true whenever $\mathbf{t}_0 \neq 0$, and moreover, $\mathbf{t} \mapsto \mathbf{t}^{-1}$ has the explicit formula

$$\mathbf{t}^{-1} = \frac{1}{\mathbf{t}_0} \left\{ \sum_{m \geq 0} \left(1 - \frac{1}{\mathbf{t}_0} \mathbf{t}\right)^{\otimes m} \right\}.$$

The Feature map $\Phi : \text{Tracks} \rightarrow \mathcal{H}$.

Definition 3.1. For $\mathbf{x} \in \text{Paths}$, $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^n$ define

$$\int d\mathbf{x}^{\otimes m} := \int_{0 \leq t_1 < \dots < t_m \leq T} d\mathbf{x}(t_1) \otimes \dots \otimes d\mathbf{x}(t_m) = \int \dot{\mathbf{x}}(t_1) \otimes \dots \otimes \dot{\mathbf{x}}(t_m) dt_1 \dots dt_m.$$

┘

It is known that if two paths $\mathbf{x}, \mathbf{y} \in \text{Paths}$ are tree-like equivalent, $\mathbf{x} \sim \mathbf{y}$, then $\int d\mathbf{x}^{\otimes m} = \int d\mathbf{y}^{\otimes m}$ for every $m \geq 0$, see [46]. In fact, for the case of reparametrization $\mathbf{x}(t) = \mathbf{y}(\tau(t))$ this follows immediately from the change of variables formula. With slight abuse of notation we now define $\int d\mathbf{x}^{\otimes m}$ for $\mathbf{x} \in \text{Tracks}$.

Definition 3.2. For $\mathbf{x} \in \text{Tracks}$, define

$$\int d\mathbf{x}^{\otimes m} := \int d\mathbf{x}_{\text{path}}^{\otimes m}$$

where $\mathbf{x}_{\text{path}} \in \text{Paths}$ is in the equivalence class of \mathbf{x} and $\int d\mathbf{x}_{\text{path}}^{\otimes m}$ is as in Definition 3.1. ┘

By [46] $\int d\mathbf{x}^{\otimes m}$ for $\mathbf{x} \in \text{Tracks}$ is well-defined in the sense that the choice of \mathbf{x}_{path} does not matter. We refer to the resulting map as the signature map (this is also known as the Chen–Fliess series or chronological exponential).

Definition 3.3. We call

$$\Phi : \text{Tracks} \rightarrow \mathcal{H}, \quad \mathbf{x} \mapsto \left(\int d\mathbf{x}^{\otimes m} \right)_{m \geq 0}$$

the signature map. ┘

A well-known key property of the map Φ is that concatenation and reversal in Tracks correspond to multiplication and inversion in \mathcal{H} . Further, the map Φ is universal (up to fixing the starting point of the track, which is why we fix the starting point $a \in \mathbb{R}^n$ and restrict to the domain $\bigcup_{b \in \mathbb{R}^n} \text{Tracks}(a, b)$) in the sense that it linearizes continuous functions on Tracks. We summarize all this in Theorem 3.4 below.

Theorem 3.4. For every $a \in \mathbb{R}^n$ the map

$$\Phi : \bigcup_{b \in \mathbb{R}^n} \text{Tracks}(a, b) \rightarrow \mathcal{H}, \quad \mathbf{x} \mapsto \left(\int d\mathbf{x}^{\otimes m} \right)_{m \geq 0}$$

is injective and

1. $\Phi(\mathbf{x} \star \mathbf{y}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})$
2. $\Phi(\overleftarrow{\mathbf{x}}) = \Phi(\mathbf{x})^{-1}$
3. for every $f \in C(\text{Tracks}, \mathbb{R})$, $\epsilon > 0$ there exists a linear functional $\ell \in \mathcal{H}^*$ such that

$$|f(\mathbf{x}) - \langle \ell, \Phi(\mathbf{x}) \rangle| < \epsilon$$

uniformly in \mathbf{x} on compacts.

Proof. This is a folk theorem in algebraic topology and control theory; see [4] and [11]. What is less standard is that we use the treelike equivalence from \sim from [46]. \square

Remark 3.5. The space $\mathcal{H} \equiv \prod_{m \geq 0} (\mathbb{R}^d)^{\otimes m}$ is graded by the tensor degree m , and $\Phi(\mathbf{x}) \equiv (\int dx^{\otimes m})_{m \geq 0}$ decays exponentially fast in m , that is

$$\left\| \int dx^{\otimes m} \right\| \leq \frac{\|\mathbf{x}\|^m}{m!}$$

(on the right hand side $\|\bullet\|$ denotes the bounded variation (semi-)norm, on the left-hand side it denotes the norm on $(\mathbb{R}^d)^{\otimes m}$). Hence, in practice one only needs to compute the first m iterated integrals of $\Phi(\mathbf{x})$. For piecewise linear tracks $\mathbf{x} \in \text{Tracks}$ – which is how we identify time series – the first m entries of the map $\Phi(\mathbf{x})$ can be computed in $O(Ld^M)$ computational steps: if a track is given by piecewise linear segments $v_1, \dots, v_L \in \mathbb{R}^n$ then

$$\left(\int d\mathbf{x}^{\otimes m} \right)_{m \in \{0, 1, \dots, M\}} = \exp(v_1) \cdots \exp(v_L),$$

where $\exp(v) := (1, v, \frac{v^{\otimes 2}}{2!}, \dots) \in \mathcal{H}$. Hence, for a low dimensional state space, the map $\Phi(\mathbf{x})$ can be approximately in time that scales linearly in the length of the path. \lrcorner

The Antipode in \mathcal{H} . The two operations of addition $\mathbf{s} + \mathbf{t}$ and multiplication $\mathbf{s} \cdot \mathbf{t}$ turn \mathcal{H} into a (non-commutative) algebra $(\mathcal{H}, +, \cdot)$. However, \mathcal{H} comes with a bit more structure, namely the so-called *antipode* map

$$\alpha : \mathcal{H} \rightarrow \mathcal{H}$$

which is defined as the linear function given by linear extensions of the map

$$(\mathbb{R}^n)^{\otimes m} \rightarrow (\mathbb{R}^n)^{\otimes m}, \quad v_1 \otimes \cdots \otimes v_m \mapsto (-1)^m v_m \otimes \cdots \otimes v_1.$$

There is an important subset G of \mathcal{H} defined by the property

$$G := \{g \in \mathcal{H} : \alpha(g) = g^{-1}\}.$$

It turns out that G in fact forms a group and will play an important role in our Bayes acts for the simple fact that the feature map Φ takes values in G . We summarize this along with some facts about α that we use later in the following lemma.

Lemma 3.6. *Let α be the antipode on \mathcal{H} . Then*

1. $\alpha^2 = 1$,
2. $\alpha(\mathbf{s} \cdot \mathbf{t}) = \alpha(\mathbf{t})\alpha(\mathbf{s})$,
3. If $\mathbf{x} \in \text{Tracks}$, then $\Phi(\mathbf{x}) \in G$,
4. For a power series $p(\mathbf{t}) = \sum_{m \geq 0} p_m \mathbf{t}^{\otimes m}$, it holds that $p \circ \alpha(\mathbf{t}) = \alpha \circ p(\mathbf{t})$ for any $\mathbf{t} \in \mathcal{H}$,

5. If $\mathbf{t} \in \mathcal{H}$ is invertible, then $\alpha(\mathbf{t}^{-1}) = \alpha(\mathbf{t})^{-1}$,

6. Let $f : \mathcal{H} \rightarrow \mathbb{R}$ have the form

$$f(\mathbf{t}) = \sum_m \sum_{i_1 \cdots i_m} f_m(\mathbf{t}_m^{i_1 \cdots i_m})$$

where we identify the degree- m tensor $\mathbf{t}_m \in (\mathbb{R}^d)^{\otimes m}$ of $\mathbf{t} = (\mathbf{t}_m)_{m \geq 0} \in \mathcal{H}$ with its coordinates $\mathbf{t}_m \simeq (\mathbf{t}_m^{i_1, \dots, i_m})_{i_1, \dots, i_m = 1, \dots, d}$. If $f_m(x) = f_m(-x)$ for every $m \geq 0$ and $x \in \mathbb{R}$, then

$$f(\mathbf{t}) = f(\alpha(\mathbf{t})).$$

Proof. Items 1 and 2 follow from the definition. Item 3 is well known, see for instance [47, Section 5]. To see item 4 note that

$$p \circ \alpha(x) = \sum_{n \geq 0} p_n(\alpha x)^{\otimes n} = \sum_{n \geq 0} p_n \alpha(x^{\otimes n}) = \alpha\left(\sum_{n \geq 0} p_n x^{\otimes n}\right) = \alpha \circ p(x),$$

by Item 2 and linearity. Item 5 follows since the inverse map $\mathbf{t} \mapsto \mathbf{t}^{-1}$ has the power series expansion

$$\mathbf{t}^{-1} = \frac{1}{\mathbf{t}_0} \left\{ \sum_{n \geq 0} \left(1 - \frac{1}{\mathbf{t}_0} \mathbf{t}\right)^{\otimes n} \right\},$$

see [48, Lemma 7.16] which together with Item 4 shows the claim. For item 6, we note that

$$f(\alpha x) = \sum_n \sum_{i_1 \cdots i_n} f_n((-1)^n x^{i_n \cdots i_1}) = \sum_n \sum_{i_1 \cdots i_n} f_n(x^{i_n \cdots i_1}) = f(x).$$

□

A simple example of a function that satisfies the requirements Item 6 in Lemma 3.6 is the sum of squares

$$L(\mathbf{t}) = \sum_m \sum_{i_1 \cdots i_m} |\mathbf{t}_m^{i_1 \cdots i_m}|^2$$

Which will be used to construct a loss function for tracks in Section 6.

From Discrete Time to Continuous Time. This section has so far focused on paths [resp. tracks], that evolves [resp. equivalence classes of evolutions] in continuous time $\mathbf{x} : [0, T] \rightarrow \mathbb{R}^n$. However, in practice one typically has only access to a discrete time observations $\mathbf{x}(t_1), \dots, \mathbf{x}(t_L) \in \mathbb{R}^n$ along some grid $0 \leq t_1 < \dots < t_L \leq T$, that is a time series. But any TS can be identified as the piecewise linear path

$$t \mapsto \mathbf{x}(t_i) + \frac{t_i - t}{t_{i+1} - t_i} (\mathbf{x}(t_{i+1}) - \mathbf{x}(t_i)) \text{ for } t \in [t_i, t_{i+1})$$

and hence also as an element of Tracks after forgetting the parametrisation. Working in continuous time when the original data is discrete might look cumbersome and unnecessary at first sight but it has several advantages. Firstly, all TS are embedded into the same space Paths respectively Tracks, even if the sample grid $t_1 < \dots < t_L$ varies from TS to TS, which would not be the case if one identifies TS as vectors. Secondly, this automatically ensures consistency in terms of high-frequency limits when the grid gets finer, that is $\sup |t_{i+1}^m - t_i^m| \rightarrow 0$ as $m \rightarrow \infty$ for a sequence $(t_i^m)_{i=1, \dots, L_m}$. Finally, many popular models are naturally formulated in continuous time rather than discrete time.

From Tracks to Paths. Our guiding philosophy is that the fundamental object is the set Tracks rather than the set Paths since the former allows to ignore the time parametrization; note that the set of time parametrisations is infinite-dimensional since every continuous function $\tau : [0, T] \rightarrow [0, T']$ can be used to reparametrize a path $t \mapsto \mathbf{x}(t)$ to $t \mapsto \mathbf{x}(\tau(t))$, hence working with Tracks factors out an infinite-dimensional class of invariances. Nevertheless, for certain applications the parametrization matters – at least to a certain degree. However, this can be easily addressed by adding time as an additional coordinate: to emphasize the dimension n of the state space \mathbb{R}^n in which the paths evolve we write Paths_n (instead of just Paths that we used until now); similarly Tracks_n for the set of equivalence classes of Paths_n . Given $\mathbf{x} \in \text{Paths}_n$ we embed

$$\text{Paths}_n \hookrightarrow \text{Paths}_{n+1}, \quad \mathbf{x} \mapsto (t \mapsto (t, \mathbf{x}(t))).$$

That is a path evolving in \mathbb{R}^n is turned into a path in \mathbb{R}^{n+1} by simply adding an additional coordinate that is time itself. This makes the parametrization part of the “shape” of the trajectory which in turn is exactly the information that distinguishes tracks, hence

$$\text{Paths}_n \hookrightarrow \text{Tracks}_{n+1}.$$

This injection shows that any scoring rule for tracks induces a scoring rule for paths.

4 Scoring Rules For Tracks and Paths

Motivated by the toy example in Section 2 with the moment feature map φ for data in \mathbb{R}^n , we now follow the analogous reasoning on the non-Euclidean space of tracks by using the feature map Φ instead of φ , note that neither the domain nor the image of Φ is a linear space as its image is the group G that is embedded into the linear space \mathcal{H} . Recall that in Section 3 we have seen that it is exactly the group structure that captures the structure of space of tracks of concatenation and reversal. This motivates us to

replace the additive inverse in $\varphi(X) - m$ by the group inverse to get $\Phi(\mathbf{X})m^{-1}$.

Our first main result Theorem 2.1 shows that this indeed leads to a proper scoring rule on the space of tracks and operations on Tracks turn into algebraic operations in decision space. Consequences of this result are Proposition 4.3 and Corollary 4.4 which show how the associated entropy on the space of tracks is invariant to time-reversal and behaves under conditioning on the past.

A Scoring Rule for Tracks. We need to introduce an additional space \mathcal{H}^* wedged between \mathcal{H} and G defined as the space of all elements $\mathbf{t} \in \mathcal{H}$ starting with a one, formally

$$\mathcal{H}^* = \{\mathbf{t} \in \mathcal{H} : \mathbf{t}_0 = 1\}.$$

Unlike \mathcal{H} , \mathcal{H}^* is not a vector space or a Hopf algebra, but it is a group like G while also being convex as a subset of \mathcal{H} in addition to being topologically closed – unlike the set of invertible elements of \mathcal{H} . We have the following sequence of inclusions

$$G \subseteq \mathcal{H}^* \subseteq \text{invertible elements of } \mathcal{H} \subseteq \mathcal{H}.$$

Definition 4.1. Let $L : \mathcal{H} \rightarrow \mathbb{R}$ be convex with a unique minimum at the unit $(1, 0, 0, \dots)$ of G . Define the left loss function as

$$\mathcal{L}^{\text{left}} : \text{Tracks} \times \mathcal{H}^* \rightarrow [0, \infty), \quad (\mathbf{x}, m) \mapsto L(m^{-1}\Phi(\mathbf{x})).$$

Applying step (I) from the scoring rule framework of Section 2, the left Bayes' act is defined as

$$a_\mu^{\text{left}} := \operatorname{argmin}_{m \in \mathcal{H}_n^*} \mathbb{E}[\mathcal{L}^{\text{left}}(m, \mathbf{X})].$$

Applying step (II) yields the proper scoring rule

$$s^{\text{left}}(\mathbf{x}, \mu) := \mathbb{E}[\mathcal{L}^{\text{left}}(a_\mu^{\text{left}}, \mathbf{x})].$$

Applying step (III) yields the (generalised) entropy, divergence, and mutual information

$$\begin{aligned} H^{\text{left}}(\mu) &:= \mathbb{E}_{\mathbf{X} \sim \mu} \mathcal{L}^{\text{left}}(a_\mu^{\text{left}}, \mathbf{X}) \\ d^{\text{left}}(\nu, \mu) &:= \mathbb{E}_{\mathbf{X} \sim \nu} \left[\mathcal{L}^{\text{left}}(a_\mu^{\text{left}}, \mathbf{X}) - \mathcal{L}^{\text{left}}(a_\nu^{\text{left}}, \mathbf{X}) \right] \\ I^{\text{left}}(\mu, \nu) &:= H(\mu) - \mathbb{E}_{U \sim \nu} [H(\mu|U)] \end{aligned}$$

on the output space $\mathcal{X} = \text{Tracks}$ and the action space $\mathcal{A} = \mathcal{H}^*$. Analogously we define the right loss function $\mathcal{L}^{\text{right}}(\mathbf{x}, m) := L(\Phi(\mathbf{x})m^{-1})$, right Bayes act a_μ^{right} and right scoring rule $s^{\text{right}}(\mathbf{x}, \mu)$ as well as right entropy, divergence, and mutual information. \lrcorner

The scoring rule framework of Definition 4.1 turns operations on tracks into algebraic operations in the decision space.

Theorem 4.2. *Let the output space be $\mathcal{X} = \text{Tracks}$, the action space $\mathcal{A} = \mathcal{H}^*$ and*

$$\mathcal{L}^{\text{left}} : \text{Tracks} \times \mathcal{H}^* \rightarrow [0, \infty) \text{ resp. } \mathcal{L}^{\text{right}} : \text{Tracks} \times \mathcal{H}^* \rightarrow [0, \infty)$$

the loss functions from Definition 4.1. The following properties hold

1. *If L is coercive, that is $L(\mathbf{t}) \rightarrow \infty$ whenever $\|\mathbf{t}\| \rightarrow \infty$, then for any Borel measure μ such that L is μ -integrable, both a_μ^{left} and a_μ^{right} exist. If L is strictly convex, then they are unique.*

2. If $\mu = \delta_{\mathbf{x}}$ then $a_{\mu}^{right} = a_{\mu}^{left} = \Phi(\mathbf{x})$

3. The Bayes' acts satisfy

$$\begin{aligned} a_{\nu \star \mu | \nu}^{left} &= \nu(\Phi) a_{\mu | \nu}^{left} \\ a_{\mu \star \nu | \nu}^{right} &= a_{\mu | \nu}^{right} \nu(\Phi) \end{aligned}$$

where $\nu(\Phi)$ denotes the pushforward measure of ν under Φ and by $\mu \star \nu | \nu$ we denote the law of $\mathbf{X} \star \mathbf{Y} | \mathbf{Y}$ where $\text{Law}(\mathbf{X}) = \mu$ and $\text{Law}(\mathbf{Y}) = \nu$.

4. If L satisfies $L(\mathbf{t}) = L(\alpha(\mathbf{t}))$, then

$$a_{\overleftarrow{\mu}}^{right} = \alpha(a_{\mu}^{left}), \quad a_{\overleftarrow{\mu}}^{left} = \alpha(a_{\mu}^{right})$$

where $\overleftarrow{\mu}$ denotes the measure μ given by running samples from μ backwards in time²

Proof. We give the proofs for the right Bayes' act as the proofs for the left Bayes' act is similar.

(1) We equip \mathcal{H}_n with its ℓ^2 norm,

$$\|\mathbf{t}\| = \sqrt{\sum_n |\mathbf{t}_n|^2}$$

which makes it into a separable Hilbert Space.

Fix some measure μ on \mathcal{X} and define the map $\psi : G_n \rightarrow \mathbb{R}$ by

$$\psi(x) = \mathbb{E}_{\mu} L(\Phi(\mathbf{X})x).$$

We want to show that ψ is convex, coercive and lower semicontinuous on $(\mathcal{H}_n, \|\cdot\|)$, as this guarantees the existence of a minimiser. This is because the unit ball of $(\mathcal{H}_n, \|\cdot\|)$ is weakly compact, hence we could choose some weakly compact and convex set C such that $\psi(x) > M$ outside of C , and since ψ is lower semicontinuous it is also weakly lower semicontinuous, and therefore since it is convex it achieves a minimum on C which must be a global minimum. Note that its minimiser must be $(a_{\mu}^{right})^{-1}$. It follows that if ψ is strictly convex, then the minimiser is unique.

Note that if L is (strictly) convex, then

$$\psi\left(\frac{1}{2}x + \frac{1}{2}y\right) = \mathbb{E}_{\mu} L\left(\frac{1}{2}\Phi(\mathbf{X})x + \frac{1}{2}\Phi(\mathbf{X})y\right) \leq (\leq) \frac{1}{2}\mathbb{E}_{\mu} L(\Phi(\mathbf{X})x) + \frac{1}{2}\mathbb{E}_{\mu} L(\Phi(\mathbf{X})y) = \frac{1}{2}\psi(x) + \frac{1}{2}\psi(y)$$

hence ψ is also (strictly) convex.

Note that $(\mathcal{H}_n, \|\cdot\|)$ is a Banach algebra, that is $\|\mathbf{t} \cdot \mathbf{s}\| \leq \|\mathbf{t}\| \cdot \|\mathbf{s}\|$. By taking multiplicative inverses, this implies that

$$\|\mathbf{t} \cdot \mathbf{s}\| \geq \frac{1}{\|\mathbf{s}^{-1}\|} \|\mathbf{t}\|.$$

²Formally $\mathbf{X} \sim \mu$ then $\overleftarrow{\mu}$ is defined as the law of $\overleftarrow{\mathbf{X}}$.

for any invertible element \mathbf{s} . As μ is Borel, and \mathcal{H}_n is a Polish space, μ is a Radon measure by [49, Theorem 7.1.7] and we may choose a compact set C such that $\mathcal{P}(\mathbf{X} \notin C) \leq \varepsilon$, and define M to be $\sup_{X \in C} \|\Phi(X)^{-1}\|$. Then on C , $\|\Phi(X)x\| \geq \frac{1}{M}\|x\|$, and since

$$\psi(x) = \mathbb{E}_\mu L(\Phi(\mathbf{X})x) \geq \mathbb{E}_{\mu|_C} L(\Phi(\mathbf{X})x)$$

and L is coercive, so is ψ .

To see that ψ is lower semicontinuous, note that for a sequence $x_k \rightarrow x$

$$\liminf_k \psi(x_k) = \liminf_k \mathbb{E}_\mu L(\Phi(\mathbf{X})x_k) \geq \mathbb{E}_\mu L(\Phi(\mathbf{X})x) = \psi(x)$$

by Fatous Lemma, the assertion follows.

(2) Since L is minimised at the unit it is clear that for $\mu = \delta_{\mathbf{x}}$, $a_\mu^{\text{left}} = a_\mu^{\text{right}} = \Phi(\mathbf{x})$ is optimal since $(a_\mu^{\text{left}})^{-1}\Phi(\mathbf{x}) = \Phi(\mathbf{x})(a_\mu^{\text{right}})^{-1} = 1$.

(3) For $\mathbf{Y} \sim \nu$ we have

$$\begin{aligned} a_{\mu \star \nu | \nu}^{\text{right}} &:= \operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} [L(\Phi(\mathbf{X} \star \mathbf{Y})m^{-1}) | \mathbf{Y}] = \\ &\operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} [L(\Phi(\mathbf{X})\Phi(\mathbf{Y})m^{-1}) | \mathbf{Y}] = \\ &\operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} [L(\Phi(\mathbf{X})m^{-1}) | \mathbf{Y}] := a_{\mu | \mathbf{Y}}^{\text{right}} \Phi(\mathbf{Y}). \end{aligned}$$

(4) If L satisfies $L(\mathbf{t}) = L(\alpha(\mathbf{t}))$, then

$$\begin{aligned} a_{\overleftarrow{\mu}}^{\text{right}} &:= \operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \overleftarrow{\mu}} L(\Phi(\mathbf{X})m^{-1}) = \\ &\operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} L(\alpha(\Phi(\mathbf{X}))m^{-1}) = \\ &\operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} L(\alpha(m^{-1})\Phi(\mathbf{X})) := \alpha(a_\mu^{\text{left}}). \end{aligned}$$

□

Entropy, Divergence, and Mutual Information on the Space of Tracks. We now focus on the (generalized) entropy, divergence and mutual information for probability measures on tracks that results from Definition 4.1.

Proposition 4.3. *For any two probability measures μ and ν on Tracks it holds that*

1. $H^{\text{left}}(\mu \star \nu | \mu) = H^{\text{left}}(\nu | \mu)$ and $H^{\text{right}}(\mu \star \nu | \nu) = H^{\text{right}}(\mu | \nu)$
2. If L satisfies $L(\mathbf{t}) = L(\alpha(\mathbf{t}))$, then

$$\begin{aligned} H^{\text{right}}(\mu) &= H^{\text{left}}(\overleftarrow{\mu}), \\ d^{\text{right}}(\nu, \mu) &= d^{\text{left}}(\overleftarrow{\nu}, \overleftarrow{\mu}), \\ I^{\text{right}}(\mu, \nu) &= I^{\text{left}}(\overleftarrow{\mu}, \overleftarrow{\nu}) \end{aligned}$$

Proof. (1)

$$\begin{aligned} H^{\text{left}}(\mu \star \nu | \mu) &= \mathbb{E}_{X \sim \mu \star \nu | \mu} L((a_{\mu \star \nu | \mu}^{\text{left}})^{-1} \Phi(X)) \\ &= \mathbb{E}_{X \sim \nu | \mu} L((a_{\nu | \mu}^{\text{left}})^{-1} \nu(\Phi)^{-1} \nu(\Phi) \Phi(X)) = H^{\text{left}}(\nu | \mu). \end{aligned}$$

(2)

$$\mathbb{E}_{X \sim \nu} L(\Phi(X)(a_{\mu}^{\text{right}})^{-1}) = \mathbb{E}_{X \sim \nu} L(\alpha(a_{\mu}^{\text{right}})^{-1} \alpha \Phi(X)) = \mathbb{E}_{X \sim \overleftarrow{\nu}} L((a_{\overleftarrow{\nu}}^{\text{left}})^{-1} \alpha \Phi(X)).$$

The other equalities follow. \square

Corollary 4.4. *If L satisfies $L(\mathbf{t}) = L(\alpha(\mathbf{t}))$ and a measure μ is reversible, that is, μ and $\overleftarrow{\mu}$ are equal up to their starting distribution, then*

$$H^{\text{right}}(\mu) = H^{\text{left}}(\mu).$$

In the experiments we will simulate sample paths from Brownian motion and since this is a reversible process it will not matter if we use the left- or right entropy.

Remark 4.5. An alternative to using the group structure in Definition 4.1 of the Bayes act is to use that the group G is embedded into the linear space \mathcal{H} and use this linear structure. That is, we define a Bayes act as $a_{\mu} := \operatorname{argmin}_{m \in \mathbb{T}(\mathbb{R}^n)} \mathbb{E}_{\mathbf{X} \sim \mu} [L(\Phi(\mathbf{X}) - m)^2]$. It is easy to show that this gives a proper scoring rule and that $a_{\mu} = \mathbb{E}[\Phi(\mathbf{X})]$. However, this scoring rule relies on the embedding of the group into its ambient vector space and does not account for or respect the group structure. Moreover, the resulting divergence and entropy reduce to just the Euclidean distance and usual variance. The same remark extends to (signature) kernel based scoring, where linear methods are used in an RKHS; see the discussion about non-kernel based scoring in the introduction. \lrcorner

Remark 4.6. We identify a stochastic process as a path- or sequence-valued random variable, possibly even ignoring its parametrization. However, for some applications the filtration of a stochastic process matters and one could ask to extend the scoring rule framework to this. A kernel that captures the filtration was introduced in [50] and a kernel algorithm and new applications given in [19]. To get a non-kernel scoring one could try to replace Φ in Definition (4.1) by the higher-rank signature from [50]. \lrcorner

5 Gradient Descent on the Space of Tracks

Given a smooth function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the simplest update rule for gradient descent is

$$x_{i+1} = x_i - \eta \nabla f(x_i) \tag{5}$$

and under additional regularity of f , the resulting sequence $(x_i)_i \subset \mathbb{R}^n$ converges to a minimizer of f . Our interest lies minimizing functions $f : \text{Tracks} \rightarrow \mathbb{R}$. In accordance with our guiding theme we do not identify these domains as linear spaces where classical gradient descent can be applied. However, we have seen that Φ provides an isomorphism

between the space of tracks and the free group G (up to forgetting the starting point of the track)

$$\text{Tracks} \simeq G$$

Hence, the minimization problem of a function f of tracks can be re-formulated as a minimization problem of a function $F = F(g) := f \circ \Phi^{-1}(g)$ on the free group G . That is, the general problem we try to solve is to find

$$\operatorname{argmin}_{g \in G} F(g) \text{ for } F : G \rightarrow \mathbb{R}$$

for any F in class of sufficiently “smooth” real-valued functions on G .

There have been many attempts to generalize gradient descent to non-linear domains. Arguably the the case of Riemannian manifolds [51] is the most well-developed among these. However, the group G does not come with a Riemannian structure (to wit, only a Sub-Riemannian structure [52]). We follow here a somewhat different approach inspired by work of Pierre Pansu [53] that directly uses the group structure to define gradients. We show that this gradient in turn allows us to give a straightforward generalization of the gradient update rule (5) from \mathbb{R}^n to G , so that the resulting sequence $(g_i)_i \subset G$ converges to the minimizer.

Pansu Derivatives. The derivative $Df(x)(\cdot)$ of a function

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

at a point x is a linear functional of \mathbb{R}^n that can be defined as the limit

$$Df(x)(h) := \lim_{h \rightarrow 0} \frac{f(x + \lambda h) - f(x)}{\lambda}.$$

Identifying \mathbb{R}^n as a the additive group $(\mathbb{R}^n, +)$ one can regard the difference quotient that appears in the limit as applying to the group operation to x and λh . Hence, if we have also have a generalization of the multiplication with the scalar λ , then the above difference quotient makes sense for other groups than the additive group $(\mathbb{R}^n, +)$. To formalize scalar multiplication, it turns out that the right notion is that of a Carnot group: a Carnot group is a Lie group C that carries a left-invariant geodesic distance $\operatorname{dist} : C \times C \rightarrow [0, \infty)$ and for each $\lambda > 0$ a bijection

$$\delta_\lambda : C \rightarrow C \text{ such that } \operatorname{dist}(\delta_\lambda(g), \delta_\lambda(h)) = \lambda \operatorname{dist}(g, h).$$

However, our focus is on the free group G and here the scaling δ_λ by a scalar $\lambda > 0$ has an explicit form

Proposition 5.1. *The group $G \subset \mathcal{H}$ equipped with the geodesic distance and*

$$\delta : \lambda \times G \rightarrow G, \quad \delta_\lambda g = \delta_\lambda(1, g_1, g_2, g_3, \dots) = (1, \lambda g_1, \lambda^2 g_2, \lambda^3 g_3, \dots).$$

forms a limit of Carnot groups.

We now have all we need to define the (Pansu) derivative. We denote by G^* the topological dual of G .

Definition 5.2. Let $f : G \rightarrow \mathbb{R}$. We define $Df : G \rightarrow G^*$ as

$$Df(g)h = \lim_{\lambda \downarrow 0} \frac{f(g\delta_\lambda h) - f(g)}{\lambda}$$

whenever this limit exists and call $Df(g)h$ the Pansu derivative of f at g in direction h . If f has a Pansu derivative for all $g \in G$ then we say that f is Pansu differentiable. Analogous we define the spaces $C^k(G, \mathbb{R})$ of k -times Pansu differentiable functions. \lrcorner

The Pansu derivative behaves very similar to the classic linear gradient. For example, for the proof of convergence of gradient descent on G we make use of the following ‘‘Taylor expansion’’.

Lemma 5.3. *If $f \in C^3(G, \mathbb{R})$, then*

$$f(gh) = f(g) + Df(g)h_1 + \frac{1}{2}D^2f(g)h_2 + O(\|h_3\|)$$

Proof. Consider the function $A : \mathbb{R} \rightarrow \mathbb{R}$

$$A(\lambda) = f(g\delta_\lambda h)$$

then A is C^3 and by a (classical linear) Taylor expansion we may write

$$A(1) = A(0) + \dot{A}(0) + \frac{1}{2}\ddot{A}(0) + O(\|A^{(3)}(0)\|)$$

which translates into the asserted equation since $h^{\otimes n}$ is contained in h_n . \square

Remark 5.4. A popular approach to differentiating functions of paths is to use a Fréchet derivative as in Malliavin calculus, i.e. one identifies the space of paths as a linear space, see [54]. However, the above Pansu derivative is of a very different nature and – by construction – respects the non-Euclidean structure of paths resp. tracks. \lrcorner

Gradient Descent on G . The idea of gradient descent to take a step in a direction that minimises f in a neighbourhood B_h . In our (Lie) group G , a natural choice is to choose some vector $v \in \mathbb{R}^n$ and use an exponential neighbourhood $ge^{\eta v}$ of $g \in G$ where \exp denotes the exponential from Lie algebra to Lie group. Hence, the question becomes how choose v to minimise $f(ge^{\eta v})$. To do so, note that

$$f(ge^{\eta v}) = f(g\delta_\eta e^v)$$

whenever $v \in \mathbb{R}^n$. Now using Lemma 5.3 shows

$$f(g\delta_\eta e^v) = f(g) + Df(g)e_1^{\eta v} + \eta^2 = f(g) + \eta Df(g) \cdot v + \eta^2.$$

This suggests that the direction of steepest descent in the exponential neighbourhood is indeed given by the Pansu derivative (henceforth, and with slight abuse of notation, we

identify the resulting element of G^* as an element of the Hilbert space \mathcal{H} since G and G^* both embed into \mathcal{H}),

$$v = -Df(g).$$

This leads to the following geometric descent rule

$$g_{i+1} = g_i e^{-\eta Df(g_i)}. \quad (6)$$

As for classic gradient descent, we need a notion of convexity to guarantee convergence to a minimum.

Definition 5.5. We say that a function $f : G \rightarrow \mathbb{R}$ is geometrically convex if

$$f(g\delta_\lambda(g^{-1}h)) \leq (1 - \lambda)f(g) + \lambda f(h)$$

for any $0 \leq \lambda \leq 1$. □

As in the linear case, one can show that if $f : G \rightarrow \mathbb{R}$ is geometrically convex and C^2 , then $D^2f(g)$ is positive definite everywhere. However, D^2f is not symmetric in general unlike in the linear case. Putting everything together allows us to mimic the convergence proof of gradient descent in linear spaces which ultimately justifies the above informal derivation of the update rule.

Theorem 5.6.

1. *The geometric update rule is transitive, that is, one may go from any point in the group to any other point using updates of the form (6).*
2. *If f is geometrically convex and bounded from below with bounded second Pansu derivatives, then for η sufficiently small the sequence in Equation (6) converges to a minimum*

Proof. Item 1 follows from Chow's theorem [55] which states that G is generated by simple exponentials. For Item 2 let L be a bound on the derivatives of f . We may write

$$f(gh) \leq f(g) + Df(g)h_1 + \frac{1}{2}L\|h_2\|.$$

Hence, if $g_{n+1} = g_n e^{-\eta Df(g_n)}$, then

$$f(g_{n+1}) \leq f(g_n) - \eta\|Df(g_n)\|^2 + \frac{1}{2}\eta^2 L\|Df(g_n)\|^2 = f(g_n) + \left(\frac{L}{2}\eta^2 - \eta\right)\|Df(g_n)\|^2$$

which is smaller than $f(g_n)$ for η small enough whenever $Df(g_n) \neq 0$. Hence this is a strictly decreasing sequence and converges to a minimum. □

Theorem 5.6 gives the analogous convergence guarantees as regular gradient descent and is simple to implement.

6 Experiments

In all our experiments we take as loss function L the squared norm,

$$L(\mathbf{t}) = \|\mathbf{t}\|^2 = \sum_{m=1}^M \|\mathbf{t}_m\|^2.$$

which is symmetric by Lemma 3.6. Note that the function L is smooth as a sum of squares, and since Φ is also smooth, since it is polynomial in the increments of its input when computed up to some fixed degree M , we can easily compute d, H , and I by automatic differentiation. For the computation of Φ we use the signatory [56] package which allows for fast and easy computations.

Remark 6.1. In the experiments we will simulate sample paths from Brownian motion. As this is a reversible process it does not matter if entropy is computed from the left or the right by Corollary 4.4, and we will normally compute it from the right. \lrcorner

6.1 Comparing a warped time-series to itself

One natural question to ask is how well the regularised versions of DTW that allow for differentiation are able to incorporate the parametrisation invariance. One drawback of these regularisation of DTW is that it typically leads to a trade-off between the smoothness and parametrization invariance.

Given two TS \mathbf{x} and \mathbf{y} we then use sDTW $d_{\text{sDTW}}(\mathbf{x}, \mathbf{y})$ which is the regularised version of DTW introduced in [8] which is differentiable and to a certain degree parametrisation invariant quantity. The scoring rule framework that we presented in the previous sections provides divergence not only between TS but probability measures on TS,

$$d(\nu, \mu) \equiv \mathbb{E}_{X \sim \nu}[S(X, \mu)] - H(\nu),$$

but as a special case, we can restrict this to point measures to get a “divergence” between two TS akin to sDTW which reduces to the formula

$$(\mathbf{x}, \mathbf{y}) \mapsto d(\delta_{\mathbf{x}}, \delta_{\mathbf{y}}) = L(\Phi(\mathbf{x})\Phi(\mathbf{y})^{-1}). \quad (7)$$

To study the different behaviour between these two “divergences” – sDTW and (7) – we generated samples paths from a Brownian motion to get \mathbf{x} and warped it a time change $\varphi : [0, T] \rightarrow [0, T]$, $\varphi(t) = T(t/T)^p$ to get \mathbf{y} where $p \in [1, \infty)$ is a parameter that determines the severity of the warping. Finally, we sampled on a discrete time grid with resolution 10^{-2} . The results are shown in Figure 1. The geometric score stays close to 0 regardless of the value of p but the sDTW divergence will increase with p and the rate of increase depends on the γ parameter. As γ tends to 0 the sDTW score will tend to the geometric score, but it will lose its smoothness while doing so. It is worth noting that $\gamma = 1$ is considered a default value. To achieve the same level of invariance enjoyed by the the divergence (7) in a DTW paradigm on can use the classical non-smooth DTW algorithms which cannot be updated using gradient descent. In contrast, the signature divergence (7) is always differentiable and parametrization invariant. Further, the signature divergence

is more general in the sense that it is not just a divergence between TS but between probability measures on TS which allows in principle for many other applications such as variational inference.

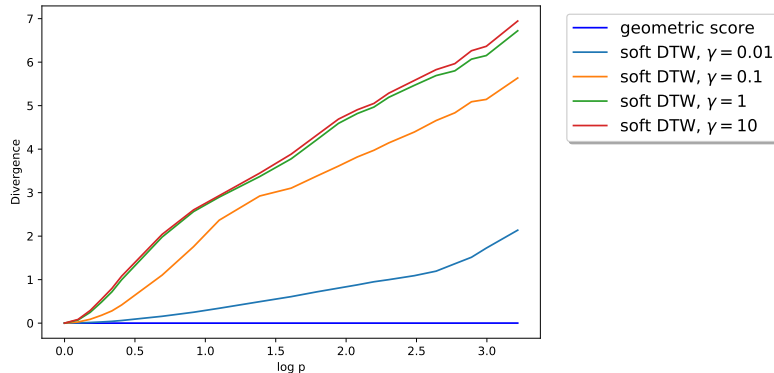


Figure 1: Geometric scoring compared to sDTW for different values of γ plotted against the logarithm of p for p between 1 and 25.

The implementation of sDTW is taken from the excellent Python package from [8].

6.2 Mutual information

Recall that the mutual information between two probability measures μ and ν is defined as

$$I(\mu, \nu) := H(\mu) - \mathbb{E}_{U \sim \nu}[H(\mu|U)]. \quad (8)$$

and provides a dependency measure between μ and ν . Closest related to independence on paths are the *signature cumulants* [57] which have been proven to be useful in applications [58, 59]. However, signature cumulants only compare probability measures on paths [tracks] with other probability measures on paths [tracks]. In contrast, the mutual information (8) allows to measure dependency between a probability measure μ on paths [tracks] and a probability measure ν on an arbitrary measurable space; in particular, this allows to measure dependence relations between TS and scalar-valued random variables.

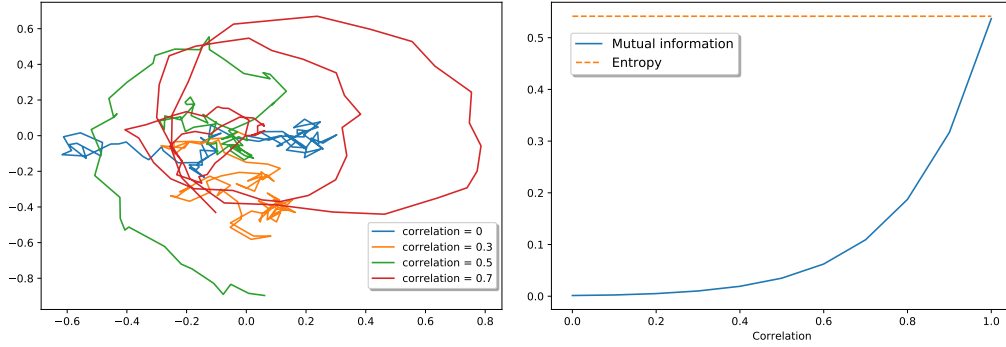


Figure 2: The left hand side shows sample trajectories from \mathbf{Y} for different correlation values ρ . The right plot shows the mutual information between \mathbf{Y} and the rotational speed ω as a function of ρ .

To demonstrate this, we consider two examples

Mutual information between a stochastic process and a scalar. We consider the stochastic process \mathbf{y} defined as

$$\mathbf{y}_t = \rho t \begin{pmatrix} \cos \omega t \\ \sin \omega t \end{pmatrix} + \sqrt{1 - \rho^2} \mathbf{x}_t$$

where ω is sampled from the uniform distribution on $[0, 8\pi)$ and \mathbf{x} is a standard 2-dimensional Brownian motion independent from ω . \mathbf{x} and \mathbf{y} are generated on the interval $[0, 1]$ on a grid with resolution 10^{-2} . The right plot in Figure 2 shows some sample trajectories of the process \mathbf{y} . The left plot in Figure 2 shows the mutual information between (the law of) ω and \mathbf{y} for various values of ρ . As expected, the mutual information increases monotone with ρ and is bounded by the entropy.

Mutual Information between two unparametrized stochastic processes. We consider two independent a Brownian motions \mathbf{x} and \mathbf{y} as before, and a random parametrisation $\phi(t) = T(t/T)^p$ where $p \in [1, 10]$ is uniformly distributed independent of \mathbf{x} . The stochastic process \mathbf{z} is defined as

$$\mathbf{z}_t = \rho \mathbf{x}_{\phi(t)} + \sqrt{1 - \rho^2} \mathbf{y}_t.$$

We then compute the mutual information between \mathbf{z} and \mathbf{x} as a function of ρ . The results are shown in Figure 3.

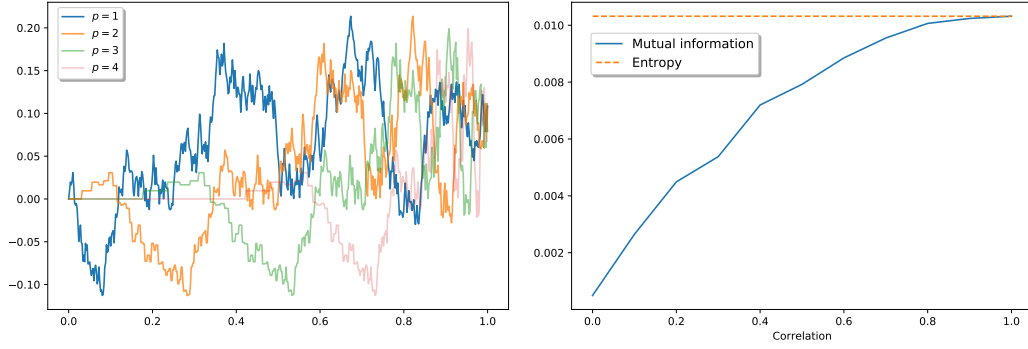


Figure 3: The left plot shows one-dimensional projections of sample paths of \mathbf{x} for different warping parameters ρ . The right plot shows the mutual information between \mathbf{x} and \mathbf{z} as a function of ρ .

Acknowledgements

PB is supported by the Engineering and Physical Sciences Research Council [EP/R513295/1]. HO is supported by the EPSRC grant “DATASIG”, the Turing Institute, and the Oxford-Man Institute of Quantitative Finance.

A Tree-like equivalence of paths

Informally, a path $x : [0, T] \rightarrow \mathbb{R}^n$ is *tree-like* if the set it traces out in \mathbb{R}^n looks like a tree. The formal definition is

Definition A.1 ([60]). A continuous path $x : [0, T] \rightarrow E$ is said to be *tree-like* if there exists an \mathbb{R} -tree τ , a continuous map $\varphi : [0, T] \rightarrow \tau$ and a map $\psi : \tau \rightarrow E$ such that $\varphi(0) = \varphi(T)$ and $x = \psi \circ \varphi$. Two paths \mathbf{x}, \mathbf{y} are *tree-like equivalent* if $\mathbf{x} \star \mathbf{y}^{-1}$ is tree-like and we denote this relation with $\mathbf{x} \sim \mathbf{y}$. \lrcorner

That is, if \mathbf{x} and \mathbf{y} follow the same trajectory in \mathbb{R}^n up to tree-like excursions then $\mathbf{x} \sim \mathbf{y}$. In particular, if \mathbf{y} is just \mathbf{x} under time-reparametrization, $\mathbf{y}(t) = \mathbf{x}(\tau(t))$ for an increasing τ , then $\mathbf{x} \sim \mathbf{y}$. See also [57, Appendix B] for more details. It is easy to check that \sim is an equivalence relation, hence we can define

$$\text{Tracks}(a, b) := \text{Paths}(a, b) / \sim, \text{ and } \text{Tracks} := \text{Paths} / \sim,$$

References

- [1] Glenn W Brier et al. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.
- [2] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [3] Philip Dawid. The geometry of proper scoring rules. *Annals of the Institute of Statistical Mathematics*, 59(1):77–93, 2007.
- [4] Kuo-Tsai Chen. Integration of paths—a faithful representation of paths by non-commutative formal power series. *Trans. Amer. Math. Soc.*, 89:395–407, 1958.
- [5] Hiroaki Sakoe and Seibi Chiba. A dynamic programming approach to continuous speech recognition. In *Proceedings of the Seventh International Congress on Acoustics, Budapest*, volume 3, pages 65–69, Budapest, 1971. Akadémiai Kiadó.
- [6] H. Sakoe and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26(1):43–49, 1978.
- [7] Marco Cuturi and Mathieu Blondel. Soft-dtw: a differentiable loss function for time-series. In *ICML*, 2017.
- [8] Mathieu Blondel, A. Mensch, and Jean-Philippe Vert. Differentiable divergences between time series. In *AISTATS*, 2021.
- [9] T. Gneiting and A. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102:359 – 378, 2007.

- [10] F. J Király and H. Oberhauser. Kernels for sequentially ordered data. *Journal of Machine Learning Research*, 20(31):1–45, 2019.
- [11] M. Fliess. Fonctionnelles causales non linéaires et indéterminées non commutatives. *Bull. Soc. Math. France*, 109(1):3–40, 1981.
- [12] Roger W. Brockett. Volterra series and geometric control theory. *Automatica*, 12(2):167 – 176, 1976.
- [13] Terry J. Lyons, Michael Caruana, and Thierry Lévy. *Differential equations driven by rough paths*, volume 1908 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007. Lectures from the 34th Summer School on Probability Theory held in Saint-Flour, July 6–24, 2004, With an introduction concerning the Summer School by Jean Picard.
- [14] Terry Lyons and Zhongmin Qian. *System Control and Rough Paths*. Oxford University Press, 2002. Oxford Mathematical Monographs.
- [15] Peter K Friz and Martin Hairer. *A course on rough paths: with an introduction to regularity structures*. Springer, 2014.
- [16] Anastasia Papavasiliou and Christophe Ladroue. Parameter estimation for rough differential equations. *The Annals of Statistics*, 39(4):2047–2073, 2011.
- [17] Ilya Chevyrev and Harald Oberhauser. Signature moments to characterize laws of stochastic processes. *arXiv e-prints*, page arXiv:1810.10971, Oct 2018.
- [18] Maud Lemerrier, Cristopher Salvi, Theodoros Damoulas, Edwin V. Bonilla, and Terry J. Lyons. Distribution regression for continuous-time processes via the expected signature. *CoRR*, abs/2006.05805, 2020.
- [19] Cristopher Salvi, Maud Lemerrier, Chong Liu, Blanka Horvath, Theodoros Damoulas, and Terry Lyons. Higher order kernel mean embeddings to capture filtrations of stochastic processes. *arXiv preprint arXiv:2109.03582*, 2021.
- [20] Adeline Fermanian, Pierre Marion, Jean-Philippe Vert, and Gérard Biau. Framing rnn as a kernel method: A neural ode approach, 2021.
- [21] Csaba Toth and Harald Oberhauser. Bayesian learning from sequential data using gaussian processes with signature covariances. In *International Conference on Machine Learning*, pages 9548–9560. PMLR, 2020.
- [22] Joel Dyer, Patrick Cannon, and Sebastian M Schmon. Approximate bayesian computation with path signatures, 2021.
- [23] Maud Lemerrier, Cristopher Salvi, Thomas Cass, Edwin V. Bonilla, Theodoros Damoulas, and Terry Lyons. Siggpde: Scaling sparse gaussian processes on sequential data, 2021.
- [24] Hao Ni, Lukasz Szpruch, Magnus Wiese, Shujian Liao, and Baoren Xiao. Conditional sig-wasserstein gans for time series generation, 2020.

- [25] Hans Bühler, Blanka Horvath, Terry Lyons, Imanol Perez Arribas, and Ben Wood. A data-driven market simulator for small data environments, 2020.
- [26] Patrick Kidger, James Foster, Xuechen Li, Harald Oberhauser, and Terry Lyons. Neural sdes as infinite-dimensional gans, 2021.
- [27] Chad Giusti and Darrick Lee. Signatures, lipschitz-free spaces, and paths of persistence diagrams, 2021.
- [28] Darrick Lee and Robert Ghrist. Path signatures on lie groups, 2020.
- [29] I. Chevyrev, V. Nanda, and H. Oberhauser. Persistence paths and signature features in topological data analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(1):192–202, Jan 2020.
- [30] Joscha Diehl, Kurusch Ebrahimi-Fard, and Nikolas Tapia. Generalized iterated-sums signatures, 2020.
- [31] Joscha Diehl, Kurusch Ebrahimi-Fard, and Nikolas Tapia. Time-warping invariants of multidimensional time series. *Acta Applicandae Mathematicae*, 170(1):265–290, May 2020.
- [32] Csaba Toth, Patric Bonnier, and Harald Oberhauser. Seq2tens: An efficient representation of sequences by low-rank tensor projections. In *International Conference on Learning Representations*, 2021.
- [33] Tobias Fissler and J. Ziegel. Order-sensitivity and equivariance of scoring functions. 2017.
- [34] David Bolin and Finn Lindgren. Excursion and contour uncertainty regions for latent gaussian models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(1):85–106, 2015.
- [35] Ilya Molchanov and Ilya S Molchanov. *Theory of random sets*, volume 87. Springer, 2005.
- [36] Tobias Fissler, Rafael Frongillo, Jana Hlavinová, and Birgit Rudloff. Forecast evaluation of quantiles, prediction intervals, and other set-valued functionals. *Electronic Journal of Statistics*, 15(1):1034–1084, 2021.
- [37] A. Dawid and M. Musio. Theory and applications of proper scoring rules. *METRON*, 72:169–183, 2014.
- [38] Peter D. Grunwald and A. Dawid. Game theory, maximum entropy, minimum discrepancy and robust bayesian decision theory. *Annals of Statistics*, 32:1367–1433, 2004.
- [39] F. Gressmann, F. J. Király, B. Mateen, and H. Oberhauser. Probabilistic supervised learning. *ArXiv e-prints*, January 2018.

- [40] I. J. Good. Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107–114, 1952.
- [41] Leonard J Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801, 1971.
- [42] J. Abernethy and Rafael M. Frongillo. A characterization of scoring rules for linear properties. In *COLT*, 2012.
- [43] Ingo Steinwart, Chloé Pasin, R. C. Williamson, and Siyu Zhang. Elicitation and identification of properties. In *COLT*, 2014.
- [44] Rafael M. Frongillo and Ian A. Kash. Vector-valued property elicitation. In *COLT*, 2015.
- [45] Christophe Reutenauer. *Free Lie algebras*. The Clarendon Press Oxford University Press, New York, 1993. Oxford Science Publications.
- [46] Ben Hambly and Terry Lyons. Uniqueness for the signature of a path of bounded variation and the reduced path group. *Ann. of Math. (2)*, 171(1):109–167, 2010.
- [47] Ilya Chevyrev, Terry Lyons, et al. Characteristic functions of measures on geometric rough paths. *The Annals of Probability*, 44(6):4049–4082, 2016.
- [48] Peter K. Friz and Nicolas B. Victoir. *Multidimensional stochastic processes as rough paths: theory and applications*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, Cambridge, 2010.
- [49] V. I Bogachev. *Measure theory. Vol. I, II*. Springer, 2007.
- [50] Patric Bonnier, Chong Liu, and Harald Oberhauser. Adapted topologies and higher rank signatures, 2021.
- [51] Silvere Bonnabel. Stochastic gradient descent on riemannian manifolds. *IEEE Transactions on Automatic Control*, 58(9):2217–2229, September 2013.
- [52] Richard Montgomery. *A tour of subriemannian geometries, their geodesics and applications*, volume 91 of *Mathematical Surveys and Monographs*. American Mathematical Society, Providence, RI, 2002.
- [53] Pierre Pansu. Métriques de Carnot-Carathéodory et Quasiisométries des Espaces Symétriques de rang un. *Annals of Mathematics*, 129(1):1–60, 1989.
- [54] Thomas Cass and Peter Friz. Malliavin calculus and rough paths. *Bulletin des Sciences Mathématiques*, 135(6):542–556, 2011. Special issue in memory of Paul Malliavin.
- [55] Wei-Liang Chow. Über Systeme von linearen partiellen Differentialgleichungen erster Ordnung. *Mathematische Annalen*, 117-117(1):98–105, December 1940.

- [56] Patrick Kidger and Terry Lyons. Signatory: differentiable computations of the signature and logsignature transforms, on both CPU and GPU. In *International Conference on Learning Representations*, 2021. <https://github.com/patrick-kidger/signatory>.
- [57] Patric Bonnier and Harald Oberhauser. Signature cumulants, ordered partitions, and independence of stochastic processes. *Bernoulli*, 26(4), November 2020.
- [58] Alexander Schell and Harald Oberhauser. Nonlinear independent component analysis for continuous-time signals, 2021.
- [59] Peter K Friz, Paul Hager, and Nikolas Tapia. Unified signature cumulants and generalized magnus expansions. *arXiv preprint arXiv:2102.03345*, 2021.
- [60] Ben Hambly and Terry Lyons. Uniqueness for the Signature of a Path of Bounded Variation. Preprint, 2002.