

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich



Nonparametric Regression of Stochastic Processes via Signatures

A. Schell and R. Alaifari

Research Report No. 2023-45

December 2023 Latest revision: May 2024

Seminar für Angewandte Mathematik Eidgenössische Technische Hochschule CH-8092 Zürich Switzerland

Nonparametric Regression of Stochastic Processes via Signatures

Alexander Schell^{*} Rima Alaifari[†]

Seminar for Applied Mathematics, ETH Zurich

31 May 2024

Abstract

Nonparametric regression of stochastic processes estimates statistical relationships between multidimensional, time-dependent data without relying on specific parametric assumptions. We propose a novel approach to this classical estimation problem by using the signature transform from rough path theory to encode the information of a stochastic process into a sequence of iterated integrals, capturing its statistical properties in a time-global and hierarchical manner. Viewing statistical regression as an operator learning problem, this signature-based discretisation allows us to characterise the conditional statistical dependence of a stochastic process on another stochastic process as the solution to a convex semi-infinite linear least squares problem. This result is based on a functional monotone class argument involving the bounded signature of the conditioning process and allows for the efficient and provably consistent nonparametric estimation of regression functions and conditional distributions for very general classes of jointly distributed stochastic processes as solutions to convex optimisation problems. The structural insights of this approach are summarised in two universally consistent regression estimators that are computable with practical algorithms and supported by broad theoretical guarantees.

Keywords: stochastic processes, nonparametric and functional regression, conditional mean embedding, conditional expectation, signature features, kernel ridge regression, function approximation

1 Introduction

Modelling and inferring meaningful patterns and relationships within complex, high-dimensional data is a central challenge in modern machine learning and statistics. At the core of this challenge lies the task of capturing statistical dependencies between data sets through conditional distributions and the structured operationalisation of statistical conditioning. These are fundamental aspects of probabilistic modelling that are crucial for a growing range of inference techniques and applications. One of the most important paradigms for this is nonparametric regression, which is predicated on the idea that the relationships between variables can be captured with minimal assumptions about their functional form. Unlike parametric methods, which specify a fixed structure for the relationship (e.g., linear or polynomial), nonparametric approaches allow the data to speak for itself, revealing the underlying patterns and relationships more flexibly.

Conceptually, nonparametric regression involves estimating a regression function that describes the conditional expectation of a response variable given one or more predictor variables. The challenge intensifies when the predictors or responses are observed as part of a stochastic process, since

^{*}alexander.schell@math.ethz.ch [†]rima.alaifari@math.ethz.ch

one must then account for temporal dependencies and potential non-stationarities inherent in the data. Applications of such nonparametric regression methods for stochastic processes are extensive and profound. In finance and econometrics, these methods can be used to model asset prices, interest rates, or volatility, capturing subtle dependencies and predicting future trends [54]. They also permeate risk assessments in financial markets and portfolios [2], or stochastic filtering and the optimisation of control systems in engineering [3, 38]. In environmental science, they help in understanding climate patterns, pollution levels, or ecological dynamics, accommodating complex interactions and temporal variations. In biomedical engineering, they are employed to analyse dependencies between physiological signals, such as heart rate variability or brain activity, where the underlying biological processes are inherently stochastic and nonlinear. Stochastic process regression is also applied in computer vision [1] and molecular dynamics [30], and it is central to survival analysis [27] and causal models [44], time series analysis and forecasting [39, 45, 55], Bayesian inversion and inference [37, 50] and statistical machine learning [21] broadly, to name just a few classical examples. Most recently, nonparametric stochastic process regression has emerged as the key concept behind large language models [62], which are essentially high-dimensional statistical regression models that approximate conditional distributions on sequential data [53, 59].

A particular challenge in computing conditional expectations and their derived statistics, present in most of the above examples and especially in sequential or language-based machine learning, is to efficiently account for potential time dependencies in the conditioning variable, i.e. the conditioning on complex multidimensional stochastic processes. This problem was first systematically considered in classical probability, where the traditional model classes of martingales and Markov processes were conceived to elegantly circumvent subtler issues of time-dependent conditioning. However, these classical 'convenience' assumptions have clear limitations (e.g. [6, 32]), making it worthwhile to revisit the general problem of time-dependent regression and conditioning with modern tools from stochastic analysis. The present work aims to contribute to this endeavour.

Specifically, this paper addresses the following central inference questions, which the above examples suggest are of significant practical relevance: Given two multidimensional stochastic processes X and Y in discrete- or continuous time and some random vector Z in Euclidean space,

how can the statistics
$$f \mapsto \mathbb{E}[f(Y) | X]$$
 and $\mathbb{E}[Z | X]$ be efficiently estimated? (1)

(The conditional law $\mathbb{P}(Y \in \cdot | X)$ is included by letting f run over the indicator functions on Y.) The approximation of such conditional expectations defines the problem of statistical regression. This has been well-addressed for time-independent X or if the temporality of $X = (X_t)$ conforms to certain parametric assumptions (e.g. [10, 5, 61] and the references therein). However, to the best of our knowledge, there are currently no rigorous nonparametric solutions to (1) for the case of general (jointly distributed) stochastic processes $X = (X_t)$ and $Y = (Y_t)$.

The present work attempts to close this gap by using tools from rough path theory, and in particular the concept of bounded (or "robust") signatures recently introduced in [12], to structure the statistical information given by the predictor process X: A bounded signature, ϕ , is an algebraically structured and bounded Hilbert-valued coordinate map over the space of sufficiently continuous paths, which includes the sample realisations of X. Using a functional monotone class argument, we show that for a linear subset \mathfrak{L} of the Hilbert co-domain $(\mathcal{H}, \langle \cdot, \cdot \rangle)$ of ϕ , the coordinate functions $\{\langle \ell, \phi(X) \rangle \mid \ell \in \mathfrak{L}\}$ form an L^2 -dense subspace of all square-integrable X-measurable random variables. As a basic step toward answering (1), this then implies the variational characterisation

$$\mathbb{E}[Z \mid X] = \lim_{k \to \infty} \left(\left\langle \alpha_{1k}, \phi(X) \right\rangle, \cdots, \left\langle \alpha_{mk}, \phi(X) \right\rangle \right), \tag{2}$$

with m the dimension of $Z \equiv (Z^1, \dots, Z^m)$ and for any minimizing sequence $(\alpha_{1k}, \dots, \alpha_{mk})_k$ of

$$\inf_{(\tilde{\alpha}_1,\dots,\tilde{\alpha}_m)\in\mathfrak{L}^m}\sum_{i=1}^m \mathbb{E} \left| Z^i - \langle \tilde{\alpha}_i, \phi(X) \rangle \right|^2.$$
(3)

The parameter optimization (3) that underlies the above approximation is notably convex, and the convergence (2) holds at least in L^2 and almost surely if the sequence (α_{ik}) is fast enough. A similar characterisation can be found for the conditional expectation given X of the bounded signatures of Y itself, leading to the variational identities (convergent in L^2 , and under conditions almost surely)

$$\mathbb{E}[f(Y) \mid X] = \lim_{l \to \infty} \lim_{k \to \infty} \left\langle \ell_l^f, \psi_{\alpha_k}(X) \right\rangle \tag{4}$$

for any measurable function f of Y such that f(Y) square-integrable. In this case, $(\psi_{\alpha_k})_k$ is an L^2 dense adaptive system of ϕ -derived model functions on sample paths that conceptually resembles (and effectively replaces) an 'algebraically structured and convexly parametrisable neural network', and $(\ell_l^f)_l$ is a dualised signature-encoding of the argument function f. Optimally parametrised instances of (ψ_{α_k}) and the function encodings (ℓ_l^f) can both be computed explicitly as solutions to two separate, well-structured and efficiently approximable convex optimisation problems.

In addition to being fully convex-optimisable, the representations (2) and (4) hold nonparametrically for any response-transforming nonlinearity f without needing any assumptions on the relationship between the marginals in (X, Z) or (X, Y), such as continuity or specific statistical dependence structures. The flexible and comparatively simple algorithmic premise behind these identities also makes the regression statistics in (1) amenable to nonparametric statistical estimation within the well-established empirical framework of kernel ridge regression. Consequently, this approach provides a theoretically and practically attractive solution to the inference problem (1).

This paper is structured as follows. Section 2 reviews the fundamental concepts behind nonparametric regression analysis and statistical conditioning and explains their use for modelling and analysing dependencies between statistically associated datasets (Sections 2.1 and 2.2). It then sets the stage for (1) by reformulating this question as an equivalent operator learning problem using the perspective of conditional mean embedding (Section 2.3). Leaving the general setting behind, Section 3 delves into the specifics of nonparametric regression of time-dependent data, covering structural basics for the associated regression spaces and time-dependent random variables (Section 3.1). The signature transform is introduced as a key tool for representing such time-dependent data (Section 3.2), complemented by a brief discussion of its basic properties and how it can serve as a bounded global coordinate map on path spaces (Section 3.3). A derivation of the variational identities (2) and (4) is provided in Section 4, which presents the main theoretical contributions of the paper. We show how the bounded signature can be used to discretise and exhaust the information of the conditioning process (Section 4.1), and utilise this discretisation to characterise the conditional expectations in (1) as solutions of convex semi-infinite linear least squares problems (Sections 4.2 and 4.3). These characterisations are then promoted to practical estimators for the nonparametric regression statistics (1) in Section 5, which builds on the theoretical insights of the previous sections to achieve an efficient data-based approximation of the proposed representations (4) and (2). Upon embedding our signature-based regression architectures in the context of vector-valued reproducing kernel Hilbert spaces and providing basic support theory for a subsequent analysis of statistical convergence properties (Section 5.1 and 5.2), we present explicitly computable nonparametric regression estimators for (1) (Sections 5.3 and 5.4) and provide an analysis of their statistical approximation properties, including convergence rates and error bounds (Section 5.5). Numerical experiments and practical example applications of our method are to be included in the forhtcoming arXiv version of this paper. Most technical proofs are provided in the Appendix. Closing with a note on existing literature, we remark that prior approximations in a manner similar to (2) were first explored in [31, 34], though the respective aspects of these works are mostly empirical and based on rather strict assumptions on (X, Z). In a spirit related to ours but with no immediate mathematical connection, [13] study nowcasting using linear regression on signatures. Finally, we note that the observation that the robust signature algebra is dense in L^p , which is crucial for (2), was made independently of us (and by other mathematical means) in the recent preprint [4], where this idea is applied to different consequences in the realm of optimal stopping.

2 Perspectives on Regression and Statistical Conditioning

This section provides a brief review of probabilistic regression and statistical conditioning, introducing basic concepts (Sections 2.1 and 2.2) as well as more specialised mathematical perspectives (Section 2.3) to structure and operationalise the underlying statistical theory. The aim is to lay bare the essential probabilistic structure behind regression-based statistical inference and learning on Polish spaces, so that, in subsequent sections, we can seamlessly extend this foundational theory in the particular setting of nonparametric regression for time-dependent multidimensional data.

Specifically, Section 2.1 motivates the framework of statistical regression analysis as a principled approach to model and analyse complex, non-functional statistical dependencies between two sets of data residing in potentially infinite-dimensional spaces. The classical idea is to view these data as samples from a pair of jointly distributed random variables and decompose their joint probability distribution into a family of conditional distributions with respect to one of their marginals. This is formalised through the well-known concept of disintegration, which Section 2.2 briefly recalls and contextualises. With these basics in hand, Section 2.3 draws on the perspective of conditional mean embedding to reformulate the problem of nonparametric (and, in general, nonlinear) regression as an equivalent, yet more tangible problem of approximating, from finite data, a bounded linear 'regression' operator (Proposition 2.4 and Definition 2.6; also Lemma 2.9). We show that this problem can be solved under a 'sufficiently nonlinear' Hilbert-valued coordinatisation (15) over the space of all response data (Proposition 2.4), and illustrate how this approach yields a comprehensive statistical framework for stochastic process regression and in the context of large language models.

2.1 Regression Analysis: Modelling Statistical Dependencies in Data

Extracting and analysing statistical patterns and dependencies from complex, high-dimensional data through conditional distributions and statistical conditioning is a fundamental problem central various inference techniques and applications.

One of the most popular approaches to this challenge is regression analysis, which aims to identify the relationship between two given data sets (x_i) and (y_i) . Here, each x_i represents an input and each y_i represents a corresponding output, and these data points reside in potentially infinitedimensional spaces \mathcal{X} and \mathcal{Y} , respectively. Classical regression analysis, and most physical models alike, aims to extract from the pairs (x_i, y_i) an essentially functional relationship of the form:

$$\mathfrak{R}^{f} = \{(x, f(x)) \mid x \in \mathcal{X}\} \quad \text{or, more broadly,} \quad \mathfrak{R} = \{(x, y) \in \mathcal{X} \times \mathcal{Y} \mid y \approx f(x)\}, \qquad (5)$$

assuming the existence of a function $f : \mathcal{X} \to \mathcal{Y}$ that maps inputs $\{x\}$ to outputs $\{y\}$. This f then encapsulates the 'true' link between the data (x_i) and (y_i) . It is usually approximated by a proxy $f_{\theta_{\star}}$ selected from a predefined family of model functions (f_{θ}) to best fit the observed pairs (x_i, y_i) .

However, many real-world scenarios defy this simplified 'functional' assumption, presenting cases where an input x in \mathcal{X} does not correspond to a singular deterministic output f(x) in \mathcal{Y} , but to a range \mathfrak{R}_x^{\sim} of possible outcomes associated with x. For instance, in machine translation, sentences x from a source language \mathcal{X} are to be translated into sentences y in a target language \mathcal{Y} . This defines the complex relation

$$\mathfrak{R}^{\sim} \coloneqq \{ (x, y) \in \mathcal{X} \times \mathcal{Y} \mid y \text{ is a valid translation of } x \}, \tag{6}$$

which typically extends beyond a simple functional relation $(y \approx f(x))$ since a given sentence x in \mathcal{X} can translate into multiple valid sentences y in \mathcal{Y} . Such a non-functional, multi-valued relationship between x and y is termed an *associative relation*, or an *association*, between x and y.

To handle these complexities, classical regression analysis can be extended into a rigorous probabilistic framework. Instead of describing associations through a deterministic function from \mathcal{X} to \mathcal{Y} , a more appropriate approach is to use a (set-valued) map from \mathcal{X} to $2^{\mathcal{Y}}$, sending x to $\mathfrak{R}_x^{\sim} := \{y \in \mathcal{Y} \mid (x, y) \in \mathfrak{R}\}$. Further complexity is captured by assigning to each $x \in \mathcal{X}$ a probability measure μ_x over \mathfrak{R}_x^{\sim} , reflecting the likelihood of different outcomes associated with x. One way to implement this is by defining a probability measure μ on $\mathcal{X} \times \mathcal{Y}$, say $\mu = \mathbb{P}_{(X,Y)}$ for an $\mathcal{X} \times \mathcal{Y}$ -valued random variable (X, Y) that models the data. (The marginal variables X and Y are called regressor and regressand, respectively.) The desired plausibility measures $(\mu_x)_{x \in \mathcal{X}}$ are then obtained as the conditional probabilities of Y given X = x. This probabilistic approach provides a well-known framework to describe associative relationships between data (x) and (y) through an analysis of the statistical dependency between the random variables X and Y that model the data.

The remainder of Section 2 reviews the necessary mathematics to develop this perspective into a coherent statistical theory. The focus of our discussion will be on how to consistently approximate the probabilistic model $(\mu_x)_{x \in \mathcal{X}}$ of an associative relationship between data from a finite number of samples. We henceforth assume that \mathcal{X} and \mathcal{Y} are Polish spaces, unless otherwise stated.

2.2 Capturing Statistical Dependencies with Probability Kernels

As discussed above, the main idea behind the probabilistic modelling of an associative relation in $\mathcal{X} \times \mathcal{Y}$ (such as (6)) is to reframe this relation as a functional relation within $\mathcal{X} \times \mathcal{M}_1(\mathcal{Y})$:

By assuming that the statistically dependent ('associated') data (x_i) and (y_i) are sampled from a joint distribution μ on $\mathcal{X} \times \mathcal{Y}$, say $(x_i, y_i) \sim \mu$, we aim to decompose μ into a measure-valued function $x \mapsto \mu_x$ from \mathcal{X} to $\mathcal{M}_1(\mathcal{Y})$, such that for each $x \in \mathcal{X}$, the measure $\mu_x \in \mathcal{M}_1(\mathcal{Y})$ represents the conditional law of data in \mathcal{Y} given x.

This concept is rigorously defined through the classical notion of a probability kernel. This is a map $\kappa : \mathcal{X} \times \mathcal{B}(\mathcal{Y}) \to [0, 1]$ with $(\kappa(x, \cdot) \mid x \in \mathcal{X}) \subseteq \mathcal{M}_1(\mathcal{Y})$ such that for each set $B \in \mathcal{B}(\mathcal{Y})$,

$$\mathcal{X} \ni x \longmapsto \kappa(x, B) \eqqcolon \kappa_x(B)$$
 is $(\mathcal{B}(\mathcal{X}), \mathcal{B}([0, 1]))$ -measurable.¹

Following this, recall that every such kernel κ can be 'fused' with any $v \in \mathcal{M}_1(\mathcal{X})$ via the coupling

$$\kappa \otimes \upsilon : \mathcal{B}(\mathcal{X} \times \mathcal{Y}) \ni A \longmapsto \int_{\mathcal{X}} \left[\int_{\mathcal{Y}} \mathbb{1}_A(x, y) \, \kappa_x(\mathrm{d}y) \right] \upsilon(\mathrm{d}x) \, \in \, [0, 1], \tag{7}$$

and the resulting map (7) is then again a Borel probability measure on $\mathcal{X} \times \mathcal{Y}$. Conversely, any given measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ can be 'divided' by each of its marginals, and the resulting 'quotient' is essentially unique. This is formalised by the concept of *disintegration*, which asserts that there is

a 'unique' probability kernel $\mu_{\mathcal{Y}|\mathcal{X}} \equiv (\mu_{\mathcal{Y}|x} \mid x \in \mathcal{X})$ such that $\mu = \mu_{\mathcal{Y}|\mathcal{X}} \otimes \mu_{\mathcal{X}}$, (8)

where 'unique' means 'unique up to inequality on a $\mu_{\mathcal{X}}$ -nullset': for any two kernels $\kappa, \tilde{\kappa} \subseteq \mathcal{M}_1(\mathcal{Y})$, the identity $\kappa \otimes \mu_{\mathcal{X}} = \tilde{\kappa} \otimes \mu_{\mathcal{X}}$ implies that $\kappa_x = \tilde{\kappa}_x \ \mu_{\mathcal{X}}$ -almost everywhere (see e.g. [23, Chap. 8]).



Figure 1: Visualisation of increasingly non-functional relations and their probabilistic description. The left panel shows a functional relation $f: x \mapsto y$, as in (5). The middle panel depicts a (slightly non-functional) associative relation where an input x maps to a whole set \Re_x of possible outputs, together with an associated probability measure μ_x capturing the likelihood of different outcomes within \Re_x . The right panel shows a fully associative relation where X and Y exhibit complex dependencies; it also demonstrates the concept of disintegrations (8), showing the decomposition of the joint distribution $\mu = \mathbb{P}_{(X,Y)}$ into the conditional distributions $\mu_{\mathcal{Y}|\mathcal{X}} \equiv (\mu_x)_{x \in \mathcal{X}}$ wrt. the marginal $\mu_{\mathcal{X}} = \mathbb{P}_X$. Altogether, these panels illustrate the transition from classical regression to more complex probabilistic models and the role of disintegrations (12) in capturing statistical dependencies.

Consequently, for a fixed \mathcal{X} -marginal $\xi \in \mathcal{M}_1(\mathcal{X})$, the assignment $\mu \mapsto \mu_{\mathcal{Y}|\mathcal{X}}$ in (8) defines a bijection

$$\mathfrak{c}_{\xi} : \mathcal{M}_{1}^{\xi}(\mathcal{X} \times \mathcal{Y}) \longrightarrow L^{0}(\xi; \mathcal{M}_{1}(\mathcal{Y})), \quad \mu \mapsto \mu_{\mathcal{Y}|\mathcal{X}}, \tag{9}$$

with inverse $\mathbf{c}_{\xi}^{-1}(\cdot) = (\cdot) \otimes \xi$. This conditioning map (9) provides the desired 'functionalisation' of any type of statistical dependence between the marginals of $\{(x_i, y_i)\}$. The map (9) thus captures every kind of associative relation on $\mathcal{X} \times \mathcal{Y}$ by interpolating between two statistical extreme cases: **Example 2.1** (Functional Dependence and Independence). Given marginals $\xi \in \mathcal{M}_1(\mathcal{X})$ and $\mu_{\mathcal{Y}} \in \mathcal{M}_1(\mathcal{Y})$ —read as the distributions of (x_i) and (y_i) , resp.—and a function $f : \mathcal{X} \to \mathcal{Y}$, consider

$$\mu_{\mathcal{Y}|\mathcal{X}}^f \coloneqq (\delta_{f(x)} \mid x \in \mathcal{X}) \quad \text{and} \quad \mu_{\mathcal{Y}|\mathcal{X}}^\perp \coloneqq (\mu_{\mathcal{Y}} \mid x \in \mathcal{X}).$$

These measure-valued functions (in fact kernels, i.e. in $L^0(\xi; \mathcal{M}_1(\mathcal{Y}))$, if f is Borel-measurable) identify associations $\mu_f \coloneqq \mu_{\mathcal{Y}|\mathcal{X}}^f \otimes \xi \ (=\mathfrak{c}_{\xi}^{-1}(\mu_{\mathcal{Y}|\mathcal{X}}^f))$ and $\mu_{\perp} \coloneqq \mu_{\mathcal{Y}|\mathcal{X}}^{\perp} \otimes \xi = \mu_{\mathcal{Y}} \otimes \xi \ (=\mathfrak{c}_{\xi}^{-1}(\mu_{\mathcal{Y}|\mathcal{X}}^{\perp}))$ that characterize functional dependence (5) and statistical independence $((x_i) \perp (y_i))$, respectively.

The concept of a conditioning map (9) allows us to express the task of learning associative relations between data (from finitely many data samples) as follows:

For
$$\mu \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$$
, approximate the kernel $\mathfrak{c}_{\xi}(\mu)$ from samples of μ . (10)

Before we transition from the informal description (10) (the term 'approximate' is made precise later) to a precise and actionable problem formulation, let us first establish some basic terminology.

Definition 2.2. For \mathcal{X}, \mathcal{Y} Polish spaces and $\xi \in \mathcal{M}_1(\mathcal{X})$, the disintegration $\mathfrak{c}_{\xi}(\mu)$ of any law $\mu \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$ is termed the conditional dependence of μ wrt. the source law ξ .

It is often convenient, not least for notational reasons, to view a measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ as the law of some $(\mathcal{X} \times \mathcal{Y})$ -valued random variable. That is, to suppose that

$$\mu = \mathbb{P}_{(X,Y)}$$
 for some $X: \Omega \to \mathcal{X}$ and $Y: \Omega \to \mathcal{Y}$, (11)

where both X and Y are some Borel-measurable maps over a joint probability space $(\Omega, \mathscr{F}, \mathbb{P})$; the data $\{(x_i, y_i)\}$ are then sample realisations of (X, Y). In common regression terminology, Y is called the 'regressand' and X is called the 'regressor'. In this setting, the marginals of μ are $\mu_{\mathcal{X}} = \mathbb{P}_X$ and $\mu_{\mathcal{Y}} = \mathbb{P}_Y$, and it is customary to denote the conditional kernel $\mu_{\mathcal{Y}|\mathcal{X}}$ in (8) as

$$\mathfrak{c}_X(\mu) \equiv \left(\mu_{\mathcal{Y}|x}(\,\cdot\,) \,:\, x \in \mathcal{X}\right) \eqqcolon \left(\mathbb{P}(Y \in (\,\cdot\,) \mid X = x) \,:\, x \in \mathcal{X}\right). \tag{12}$$

A solution to (10) then provides a data-based approximation ('estimator') of the statistics

$$\mathbb{P}(Y \in A \mid X = x) = \int_{A} d\mu_{\mathcal{Y}|x}, \quad \mathbb{P}(X \in A, Y \in B) = \int_{A} \int_{B} d\mu_{\mathcal{Y}|x} \mathbb{P}_{X}(dx) \quad \text{and}$$

$$\mathbb{E}[f(Y) \mid X] = \int_{\mathcal{Y}} f d\mu_{\mathcal{Y}|X}, \quad \text{for any} \quad A \times B \in \mathcal{B}(\mathcal{X} \times \mathcal{Y}) \quad \text{and} \quad f \in \mathcal{L}^{2}(\mathbb{P}_{Y}),$$
(13)

where $\mu_{\mathcal{Y}|X} \coloneqq \mathfrak{c}_{\xi}(\mu) \circ X$. These are of central importance to various applications.

Remark 2.3. Any conditional dependence (CD) can be lifted to the (conditional) law of a canonical $(\mathcal{X} \times \mathcal{Y})$ -valued random variable, and the converse is also true, giving a one-to-one correspondence between CDs and paired random variables, see Remark B.1. Thus, the 'lifting assumption' (11) entails no loss of generality and is, in fact, not an assumption but rather a provable statement. \blacklozenge

2.3 Nonparametric Regression as an Operator Learning Problem

Given a law $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ with marginal $\xi := \mu_{\mathcal{X}}$, the main challenge we want to address is how to efficiently approximate the conditional dependence $\mathfrak{c}_{\xi}(\mu)$ from a finite number of μ -drawn samples. In fact, we will tackle the task (10) in the setting of time-dependent data, where \mathcal{X} and \mathcal{Y} are infinite-dimensional function spaces consisting of vector-valued paths (Sect. 3). This sequential setting poses several statistical and analytical challenges, not least because the resulting relation spaces $\mathcal{X} \times \mathcal{Y}$ are generally not locally compact. Our approach to handle this is to work with sufficiently nonlinear 'coordinate maps' on \mathcal{X} and \mathcal{Y} that ultimately allow us to estimate (13) with bespoke approximation architectures grounded in rigorous statistical theory (see Sections 4 and 5).

But first, we must ask how, and in what precise sense, we can approximate a conditional dependence (8)—an uncountable family of measures—from finite data. An established² approach is to use functional analysis and view such families as explicitly representable bounded linear operators.

Proposition 2.4 (Conditional Mean Embedding). For any law $\mu \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$, the conditional dependence $\mathfrak{c}_{\xi}(\mu) \equiv (\mu_x)_{x \in \mathcal{X}} \subseteq \mathcal{M}_1(\mathcal{Y})$ of μ wrt. ξ can be identified with a bounded linear operator

$$\mathfrak{c}_{\xi}(\mu) : L^{2}(\mu_{\mathcal{Y}}) \longrightarrow L^{2}(\xi), \quad f \longmapsto \mu_{\cdot}(f) \coloneqq \left[x \mapsto \int_{\mathcal{Y}} f \, \mathrm{d}\mu_{x} \right], \tag{14}$$

Moreover, given a separable Hilbert space $(\mathcal{H}_{\mathcal{Y}}, \langle \cdot, \cdot \rangle)$ together with a Borel-measurable map

 $q: \mathcal{Y} \to \mathcal{H}_{\mathcal{Y}} \quad such \ that \quad \left\{ \langle \ell, q(\cdot) \rangle \mid \ell \in \mathcal{H}_{\mathcal{Y}} \right\} \ is \ L^2 \text{-dense} \ in \ L^2(\mu_{\mathcal{Y}}), \tag{15}$

 $^{^2}$ See, for instance, [41, 33] and the references therein.

then the (so-called) regression operator (14) and, thus, the conditional dependence $\mathfrak{c}_{\xi}(\mu)$ itself, can be identified with the function

$$\mu_{\cdot}(q) \in L^{2}(\xi; \mathcal{H}_{\mathcal{Y}}) \quad given \ by \quad x \mapsto \mu_{x}(q) \coloneqq \int_{\mathcal{Y}} q(y) \ \mu_{x}(\mathrm{d}y). \tag{16}$$

Specifically then, for each $f \in L^2(\mu_{\mathcal{Y}})$ there is a sequence $(\ell_j^{(f)})_j \subset \mathcal{H}_{\mathcal{Y}}$ such that

$$\mathfrak{c}_{\xi}(\mu)(f) = \lim_{j \to \infty} \left\langle \ell_j^{(f)}, \mu_{\cdot}(q) \right\rangle \quad in \quad L^2(\xi).$$
(17)

(For functions f such that $f = \langle \ell, q(\cdot) \rangle$ for some $\ell \in \mathcal{H}_{\mathcal{Y}}$, we can choose $\ell_j^{(f)} \coloneqq \ell$ for each j.)

Proof. Mostly straightforward consequences of the definitions, but see Appendix B.1.2.

Remark 2.5 (Parametrizing Conditional Dependencies). The representation (17) based on (15) can be seen as an 'infinite-dimensional parametrization' of the regression operator (14) representing a conditional dependence $\mathfrak{c}_{\xi}(\mu)$; see also Remark B.2. The parameter domain for this representation can be chosen as the L^2 -closure of a separable Hilbert space. For conditional dependencies between time-dependent data, this domain has a highly-structured and conveniently explicit form; see Proposition 4.10 and Lemma 3.12. These structural benefits, as shown in Section 5, facilitate translating the operator representation (17) into an efficient scheme for estimating $\mathfrak{c}_{\xi}(\mu)$.

The functional analytic viewpoint of Proposition 2.4—specifically, the one-to-one correspondence between (9) and (14)—facilitates a rigorous and straightforward formulation of (10) in terms of operator learning. For any fixed pair \mathcal{X}, \mathcal{Y} of Polish spaces, denote $\mathcal{Z} \coloneqq \mathcal{X} \times \mathcal{Y}$ and let $\mathscr{Z} \coloneqq \mathcal{Z}^{\mathbb{N}}$.

Definition 2.6 (Learning CDs). Let $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$. A sequence of operator-valued maps

$$\hat{T}_n : \mathscr{Z} \longrightarrow \left[\mathcal{L}^2(\mu_{\mathcal{Y}}) \to \mathcal{L}^2(\mu_{\mathcal{X}}) \right] \qquad (n \in \mathbb{N})$$
 (18)

will be called a consistent estimator of $\mathfrak{c}_{\mu\chi}(\mu)$ if, for any iid samples $Z_1, Z_2, \ldots \sim \mu \eqqcolon \mathbb{P}$, we have

$$\hat{T}_n((Z_j)) \xrightarrow[n \to \infty]{\mathbb{P}} \mathfrak{c}_{\mu_{\mathcal{X}}}(\mu) \quad \text{strongly, that is:} \quad \lim_{n \to \infty} \mathbb{P}\Big(\big\| \mu_{\cdot}(f) - \hat{T}_n((Z_j))(f) \big\|_{L^2(\mu_{\mathcal{X}})} \ge \varepsilon \Big) = 0 \quad (19)$$

for each $f \in \mathcal{L}^2(\mu_{\mathcal{Y}})$, which is required to hold for any $\varepsilon > 0$. We call universally consistent a sequence (\hat{T}_n) for which both (18) and (19) hold for every measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$.

Remark 2.7 (Stochastic Process Regression). The main goal of this paper is to propose a universally consistent estimator of conditional dependencies on spaces \mathcal{X}, \mathcal{Y} of multivariate time-dependent data. In other words (see Remark 2.3), we aim to estimate the conditional dependencies between coupled multidimensional (discrete- or continuous-time) stochastic processes X and Y. In this context, we recognize in (14) the familiar expressions

$$\mu_{\cdot}(f) : x \mapsto \mu_{x}(f) \equiv \mathbb{E}[f(Y) \mid X = x] \quad \text{and} \quad \mu_{X}(f) \coloneqq \mu_{\cdot}(f) \circ X \equiv \mathbb{E}[f(Y) \mid X]$$
(20)

(recall (13)), or in other words, $\mu_{\cdot}(f)$ is precisely the regression function of f(Y) on \mathcal{X} . In particular,

$$\mathbb{P}(Y \in (\cdot) \mid X = x) \equiv \mu_x(\mathbb{1}_{(\cdot)}) : \mathcal{B}(\mathcal{Y}) \ni A \longmapsto \mu_x(\mathbb{1}_A) \left(=\mu_x(A)\right) \in [0, 1]$$
(21)

is the conditional distribution of Y given X = x, cf. (12), where $\mathbb{P}(Y \in (\cdot) | X) \coloneqq \mu_X(\mathbb{1}_{(\cdot)})$ as usual. If the space \mathcal{Y} of possible outcomes is finite, then we have the explicit embedding

$$\mathbb{P}(Y \in A \mid X) = \sum_{y \in \mathcal{Y}} \mathbb{1}_A(y) \mu_X(\{y\}) \stackrel{(15)}{=} \lim_{j \to \infty} \sum_{y \in \mathcal{A}} \left\langle \ell_j^{(y)}, \mathbb{E}[q(Y) \mid X] \right\rangle \qquad \left(A \in \mathcal{B}(\mathcal{Y})\right)$$
(22)

of the conditional law $\mathbb{P}(Y \in \cdot | X)$, where $\{(\ell_j^{(y)}) | y \in \mathcal{Y}\} \subset \mathcal{H}_{\mathcal{Y}}$ is some (deterministic) family of extractors obtainable as minimizing sequences of the regression problem $\inf\{\|\mathbb{1}_{\{y\}} - \langle \ell, q \rangle\|_{L^2(\mathbb{P}_Y)}^2 | \ell \in \mathcal{H}_{\mathcal{Y}}\}$ (Section 4.3). A consistent estimator (\hat{T}_n) , as in (18) and (19), then yields an approximation

$$\lim_{n \to \infty} \mathbb{P} \Big(\sup_{A \in \mathcal{B}(\mathcal{Y})} \left| \mathbb{P}(Y \in A \,|\, X) - \hat{T}_n(A) \right| \ge \varepsilon \Big) = 0$$
(23)

for any given $\varepsilon > 0$, where $\hat{T}_n(A) := \hat{T}_n(\mathbb{1}_A)$ (see Section 5.5). The representation (17) (if it exists) suggests a general blueprint for such estimators (\hat{T}_n) , namely a 'composite architecture' of the form

$$\hat{T}_n(f) = \left\langle \hat{\ell}_n^{(f)}, \hat{\Xi}_n^q \right\rangle \tag{24}$$

for each $f \in \mathcal{L}^2(\mathbb{P}_Y)$, where $\hat{\ell}_n^{(f)}$ and $\hat{\Xi}_n^q$ are estimators of the components $(\ell_n^{(f)})$ and $\mu_{\cdot}(q) \equiv \mathbb{E}[q(Y) | X = \cdot]$ from (17), respectively. This paper presents two universally consistent instances of this architecture and analyses their statistical approximation properties (19) (see Sects. 5.3–5.5).

A prominent recent application of this regression framework for stochastic processes is found in computational linguistics:

Example 2.8 (Large Language Models). A particular use case of the above regression framework is the efficient representation and estimation of conditional distributions (21) of a [random vector or] stochastic process Y given a stochastic process X, as outlined in (22) and (23). This is the central task of Large Language Models (LLMs) such as Generative Pre-trained Transformers (GPT), whose goal it is to predict—that is, to sample from an estimated conditional distribution—a contextually appropriate response $y \sim Y$ (potential text completions) given a sequential input $x \sim X$ (supplied text prompts). In this context, \mathcal{X} is the space of all possible inputs, such as appropriately vectorized text prompts of a given length, and \mathcal{Y} is the (finite, limited by the vocabulary size) set of all possible token sequences of a given length that make up the model's potential responses. An LLM learns from data a conditional dependence $\mathfrak{c}_{\xi}(\mu)$ as in (12), for μ the joint distribution of (X, Y) and ξ the law of X, and this $\mathfrak{c}_{\mathfrak{c}}(\mu)$ captures the complex statistical relationship between input sequences (X) and potential outcomes (Y). This learning is operationalised (as expressed by (19) and (23)) through an estimation of the conditional distribution $\mathbb{P}(Y \in (\cdot) \mid X)$, and in practice translates to computing the (empirical) likelihood of different possible model responses given a particular user input. (Given a prompt, the model calculates the probability of various next tokens based on the learned conditional distribution.) The necessary computations are performed by estimators T_n of the general form (18), which in contemporary models are typically implemented as sequence-processing neural networks with (empirically fine-tuned) architectures like transformers. These estimators serve to translate sampled text data into approximations of the conditional mean embeddings (14) and (17) (or other suitable representations of the probability kernel $\mathfrak{c}_{\mathfrak{c}}(\mu)$ that functionally express the associative relationship between X and Y. This yields a statistical approximation of the intricate relationship between input and output text data, underpinning the generative capabilities of LLMs and providing us with probabilistically founded function approximators capable of emulating the nuances and complexities of natural language. Applied to this context, the probabilistic concepts and novel statistical estimation methods in this paper may bring practical improvements and contribute to our understanding of the mathematical foundations of LLMs and statistical language models at large. For more information on the statistics of LLMs, see [?] and the references therein.

A mathematically insightful interpretation of the regression operator $\mathfrak{c}_{\xi}(\mu)$ in (14) is to view it as providing, in an L^2 -variational sense, the optimal transfer of information from \mathcal{Y} to \mathcal{X} :

Lemma 2.9. The operator (14) is isometrically isomorphic to the restricted projection operator

$$P_{\mathcal{Y}|\mathcal{X}} : L^2(\mathcal{X} \times \mathcal{Y}, \mathcal{X} \times \mathcal{B}(\mathcal{Y}), \mu) \longrightarrow L^2(\mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{X}) \times \mathcal{Y}, \mu), \quad \tilde{f} \longmapsto P\tilde{f}, \tag{25}$$

for the orthogonal projection $P: L^2(\mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{X} \times \mathcal{Y}), \mu) \twoheadrightarrow L^2(\mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{X}) \times \mathcal{Y}, \mu).$

Proof. See Appendix B.1.3. (Here, $\mathcal{X} \times \mathcal{B}(\mathcal{Y}) \equiv \{\mathcal{X} \times A \mid A \in \mathcal{B}(\mathcal{Y})\}$, and likewise for $\mathcal{B}(\mathcal{X}) \times \mathcal{Y}$.)

This paper combines the classical variational perspective of Lemma 2.9 with a specific 'feature-based parametrisation' of the form (17) to develop two universally consistent estimators for conditional dependencies between jointly distributed, multidimensional stochastic processes. These estimators, of the type (24), are derived from an ensemble of core rough path theory objects and enjoy a clear mathematical and conceptual interpretation (Section 4). They are also computationally efficient and supported by broad theoretical guarantees, including explicit convergence rates (Section 5).

3 Time-Dependent Data and the Signature Transform

So far, we have not specified the regressor-regressand relation $X \sim Y$ beyond its inclusion in the product $\mathcal{X} \times \mathcal{Y}$ for general Polish spaces \mathcal{X} and \mathcal{Y} . In this section, we sharpen this assumption and focus on learning statistical relationships between multidimensional time-dependent data. Specifically, we consider conditional dependencies (8) pertaining to (sufficiently regular) subspaces

$$\mathcal{X}$$
 of $C([0,1];\mathbb{R}^{d_1})$ and \mathcal{Y} of $C([0,1];\mathbb{R}^{d_2})$, for any $d_1, d_2 \in \mathbb{N}$. (26)

This setting will allow for an effective and rigorous description of time-dependent statistical relationships betweens multidimensional covariates (X, Y) that are genuinely probabilistic, or modelled as such due to their inherent complexity, and exhibit nuanced temporal and spatial dynamics.

Suitably large data spaces \mathcal{X} and \mathcal{Y} as in (26) are obtained in Section 3.1 as $\|\cdot\|_{\infty}$ -dense linear Polish subspaces of continuous paths by imposing some natural continuity assumptions on the points in $C([0,1];\mathbb{R}^d)$. These data spaces \mathcal{X} and \mathcal{Y} , when adequately normed, are effectively (topological) 'Hilbert manifolds' as they admit single, globally-defined coordinate charts:

$$q_{\mathcal{X}} : \mathcal{X} \longrightarrow \mathcal{H}_{\mathcal{X}} \quad \text{and} \quad q_{\mathcal{Y}} : \mathcal{Y} \longrightarrow \mathcal{H}_{\mathcal{Y}}$$
(27)

mapping the original data spaces \mathcal{X} and \mathcal{Y} into structurally well-behaved Hilbert spaces $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$. In the language of RKHS, the functions (27), called '(bounded) signature transforms' and defined in Sections 3.2 and 3.3, can be viewed as highly structured feature maps for multidimensional sequential data. Serving the role of (15), these maps integrate seamlessly into the context of (10), allowing for a fruitful extension of the general approximation approach (17) (see Sections 4 and 5).

3.1 Spaces of Sequential Data and Stochastic Processes

We specialise the general setting of Section 2 to the case of time-dependent data, and ensure the measure-theoretic adequacy of the resulting state spaces \mathcal{X} and \mathcal{Y} for immediate integration into the framework of Section 2.3. To this end, we define a time-dependent, or 'sequential', datum $z = (z_t)$ as a continuously ordered family of vectors,

$$z \coloneqq (z_t \mid t \in I) \equiv (z_t)_{t \in I} \quad \text{with} \quad z_t \equiv (z_t^1, \dots, z_t^d)^{\mathsf{T}} \in \mathbb{R}^d,$$
(28)

that is, a continuous map $z: I \to \mathbb{R}^d$ over some fixed compact interval $I \subset \mathbb{R}$, usually I = [0, 1]. The probabilistic description of time-dependent information (28) is largely a 'macroscopic' endeavour, and so we view the objects (28) primarily as (highly-structured) points of the so-called path spaces

$$C(I;Z) \coloneqq \{ z \equiv (z_t) \in Z^I \mid I \ni t \mapsto z_t \text{ is continuous} \}, \text{ normed by } \|z\|_{\infty} \coloneqq \sup_{t \in I} |z_t|, (29)$$

where $(Z, |\cdot|)$ is a given normed space, usually $Z = \mathbb{R}^d$ for some $d \in \mathbb{N}$; we write $\mathcal{C}_d \coloneqq C([0, 1]; \mathbb{R}^d)$.

Following on from the general spaces (29), choose I = [0, 1] wlog and consider the Banach space

$$(\mathcal{C}_d \coloneqq C([0,1]; \mathbb{R}^d); \|\cdot\|_{\infty})$$
 of all continuous paths $z \equiv (z_t) : [0,1] \to \mathbb{R}^d$.

For simplicity, this paper operates on the 'smooth core' \mathcal{C}^1_d of absolutely continuous paths in \mathcal{C}_d ,

$$\mathcal{C}_{d}^{1} \coloneqq \left\{ z \in \mathcal{C}_{d} \; \middle| \; \exists \, \dot{z} \in L^{1}([0,1];\mathbb{R}^{d}) : \, z_{\cdot} = z_{0} + \int_{0}^{\cdot} \dot{z}_{s} \, \mathrm{d}s \right\}, \quad \text{norm} \quad \|z\|_{1-\text{var}} \coloneqq |z_{0}| + \|\dot{z}\|_{L^{1}}, \quad (30)$$

the so-called 1-variation norm. We note that the dense subspace \mathcal{C}_d^1 of \mathcal{C}_d is chosen for technical convenience. However, as usual in rough path theory, this choice incurs essentially no loss of generality, as everything that follows can be extended canonically to spaces of rougher paths [18].

Remark 3.1 (Discrete-Time Data). Note that the above setting covers both continuous- and discrete-time data. Discrete-time data can be naturally embedded in the path spaces (30) through order-preserving piecewise-linear interpolation of discretely observed sequential information. For details on this embedding, see e.g. [49, Section B.2].

3.1.1 Measure-Theoretical Preliminaries for Path Spaces

Next, let us note the following convenient topological and measure-theoretical properties.

Lemma 3.2. The space C_d^1 is a Borel subset of $(C_d, \|\cdot\|_{\infty})$ and a separable Banach space wrt. the norm $\|\cdot\|_{1-\operatorname{var}}$, and the spaces $(\mathcal{C}^1_d, \|\cdot\|_{\infty})$ and $(\mathcal{C}^1_d, \|\cdot\|_{1-\operatorname{var}})$ have the same Borel σ -algebra.

Proof. The first part of the lemma is well-known, and the second part is shown in Appdx. B.2.2.

Specifying the sequential setup (26), fix any predictor- and response-dimension $d_{\mathcal{X}}, d_{\mathcal{Y}} \in \mathbb{N}$ and set

$$\mathcal{X} \coloneqq \left(\mathcal{C}^{1}_{d_{\mathcal{X}}}, \|\cdot\|_{1-\operatorname{var}}\right) \quad \text{and} \quad \mathcal{Y} \coloneqq \left(\mathcal{C}^{1}_{d_{\mathcal{Y}}}, \|\cdot\|_{1-\operatorname{var}}\right).$$
(31)

For a topology on the Cartesian product $\mathcal{X} \times \mathcal{Y}$ hosting all relations between the points in (31), let

$$\|(x,y)\|_{1-\operatorname{var}} \coloneqq |(x_0,y_0)| + \sup_{\mathcal{D}} \sum_{(t_{\nu})\in\mathcal{D}} |(x_{t_{\nu}},y_{t_{\nu}}) - (x_{t_{\nu-1}},y_{t_{\nu-1}})| \quad \text{and} \quad \|(x,y)\|_{\infty} \coloneqq \sup_{t\in[0,1]} |(x_t,y_t)|,$$

where the first supremum runs over the set \mathcal{D} of all (finite) dissections (t_{ν}) of the interval [0, 1].

Remark 3.3. Lemma 3.4 below asserts that fusing (31) to the product $(\mathcal{X} \times \mathcal{Y}; \|\cdot\|_{1-\text{var}})$ bears no measure-theoretic complications, confirming that the sequential setting (31) integrates smoothly into the framework discussed in Section 2. Alongside Lemma 3.2, this lemma further ensures that the choice of path-space topology—whether $\|\cdot\|_{1-var}$ or $\|\cdot\|_{\infty}$ —is inconsequential when working with the Borel measure spaces $\mathcal{M}_1(\mathcal{X})$, $\mathcal{M}_1(\mathcal{Y})$, or $\mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$, since both topologies induce identical σ -algebras on the respective data spaces \mathcal{X}, \mathcal{Y} and $\mathcal{X} \times \mathcal{Y}$. ٠

Lemma 3.4. The space $\mathcal{Z} \coloneqq (\mathcal{X} \times \mathcal{Y}, \|\cdot\|_{1 \text{-var}})$ is Polish, and its Borel- σ -algebra $\mathcal{B}(\mathcal{Z}) \equiv \mathcal{B}(\mathcal{Z}, \|\cdot\|_{1 \text{-var}})$ $\|_{1-\operatorname{var}}$ coincides with $\mathcal{B}(\mathcal{Z}, \|\cdot\|_{\infty})$. For maps $X : (\Omega, \mathscr{F}) \to \mathcal{X}$ and $Y : (\Omega, \mathscr{F}) \to \mathcal{Y}$ on a measurable space (Ω, \mathscr{F}) , the joint process (X, Y) is $(\mathscr{F}, \mathcal{B}(\mathcal{Z}))$ -measurable iff X and Y are Borel-measurable. \square

Proof. See Appendix B.2.3.

As for the lift of a measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ to the law of an $\mathcal{X} \times \mathcal{Y}$ -valued pair (X, Y) of random variables, as described in (11), we as usual define a (continuous-time) stochastic process as any map

$$S: \Omega \to C(I; Z)$$
 such that $\omega \mapsto S(\omega) \equiv (S_t(\omega))_{t \in I}$ is $(\mathscr{F}, \mathcal{B}(C(I; Z)))$ -measurable, (32)

where the above is defined over some probability space $(\Omega, \mathscr{F}, \mathbb{P})$ (usually left implicit).

Some useful technical details on stochastic processes (32) are collected in Remark B.3.

3.1.2 Example Applications: Stochastic Process Prediction and Classification

As a brief preview of potential future applications, note that the standard regression task of forecasting a stochastic process's future from its past fits seamlessly into the general framework (11):

Any process $S = (S_t(\omega)) : I \times \Omega \to Z$ as in (32) generates two types of sub- σ -algebras of \mathscr{F} ,

$$\sigma(S) \coloneqq \sigma(S^{-1}(B) \mid B \in \mathcal{B}(C(I;Z))) \quad \text{and} \quad \sigma_{t]}(S) \coloneqq \sigma(S_s \mid s \in I : s \le t), \quad (33)$$

for any fixed $t \in I$. The following well-known observation reminds us that (33) are closely related. More specifically, Lemma 3.5 ensures that the common regression application of predicting the future evolution of a stochastic process (S) based on its history (i.e., $\sigma(S_s \mid s \leq t_0)$, for some $t_0 \in I$) is an immediate special instance of the general framework (11): Simply choose $(X, Y) := (S^{\wedge t_0}, S)$.

Lemma 3.5 (Prediction). Let $S = (S_t)_{t \in I}$ be as in (32), and $t_0 \in I$. Then we have $\sigma_{t_0}(S) = \sigma(S^{\wedge t_0})$ for the stochastic process $S^{\wedge t_0} \coloneqq (S_{t \wedge t_0})_{t \in I}$.

Proof. Directly by definition of $\mathcal{B}(C(I; Z))$ and (33), see e.g. [24, Problem 2.4.4.2 (p. 60)].

A second major machine learning task that can be subsumed under the above framework (11) & (10) is the (probabilistic) classification of time-series/stochastic processes: Suppose that any given datum $x \in \mathcal{X}$, modelled as a sample of some abstract process (32) (that is, $x = S(\omega)$ for an $\omega \in \Omega$) is to be assigned a class label $c_x \in \{c_1, \ldots, c_k\} \subset \mathbb{R}$. To ensure this assignment is well-defined, let us call the process S classifiable (under k-many labels) if there exists a measurable partition $\Omega = \bigsqcup_{i=1}^{k} E_i$ (with $E_i \in \mathscr{F}$ for all $i \in [k]$, and $E_j \cap E_k = \emptyset$ if $j \neq k$) such that: for each $\omega \in \Omega$ there is $i = i_\omega \in [k]$ with $S^{-1}(S(\omega)) \subseteq E_i$. We then define $C \coloneqq \sum_{i=1}^{k} c_i \mathbb{1}_{E_i} : \Omega \to \mathbb{R}$ and call this the class variable associated to S. The following lemma ensures that classification can be treated as a special instance³ of the regression framework (11) & (10) (recalling (13)): Simply choose $(X, Y) \coloneqq (S, C)$.

Lemma 3.6 (Classification). Let S be classifiable and C its associated class variable. Then the class labelling $\lambda : \mathcal{X} \ni x \mapsto c_x \in \{c_i\}$ is well-defined on $S(\Omega)$, the conditional dependence of $\mu := \mathbb{P}_{(S,C)}$ wrt. \mathbb{P}_S is $\mathfrak{c}_S(\mu) = (\delta_{\lambda(x)} \mid x \in S(\Omega))$, and we have $\lambda = \arg \max_{c \in C(\Omega)} \mathbb{P}(C = c \mid S = \cdot) (\mathbb{P}_S \text{-a.e.})$.

Proof. Let there be k-many class labels, say $c_1, \ldots, c_k \in \mathbb{R}$. Since S is classifiable, we know that for each $x \in S(\Omega)$ there is a unique $i_x \in [k]$ with $S^{-1}(\{x\}) \subseteq E_{i_x}$. This defines a function $\varphi: S(\Omega) \to [k], x \mapsto i_x$, which in turn determines the class labelling λ to $\lambda: S(\Omega) \ni x \mapsto c_{\tau(\varphi(x))}$, for $\tau: [k] \xrightarrow{\sim} [k]$ some fixed permutation; this proves the first claim. Next, note that, by construction, $C = \lambda(S)$ pointwise on Ω (assuming $\tau = \text{id wlog}$). This shows $\mathfrak{c}_S(\mu) = (\delta_{\lambda(x)})$, cf. Example 2.1, and also: $\mathbb{P}(C = c \mid S) = \mathbb{E}[\mathbb{1}_{\{c\}}(\lambda(S)) \mid S] = \mathbb{1}_{\{c\}}(\lambda(S)) = \delta_{c,\lambda(S)}$, which proves the last claim. \Box

3.2 The Signature Representation of Time-Dependent Data

We proceed to define the (bounded) signature transforms (27) over paths [9, 35, 12], which are coordinate charts that concisely and efficiently embed the previously introduced data domains (31) into (bounded subsets of) well-organized Hilbert spaces. These charts are conveniently computable, continuous, and graded injections, formulated as multivariate formal power series with coefficients given by rapidly decaying iterated integrals. As later sections will show, the rich structure of these charts, carefully discussed below, positions them as valuable tools for a precise stochastic analysis of time-dependent conditional dependencies (9) within a streamlined functional analytic context.

³ Note that the class variable C lifts to a stochastic process $C : \Omega \to \mathcal{C}_1^1$ under the trivial embedding $\mathbb{R} \hookrightarrow \mathcal{C}_1^1$ which identifies the reals with constant paths.

3.2.1 Definition of the Signature

At the most basic level, the signature transform is a faithful (i.e., one-to-one) compression that maps a path to a hierarchically graded list of countably many numerical coordinates of that path. A convenient indexing of these coordinates requires some basic 'multiindex notation':

Notation 3.7 (Words and Formal Power Series). Let $d \in \mathbb{N}$ be fixed. We denote by

$$[d]^* \coloneqq \{\emptyset, 1, 12, 21, d11, ddd1211d, \ldots\}$$

the free monoid over the alphabet $[d] := \{1, 2, ..., d\}$, representing all finite sequences (or 'words') of zero or more elements from [d] (the 'letters'), with \emptyset symbolizing the empty word. The space $\mathbb{R}[[d]] \cong \mathbb{R}[[\mathbf{x}_1, \ldots, \mathbf{x}_d]]$) denotes the (free) algebra of all multivariate formal power series in the non-commutative variables $\mathbf{1} \cong \mathbf{x}_1, \ldots, \mathbf{m} \cong \mathbf{x}_d$), with $\mathbf{1} \coloneqq \mathbf{1} \cdot \emptyset$ its multiplicative unit. Explicitly,

$$\mathbb{R}[[d]] = \{ \boldsymbol{t} : [d]^* \to \mathbb{R} \mid \boldsymbol{t} \text{ is a map} \} \equiv \{ \sum_{w \in [d]^*} \boldsymbol{t}_w \cdot w \cong (\boldsymbol{t}_w)_{w \in [d]^*} \mid \boldsymbol{t}_w \in \mathbb{R} \} \cong \prod_{\nu=0}^{\infty} (\mathbb{R}^d)^{\otimes \nu},$$
(34)

where each word $\mathbf{i}_1 \cdots \mathbf{i}_d \in [d]^*$ is identified with its associated elementary tensor $e_1 \otimes \cdots \otimes e_d \in (\mathbb{R}^d)^{\otimes m}$, writing $(e_i)_{i \in [d]}$ for the standard basis in \mathbb{R}^d . The length |w| of a word w is defined as the number of its letters, so that any length-k word $(k \in \mathbb{N})$ is of the form

$$w = \mathbf{i}_1 \mathbf{i}_2 \cdots \mathbf{i}_{k-2} \mathbf{i}_{k-1} \mathbf{i}_k \in [d]^*, \quad \text{with} \quad \mathbf{i}_{\nu} \in \{1, \dots, \mathsf{d}\} \text{ for each } \nu \in [k].$$
(35)

All of this extends to the alphabet $[\underline{d}]^*$ with $[\underline{d}] \coloneqq \{0, 1, \dots, d\}$. For any $k \in \mathbb{N}$, we write $\Delta_k \coloneqq \{(t_{\nu}) \in [0, 1]^k \mid 0 \le t_1 \le t_2 \le \dots \le t_k \le 1\}$ for the k-dimensional standard simplex.

Take any path $z \equiv (z^1, \dots, z^d) \in \mathcal{C}_d^1$ with components $z^i \in \mathcal{C}_1^1$. For compact notation, let us for any word $w \equiv \mathbf{i}_1 \mathbf{i}_2 \cdots \mathbf{i}_k \in [d]^*$ $(k \in \mathbb{N})$ consider the *w*-indexed differential *k*-form

$$\mathrm{d} z^w \coloneqq \mathrm{d} z^{i_1} \wedge \mathrm{d} z^{i_2} \wedge \cdots \wedge \mathrm{d} z^{i_k} = \dot{z}_{t_1}^{i_1} \dot{z}_{t_2}^{i_2} \cdots \dot{z}_{t_k}^{i_k} \,\mathrm{d} t_1 \wedge \mathrm{d} t_2 \wedge \cdots \wedge \mathrm{d} t_k.$$

The following map is (essentially) a global coordinate chart that elucidates the spaces (31) of sequential data by embedding their elements into a Hilbert space where they are easier to analyse.

Definition 3.8 (Signature). The signature $\mathfrak{sig} : \mathcal{C}_d^1 \to \mathbb{R}[[d]]$ sends a path z to the formal power series

$$\mathfrak{sig}(z) \coloneqq \sum_{w \in [d]^*} \int_{\Delta_{|w|}} \mathrm{d}z^w \cdot w \cong \left(\int_{\Delta_{|w|}} \mathrm{d}z^w \ \middle| \ w \in [1, \dots, \mathbf{d}]^* \right).$$
(36)

Written out, the w-th signature coefficient of a path $z = (z_t^1, \cdots, z_t^d)_{t \in [0,1]}$, denoted $\mathfrak{sig}_w(z)$, reads

$$\int_{\Delta_{|w|}} \mathrm{d}z^{w} \equiv \int_{0}^{1} \int_{0}^{t_{k}} \int_{0}^{t_{k-1}} \cdots \int_{0}^{t_{3}} \int_{0}^{t_{2}} \mathrm{d}z_{t_{1}}^{i_{1}} \mathrm{d}z_{t_{2}}^{i_{2}} \cdots \mathrm{d}z_{t_{k-2}}^{i_{k-2}} \mathrm{d}z_{t_{k-1}}^{i_{k-1}} \mathrm{d}z_{t_{k}}^{i_{k}}, \tag{37}$$

for any word $w \in [d]^*$ of the general form (35); these are each iterated Lebesgue-Stieltjes integrals.

3.2.2 Some Basic Properties of the Signature

Transitioning from a path's trace to its graph, let us now map a path $z \equiv (z^1, \dots, z^d) \in \mathcal{C}_d$ to

$$\bar{z} \coloneqq (t, z_t)_{t \in [0,1]} \equiv \left(z_t^0, z_t^1, \cdots, z_t^d\right)_{t \in [0,1]} \in \mathcal{C}_{d+1}.$$
(38)

The mapping $\bar{\iota} : z \mapsto \bar{z}$ embeds \mathcal{C}^1_d into \mathcal{C}^1_{d+1} , and composing it with (36) results in an actual embedding of paths into the space of formal power series.



Figure 2: The signature, sig, is a global Hilbert-valued coordinate chart on (sufficiently continuous) time-dependent data that maps a multidimensional path to an ordered list of its iterated integrals.

Theorem 3.9 ([20]). The augmented signature map

$$\mathfrak{sig} \coloneqq \mathfrak{sig} \circ \overline{\iota} : \mathcal{C}_d^1 \longrightarrow \mathbb{R}[[\underline{d}]], \quad z \mapsto \mathfrak{sig}(\overline{z}), \quad is \ injective. \tag{39}$$

Proof. Immediate by [20, Theorem 4] and the fact that, for any $x, y \in C_d^1$, due to the strict monotonicity of their first component the augmented paths \bar{x} and \bar{y} are treelike equivalent iff x = y. \Box

To exploit the transform (36), or rather its augmented version (39), for our analysis, we need to enhance its co-domain with additional analytic structure. Specifically, we see next that a subspace of $\mathbb{R}[[d]]$ [resp. $\mathbb{R}[[\underline{d}]]$] which includes the range im(\mathfrak{sig}) [resp. im(\mathfrak{sig})] of the [augmented] signature, can be easily structured as a Hilbert space. This requires some preliminary basic notation:

Notation 3.10 (Gradation, Projections, and Inner Products). Any $[d]^*$ -indexed infinite tuple $a \equiv (a_w)_{w \in [d]^*} \subset \mathbb{R}$, such as (36), can be injected into the set (34) via $a = \sum_{w \in [d]^*} t_a(w) \cdot w \in \mathbb{R}[[d]]$, where $t_a(\mathfrak{i}_1 \cdots \mathfrak{i}_m) \coloneqq a_{i_1 \cdots i_m}$. Upon grouping the summands in (34) by their wordlength, we get

$$V \equiv \mathbb{R}[[d]] = \prod_{m=0}^{\infty} V_m \quad \text{with} \quad V_m \coloneqq \bigoplus_{w \in [d]^* : |w| = m} \mathbb{R}w$$
(40)

for $V_0 := \mathbb{R}$. The decomposition (40) of $\mathbb{R}[[d]]$ into an infinite product of homogeneous components (V_m) defines a gradation of V, accompanied by projections $\pi_m : V \to V_m$, $\boldsymbol{a} \mapsto \pi_m(\boldsymbol{a}) := \boldsymbol{a}_m \equiv \sum_{|w|=m}^m a_w \cdot w$; we set $\pi_{[m]} \equiv \sum_{\nu=1}^m \pi_\nu : V \longrightarrow V_{[m]} := \bigoplus_{j=0}^m V_j \subset V$. Finally, the inner product

$$\langle \cdot, \cdot \rangle : V \times V \to \overline{\mathbb{R}}, \quad (\boldsymbol{a}, \boldsymbol{b}) \mapsto \sum_{w \in [\boldsymbol{d}]^*} \langle \boldsymbol{a}, w \rangle \cdot \langle \boldsymbol{b}, w \rangle \eqqcolon \langle \boldsymbol{a}, \boldsymbol{b} \rangle,$$
(41)

is defined as the (infinite) bilinear extension of $\langle u, v \rangle \coloneqq \delta_{uv}, u, v \in [d]^*$. In other words, $[d]^*$ is an ONS wrt. $\langle \cdot, \cdot \rangle$, and we note that $\langle \cdot, \cdot \rangle = \sum_{m \ge 0} \langle \pi_m(\cdot), \pi_m(\cdot) \rangle$ pointwise on $V \times V$. **Example 3.11.** To illustrate the above notation, note for instance that, for d > 4 say,

$$3 \mathfrak{sig}_{21}(x) - 25 \mathfrak{sig}_{44}(x) = \langle 3 \cdot 2\mathbf{1} - 25 \cdot 44, \mathfrak{sig}(x) \rangle = 3 \int_0^1 \int_0^t \mathrm{d}x_s^2 \mathrm{d}x_t^1 - 25 \int_0^1 \int_0^t \mathrm{d}x_s^4 \mathrm{d}x_t^4,$$

$$\langle 10230, \underline{\mathfrak{sig}}(x) \rangle = \int_{\Delta_5} \mathrm{d}x^1 \wedge \mathrm{d}s \wedge \mathrm{d}x^2 \wedge \mathrm{d}x^3 \wedge \mathrm{d}t = \int_0^1 \int_0^{t_5} \int_0^{t_4} \int_0^{t_5} \int_0^{t_2} \dot{x}_{t_1}^1 t_2 \dot{x}_{t_3}^2 \dot{x}_{t_4}^3 t_5 \, \mathrm{d}t_1 \mathrm{d}t_2 \mathrm{d}t_3 \mathrm{d}t_4 \mathrm{d}t_5$$

and further $\pi_4(1232 + 2 \cdot 243 - 0.5 \cdot 22 - 3 \cdot 2334) = 1232 - 3 \cdot 2334$, $\pi_2(231) = 0 (\equiv 0 \cdot \emptyset)$, $\pi_{[3]}(1 - 32 + 1345 + 6 \cdot 333 + 2 \cdot 32) = 1 + 32 + 6 \cdot 333$, and $\langle 3 \cdot 11 - 2 \cdot 132, 123 - 4 \cdot 11 \rangle = -12$. \blacklozenge The structure (41) allows us to configure V as a Hilbert space, under convergence constraints:

The structure (41) above us to compute v as a finder space, under convergence constraints.

Lemma 3.12 (Hilbert Codomain of (36)). For $(V, \langle \cdot, \cdot \rangle)$ the power series algebra from above, let

$$\mathcal{H}_d \coloneqq \left\{ \boldsymbol{t} \in V \mid \|\boldsymbol{t}\| \coloneqq \sqrt{\sum_{m \ge 0} \|\pi_m(\boldsymbol{t})\|_m^2} < \infty \right\} \quad with \quad \|\cdot\|_m \coloneqq \sqrt{\langle \cdot, \cdot \rangle_m} \,, \tag{42}$$

where $\langle \cdot, \cdot \rangle_m \coloneqq \langle \pi_m(\cdot), \pi_m(\cdot) \rangle$ for each $m \ge 0$. Then $(\mathcal{H}_d, \langle \cdot, \cdot \rangle)$ is a separable Hilbert space with orthonormal basis $(w \mid w \in [d]^*)$, containing the image $\mathfrak{sig}(\mathcal{C}^1_d)$ and all of $\mathbb{R}[d]$. Moreover, the maps

$$\mathfrak{sig}: \mathcal{C}_d^1 \to \mathcal{H}_d \qquad and \qquad \mathfrak{\underline{sig}}: \mathcal{C}_d^1 \to \mathcal{H}_{\underline{d}}$$

$$\tag{43}$$

are both continuous wrt. the p-variation topology, for any $p \ge 1$.

The statements of this lemma are all well-known, but see Appendix B.2.4.

Remark 3.13. By weighting the graded components (V_m) in (40) before ℓ^2 -direct-summing them into a composite Hilbert space as in (42), we can generalise the standard space \mathcal{H}_d to a whole family of alternate sig- resp. sig-containing Hilbert codomains \mathcal{H}_d^{γ} . See Remark B.4 for details.

Geometrically, the (augmented) signature map \underline{sig} from (43) serves as a global coordinate chart⁴ for the Hilbert manifold C_d^1 . Similar to general coordinate charts that map from a less intuitive space (like (30)) into a more comprehensible 'analysis space' (such as (42)), we will use the signature representation (36) of the paired spaces (31) to facilitate the analysis of statistical dependencies between sequential data through (17) and (27) via the lens (43).

To this end, the next subsection bridges the gap between (43) and the desired transforms (27).

3.3 Bounded Signature Transforms

The primary aim of this paper is to propose a universally consistent estimator for conditional dependencies in sequential data, as detailed in Definition 2.6. Employing the composite approach (24), as suggested by (17), requires the square integrability, with respect to any measure in $\mathcal{M}_1(\mathcal{Y})$, of the (31)-tailored coordinate chart $q_{\mathcal{Y}}$ from (27) (whose prototype is (15)). The only way to ensure this is for the function $q_{\mathcal{Y}}$ to be bounded. This is achieved by the concepts of tensor normalization and 'robust signatures' introduced in [12], which this section implements to confine the signature transforms (43) to a bounded subdomain of (42) in a structure-preserving fashion. This finalizes the transition from the default charts (43) to the (15)-respecting enhancements (27).

To begin, let us state some basic facts on the growth [wrt. their gradation index m] and continuity of the coordinate representations (43). This section consistently applies Notation 3.10.

Given $\lambda \in \mathbb{R}$, the λ -dilation is the map $\delta_{\lambda} : V \ni \mathbf{t} \mapsto \sum_{m=0}^{\infty} \lambda^m \pi_m(\mathbf{t}) \in V$ (thus, $\delta_1 = \mathrm{id}_V$).

 $[\]overline{4}$ Notwithstanding the continuity of its inverse \mathfrak{sig}^{-1} , which does not concern us in this paper.

Proposition 3.14 (Signature Decay; e.g. [36, Theorem 3.7 (case p = 1)]). We have that

$$\left\|\pi_m(\mathfrak{sig}(x))\right\|_m \le \|x\|_{1-\operatorname{var}}^m/(m!\beta), \quad \text{for each } (x,m) \in \mathcal{C}^1_d \times \mathbb{N}_0, \tag{44}$$

for some [(x,m)-independent] constant $\beta > 1$; i.e., the signature decays factorially. In particular,

$$\mathfrak{sig}(\mathcal{C}_d^1) \ \subset \ \mathcal{H}_d^{\downarrow} \coloneqq \left\{ \boldsymbol{t} \in V \ \Big| \ \sum\nolimits_{m \geq 0} \| \pi_m(\boldsymbol{t}) \|_m \lambda^m < \infty, \ \forall \, \lambda > 0 \right\} \ \subseteq \ \mathcal{H}_d.$$

and consequently $\delta_{\lambda}(\mathfrak{sig}(\mathcal{C}^{1}_{d})) \subset \mathcal{H}^{\downarrow}_{d}$ for each $\lambda > 0$, since clearly $\delta_{\lambda} : \mathcal{H}^{\downarrow}_{d} \to \mathcal{H}^{\downarrow}_{d}$ for each $\lambda > 0$.

Lemma 3.15 ('Strong Continuity'). Endow the subspace $\mathcal{H}_d^{\downarrow}$ with the locally convex topology τ_{\downarrow} that is induced by the family of norms

$$\|\|\cdot\|\|_{\lambda} \, : \, \mathcal{H}_d^{\downarrow} o \mathbb{R}, \quad \boldsymbol{t} \mapsto \|\|\boldsymbol{t}\|\|_{\lambda} \, \coloneqq \, \sum_{m \ge 0} \|\pi_m(\boldsymbol{t})\|\lambda^m, \qquad \lambda > 0$$

Then $(\mathcal{H}_d^{\downarrow}, \tau_{\downarrow})$ is separable and metrizable Hausdorff and, for each $p \geq 1$, the signature transform

$$\mathfrak{sig} : \left(\mathcal{C}_d^1, \|\cdot\|_{p\text{-var}}\right) \to \left(\mathcal{H}_d^{\downarrow}, \tau_{\downarrow}\right) \quad is \ continuous.$$

$$\tag{45}$$

Proof. The topological qualities of $(\mathcal{H}_d^{\downarrow}, \tau_{\downarrow})$ are due to [11, Corollary 2.4], while the continuity assertion about the signature is contained in [11, Corollary 5.5].

Remark 3.16 (Comparison of Topologies). Clearly, for any topological space \mathcal{T} , a map $\varphi : \mathcal{H}_d^{\downarrow} \to \mathcal{T}$ is τ_{\downarrow} -continuous if φ is $\|\cdot\|_{\lambda}$ -continuous for at least one $\lambda > 0$. Moreover, $\mathcal{H}_d^{\downarrow}$ is a subspace of each \mathcal{H}_d^{γ} (Remark B.4), and the above topology τ_{\downarrow} on $\mathcal{H}_d^{\downarrow}$ is *finer* than the (127)-induced subspace topology τ_{γ} on $\mathcal{H}_d^{\downarrow}$ (Appendix B.2.5). Thus, statement (45) also holds for $\mathcal{H}_d^{\downarrow}$ replaced by \mathcal{H}_d^{γ} .

Lemma 3.17. If $\lambda_{\cdot} : (\mathcal{H}_d^{\downarrow}, \tau_{\downarrow}) \to \mathbb{R}_{>0}$ is a continuous positive scalar field, then the map

$$\Lambda : (\mathcal{H}_d^{\downarrow}, \tau_{\downarrow}) \to (\mathcal{H}_d^{\downarrow}, \|\cdot\|), \quad \boldsymbol{t} \mapsto \delta_{\lambda_{\boldsymbol{t}}} \boldsymbol{t}, \quad is \ continuous.$$
(46)

Consequently and for any fixed $p \geq 1$, the Λ -scaled augmented signature transform

$$\Lambda \circ \underline{\mathfrak{sig}} : (\mathcal{C}^1_d, \|\cdot\|_{p\text{-var}}) \to (\mathcal{H}^{\downarrow}_{\underline{d}}, \|\cdot\|), \quad x \mapsto \delta_{\lambda_{\underline{\mathfrak{sig}}(x)}}(\underline{\mathfrak{sig}}(x)), \quad is \ continuous.$$
(47)

Proof. See Appendix B.2.6

With the above observations in mind, we can now proceed to modify the Hilbert charts (43) to ensure they are integrable with respect to any (Borel) probability measure on their respective domains (31), which serve as our regression spaces. The idea, originally proposed in [12], is to continuously inject ('squish') the image $\operatorname{sig}(\mathcal{C}_d^1)$ of the chart (39) into a ball of finite radius in \mathcal{H}_d by composing the chart sig with a bounded 'squeezing dilation' of the form (46). Provided that the chosen dilation (46) is also a continuous injection, this ensures that the squished signature (47) is a Hilbert chart of the envisioned kind (27). The next definition summarizes this procedure.

Definition 3.18 (cf. [12]). We call *feature normalisation* (fN) any injective map of the form

$$\Lambda : \mathcal{H}_{\underline{d}}^{\downarrow} \to \mathcal{H}_{\underline{d}}^{R} \coloneqq \{ \boldsymbol{t} \in \mathcal{H}_{\underline{d}} \mid \|\boldsymbol{t}\| \le R \}, \quad \boldsymbol{t} \mapsto \delta_{\lambda_{\boldsymbol{t}}} \boldsymbol{t},$$
(48)

for R > 0 some fixed constant and $\lambda : (\mathcal{H}_d^{\downarrow}, \tau_{\downarrow}) \ni \mathbf{t} \mapsto \lambda_{\mathbf{t}} \in \mathbb{R}_{>0}$ continuous. Given an fN Λ , we call

$$\underline{\mathfrak{sig}}_{\Lambda} \coloneqq \Lambda \circ \underline{\mathfrak{sig}} \, : \, \mathcal{C}^1_d \longrightarrow \mathcal{H}^{R_{\Lambda}}_{\underline{d}} \tag{49}$$

a (Λ -)bounded signature transform (here, $R_{\Lambda} \coloneqq \sup\{\|\Lambda(t)\| | t \in \mathcal{H}_{d}^{\downarrow}\}$).

[12, Section 3.2] show that feature normalisations exist and can be conveniently constructed. Furthermore, [12] show that the specific form (46) of these transforms preserves certain algebraic properties of \mathfrak{sig} , which is crucial for proving that bounded signatures fulfil the universality property (15); see Section 4 below. The following result is a slight extension of [12, Theorem 21].

Proposition 3.19 (cf. [12, Theorem 21]). Let $\Lambda = \delta_{\lambda_{(\cdot)}}$ be an fN with scalar field $\lambda_{\cdot} : (\mathcal{H}_{\underline{d}}^{\downarrow}, \tau_{\downarrow}) \to \mathbb{R}_{>0}$, and set $\underline{\lambda}_{(\cdot)} \coloneqq \lambda_{\cdot} \circ \underline{\mathfrak{sig}}(\cdot)$. Further, abbreviate $\mathcal{Z} \coloneqq (\mathcal{C}_{d}^{1}, \|\cdot\|_{1-\operatorname{var}})$ and denote by

$$\mathcal{A}_{\Lambda} \coloneqq \operatorname{span}_{\mathbb{R}} \left\{ \underline{\xi}_{w}^{\Lambda} : \mathcal{Z} \ni x \mapsto \underline{\lambda}_{x}^{|w|} \int_{\Delta_{|w|}} \mathrm{d}\bar{x}^{w} \mid w \in [\underline{d}]^{*} \right\}, \quad thus \quad \underline{\mathfrak{sig}}_{\Lambda} = \sum_{w \in [\underline{d}]^{*}} \underline{\xi}_{w}^{\Lambda} \cdot w, \quad (50)$$

the linear span of the component functions $(\underline{\xi}_w^{\Lambda} \mid w \in [\underline{d}]^*)$ of the Λ -bounded signature $\underline{\mathfrak{sig}}_{\Lambda}$. Then \mathcal{A}_{Λ} is a point-separating and non-vanishing⁵ subalgebra of $C_b(\mathcal{Z})$, the algebra of all bounded continuous functions on \mathcal{Z} . Moreover, \mathcal{A}_{Λ} is dense in $(C_b(\mathcal{Z}), \tau_{\mathrm{str}}^{\mathcal{Z}})$, for $\tau_{\mathrm{str}}^{\mathcal{Z}}$ the strict topology⁶ on $C_b(\mathcal{Z})$. If instead the domain \mathcal{Z} is a bounded subset of $(\mathcal{C}_d^1, \|\cdot\|_{1-\mathrm{var}})$, then all of the above holds for $\Lambda = \mathrm{id}_{\mathcal{H}_{\bullet}^{\perp}}$.

Proof. The fact that \mathcal{A}_{Λ} is an algebra is due to the span of all iterated integrals (37) being closed under multiplication, a property preserved by the dilation structure (46). The inclusion $\mathcal{A}_{\Lambda} \subset C_b(\mathcal{Z})$ follows clearly from (48) and Lemma 3.17. That \mathcal{A}_{Λ} is point-separating is ensured by Theorem 3.9 and the injectivity of Λ , while \mathcal{A}_{Λ} being non-vanishing is evident since $\underline{\xi}_{\emptyset}^{\Lambda} = 1$ by definition of (49).

Given the preceding discussion, the statements regarding the strict topology follow immediately from [19, Thm. 3.1], cf. [12, Thm. 21]. All details are provided in Appdx. B.2.7 for completeness. \Box

Proposition 3.19, i.e. the asserted denseness of \mathcal{A}_{Λ} in $C_b(\mathcal{Z})$ with respect to the strict topology, generalises the Stone-Weierstrass theorem to the non locally-compact setting of path spaces (31). This aspect of (49) holds secondary significance for the main objective of this article and will be utilized merely to substantiate a few more peripheral corollaries in Section 4.3.

4 Signature Representation of Time-Dependent Association

The signature transform, introduced as a global Hilbert-valued coordinate map on spaces of sufficiently continuous paths (Section 3.2), is a powerful and well-structured tool for analysing timedependent data, and it retains this property when injected into bounded subsets (Section 3.3). With this data-analytic concept in hand, we now return to our main task (10) of estimating statistical dependencies, represented as regression operators (Proposition 2.4), between sets (26) of time-dependent multidimensional data. Our strategy for achieving this is as follows:

The statistical approximation of conditional dependencies (8) relies on the σ -algebra generated by the regressor, which serves as an information reservoir for the approximants of (14). Section 4.1 uses the bounded signature transform to discretise this reservoir (Lemma 4.3) and exhaust its operationalisable information with a system of simple signature-based functions (Proposition 4.5 and Corollary 4.6). These results provide us with a tailor-made architecture (58) for the convex approximation of statistical dependencies, with the gradation index m in (40) of the signature serving as the central control parameter for the resolution of this approximation. This architecture is put to use in Section 4.2, where the conditional expectation of the bounded signature of a process Y given another process X, as well as the conditional expectation of a random vector Z given X, are characterised as the solutions of conveniently approximable, convex semi-infinite linear least squares problems (Theorem 4.8 and Corollary 4.9). These solutions yield 'variational signature representations' of the respective conditional expectations, which are then extended in Section 4.3 to asymptotic variational representations (17) of the regression operator (14) (Proposition 4.10).

⁵ A family $\mathcal{A} \subseteq C(\mathcal{Z})$ will be called *non-vanishing* if: $\forall x \in \mathcal{Z}$ there exists $\varphi \in \mathcal{A}$ such that $\varphi(x) \neq 0$. ⁶ See e.g. [19] or, for internal reference, Definition B.5.

In combination, this provides a convex nonparametric regression framework for stochastic processes that makes essentially no assumptions about the joint distribution of regressor and regressand, and allows for a comprehensive statistical estimation theory and error control (Section 5).

4.1 A Dense Discretisation of Time-Dependent Information

To begin, we re-instantiate the sequential setting of Section 3.1, meaning that from here onwards:

$$\mathcal{X} \coloneqq \left(\mathcal{C}^{1}_{d_{\mathcal{X}}}, \|\cdot\|_{1-\operatorname{var}}\right) \quad \text{and} \quad \mathcal{Y} \coloneqq \left(\mathcal{C}^{1}_{d_{\mathcal{Y}}}, \|\cdot\|_{1-\operatorname{var}}\right) \quad \text{for some} \quad d_{\mathcal{X}}, d_{\mathcal{Y}} \in \mathbb{N}, \tag{51}$$

denoting $(\underline{d}, \underline{\tilde{d}}) \coloneqq (\underline{d}_{\mathcal{X}}, \underline{d}_{\mathcal{Y}})$ for brevity. Then, any probability measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ admits a lift $\mu = \mathbb{P}_{(X,Y)}$ to the law of some pair of jointly distributed stochastic processes

$$X : \Omega \to \mathcal{X} \quad \text{and} \quad Y : \Omega \to \mathcal{Y}$$
 (52)

over some (complete) probability space $(\Omega, \mathscr{F}, \mathbb{P})$, see (11) and Remark 2.3. These X and Y induce

sub-
$$\sigma$$
-algebras $\Sigma_X \coloneqq \sigma(X)$ and $\Sigma_Y \coloneqq \sigma(Y)$ of \mathscr{F}

Remark 4.1. Let us include a few basic remarks on the above σ -algebras and measurability.

- (i) We recall the set system Σ_X to serve as an information base for the best-approximation (in the Bochner-L²-distance) Y_X =: E[Y | X] of the process Y from within all Y-valued measurable functions of X, cf. also Lemma 2.9. Recall that this best-approximation E[Y | X] exists if Y is Bochner-integrable (which we don't need to assume for our actual purposes), and it is unique up to P-almost sure equality; see e.g. [48, Theorem II.2.1]. For the 'canonical' probability space (Ω, 𝔅, P) := (𝔅 × 𝔅, 𝔅(𝔅 × 𝔅), μ) (cf. Remark B.1), the sub-σ-algebras Σ_X and Σ_Y of 𝔅 correspond to the sub-σ-algebras 𝔅(𝔅) × 𝔅 and 𝔅 × 𝔅(𝔅) of 𝔅(𝔅 × 𝔅), respectively.
- (ii) Throughout, any random variable is understood to be Borel-measurable. In particular: for any Hilbert space $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$ and function $\mathbb{Z} : \Omega \to \mathcal{H}$, we call \mathbb{Z} an \mathcal{H} -valued random variable $(\Leftrightarrow: \mathbb{Z} \in \mathcal{L}^0(\mathbb{P}; \mathcal{H}))$ iff \mathbb{Z} is $(\mathscr{F}, \sigma(\|\cdot\|_{\mathcal{H}}))$ -measurable. The latter measurability is equivalent to \mathbb{Z} being Bochner-measurable thanks to [58, Prop. 1.8], provided that \mathcal{H} is separable.

A best-approximation Y_X of the full-process Y given X is generally difficult to come by, especially due to the 'analytical intractability' of the path space \mathcal{Y} . Luckily, such an approximation Y_X is also not directly required: Instead of their Banach-valued default representation (52), it suffices for (10) to examine the covariates X and Y through the lens of their bounded Hilbert-coordinates (49).

So for any two feature normalisations $\Xi: \mathcal{H}_{\mathcal{X}}^{\downarrow} \to \mathcal{H}_{\mathcal{X}}$ and $\Lambda: \mathcal{H}_{\mathcal{Y}}^{\downarrow} \to \mathcal{H}_{\mathcal{Y}}$, see (48), we denote by

$$X^{\Xi} \coloneqq \underline{\mathfrak{sig}}_{\Xi}(X) \quad \text{and} \quad Y^{\Lambda} \coloneqq \underline{\mathfrak{sig}}_{\Lambda}(Y)$$
 (53)

the respective Hilbert-valued representations of X and Y. We then know from (49) that X_{Ξ} and \mathbb{Y}_{Λ} are bounded Hilbert-valued random variables, specifically:

$$\mathbb{X}^{\Xi} \in \mathcal{L}^{\infty}(\mathbb{P}; \mathcal{H}_{\mathcal{X}}) \coloneqq \mathcal{L}^{\infty}(\Omega, \mathscr{F}, \mathbb{P}; \mathcal{H}_{\mathcal{X}}) \qquad \text{and} \qquad \mathbb{Y}^{\Lambda} \in \mathcal{L}^{\infty}(\mathbb{P}; \mathcal{H}_{\mathcal{Y}})$$

with $\mathcal{H}_{\mathcal{X}} \coloneqq \mathcal{H}_{\underline{d}}$ and $\mathcal{H}_{\mathcal{Y}} \coloneqq \mathcal{H}_{\underline{\tilde{d}}}$ the Hilbert spaces of square-summable power series from (42). Differentiating degrees of integrability $1 \leq p < \infty$, we also introduce the (Bochner-) L^p -spaces

$$L^{p}(\mathbb{P};\mathcal{H}) \coloneqq \left\{ \mathbb{Z} \in L^{0}(\mathbb{P};\mathcal{H}) \mid \|\mathbb{Z}\|_{L^{p}(\mathcal{H})} < \infty \right\} \quad \text{for} \quad \|\mathbb{Z}\|_{L^{p}(\mathcal{H})} \coloneqq \left(\int_{\Omega} \|\mathbb{Z}(\omega)\|_{\mathcal{H}}^{p} \mathrm{d}\mathbb{P}\right)^{1/p} = \mathbb{E}\left[\|\mathbb{Z}\|_{\mathcal{H}}^{p}\right]^{\frac{1}{p}}.$$
(54)

Recall that $(L^p(\mathbb{P};\mathcal{H}), \|\cdot\|_{L^p(\mathcal{H})})$ is Banach for $1 \leq p < \infty$, and Hilbert for p = 2 with inner product

$$\langle \mathbb{Z}_1, \mathbb{Z}_2 \rangle_{L^2(\mathcal{H})} \coloneqq \int_{\Omega} \langle \mathbb{Z}_1(\omega), \mathbb{Z}_2(\omega) \rangle_{\mathcal{H}} d\mathbb{P} = \mathbb{E}[\langle \mathbb{Z}_1, \mathbb{Z}_2 \rangle_{\mathcal{H}}].$$
 (55)

The boundedness of (53) ensures that there will be no integrability concerns for us here.

Lemma 4.2. We have that $X^{\Xi} \in L^p(\mathbb{P}; \mathcal{H}_{\mathcal{X}})$ and $Y^{\Lambda} \in L^p(\mathbb{P}; \mathcal{H}_{\mathcal{Y}})$, for each $p \geq 1$.

Proof. The inclusion $\mathbb{Y}^{\Lambda} \in \mathcal{L}^0(\mathbb{P}; \mathcal{H}_{\mathcal{Y}})$ is due to (47) and the fact that in our setting, Bochner- and Borel-measurability coincide thanks to the separability of $\mathcal{H}_{\mathcal{Y}}$, see Remark 4.1 (ii). The unrestricted integrability of \mathbb{Y}^{Λ} is clear since Λ is bounded by definition (48) of a feature normalisation.

Let us also make the following simple but useful observation, proved in Appendix B.2.8.

Lemma 4.3. For each feature normalisation $\Xi : \mathcal{H}^{\downarrow}_{\mathcal{X}} \to \mathcal{H}_{\mathcal{X}}$, we have that

$$\Sigma_X = \sigma(\underline{\mathfrak{sig}}_{\Xi}(X)) = \sigma(\underline{\xi}_w^{\Xi}(X) \mid w \in [\underline{d}]^*)$$
(56)

(in the notation of Proposition 3.19).

Let us further introduce the spaces (which for p = 2 can be identified with the domain of (25))

$$L^{p}_{X}(\mathcal{H}_{\mathcal{Y}}) \coloneqq L^{p}(\mathbb{P}, \Sigma_{X}; \mathcal{H}_{\mathcal{Y}}) \equiv \{\mathbb{Z} \in L^{p}(\mathbb{P}; \mathcal{H}_{\mathcal{Y}}) \mid \mathbb{Z}^{-1}(\mathcal{B}(\mathcal{H}_{\mathcal{Y}})) \subseteq \Sigma_{X}\} \quad (p \in [1, \infty))$$
(57)

of all Σ_X -measurable mean-*p*-integrable $\mathcal{H}_{\mathcal{Y}}$ -valued random variables, and the parameter space

$$\mathfrak{L}^2_{X_{\Xi}} \coloneqq \left\{ \alpha : \left[\underline{\tilde{d}}\right]^* \to \mathbb{R}[\underline{d}], \, w \mapsto \alpha_w \mid \|\alpha\|_{\mathfrak{L}}^2 \coloneqq \sum_{w \in [\underline{\tilde{d}}]^*} \mathbb{E}\left[\langle \alpha_w, \, \mathbb{X}^{\Xi} \rangle^2\right] < \infty \right\}$$

The following proposition provides an architecture of readily adjustable model functions through which we can compute the conditional expectation $\mathbb{E}[\mathbb{Y}^{\Lambda} | X]$ as the unique solution to a convex and practically implementable 'least-squares type' optimization problem (Theorem 4.8). The proof of this proposition is based on the classical fact that if a class of bounded functions contains all indicator functions of a π -system and is closed under taking pointwise limits of uniformly bounded, monotone sequences, then this class contains all bounded, measurable functions with respect to the σ -algebra generated by the π -system. This fact is formulated as Lemma 4.4 below.

Lemma 4.4 (Functional Monotone Class). Suppose that H is a vector space of bounded realvalued functions on a measurable space \mathscr{X} such that H contains the constants and is closed under bounded monotone convergence (that is, for any increasing sequence $(\varphi_k) \subset H$ of positive, uniformly bounded functions, the (pointwise) limit $\varphi := \lim_{k \to \infty} \varphi_k$ lies in H). Let \mathfrak{C} be a subset of H which is closed under pointwise multiplication, then H contains all $\sigma(\mathfrak{C})$ -measurable bounded functions.

Proof. See, for instance, [22, Theorem A.1].

Proposition 4.5. For any given feature normalisation $\Xi : \mathcal{H}_{\mathcal{X}}^{\downarrow} \to \mathcal{H}_{\mathcal{X}}$, consider the family

$$\Psi_X \coloneqq \left\{ \psi_\alpha : \mathcal{X} \to \mathcal{H}_Y \cup \{\infty\} \ \middle| \ \alpha \equiv (\alpha_w) \in \mathfrak{L}^2_{X_\Xi} \right\}, \quad with \quad \psi_\alpha \coloneqq \sum_{w \in [\underline{\tilde{d}}]^*} \left\langle \alpha_w, \, \underline{\mathfrak{sig}}_{\Xi}(\,\cdot\,) \right\rangle \cdot w, \ (58)$$

of $\mathfrak{L}^2_{X_{\Xi}}$ -parametrised (\mathbb{P}_X -a.e.-defined) 'simple' functions. Then, the family of random variables

$$\Psi_X(X) \coloneqq \{\psi(X) \mid \psi \in \Psi_X\} \text{ is an } \|\cdot\|_{L^2(\mathcal{H}_{\mathcal{Y}})} \text{-dense subset of } L^2_X(\mathcal{H}_{\mathcal{Y}}).$$
(59)

Proof. We first show $\Psi_X(X) \subseteq L^2_X(H_{\mathcal{Y}})$, for which we fix any $\alpha \in \mathfrak{L}^2_X$. By monotone convergence,

$$\|\psi_{\alpha}(X)\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}^{2} = \sum_{w \in [\tilde{d}]^{*}} \mathbb{E}\big[|\langle \alpha_{w}, \underline{\mathfrak{sig}}_{\Xi}(X)\rangle|^{2}\big] = \|\alpha\|_{\mathfrak{L}}^{2} < \infty.$$

In particular, $\|\psi_{\alpha}(X)\| < \infty$ almost surely, whence the $(\Sigma_X, \mathcal{B}(\mathcal{H}_Y))$ -measurability of $\psi_{\alpha}(X)$ follows via Pettis measurability theorem (e.g. [58, Theorem 1.11], applicable by Remark B.3 (iv)) from the fact that: (a) the compositions $\langle w, \psi_{\alpha}(X) \rangle = \langle \alpha_w, \mathfrak{sig}_{\Xi}(X) \rangle$ are each Σ_X -measurable by Lemma 3.17 and Lemma 4.3, and (b) ($\langle w, \cdot \rangle \mid w \in [d]$) is a (Schauder) basis of the [topological] dual of \mathcal{H}_{d} . (Note that in our setting, Bochner- and Borel-measurability coincide; see Remark 4.1 (ii).)

Next we prove (59), for which we consider any fixed $p \in [1,\infty)$ (a setting that, for the sake of generality, extends beyond the requirements of (59), where only p = 2 is needed). For any $\mathbb{Z} \in L^p_X(\mathcal{H}_{\mathcal{Y}}) \eqqcolon G \text{ and any given } w \in [\tilde{d}]^*, \text{ the coordinate } \mathbb{Z}_w \coloneqq \langle w, \mathbb{Z} \rangle \text{ is in } L^p(\Omega, \Sigma_X, \mathbb{P}) \eqqcolon G_1.$

$$H := \overline{\mathcal{A}_{\Xi}(X)}^{L^{r}} \cap L^{\infty}$$
 and $\mathfrak{C} := \mathcal{A}_{\Xi}(X).$

In other words, \mathfrak{C} is the vector space of images of X under the maps from \mathcal{A}_{Ξ} , where \mathcal{A}_{Ξ} the space of all Ξ -scaled signature polynomials as defined by $(50)|_{\Lambda=\Xi}$, and H is the set of all bounded Σ_X -measurable random variables which are in the G_1 -closure of \mathfrak{C} . From Proposition 3.19 we know that \mathcal{A}_{Ξ} is a subalgebra of $C_b(\mathcal{X})$, and consequently \mathfrak{C} is a subset of H which is closed under pointwise multiplication. Next, let us show that H satisfies the hypotheses of Lemma 4.4: First, it is clear that H is a vector space (as the intersection of two vector spaces) which also contains the constants since \mathfrak{C} contains the constants. To check for the appropriate closedness of H, note that for any monotone sequence $(\mathbf{z}_k) \subset H$ such that $\sup_k |\mathbf{z}_k| \leq C$ for some C > 0 and $\mathbf{z} := \lim_{k \to \infty} \mathbf{z}_k$ pointwise, we have $\mathbf{z}_k \to \mathbf{z}$ in L^p by dominated convergence, implying $\mathbf{z} \in H$ as required.

The above pair (H, \mathfrak{C}) thus qualifies for the application of Lemma 4.4, which yields that H contains all bounded $\sigma(\mathfrak{C})$ -measurable functions. But since $\mathfrak{C} = \operatorname{span}_{\mathbb{R}}\{\underline{\xi}_{w}^{\Xi}(X) \mid w \in [\underline{\tilde{d}}]\}$ and hence $\sigma(\mathfrak{C}) = \sigma(\underline{\xi}_{w}^{\Xi}(X) \mid w \in [\underline{\tilde{d}}])$, Lemma 4.3 implies $\sigma(\mathfrak{C}) = \Sigma_{X}$, which shows that in fact we proved

$$G_1 \cap L^{\infty} \subseteq H \subseteq \overline{\mathfrak{C}}^{L^p}.$$
 (60)

Before using this observation to prove (59), note that for the above w-coordinate \mathbb{Z}_w we have

$$\mathbb{Z}_w = L^p \lim_{n \to \infty} \mathbb{Z}_w^{\langle n \rangle} \quad \text{for the truncations} \quad \mathbb{Z}_w^{\langle n \rangle} \coloneqq \max\left(-n, \min(\mathbb{Z}_w, n)\right) \in G_1 \cap L^\infty$$

(the L^p -convergence holds by dominated convergence). But since $(\mathbb{Z}_w^{\langle n \rangle}) \subset \overline{\mathfrak{C}}^{L^p}$ by (60), we find

$$\mathbb{Z}_w \in \overline{\mathfrak{C}}^{L^p}$$
, that is: $\mathbb{Z}_w = L^p \lim_{j \to \infty} \langle \alpha_{w,j}, \underline{\mathfrak{sig}}_{\Xi}(X) \rangle$ for some $(\alpha_{w,j})_j \subset \mathbb{R}[\underline{d}].$ (61)

Since (61) holds for p = 2 and for all $w \in [\underline{\tilde{d}}]^*$, the conclusion (59) is now within very close reach: Fix any $\varepsilon > 0$. Abbreviating $\varphi_{\wp} \coloneqq \langle \wp, \mathfrak{sig}_{\Xi}(X) \rangle$ for $\wp \in \mathbb{R}[\underline{d}]$, choose some

$$\alpha_w^{\star} \in \mathbb{R}[\underline{d}] \quad \text{such that} \quad \|\mathbb{Z}_w - \varphi_{\alpha_w^{\star}}\|_{G_1}^2 \le \varepsilon^2 (2d_{\mathcal{Y}} + 2)^{-|w|}/2 \qquad \left(w \in [\underline{\tilde{d}}]^*\right). \tag{62}$$

For the coefficient vector $\alpha^* \coloneqq (\alpha^*_w)_{w \in [\underline{\tilde{d}}]^*}$, we then obtain

$$\|\alpha^{\star}\|_{\mathfrak{L}}^{2} = \sum_{w \in [\underline{\tilde{d}}]^{\star}} \mathbb{E}\left[\varphi_{\alpha_{w}^{\star}}^{2}\right] \leq 2\sum_{m=0}^{\infty} \sum_{|w|=m} \|\mathbb{Z}_{w} - \varphi_{\alpha_{w}^{\star}}\|_{G_{1}}^{2} + \|\mathbb{Z}_{w}\|_{G_{1}}^{2} \leq \varepsilon^{2} \sum_{m=0}^{\infty} 2^{-m} + 2\|\mathbb{Z}\|_{G}^{2} < \infty, \quad (63)$$

where the penultimate inequality is due to there being $\sharp\{w \in [\underline{\tilde{d}}]^* \mid |w| = m\} = (d_{\mathcal{Y}} + 1)^m$ many length-*m* words in the monoid $[\underline{\tilde{d}}]^*$. So $\alpha^* \in \mathfrak{L}^2_{X_{\Xi}}$ and thus $\psi_{\alpha^*} \in \Psi_X(X)$, and from (63) we get

$$\|\mathbb{Z}-\psi_{\alpha^{\star}}\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}^{2} = \sum_{m=0}^{\infty} \sum_{|w|=m} \|\mathbb{Z}_{w}-\varphi_{\alpha^{\star}_{w}}\|_{G_{1}}^{2} \leq \varepsilon^{2}.$$

Since both $\mathbb{Z} \in L^2_X(\mathcal{H}_Y)$ and $\varepsilon > 0$ were arbitrary, the claim (59) is established.

The next statement, which is of independent interest, is an immediate corollary of the above proof.

Corollary 4.6. In the setting and notation of Proposition 3.19, we have for any $v \in \mathcal{M}_1(\mathcal{Z})$ that

$$\mathcal{A}_{\Lambda}$$
 is an $\|\cdot\|_{L^{p}(v)}$ -dense subset of $L^{p}(v)$, for each $p \in [1, \infty)$.

Proof. This is precisely statement (61), upon replacing (Ξ, \mathbb{P}_X) by (Λ, v) (which does not impact the proof).

Corollary 4.7. Let $k \in \mathbb{N}$ and $v \in \mathcal{M}_1(\mathcal{X})$, and $\Xi : \mathcal{H}_{\mathcal{X}}^{\downarrow} \to \mathcal{H}_{\mathcal{X}}$ be any feature normalisation. Then for the family of vector-valued (bounded) signature polynomials

$$\psi_{\alpha}^{[k]} \coloneqq \sum_{i=1}^{k} \langle \alpha_{i}, \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \cdot e_{i} : \mathcal{X} \longrightarrow \mathbb{R}^{k}, \qquad \alpha \equiv (\alpha_{1}, \dots, \alpha_{k}) \in \mathfrak{L}_{k}^{2} \coloneqq \left(\mathbb{R}[\underline{d}]^{*}\right)^{\times k},$$

we have that $\{\psi_{\alpha}^{[k]} \mid \alpha \in \mathfrak{L}_k^2\}$ is an $\|\cdot\|_{L^2(\upsilon;\mathbb{R}^k)}$ -dense subset of $L^2(\upsilon;\mathbb{R}^k)$.

Proof. Immediate by Corollary 4.6 upon recalling the definition in (41) of $\langle \cdot, \cdot \rangle \equiv \langle \cdot, \cdot \rangle_{\mathcal{H}_{\mathcal{X}}}$.

4.2 The Conditional Expected Signature of Y given X

Merging Proposition 4.5 with the classical perspective on conditional expectations as L^2 -projections (Lemma 2.9) allows us to compute the conditional expected signature $\mathbb{E}[\mathbb{Y}^{\Lambda} | X]$ as the solution to a conveniently posed convex optimization problem, as Theorem 4.8 below shows.

In order to align this result with the overarching goal of approximating the regression operator (14) on time-dependent data (51) (Section 2.3), let us recall that, in the notation of Proposition 2.4:

$$\mathbb{E}[\mathbb{Y}^{\Lambda} \mid X] = \mu_{\cdot}(q) \quad \text{for} \quad \mu \coloneqq \mathbb{P}_{(X,Y)} \quad \text{and} \quad q \coloneqq \underline{\mathfrak{sig}}_{\Lambda}.$$
(64)

The following theorem thus provides a principled approximation of the inner component (16) for the proposed regression strategy that combines (17) and (64).

Theorem 4.8. Adopting the setting and notation of Proposition 4.5, we have that

$$\mathbb{E}\left[\mathbb{Y}^{\Lambda} \,\middle|\, X\right] = \lim_{k \to \infty} \psi_{\alpha_k}(X) \quad in \ L^2_X(\mathcal{H}_{\mathcal{Y}}) \tag{65}$$

for any minimizing sequence $(\alpha_k) \subset \mathfrak{L}^2_{X_{\Xi}}$ of the convex linear least squares problem

$$\inf_{\alpha \in \mathfrak{L}^{2}_{X_{\Xi}}} \mathbb{E} \Big[\| \mathbb{Y}^{\Lambda} - \psi_{\alpha}(X) \|^{2} \Big].$$
(66)

Denoting $\Phi(\mathbb{Z}) := \mathbb{E}[\|\mathbb{Y}^{\Lambda} - \mathbb{Z}\|^2]$ and $\gamma := \inf_{\alpha \in \mathfrak{L}^2_{X_{\Xi}}} \Phi(\psi_{\alpha}(X))$, the convergence (65) holds \mathbb{P} -almost surely if (α_k) runs 'fast enough' in the sense that it satisfies $\sum_{k=0}^{\infty} (\Phi(\psi_{\alpha_k}(X)) - \gamma)^{1/2} < \infty$.

Sketch of Proof. We combine Proposition 4.5 with an $\mathcal{H}_{\mathcal{Y}}$ -valued version of Lemma 2.9:

Since $\mathbb{Y}_X^{\Lambda} := \mathbb{E}[\mathbb{Y}^{\Lambda} \mid X]$ is the orthogonal projection of $\mathfrak{sig}_{\Lambda}(Y)$ with respect to the decomposition

$$L^2(\mathbb{P};\mathcal{H}_{\mathcal{Y}}) = L^2_X(\mathcal{H}_{\mathcal{Y}}) \oplus L^2_X(\mathcal{H}_{\mathcal{Y}})^{\perp},$$

cf. (Rem. 4.1 (i) and) Lemma 2.9, the Hilbert projection theorem yields the variational characterization

$$\mathbb{Y}_{X}^{\Lambda} = \underset{\mathbb{Z} \in L_{X}^{2}(\mathcal{H}_{\mathcal{Y}})}{\operatorname{arg\,min}} \mathbb{E} \| \mathbb{Y}^{\Lambda} - \mathbb{Z} \|^{2}, \quad \text{and hence} \quad \Phi(\mathbb{Y}_{X}^{\Lambda}) = \underset{\alpha \in \mathfrak{L}_{X_{\Xi}}^{2}}{\operatorname{inf}} \Phi(\psi_{\alpha}(X)), \quad (67)$$

of \mathbb{Y}_X^{Λ} as the L^2 -proximum of \mathbb{Y}^{Λ} within the closed subspace $L_X^2(\mathcal{H}_{\mathcal{Y}})$. Combining this characterisation with Proposition 4.5 leads to the second identity in (67), which in turn implies the desired convergence (65), including the stated almost sure version, via the completeness of $L_X^2(\mathcal{H}_{\mathcal{Y}})$ and by applying the parallelogram law for the norm $\|\cdot\|_{L^2_{\mathcal{Y}}(\mathcal{H}_{\mathcal{Y}})}$. See Appendix B.2.9 for the details. \Box

The next result, which is a corollary to the proof of Theorem 4.8, presents a way to directly compute the conditional expectation (20) for a *fixed* function $f \in \mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k)$. Ultimately, this computation can be achieved from samples of the pair (X, Z) := (X, f(Y)), see Theorem 5.9.

Corollary 4.9. Let $k \in \mathbb{N}$. For any random vector $Z \in L^2(\Omega, \mathscr{F}, \mathbb{P}; \mathbb{R}^k)$ we have that

$$\mathbb{E}[Z \mid X] = \lim_{j \to \infty} \psi_{\alpha_j}^{[k]}(X) \quad in \ L^2_X(\mathbb{R}^k)$$
(68)

for each minimizing sequence $(\alpha_j) \subset \mathfrak{L}^2_k$ of the convex semi-infinite linear least-squares problem

$$\inf_{\alpha \in \mathfrak{L}^2_k} \mathbb{E}\Big[\big| Z - \psi^{[k]}_{\alpha}(X) \big|^2 \Big].$$
(69)

Denoting $\Phi_k(\mathcal{W}) \coloneqq \mathbb{E}[||Z - \mathcal{W}||^2]$ and $\eta \coloneqq \inf_{\alpha \in \mathfrak{L}^2_k} \Phi_m(\psi_{\alpha}^{[k]}(X))$, the convergence (68) holds \mathbb{P} -a.s. if the sequence (α_j) runs 'fast enough' in the sense that it satisfies $\sum_{j=0}^{\infty} \left(\Phi_k(\psi_{\alpha_j}^{[k]}(X)) - \eta\right)^{1/2} < \infty$.

Proof. The proof of Theorem 4.8 stays valid up to the first identity in (67) if we replace $(\mathbb{Y}^{\Lambda}, L^2_X(\mathcal{H}_{\mathcal{Y}}))$ with $(Z, L^2_X(\mathbb{R}^k))$, where $L^2_X(\mathbb{R}^k) \equiv \{\mathcal{W} \in L^2(\mathbb{R}^k) \mid \mathcal{W} \text{ is } (\Sigma_X, \mathcal{B}(\mathbb{R}^k))\text{-measurable}\}$. In particular

$$\mathbb{E}[Z \mid X] = \underset{\mathcal{W} \in L^2_X(\mathbb{R}^k)}{\operatorname{arg\,min}} \mathbb{E}|Z - \mathcal{W}|^2, \tag{70}$$

and since $\{\psi_{\alpha}^{[m]}(X) \mid \alpha \in \mathfrak{L}^2_k\}$ is $\|\cdot\|_{L^2_X(\mathbb{R}^k)}$ -dense in $L^2_X(\mathbb{R}^k)$ by Corollary 4.7, we see Corollary 4.9 follow from (70) in the same way Theorem 4.8 follows from the first identity in (67).

Provided that (samples of) the association (X, Z) := (X, f(Y)) can be directly observed, the approach of Corollary 4.9 offers a more direct alternative for computing $\mathbb{E}[f(Y) | X]$ than the compound strategy (17) based on Theorem 4.8 (via (64)). However, there are scenarios where the latter method proves advantageous, such as when the pairs (X, Y) and (Y, f(Y)) are observed independently and the evaluation $Y \mapsto f(Y)$ is challenging or infeasible. The compound strategy (17) is also beneficial in cases requiring rapid online evaluations of the regression operator (14) across a wide range of different argument functions (e.g. with LLMs, where individual f may be indicators over the model's vocabulary, cf. (21) and Example 2.8). A further comparison of the advantages and disadvantages of each of these estimation methods is given in Remark B.8.

4.3 Signature-Based Regression of Stochastic Processes

The conditional expectation $\mathbb{E}[\mathbb{Y}^{\Lambda} | X]$ of the bounded signature coordinates \mathbb{Y}^{Λ} (see (53)) of a process Y given another process X can be variationally characterized as the solution to a convex least squares problem, see (65) and (66) in Theorem 4.8. This characterisation can be extended to a convex nonparametric estimation method for stochastic process regression (Remark 2.7) through the operator-learning scheme defined by (17) and (64). The central objective of this regression method is the consistent and efficient approximation of X-conditional functionals of Y, see (20), and this subsection presents the probabilistic foundations of this regression approach.

Note that, since the conditional expectation operators in (20) are linear in f, the following results all trivially extend to vector-valued functions of Y.

Proposition 4.10. For any stochastic processes X and Y as in (52) and any $f \in L^2(\mathbb{P}_Y)$, we have

$$\mathbb{E}[f(Y) \mid X] = \lim_{l \to \infty} \left\langle \mathbb{E}[\mathbb{Y}^{\Lambda} \mid X], \, \ell_f^{(l)} \right\rangle \quad in \ L^2_X(\mathbb{R}) \tag{71}$$

for each minimizing sequence $(\ell_f^{(l)} | l \in \mathbb{N})$ of the convex optimization problem

$$\inf_{\ell \in \mathbb{R}[\tilde{d}]} \int_{\mathcal{Y}} \left(f(y) - \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle \right)^2 \mathbb{P}_{Y}(\mathrm{d}y).$$
(72)

In particular, for $(\ell_f^{(l)})$ as above and any sequence $(\psi_{\alpha_k}(X))$ as in (65) or (168), we have

$$\mathbb{E}[f(Y) \mid X] = \lim_{l \to \infty} \lim_{k \to \infty} \left\langle \psi_{\alpha_k}(X), \, \ell_f^{(l)} \right\rangle \quad in \ L^2_X(\mathbb{R}).$$
(73)

Proof. Fix any $f \in L^2(\mathbb{P}_Y)$. By virtue of Corollary 4.6, there is a sequence $(\ell_l) \subset \mathbb{R}[\tilde{d}]$ such that $\int_{\mathcal{Y}} (f(y) - \langle \ell_l, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle)^2 \mathbb{P}_Y(\mathrm{d}y) = \|f(Y) - \langle \ell_l, \underline{\mathfrak{sig}}_{\Lambda}(Y) \rangle\|_{L^2(\Sigma_Y)}^2 \to 0$ as $l \to \infty$, so the infimum (72) is zero. Hence for any minimizing sequence $(\ell_f^{(l)}) \subset \mathbb{R}[\tilde{d}]$ of (72),

$$\begin{aligned} \left\| \mathbb{E}[f(Y) \mid X] - \mathbb{E}[\langle \ell_f^{(l)}, \mathbb{Y}^{\Lambda} \rangle \mid X] \right\|_{L^2(\Sigma_X)} &= \left\| \mathbb{E}[f(Y) - \langle \ell_f^{(l)}, \mathbb{Y}^{\Lambda} \rangle \mid X] \right\|_{L^2(\Sigma_X)} \\ &\leq \left\| f(Y) - \langle \ell_f^{(l)}, \mathbb{Y}^{\Lambda} \rangle \right\|_{L^2(\Sigma_Y)} \longrightarrow 0 \quad \text{as } l \to \infty, \end{aligned}$$
(74)

where the last line is due to Jensen's inequality and the tower property of conditional expectations.

Now inherited from the analogous 'commuting' property of Bochner integrals, we have that

$$\mathbb{E}[\langle \ell_f^{(l)}, \mathbb{Y}^{\Lambda} \rangle \,|\, X] = \left\langle \ell_f^{(l)}, \mathbb{E}[\mathbb{Y}^{\Lambda} \,|\, X] \right\rangle \qquad (\forall \, l \in \mathbb{N})$$
(75)

with probability one, see for instance [48, Theorem II.2.3]. Combining (74) and (75) proves (71).

The convergence (73) is clear from (71) and (65). Indeed: Since, for each $\ell \in \mathbb{R}[d]$, $\eta_k(\ell) := \langle \psi_{\alpha_k}(X), \ell \rangle \in L^2_X(\mathbb{R})$ is merely a (finite) linear combination of $[L^2_X(\mathbb{R})$ -valued] coefficients of $\psi_{\alpha_k}(X)$, the $L^2_X(\mathbb{R})$ -convergence $\langle \mathbb{E}[\mathbb{Y}^{\Lambda} | X], \ell \rangle = \lim_{k \to \infty} \eta_k(\ell)$ is readily implied by (65) [cf. (54)]. \Box

Almost-sure versions of (71) and (73) hold, e.g., if Y is compactly supported and f is continuous.

Corollary 4.11. Let X and Y be as in (52) but with $\mathfrak{D}_Y \coloneqq \operatorname{supp}(\mathbb{P}_Y)$ compact. Then for any function $f \in C(\mathfrak{D}_Y)$ and any minimizing sequence $(\ell_f^{(l)})_l$ of the optimisation problem

$$\inf_{\ell \in \mathbb{R}[\tilde{d}]} \left\| f - \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(\cdot) \rangle \right\|_{\infty;\mathfrak{D}_{Y}}$$
(76)

and any approximating sequence $(\psi_{\alpha_k}(X))$ as in (65) or (168), we have that

$$\lim_{l \to \infty} \left\langle \mathbb{E} \left[\mathbb{Y}^{\Lambda} \, \big| \, X \right], \, \ell_f^{(l)} \right\rangle \stackrel{\alpha}{=} \mathbb{E} \left[f(Y) \, \big| \, X \right] \stackrel{\beta}{=} \lim_{l \to \infty} \lim_{k \to \infty} \left\langle \psi_{\alpha_k}(X), \, \ell_f^{(l)} \right\rangle, \tag{77}$$

where the convergence in (α) holds \mathbb{P} -a.s. and the convergence in (β) holds wrt. $\|\cdot\|_{L^2_X(\mathbb{R})}$; if in addition $(\psi_{\alpha_k}(X))$ is such that (65) or (169) converges almost surely, then (β) also holds \mathbb{P} -a.s.

Proof. See Appendix B.2.10.

Applied to the special cases $f = \mathbb{1}_A$ for arbitrary Borel sets $A \subseteq \mathcal{Y}$, Proposition 4.10 yields the following variational representation of the conditional law $\mathbb{P}(Y \in \cdot | X)$.

Corollary 4.12. For any stochastic processes X and Y as in (52) and any $A \in \mathcal{B}(\mathcal{Y})$, we have

$$\mathbb{P}(Y \in A \mid X) = \lim_{l \to \infty} \lim_{k \to \infty} \left\langle \tilde{\psi}_{\alpha_k}(X), \, \ell_A^{(l)} \right\rangle \quad in \ L^2_X(\mathbb{R})$$

for each sequence $(\tilde{\psi}_{\alpha_k}(X))$ as in (65) or (168) and any minimizing sequence $(\ell_A^{(l)} | l \in \mathbb{N})$ of

$$\inf_{\ell \in \mathbb{R}[\hat{d}]} \int_{\mathcal{Y}} \left(\mathbb{1}_A(y) - \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle \right)^2 \mathbb{P}_Y(\mathrm{d} y).$$

For an almost sure variant of this result, note that if $A \subseteq \mathcal{Y}$ is open, then $\mathbb{1}_A$ admits a monotone pointwise approximation by a sequence of nonnegative uniformly bounded functions in $C_b(\mathcal{Y})$:

$$\mathbb{1}_A \uparrow h_A^{(\nu)} \quad (\nu \to \infty) \quad \text{pointwise, e.g. } h_A^{(\nu)} : y \mapsto \frac{\mathrm{d}(y, \mathcal{Y} \setminus A)}{\mathrm{d}(y, \mathcal{Y} \setminus A) + \mathrm{d}(y, F_\nu(A))} \in C_b(\mathcal{Y}), \quad (78)$$

where $d(y, C) \coloneqq \inf_{z \in C} \|y - z\|_{1-\text{var}}$ (for $C \subset \mathcal{Y}$) and $F_{\nu}(A) \coloneqq \{y \in \mathcal{Y} \mid d(y, \mathcal{Y} \setminus A) \ge \nu^{-1}\}$.

Corollary 4.13. Adopt the setting and notation of Corollary 4.13 and suppose in addition that the support $\mathfrak{D}_Y \coloneqq \operatorname{supp}(\mathbb{P}_Y)$ is compact and $A \subseteq \mathcal{Y}$ is open. Then

$$\lim_{\nu \to \infty} \lim_{\mu \to \infty} \left\langle \mathbb{E} \left[\mathbb{Y}^{\Lambda} \, \big| \, X \right], \, \ell_{A_{\nu}}^{(\mu)} \right\rangle \stackrel{\alpha}{=} \mathbb{P} (Y \in A \mid X) \stackrel{\beta}{=} \lim_{\nu \to \infty} \lim_{\mu \to \infty} \lim_{k \to \infty} \left\langle \tilde{\psi}_{\alpha_{k}}(X), \, \ell_{A_{\nu}}^{(\mu)} \right\rangle,$$

which, for any sequence $(h_A^{(\nu)})_{\nu} \subset C_b(\mathcal{Y})$ as in (78), holds for any $(\ell_{A_{\nu}}^{(\mu)} \mid \mu, \nu \in \mathbb{N})$ such that:

$$(\ell_{A_{\nu}}^{(\mu)})_{\mu} \quad \text{is minimizing sequence of} \quad \inf_{\ell \in \mathbb{R}[\tilde{d}]} \|h_{A}^{(\nu)} - \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(\cdot) \rangle \|_{\infty;\mathfrak{D}_{Y}}, \quad \text{for each } \nu \in \mathbb{N};$$

under these assumptions, the convergence (α) holds \mathbb{P} -a.s. and (β) holds wrt. $L^2_X(\mathbb{R})$ -convergence, and (β) holds \mathbb{P} -a.s. if ($\tilde{\psi}_{\alpha_k}(X)$) is chosen such that (65) or (168) converges almost surely.

Since $\mathbb{P}(Y \in A \mid X) = 1 - \mathbb{P}(Y \in A^c \mid X)$, analogous statements hold if A is closed.

Proof. From (78) and the definition of conditional probability, we have

$$\mathbb{P}(Y \in A \mid X) = \mathbb{E}\left[\mathbb{1}_A(Y) \mid X\right] = \lim_{\nu \to \infty} \mathbb{E}\left[h_A^{(\nu)}(Y) \mid X\right] \quad \mathbb{P}\text{-a.s.}$$
(79)

via the conditional monotone convergence theorem. The remainder of the corollary then follows from (79) by applying Proposition 4.10 to the continuous functions $f = h_A^{(\nu)} \Big|_{\mathcal{D}_{\mathcal{N}}}$ for each $\nu \in \mathbb{N}$. \Box

5 Signature-Based Regression: Estimation and Consistency

This section extends the probabilistic results of Section 4.3 to a practically applicable, convex, and nonparametric statistical method to perform regression on and between stochastic processes. Specifically, we derive two signature-based regression estimators (Algorithms I and II, cf. Definition 2.6) whose universal consistency guarantees come with explicit convergence rates (Theorem 5.9).

These estimators rely on a data-based discretization of the optimisation problems (66) and (69) and of the associated approximations (72) and (73). The initial hurdle is that the regressors $\Psi_X \equiv \{\psi_\alpha \mid \alpha \in \mathcal{L}_X^2\}$, while defined only almost everywhere wrt. \mathbb{P}_X in (58), require pointwise evaluation when applied to specific data points. This is addressed in Section 5.1 by embedding Proposition 4.5 within the context of vector-valued reproducing kernel Hilbert spaces. In this framework, Section 5.2 adapts the tools and ideas from Section 4 to prepare for a regularised-least-squares-type empirical estimation analysis. This leads to several auxiliary results that structure the algorithmisation of Proposition 4.10 and culminate in Section 5.3, where these notions are combined to a first estimator for the regression operator $f \mapsto \mu_{-}(f)$ in (20). Section 5.4 specialises this approach to finite-dimensional regressands (Corollary 4.9) and proposes a second, specialised estimator for this scenario. Finally, Section 5.5 synthesizes all previous ideas and results to demonstrate the statistical consistency of the proposed estimators, convergence rates and error bounds (Theorem 5.9).

5.1 Signature Regressors form a Dense Vector-Valued RKHS

To underpin the approximation strategies of Section 4 with a clean and robust estimation theory, it is beneficial to transfer the insights of Proposition 4.5 to the framework of vector-valued reproducing kernel Hilbert spaces (vRKHS), cf. [42]. To this end, we modify the function space in (58) to obtain a dense subsystem of functions in $L_X^2(\mathcal{H}_{\mathcal{Y}})$ with vRKHS structure, thus ensuring that these modified signature-regressors can be evaluated pointwise. This adaptation is crucial for empirically approximating (via (68) and (71)) the probabilistic quantities in (13) from discrete-data points.

(The theory of vRKHS was first developed in [52], and we refer to [7] for a modern exposition.)

Proposition 5.1. Let Ξ be a feature normalisation on $\mathcal{H}_{\mathcal{X}}$, and denote by $(\mathcal{H}_{\kappa}, \|\cdot\|_{\kappa}) \subset \mathbb{R}^{\mathcal{X}}$ the RKHS with reproducing kernel $\kappa : (x, y) \mapsto \langle \mathfrak{sig}_{\Xi}(x), \mathfrak{sig}_{\Xi}(y) \rangle$. Then the space of functions

$$\mathfrak{H}_{\Xi} \coloneqq \left\{ f : \mathcal{X} \to \mathcal{H}_{\mathcal{Y}} \mid \left\{ \langle f, w \rangle \mid w \in [\tilde{d}]^* \right\} \subset \mathcal{H}_{\kappa} \quad and \quad \sum_{w \in [\tilde{d}]^*} \| \langle f, w \rangle \|_{\kappa}^2 < \infty \right\}$$
(80)

is a separable vRKHS with reproducing kernel $K: \mathcal{X}^{\times 2} \to \mathfrak{L}(\mathcal{H}_{\mathcal{Y}}), (x, y) \mapsto \kappa(x, y) \mathrm{Id}_{\mathcal{H}_{\mathcal{Y}}}$. Moreover,

$$\mathfrak{H}_{\Xi} \quad is \ dense \ in \quad \left(L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}}), \|\cdot\|_{L^2(\mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})}\right), \tag{81}$$

and $\mathfrak{H}_{\Xi}^{[m]} = \{\tilde{f} \equiv \pi_{[m]} \circ f \mid f \in \mathfrak{H}_{\Xi}\}\)$ is dense in $\left(L^2(\mathcal{X}, \mathbb{P}_X; \pi_{[m]}(\mathcal{H}_{\mathcal{Y}})), \|\cdot\|_{L^2(\mathbb{P}_{\mathcal{X}}; \pi_{[m]}(\mathcal{H}_{\mathcal{Y}}))}\right)\)$ for each $m \geq 0$, where $\mathfrak{H}_{\Xi}^{[0]} \cong \mathcal{H}_{\kappa}$. Recalling Proposition 4.5 and introducing the space

$$\Psi_X^{\star} \coloneqq \left\{ \sum_{w \in [\tilde{d}]^*} \langle u_w, \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle w \, \middle| \, (u_w) \subset \mathcal{H}_{\mathcal{X}} \, : \, \sum_{w \in [\tilde{d}]^*} \mathbb{E} \left[\langle u_w, X^{\Xi} \rangle^2 \right] < \infty \right\}$$

we further have the $\|\cdot\|_{L^2(\mathcal{X},\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}})}$ -dense—and, in the case of \mathfrak{H}_{Ξ} , also continuous—embeddings

$$\mathfrak{H}_{\Xi}, \Psi_X \subseteq \Psi_X^{\star} \subseteq L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}}).$$

$$(82)$$

Proof. See Appendix B.3.1.

5.2 Preliminary Results via Hilbert-Valued Regularised Least Squares

Proposition 5.1 together with Proposition 4.10 and Corollary 4.9, allows us to leverage the consistency theory of 'vectorial' regularised least squares routines, as outlined in [42], to derive estimators that are universally consistent in the sense of Definition 2.6. This section provides some auxiliary results to support this.

For this, we let $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ be a Borel probability measure with marginals $\mu_{\mathcal{X}} \coloneqq \mu \circ \pi_{\mathcal{X}}^{-1}$ and $\mu_{\mathcal{Y}} \coloneqq \mu \circ \pi_{\mathcal{Y}}^{-1}$, and appeal to an M-estimation scheme given by the following family of objectives:

$$\begin{split} \phi_{\mu}^{(m)} &: L^{2}\left(\mu_{\mathcal{X}}; \mathcal{H}_{\mathcal{Y}}^{[m]}\right) \longrightarrow \overline{\mathbb{R}}_{+}, \quad g \mapsto \phi_{\mu}^{(m)}(g) \coloneqq \int_{\mathcal{X} \times \mathcal{Y}} \left\|\pi_{[m]}(\underline{\mathfrak{sig}}_{\Lambda}(y)) - g(x)\right\|^{2} \mathrm{d}\mu(x, y), \quad (m \in \overline{\mathbb{N}}), \\ \phi_{\mu|\lambda}^{(m)} &: \mathfrak{H}_{\Xi}^{[m]} \longrightarrow \overline{\mathbb{R}}_{+}, \qquad g \mapsto \phi_{\mu|\lambda}^{(m)}(g) \coloneqq \phi_{\mu}^{(m)}(g) + \lambda \|g\|_{\mathfrak{H}_{\Xi}}^{2}, \qquad (\lambda > 0), \end{split}$$

where $\mathcal{H}_{\mathcal{Y}}^{[m]} \coloneqq \pi_{[m]}(\mathcal{H}_{\mathcal{Y}})$ and $\mathfrak{H}_{\Xi}^{[m]} \coloneqq \pi_{[m]}(\mathfrak{H}_{\Xi})$ for $\pi_{[m]}$ as in Section 3.10 and $\pi_{[\infty]} \coloneqq \mathrm{id}_{\mathcal{H}_{\mathcal{Y}}}$. For any given $f \in \mathcal{L}^2(\mu_{\mathcal{Y}})$, we additionally consider

$$\begin{split} \phi^{f}_{\mu_{\mathcal{Y}}} &: \mathcal{H}_{\mathcal{Y}} \longrightarrow \mathbb{R}_{+}, \quad \ell \mapsto \phi^{f}_{\mu_{\mathcal{Y}}}(\ell) \coloneqq \int_{\mathcal{Y}} \left| f(y) - \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle \right|^{2} \mu_{\mathcal{Y}}(\mathrm{d}y), \\ \phi^{f}_{\mu_{\mathcal{Y}}|\lambda} &: \mathcal{H}_{\mathcal{Y}} \longrightarrow \mathbb{R}_{+}, \quad \ell \mapsto \phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\ell) \coloneqq \phi^{f}_{\mu_{\mathcal{Y}}}(\ell) + \lambda \left\| \langle \ell, \underline{\mathfrak{sig}}_{\Lambda}(\cdot) \rangle \right\|_{\kappa}^{2}. \end{split}$$
(83)

Finally, for non-zero $(f, \Upsilon) \in L^2(\mu_{\mathcal{Y}}) \times L^2(\mu_{\mathcal{X}}; \mathcal{H}_{\mathcal{Y}})$ we define

$$\vartheta_f(\tilde{\epsilon}) \coloneqq \frac{\tilde{\epsilon}^2}{4} \Big(\inf \left\{ \|g\|_{\kappa}^2 \, \big| \, g \in \mathcal{H}_{\kappa} \colon \|f - g\|_{L^2(\mathbb{P}_Y)}^2 \le \tilde{\epsilon}^2/2 \right\} \Big)^{-1} \quad \text{and} \tag{84}$$

$$\tilde{\vartheta}_{\Upsilon}(\tilde{\epsilon}) \coloneqq \frac{\tilde{\epsilon}^2}{4} \left(\inf \left\{ \|g\|_{\mathfrak{H}_{\Xi}}^2 \, \big| \, g \in \mathfrak{H}_{\Xi}^{[m]} : \, \|\Upsilon - g\|_{L^2(\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}})}^2 \le \tilde{\epsilon}^2/2 \right\} \right)^{-1},\tag{85}$$

where we set $\vartheta_f(\tilde{\epsilon}) \coloneqq \infty$ [resp. $\tilde{\vartheta}_{\Upsilon}(\tilde{\epsilon}) \coloneqq \infty$] if the infimum in (84) [resp. in (85)] is zero.

(Note $\tilde{\epsilon} \mapsto \vartheta_f(\tilde{\epsilon})$ and $\tilde{\epsilon} \mapsto \tilde{\vartheta}^{(m)}_{\Upsilon}(\tilde{\epsilon})$ are well-defined maps from $(0, \infty)$ to $(0, \infty]$ by Proposition 5.1.)

The following technical results elucidate basic properties and interrelations of the above auxiliary functions, preparing them for use in deriving our announced regression estimators.

Lemma 5.2. For each $f \in \mathcal{L}^2(\mu_{\mathcal{Y}})$ and any $\lambda > 0$ and $m \in \overline{\mathbb{N}}$, both minimizers

$$\ell^{f}_{\mu_{\mathcal{Y}},\lambda} \in \operatorname*{arg\,min}_{\ell \in \mathcal{H}_{\mathcal{Y}}} \phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\ell) \qquad and \qquad \Upsilon^{*}_{\mu,\lambda,m} = \operatorname*{arg\,min}_{g \in \mathfrak{H}_{\Xi}^{[m]}} \phi^{(m)}_{\mu|\lambda}(g) \tag{86}$$

exist and the minimizer $\Upsilon^*_{\mu,\lambda,m}$ is unique. Moreover, if $\mu_{\mathcal{Y}}$ is an empirical measure of the form

$$\mu_{\mathcal{Y}} = \frac{1}{|I|} \sum_{i \in I} \delta_{y_i} \qquad \text{for some} \quad \{y_i \mid i \in I\} \subset \mathcal{Y} \quad \text{with } I \text{ finite}, \tag{87}$$

then $\left(\arg\min_{\ell\in\mathcal{H}_{\mathcal{V}}}\phi^{f}_{\mu_{\mathcal{V}}\mid\lambda}\right)\cap V_{\mu_{\mathcal{V}}}\neq\emptyset$ for $V_{\mu_{\mathcal{V}}}\coloneqq \operatorname{span}\{\underline{\operatorname{sig}}_{\Lambda}(y_{i})\mid i\in I\}$; if, in addition, the points in $\{y_{i}\mid i\in I\}$ are pairwise distinct, then $\ell^{f,\star}_{\mu_{\mathcal{V}},\lambda}\coloneqq \operatorname{argmin}_{\ell\in V_{\mu_{\mathcal{V}}}}\phi^{f}_{\mu_{\mathcal{V}}\mid\lambda}(\ell)$ is unique and satisfies:

$$\|\ell^{f}_{\mu_{\mathcal{Y}},\lambda}\|_{\mathcal{H}_{\mathcal{Y}}} \ge \|\ell^{f,\star}_{\mu_{\mathcal{Y}},\lambda}\|_{\mathcal{H}_{\mathcal{Y}}} \quad for \ each \ \ \ell^{f}_{\mu_{\mathcal{Y}},\lambda} \ as \ in \ (86), \tag{88}$$

 $\begin{array}{l} \text{and also } \left\| \ell_{\mu_{\mathcal{Y}},\lambda}^{f,\star} \right\|_{\mathcal{H}_{\mathcal{Y}}} = \left\| \arg\min_{h \in \mathcal{H}_{\tilde{\kappa}}} \tilde{\phi}_{f,\lambda}(h) \right\|_{\tilde{\kappa}} \text{ for the map } \tilde{\phi}_{f,\lambda} : \mathcal{H}_{\tilde{\kappa}} \ni h \mapsto \|f-h\|_{L^{2}(\mu_{\mathcal{Y}})}^{2} + \lambda \|h\|_{\tilde{\kappa}}^{2}, \\ \text{where the domain } \mathcal{H}_{\tilde{\kappa}} \text{ is an RKHS with kernel } \tilde{\kappa} : \mathcal{Y}^{\times 2} \ni (x,y) \mapsto \langle \underline{\mathfrak{sig}}_{\Lambda}(x), \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle. \end{array}$

Proof. Delegated to Appendix B.3.2.

In the following, suppose that, for any fixed X and Y as in (52), we are given observations

$$\mathfrak{Y} \equiv \left\{ Y^{(l)} \mid l \in [n] \right\} \quad \text{and} \quad \mathfrak{Z} \equiv \left\{ (X^{(j)}, Y^{(j)}) \mid j \in [N] \right\} \qquad (n, N \in \mathbb{N}), \tag{89}$$

modelled as iid copies of Y and (X, Y), respectively.

For the associated empirical (random) measures $\hat{\mu}_{\mathfrak{Y}} \coloneqq \frac{1}{n} \sum_{l=1}^{n} \delta_{Y^{(l)}}$ and $\hat{\mu}_{\mathfrak{Z}} \coloneqq \frac{1}{N} \sum_{j=1}^{N} \delta_{(X^{(j)}, Y^{(j)})}$ and for any given $f \in \mathcal{L}^2(\mathbb{P}_Y)$, we follow the notation in (86) and choose any minimizers

$$\hat{\ell}^{f}_{n,\lambda} \coloneqq \ell^{f}_{\hat{\mu}_{\mathfrak{Y}},\lambda} \quad \text{and} \quad \hat{\Upsilon}^{*}_{N,\lambda,m} \coloneqq \Upsilon^{*}_{\hat{\mu}_{\mathfrak{Z}},\lambda,m} \quad \text{and} \quad \hat{\Upsilon}^{*}_{N,\lambda} \coloneqq \Upsilon^{*}_{\hat{\mu}_{\mathfrak{Z}},\lambda,\infty}.$$
(90)

Recall that the objects in (90) are all $(\mathfrak{Y}$ - resp. \mathfrak{Z} -dependent) random variables.

We describe the empirical minimizers in (90) and distinguish three associated types of error.

Proposition 5.3. Let \mathfrak{Y} , \mathfrak{Z} be as in (89), and consider any $f \in \mathcal{L}^2(\mathbb{P}_Y)$. Then for all $\lambda > 0$,

$$\langle \hat{\Upsilon}^*_{N,\lambda,m}, w \rangle \in \operatorname{span}\left(\langle \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \, \big| \, j \in [N] \right), \quad for \ each \ w \in [\tilde{d}]^* : \, |w| \le m \quad (m \in \overline{\mathbb{N}}).$$
(91)

Moreover, with $m_{\tilde{d}} \coloneqq \sum_{\nu=0}^{m} \tilde{d}^{\nu}$ and $q_{j}^{*} \coloneqq \langle \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle$ we have that, with probability one,

$$\sum_{l=1}^{n} \alpha_{l} \underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)}) \in \operatorname*{argmin}_{\ell \in \mathcal{H}_{\mathcal{Y}}} \phi^{f}_{\hat{\mu}_{\mathfrak{Y}}|\lambda}(\ell) \quad and \quad \hat{\Upsilon}^{\mathrm{lex}}_{N,\lambda,m} = \tilde{A}^{\top}_{*} \cdot (q_{j}^{*})_{j=1}^{N}$$
(92)

if and only if the coefficients $\alpha^* \equiv (\alpha_l)_{l=1}^n \in \mathbb{R}^n$ and $\tilde{A}_* \equiv \left(\tilde{a}_{(1)} \middle| \cdots \middle| \tilde{a}_{(m_d)} \right) \in \mathbb{R}^{N \times m_d}$ are such that

$$\left(\boldsymbol{A}_{\mathfrak{Y}}^{\top}\boldsymbol{A}_{\mathfrak{Y}}+n\lambda\boldsymbol{A}_{\mathfrak{Y}}\right)\alpha^{*}=\boldsymbol{A}_{\mathfrak{Y}}b_{\mathfrak{Y},f}\qquad and \tag{93}$$

$$(\boldsymbol{A}_{\boldsymbol{3}}^{\top}\boldsymbol{A}_{\boldsymbol{3}} + N\lambda\boldsymbol{A}_{\boldsymbol{3}})\tilde{A}_{*} = \boldsymbol{A}_{\boldsymbol{3}}\boldsymbol{B}_{\boldsymbol{3}},\tag{94}$$

where the above systems are formulated in terms of the data-based (random) matrices and vectors

$$\boldsymbol{A}_{\mathfrak{Y}} \coloneqq \left(\left\langle \underline{\mathfrak{sig}}_{\Lambda}(Y^{(k)}), \underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)}) \right\rangle \right)_{k,l=1}^{n} \quad and \quad b_{\mathfrak{Y},f} \coloneqq \left(f(Y^{(l)}) \right)_{l=1}^{n}, \tag{95}$$

$$\boldsymbol{A}_{\mathfrak{Z}} \coloneqq \left(\left\langle \underline{\mathfrak{sig}}_{\Xi}(X^{(i)}), \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}) \right\rangle \right)_{i,j=1}^{N} \quad and \quad \boldsymbol{B}_{\mathfrak{Z}} \coloneqq \left(\left\langle \underline{\mathfrak{sig}}_{\Lambda}(Y^{(i)}), \eta^{-1}(j) \right\rangle \right)_{(i,j)\in[N]\times[m_{\tilde{d}}]}; \quad (96)$$

here, $\eta: [\tilde{d}]^* \to \mathbb{N}_0$ is the length-lexicographic⁷ (or shortlex) ordering of the words in $[\tilde{d}]^*$ and

$$\hat{\Upsilon}_{N,\lambda,m}^{\text{lex}} \equiv \left(\langle \hat{\Upsilon}_{N,\lambda,m}, \eta^{-1}(1) \rangle, \dots, \langle \hat{\Upsilon}_{N,\lambda,m}, \eta^{-1}(m_{\tilde{d}}) \rangle \right) : \mathcal{X} \to \mathbb{R}^{m_{\tilde{d}}}$$
(97)

is the length-lexicographically indexed version of $\hat{\Upsilon}^*_{N,\lambda,m}$.

Proof. See Appendix B.3.3.

Remark 5.4. The matrix A_3 in (95) [resp. the matrix $A_{\mathfrak{Y}}$ in (96)] is invertible iff the observed samples $(X^{(j)})_j$ in \mathfrak{Z} [resp. the samples $(Y^{(l)})_l$ in \mathfrak{Y}] are pairwise disjoint. For a proof, see the argumentation around (151) in Appendix B.3.2.

⁷ Recall that in the length-lexicographic ordering of $[\tilde{d}]^*$, words are primarily sorted by wordlength with the shortest words first, and words of the same length are sorted into lexicographical order; see e.g. [56, p. 14].

Revisiting the notation from (86), let us now complement (90) by any fixed choice of minimizers

$$\ell^f_{\lambda} \coloneqq \ell^f_{\mathbb{P}_{Y,\lambda}}$$
 and $\Upsilon^*_{\lambda,m} \coloneqq \Upsilon^*_{\mathbb{P}_{(X,Y)},\lambda,m}$ and $\Upsilon^*_{\lambda} \coloneqq \Upsilon^*_{\mathbb{P}_{(X,Y)},\lambda,\infty}.$

To extend the uniqueness and optimality characteristic (88) to singular kernel matrices (95), let

$$\hat{\ell}_{n,\lambda,f}^{\star} \coloneqq \sum_{l=1}^{n} \alpha_{l}^{\star} \underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)}) \quad \text{for} \quad \alpha^{\star} \equiv \left(\alpha_{l}^{\star}\right)_{l=1}^{n} \quad \text{the minimum norm solution of} \quad (93); \tag{98}$$

equivalently, $\alpha^* = (\mathbf{A}_{\mathfrak{Y}}^{\top} \mathbf{A}_{\mathfrak{Y}} + n\lambda \mathbf{A}_{\mathfrak{Y}})^+ \mathbf{A}_{\mathfrak{Y}} b_{\mathfrak{Y},f}$ for $(\cdot)^+$ the Moore-Penrose pseudoinverse [43].

Lemma 5.5. In the above notation, consider for any given $f \in \mathcal{L}^2(\mathbb{P}_Y)$ the expressions

$$r_{f}^{\mathrm{I}}(\lambda) \coloneqq \left\| f - \langle \ell_{\lambda}^{f}, \underline{\mathfrak{sig}}_{\Lambda} \rangle \right\|_{L^{2}(\mathbb{P}_{Y})} \quad and \quad r_{f,\lambda}^{\mathrm{II}}(n) \coloneqq \left\| \langle \ell_{\lambda}^{f} - \hat{\ell}_{n,\lambda}^{f}, \underline{\mathfrak{sig}}_{\Lambda} \rangle \right\|_{L^{2}(\mathbb{P}_{Y})}, \tag{99}$$

referred to as the f-approximation error and the f-estimation error, respectively. Set further

$$r_{f,\lambda,n}^{\mathrm{III}}(m) \coloneqq \left\| \hat{\ell}_{n,\lambda,f}^{\star} - \pi_{[m]}(\hat{\ell}_{n,\lambda,f}^{\star}) \right\|_{\mathcal{H}_{\mathcal{Y}}},\tag{100}$$

which, for a given $m \in \overline{\mathbb{N}}$, we refer to as the f-cutoff error. Recalling that $\mathbb{E}[\mathbb{Y}^{\Lambda} | X] = \varphi_{\star}(X)$ for some unique $\varphi_{\star} \in L^{2}(\mathbb{P}_{X}; \mathcal{H}_{\mathcal{Y}})$ by the Doob-Dynkin lemma, let us finally consider the expressions

$$R_m^{\mathrm{I}}(\lambda) \coloneqq \left\| \pi_{[m]} \left(\varphi_\star - \Upsilon_\lambda^* \right) \right\|_{L^2(\mathbb{P}_X; \mathcal{H}_\mathcal{Y})} \quad and \quad R_{m,\lambda}^{\mathrm{II}}(N) \coloneqq \left\| \pi_{[m]} \left(\Upsilon_\lambda^* - \hat{\Upsilon}_{N,\lambda}^* \right) \right\|_{L^2(\mathbb{P}_X; \mathcal{H}_\mathcal{Y})}, \quad (101)$$

referred to as the inner approximation error and inner estimation error, respectively. Then we have:

(i)
$$\sup \left\{ r_f^{\mathrm{I}}(\lambda) \, \middle| \, 0 < \lambda \leq^8 \vartheta_f(\epsilon) \right\} \leq \epsilon, \quad \text{for all } \epsilon > 0;$$

(ii) it holds that $\pi_{[m]}(\varphi_{\star}) = \arg\min_{g \in L^2(\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}}^{[m]})} \phi_{\mathbb{P}_{(X,Y)}}^{(m)}(g)$ for each $m \in \overline{\mathbb{N}}$, and also:

$$R_m^{\mathrm{I}}(\lambda)^2 = \boldsymbol{\phi}_{\mathbb{P}_{(X,Y)}}^{(m)} \big(\pi_{[m]}(\Upsilon_{\lambda}^*) \big) - \boldsymbol{\phi}_{\mathbb{P}_{(X,Y)}}^{(m)} \big(\pi_{[m]}(\varphi_{\star}) \big), \quad \text{for all } \lambda > 0 \,;$$

(iii) for each
$$m \in \overline{\mathbb{N}}$$
: $\sup \left\{ R_m^{\mathrm{I}}(\lambda) \, \big| \, 0 < \lambda \leq \tilde{\vartheta}_{\pi_{[m]}(\varphi_{\star})}(\epsilon) \right\} \leq \epsilon$, for all $\epsilon > 0$;

(iv) there is an explicit monotone decreasing null sequence $(\beta_m^{(\text{cut})}) \subset \mathbb{R}_+$, depending on $(\mathfrak{Y}, \Lambda, \lambda, f)$, such that with probability one:

$$r_{f,\lambda,n}^{\mathrm{III}}(m) \leq \beta_m^{(\mathrm{cut})} \quad \text{for all } m \geq 0;$$

(v) for any $\delta > 0$, there are explicit strictly decreasing null sequences $(\beta_n^{(\text{out})}), (\beta_N^{(\text{in})}) \subset \mathbb{R}_+,$ depending on $(\Lambda, \lambda, \delta)$ and $(\Lambda, \Xi, \lambda, \delta)$ respectively, such that for each $n, N \in \mathbb{N}$,

$$\mathbb{P}\left(r_{f,\lambda}^{\mathrm{II}}(n) \ge \beta_n^{\mathrm{(out)}}\right) \le \delta \qquad and \qquad \mathbb{P}\left(R_{m,\lambda}^{\mathrm{II}}(N) \ge \beta_N^{\mathrm{(in)}}\right) \le \delta.$$
(102)

Proof. Reported in Appendix B.3.4.

Remark 5.6. Specific rate sequences $(\beta_n^{(\text{out})})$ and $(\beta_N^{(\text{in})})$ for which (102) holds are, for example,

$$\beta_n^{(\text{out})} \coloneqq \frac{c_\Lambda}{\lambda\sqrt{n\delta}} \quad \text{and} \quad \beta_N^{(\text{in})} \coloneqq \sqrt{\frac{c_\Xi c_\Lambda}{\lambda^2 N\delta}} \quad (n, N \in \mathbb{N})$$
(103)

with constants $c_Z \coloneqq \sup_{\boldsymbol{t} \in \mathcal{H}_{\boldsymbol{\xi}}} \|Z(\boldsymbol{t})\|_{\mathcal{H}_{\boldsymbol{\xi}}}^2$ for $(Z, \boldsymbol{\xi}) \in \{(\Lambda, \mathcal{Y}), (\Xi, \mathcal{X})\}$; see proof of Lemma 5.5 (v). ⁸ If $\vartheta_f(\boldsymbol{\epsilon}) = \infty$, then the constraint ' $\lambda \leq \vartheta_f(\boldsymbol{\epsilon})$ ' is to be dropped; likewise for (iii).

5.3 A Convex Estimator for Stochastic Process Regression

Let now $f = (f_1, \dots, f_k) \in \mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k)$, for any $k \in \mathbb{N}$. In order to structure the results of Sections 5.1 and 5.2 into an estimator for $\mathbb{E}[f(Y) \mid X]$, we introduce, for each $\varepsilon, \delta > 0$ and $\lambda > 0$, the moduli

$$\lambda_f^{\mathrm{I}}(\varepsilon) \coloneqq \min\{\vartheta_{f_i}(\varepsilon/(5k)) \mid i = 1, \dots, k\} \quad \text{and} \\ n_{f,\lambda}(\varepsilon, \delta) \coloneqq \min\{n \in \mathbb{N} \mid \beta_n^{(\mathrm{out})}(\delta/(2k); \lambda) \le \varepsilon/(5k)\}$$
(104)

for any fixed rate sequence $(\beta_n^{(\text{out})}) \equiv (\beta_n^{(\text{out})}(\delta;\lambda))$ as in Lemma 5.5 (v), e.g. the one in (103), and

$$m_{f,\lambda}(\varepsilon;\mathfrak{Y}) \coloneqq \min\left\{m \in \mathbb{N} \mid \|\boldsymbol{\alpha}_{\lambda}^{f}\|_{1,2}\tilde{\beta}_{m} \le \varepsilon/(5k)\right\}$$
(105)

for the norm $||A||_{1,2} \coloneqq \max_{1 \le j \le n} ||A_{:j}||_2$ and the \mathfrak{Y} -dependent sequence $(\tilde{\beta}_m)$ from (158), and with

$$\boldsymbol{\alpha}_{\lambda}^{f} \equiv \left(\alpha_{li}^{(\lambda)}\right)_{(l,i)\in[n]\times[k]} \coloneqq \left(\boldsymbol{A}_{\mathfrak{Y}}^{\top}\boldsymbol{A}_{\mathfrak{Y}} + n\lambda\boldsymbol{A}_{\mathfrak{Y}}\right)^{+}\boldsymbol{A}_{\mathfrak{Y}}B_{\mathfrak{Y},f}$$

for the data matrices $A_{\mathfrak{Y}}$ from (95) and $B_{\mathfrak{Y},f} \coloneqq (f_j(Y^{(l)}))_{(l,j)\in[n]\times[k]}$. For $m \in \mathbb{N}$, we further set

$$\lambda_{\lambda,m}^{\mathrm{I}}(\varepsilon;\mathfrak{Y}) \coloneqq \tilde{\vartheta}_{\varphi_{m}^{\star}}(\varepsilon/(5c_{\lambda,m;\mathfrak{Y}})) \quad \text{and} \\ N_{\lambda,\tilde{\lambda},m}(\varepsilon,\delta;\mathfrak{Y}) \coloneqq \min\left\{ N \in \mathbb{N} \, \big| \, \beta_{N}^{(\mathrm{in})}(\delta/2;\tilde{\lambda}) \le \varepsilon/(5c_{\lambda,m;\mathfrak{Y}}) \right\}$$
(106)

for $\varphi_m^{\star} \coloneqq \pi_{[m]}(\varphi_{\star})$ as in Lemma 5.5, with the constant $c_{\lambda,m;\mathfrak{Y}} \coloneqq \sqrt{\sum_{i=1}^k \|\pi_{[m]}(\hat{\ell}_{n,\lambda,f_i}^{\star})\|_{\mathcal{H}_{\mathcal{Y}}}^2}$ computable from (98) and for any fixed sequence $(\beta_N^{(in)}) \equiv (\beta_N^{(in)}(\tilde{\delta};\tilde{\lambda}))$ as in Lem. 5.5 (v), e.g. (103). Finally, for any dimensions $n, N \in \mathbb{N}$ of the data sets (89) and parameters $\underline{\lambda} \equiv (\lambda_1, \lambda_2) \in \mathbb{R}^2_{>0}$ and $m \in \mathbb{N}$, consider the estimator (of the regression function between f(Y) and X)

$$\hat{\mathscr{R}}^{\mathrm{I}}_{\underline{\lambda},f}[n,N,m] \coloneqq \sum_{i=1}^{k} \left[\sum_{l=1}^{n} \alpha_{li}^{\lambda_{1}} \langle \pi_{[m]}^{\mathrm{lex}}(\underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)})), \hat{\Upsilon}^{\mathrm{lex}}_{N,\lambda_{2},m}(\cdot) \rangle_{2} \right] e_{i} \, : \, \mathcal{X} \longrightarrow \mathbb{R}^{k} \tag{107}$$

with $(e_i \mid i \in [k])$ the standard basis of \mathbb{R}^k , $\langle \cdot, \cdot \rangle_2$ the Euclidean inner product on $\mathbb{R}^{m_{\tilde{d}}}$, and where

$$\hat{\Upsilon}_{N,\lambda_2,m}^{\text{lex}} = \sum_{\nu=1}^{m_{\tilde{d}}} \left[\sum_{j=1}^{N} \beta_{j\nu}^{(\lambda_2)} \langle \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \right] \tilde{e}_{\nu} \,, \qquad \text{with} \quad m_{\tilde{d}} \coloneqq \frac{\tilde{d}^{m+1}-1}{\tilde{d}-1},$$

is as defined in (92), that is with the coefficient matrix $(\beta_{j\nu}^{(\lambda_2)}) =: \beta_{\lambda_2}$ solving

$$(\boldsymbol{A}_{\boldsymbol{3}}^{\top}\boldsymbol{A}_{\boldsymbol{3}} + N\lambda_{2}\boldsymbol{A}_{\boldsymbol{3}})\boldsymbol{\beta}_{\lambda_{2}} = \boldsymbol{A}_{\boldsymbol{3}}\boldsymbol{C}_{\boldsymbol{3}} \quad \text{for} \quad \boldsymbol{C}_{\boldsymbol{3}} \coloneqq \left(\left\langle\underline{\mathfrak{sig}}_{\Lambda}(\boldsymbol{Y}^{(i)}), \eta^{-1}(j)\right\rangle\right)_{(i,j)\in[N]\times[m_{\tilde{d}}]};$$

here, the map $\pi_{[m]}^{\text{lex}} \coloneqq \eta \circ \pi_{[m]} : \mathcal{H}_{\mathcal{Y}} \to \mathbb{R}^{m_{\tilde{d}}}$ is the shortlex-ordered projection onto $\mathcal{H}_{\mathcal{Y}}^{[m]}$.

The function defined in equation (107) serves as a data-based approximation of the (regression) function $\psi_f \in L^2(\mathbb{P}_X; \mathbb{R}^k)$ such that $\psi_f(X) = \mathbb{E}[f(Y)|X]$. Specifically, we anticipate that:

$$\mathbb{E}[f(Y) \mid X] \approx \hat{\mathscr{R}}^{\mathrm{I}}_{\lambda, f}[n, N, m](X), \tag{108}$$

where this approximation is expected to improve with sufficiently large values of n, N and with sufficiently small values of λ_1, λ_2 . The precise quality of the approximation (108), and the requisite parameter choices for achieving this approximation within given error bounds, are established in Theorem 5.9 below, based on the threshold quantities defined in equations (104), (105), and (106). In summary, the above combines to the following estimator for the regression operator (118).

Algorithm I: Computing the Regression Operator
$$f \mapsto \mu_{\cdot}(f)$$
 via Signatures

- **Goal:** For processes X in $\mathbb{R}^{d_{\mathcal{X}}}$ and Y in $\mathbb{R}^{d_{\mathcal{Y}}}$, compute $\mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k) \ni f \mapsto \psi_f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}^k)$, for regression functions $\psi_f \coloneqq \mu_{\cdot}(f)$ (with $\mu \coloneqq \mathbb{P}_{(X,Y)}$) such that $\psi_f(X) = \mathbb{E}[f(Y) \mid X]$.
- **1. Input:** Samples $\mathfrak{z} \coloneqq \{(x_j, y_j) \mid j \in [N]\}$ and $\mathfrak{y}_f \coloneqq \{(y_l, f(y_l)) \mid l \in [n]\}$ of the associations (X, Y) and (Y, f(Y)), for $f \equiv (f_1, \cdots, f_k) \in \mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k)$, respectively.

Hyperparameters: tensor normalisations (Ξ, Λ) , resolution parameters $\lambda_f, \lambda_{\Upsilon}, m > 0$.

2. Compute the weights $\beta_{\lambda_{\Upsilon}} \equiv (\beta_{j\nu}^{(\lambda_{\Upsilon})}) \in \mathbb{R}^{N \times m_{d_{\mathcal{Y}}}}$ and $\alpha_{\lambda_{f}}^{f} \equiv (\alpha_{li}^{(\lambda_{f})}) \in \mathbb{R}^{n \times k}$ solving

with $m_{d_{\mathcal{Y}}} \coloneqq (\underline{d}_{\mathcal{Y}}^{m+1} - 1)/(\underline{d}_{\mathcal{Y}} - 1)$ and for the data-dependent coefficient matrices

$$\begin{aligned} \boldsymbol{A}_{\mathfrak{z}} &\coloneqq \left(\left\langle \underline{\mathfrak{sig}}_{\Xi}(x_{i}), \underline{\mathfrak{sig}}_{\Xi}(x_{j}) \right\rangle \right)_{i,j=1}^{N} \quad \text{and} \quad \boldsymbol{C}_{\mathfrak{z}} &= \left(\left\langle \underline{\mathfrak{sig}}_{\Lambda}(y_{i}), \operatorname{lex}^{-1}(j) \right\rangle \right)_{(i,j)\in[N]\times[m_{d_{\mathcal{Y}}}]}, \\ \boldsymbol{A}_{\mathfrak{y}} &\coloneqq \left(\left\langle \underline{\mathfrak{sig}}_{\Lambda}(y_{i}), \underline{\mathfrak{sig}}_{\Lambda}(y_{j}) \right\rangle \right)_{i,j=1}^{n} \quad \text{and} \quad \boldsymbol{B}_{\mathfrak{y},f} = \left(f_{j}(y_{i}) \right)_{(i,j)\in[n]\times[k]}. \end{aligned}$$

Here, lex: $[\underline{d}_{\mathcal{Y}}]^* \to \mathbb{N}_0$ is the shortlex ordering of the words in $[\underline{d}_{\mathcal{Y}}]^*$ (cf. footnote 7).

3. Compute the function, using the weights $\beta_{\lambda\gamma}$ and the standard basis (\tilde{e}_{ν}) of $\mathbb{R}^{m_{d_{\mathcal{Y}}}}$,

$$\hat{\Upsilon}_{N,\lambda_{\Upsilon},m}^{\beta_{\lambda_{\Upsilon}}} \coloneqq \sum_{\nu=1}^{m_{d_{\Upsilon}}} \left[\sum_{j=1}^{N} \beta_{j\nu}^{(\lambda_{\Upsilon})} \langle \underline{\mathfrak{sig}}_{\Xi}(x_j), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \right] \tilde{e}_{\nu}.$$
(110)

4. Compute the function, using the weights $\alpha_{\lambda_f}^f$ and the quantity $\hat{\Upsilon}_{N,\lambda_{\Upsilon},m}^{\beta_{\lambda_{\Upsilon}}}$,

$$\hat{\psi}_{f} := \sum_{i=1}^{k} \left[\sum_{l=1}^{n} \alpha_{li}^{(\lambda_{f})} \langle \pi_{[m]}^{\text{lex}}(\underline{\mathfrak{sig}}_{\Lambda}(y_{l})), \hat{\Upsilon}_{N,\lambda_{\Upsilon},m}^{\beta_{\lambda_{\Upsilon}}}(\cdot) \rangle_{2} \right] e_{i}$$
(111)

for the standard basis (e_i) of \mathbb{R}^k and with $\pi_{[m]}^{\text{lex}} : \mathcal{H}_{d_{\mathcal{Y}}} \to \mathbb{R}^{m_{d_{\mathcal{Y}}}}, (q_w) \mapsto \sum_{\nu=1}^{m_{d_{\mathcal{Y}}}} q_{\text{lex}^{-1}(\nu)} \tilde{e}_{\nu}$.

5. Output: $\hat{\psi}_f$ (estimator of the regression function ψ_f).

The fact that this estimator is a particular instance of Definition 2.6 is emphasised in Remark B.7. **Remark 5.7.** The estimator (111), as well as the estimator (117) below, relies on the (bounded) signature kernels $\kappa_{\Theta} : \mathbb{Z} \times \mathbb{Z} \ni (z_1, z_2) \mapsto \langle \underline{sig}_{\Theta}(z_1), \underline{sig}_{\Theta}(z_2) \rangle \in \mathbb{R}$, where $(\Theta, \mathbb{Z}) \in \{(\Xi, \mathcal{X}), (\Lambda, \mathcal{Y})\}$, for both its parametrisation (109) and pointwise evaluation (110). These kernels, like the (bounded) signatures themselves [12, 26], can be efficiently computed or approximated with several specialised algorithms and estimators. For an overview of related works, see e.g. [8, 29, 47] and the references therein.

5.4 Preliminary Results and Estimator for Finite-Dimensional Covariate

The scenario of Corollary 4.9, where we condition on X some finite-dimensional covariate $Z \in L^2(\mathbb{P}; \mathbb{R}^k)$ instead of the Hilbert-valued covariate $\underline{sig}_{\Lambda}(Y)$, permits structural results similar to those in Lemma 5.5 and Propositions 5.1 and 5.3. These statements, given below, can be proven entirely in parallel to the aforementioned results and thus may be regarded as corollaries to the above proofs.

In a slight modification from the setting in (89), this Z-covaried scenario operates on the given data

$$\mathfrak{W} \equiv \left\{ \left(X^{(j)}, \left(Z^1_{(j)}, \cdots, Z^k_{(j)} \right) \right) \mid j \in [M] \right\} \qquad (M \in \mathbb{N}),$$
(112)

i.e., on some fixed ensemble of M-many iid copies of (X, Z). We can then observe the following.

Corollary 5.8. Let X be as in (52), and consider the subspace $\mathfrak{h}_k \coloneqq \{h \equiv (h^1, \cdots, h^k) : \mathcal{X} \to \mathbb{R}^k \mid h^i \in \mathcal{H}_{\kappa}, \forall i \in [k]\} \subset L^2(\mathbb{P}_X; \mathbb{R}^k) \text{ together with the objectives}$

$$\begin{split} \phi_{Z} &: L^{2}(\mathbb{P}_{X}; \mathbb{R}^{k}) \longrightarrow \mathbb{R}_{+}, \quad g \mapsto \int_{\mathcal{X} \times \mathbb{R}^{k}} \left| z - g(x) \right|^{2} \mathbb{P}_{(X,Z)}(\mathrm{d}x, \mathrm{d}z), \\ \phi_{Z|\lambda} &: \mathfrak{h}_{k} \longrightarrow \mathbb{R}_{+}, \qquad g \mapsto \phi_{Z}(g) + \lambda \|g\|_{\mathfrak{h}_{k}}^{2}, \\ \hat{\phi}_{\mathfrak{W},\lambda} &: \mathfrak{h}_{k} \longrightarrow \mathbb{R}_{+}, \qquad g \mapsto \int_{\mathcal{X} \times \mathbb{R}^{k}} \left| z - g(x) \right|^{2} \hat{\mu}_{\mathfrak{W}}(\mathrm{d}x, \mathrm{d}z) + \lambda \|g\|_{\mathfrak{h}_{k}}^{2}, \end{split}$$

for $\lambda > 0$ and where $\hat{\mu}_{\mathfrak{W}} := \frac{1}{M} \sum_{j \in [M]} \delta_{(X^{(j)}, Z^{(j)})}$. Then both the minimizer $\Upsilon^*_{Z, \lambda}$ of $\phi_{Z|\lambda}$ and

$$\hat{\Upsilon}^*_{Z,\lambda;M}\coloneqq \operatorname*{arg\,min}_{h\in\mathfrak{h}_k}\hat{\phi}_{\mathfrak{W},\lambda}(h)$$
 are unique

and the latter minimizer admits the representation (with $(e_i \mid i \in [k])$ the standard basis of \mathbb{R}^k)

$$\hat{\Upsilon}^*_{Z,\lambda;M} = \sum_{i=1}^k \left[\sum_{j=1}^M \hat{\alpha}_{ji}^{(\lambda)} \langle \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \right] e_i$$
(113)

for any coefficient matrix $\hat{A}_{\star} \equiv \left(\hat{\alpha}_{ji}^{(\lambda)}\right) \in \mathbb{R}^{M \times k}$ which, abbreviating $\tilde{q}_j \coloneqq \underline{\mathfrak{sig}}_{\Xi}(X^{(j)})$, solves

$$(\boldsymbol{A}_{\mathfrak{W}}^{\top}\boldsymbol{A}_{\mathfrak{W}} + M\lambda\boldsymbol{A}_{\mathfrak{W}})\hat{A}_{*} = \boldsymbol{A}_{\mathfrak{W}}\boldsymbol{B}_{\mathfrak{W}} \quad for \quad \boldsymbol{A}_{\mathfrak{W}} \coloneqq (\langle \tilde{q}_{i}, \tilde{q}_{j} \rangle)_{i,j=1}^{M} \text{ and } \boldsymbol{B}_{\mathfrak{W}} \coloneqq (\boldsymbol{Z}_{(i)}^{j})_{(i,j)\in[M]\times[k]}.$$
(114)

Then the following holds:

- (i) the space \mathfrak{h}_k is a separable vRKHS with reproducing kernel $\mathcal{K} : \mathcal{X}^{\times 2} \to \mathfrak{L}(\mathbb{R}^k), (x, y) \mapsto \kappa(x, y) \mathrm{Id}_{k \times k}$, and also a $\|\cdot\|_{L^2(\mathbb{P}_X;\mathbb{R}^k)}$ -dense subspace of $L^2(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbb{P}_X; \mathbb{R}^k)$;
- (ii) for (the unique) $\psi_{X,Z} \in L^2(\mathbb{P}_X; \mathbb{R}^k)$ such that $\mathbb{E}[Z|X] = \psi_{X,Z}(X)$ (Doob-Dynkin) and the modulus⁹ $\hat{\vartheta}_{\psi_{X,Z}}(\tilde{\varepsilon}) \coloneqq \frac{\tilde{\epsilon}^2}{4} (\inf \{ \|h\|_{\mathfrak{h}_k}^2 | h \in \mathfrak{h}_k : \|\psi_{X,Z} h\|_{L^2(\mathbb{P}_X; \mathbb{R}^k)}^2 \leq \tilde{\epsilon}^2/2 \})^{-1}$, we have that

$$\sup\left\{\rho_{Z}^{\mathrm{I}}(\lambda) \coloneqq \left\|\psi_{X,Z} - \Upsilon_{Z,\lambda}^{*}\right\|_{L^{2}(\mathbb{P}_{X};\mathbb{R}^{k})} \left| 0 < \lambda \leq^{10} \hat{\vartheta}_{\psi_{X,Z}}(\epsilon)\right\} \leq \epsilon, \quad \text{for all } \epsilon > 0;$$

(iii) for any $\delta > 0$, there is an explicit monotone null sequence $(\hat{\beta}_M) \equiv (\hat{\beta}_M(\delta, \lambda))_{M \in \mathbb{N}} \subset \mathbb{R}_+$, depending on (Ξ, λ, δ) , such that for each $M \in \mathbb{N}$ we have

 $\mathbb{P}\big(\rho_{Z,\lambda}^{\mathrm{II}}(M) \geq \tilde{\beta}_M\big) \leq \delta \qquad for \quad \rho_{Z,\lambda}^{\mathrm{II}}(M) \coloneqq \big\|\Upsilon_{Z,\lambda}^* - \hat{\Upsilon}_{Z,\lambda;M}^*\big\|_{L^2(\mathbb{P}_X;\mathbb{R}^k)}.$

 $[\]frac{(PZ,\lambda)^{(M)} = PM = 0}{9 \text{ We set } \hat{\vartheta}_{\psi_{X,Z}}(\tilde{\epsilon}) \coloneqq \infty \text{ if the defining infimum is zero.} \quad 10 \text{ If } \hat{\vartheta}_{\psi_{X,Z}}(\epsilon) = \infty, \text{ then the constraint } \lambda \leq \hat{\vartheta}_{\psi_{X,Z}}(\epsilon), \text{ is to be dropped.}$

Proof. The assertions concerning \mathfrak{h}_k are derived entirely analogously to the corresponding properties of (80), as demonstrated in the proof of Proposition 5.1. The uniqueness of $(\hat{\Upsilon}_{Z,\lambda;M})$ is established analogously to the uniqueness of $\Upsilon^*_{\mu,\lambda,m}$ in (86), while its data-based representation in (113) is derived similarly to the $\hat{\Upsilon}^{\text{lex}}_{N,\lambda,m}$ -representation in (92). Statement (ii) is derived in the same manner as Lemma 5.5 (iii), while the corollary's statement (iii) holds, for example, for the sequence $(\tilde{\beta}_M) :=$ $(\frac{\sqrt{c \in \mathbb{H}|Z||_{L^2(\mathbb{P};\mathbb{R}^k)}}{\lambda\sqrt{M\delta}})$, as established by the same reasoning as in the proof of Lemma 5.5 (v).

Similar to how the results of Section 5.2 informed the construction of the estimator (107), setting

$$\hat{\mathscr{R}}_{\lambda,Z}^{\mathrm{II}}[M] \coloneqq \sum_{i=1}^{k} \left[\sum_{j=1}^{M} \hat{\alpha}_{ij}^{(\lambda)} \langle \underline{\mathfrak{sig}}_{\Xi}(X^{(j)}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \right] e_{i} \, : \, \mathcal{X} \to \mathbb{R}^{k} \tag{115}$$

yields an estimator of the regression function between Z and X. Through the data (112), this estimator relies on a parameter $\lambda > 0$ and coefficients $\hat{A}_{\star} = (\hat{\alpha}_{ii}^{(\lambda)}) \in \mathbb{R}^{M \times k}$ solving (114).

When applied to the vector $Z \coloneqq f(Y)$, the function (115) is designed to yield an approximation

$$\mathbb{E}[f(Y) \mid X] \approx \hat{\mathscr{R}}_{\lambda,Z}^{\mathrm{II}}[M](X), \tag{116}$$

similar to (108), where this approximation is expected to improve with sufficiently large M and with sufficiently small λ . The precise quality of the approximation (116), and the requisite parameter choices for achieving this approximation within given error bounds, are established in Theorem 5.9 below, based on the threshold quantities defined in Corollary 5.8.

The results of Section 5 combine to the following estimators for the regression operator (118).

Algorithm II: Computing $\psi_Z \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}^k)$ such that $\psi_Z(X) = \mathbb{E}[Z \mid X] \ (Z \in L^2(\mathbb{P}; \mathbb{R}^k)).$	

Goal: For a stochastic process X in \mathbb{R}^d and a random vector Z in \mathbb{R}^k , estimate a regression function $\psi_Z \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}^k)$ such that $\psi_Z(X) = \mathbb{E}[Z \mid X]$.

1. Input: Samples $\mathfrak{w} \coloneqq \left\{ \left(x_i, (z_i^1, \dots, z_i^k) \right) \mid i \in [M] \right\}$ of the association (X, Z).

Hyperparameters: tensor normalisation $\Xi: \mathcal{H}_d^{\downarrow} \to \mathcal{H}_d$, regularity parameter $\lambda > 0$.

2. Compute the weights $\hat{\boldsymbol{\alpha}}_{\lambda} \equiv (\hat{\alpha}_{ij}^{(\lambda)}) \in \mathbb{R}^{k \times M}$ by solving

$$(\boldsymbol{A}_{\boldsymbol{w}}^{\top}\boldsymbol{A}_{\boldsymbol{w}} + M\lambda\boldsymbol{A}_{\boldsymbol{w}})\hat{\boldsymbol{\alpha}}_{\lambda} = \boldsymbol{A}_{\boldsymbol{w}}\boldsymbol{B}_{\boldsymbol{w}}$$

for the data-dependent coefficient matrices

$$\boldsymbol{A}_{\boldsymbol{\mathfrak{w}}} := \left(\left\langle \underline{\mathfrak{sig}}_{\Xi}(x_i), \underline{\mathfrak{sig}}_{\Xi}(x_j) \right\rangle \right)_{i,j=1}^{M} \quad \text{and} \quad \boldsymbol{B}_{\boldsymbol{\mathfrak{w}}} = \left(z_i^j \right)_{(i,j) \in [M] \times [k]}$$

3. Compute the function, using the weights $\hat{\alpha}_{\lambda}$ and the standard basis (e_i) of \mathbb{R}^k ,

$$\hat{\psi}_{Z} := \sum_{i=1}^{k} \left[\sum_{j=1}^{M} \hat{\alpha}_{ij}^{(\lambda)} \langle \underline{\mathfrak{sig}}_{\Xi}(x_{j}), \underline{\mathfrak{sig}}_{\Xi}(\cdot) \rangle \right] e_{i}.$$
(117)

4. Output: $\hat{\psi}_Z$ (estimator of the regression function ψ_Z).

The fact that this estimator is a particular instance of Definition 2.6 is emphasised in Remark B.8.

5.5 Convex and Consistent Stochastic Process Regression

The main objective of this paper is to establish a universally consistent framework for the nonparametric regression of stochastic processes. In Section 2.3, particularly through Proposition 2.4 and Definition 2.6, we have have identified as our central target object for this the regression operator

$$\mathfrak{c}_{\xi}(\mu) : L^{2}(\mu_{\mathcal{Y}}) \longrightarrow L^{2}(\xi), \quad f \longmapsto \mu_{\cdot}(f), \tag{118}$$

defined for any measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ with \mathcal{X} -marginal $\xi \coloneqq \mu_{\mathcal{X}} \in \mathcal{M}_1(\mathcal{X})$. Recall from (20) and (21) that if μ is the joint law of two stochastic processes X and Y as in (52), i.e., $\mu = \mathbb{P}_{(X,Y)}$ and thus $\xi = \mathbb{P}_X$, then (118) characterises the conditional statistical dependence of Y given X.

Assuming $\mu = \mathbb{P}_{(X,Y)}$ entails no loss of generality (Remark 2.3), and the samples of μ then read:

$$\mathfrak{Z} \equiv \left\{ (X^{(j)}, Y^{(j)}) \mid j \in [N] \right\} \qquad (N \in \mathbb{N})$$

$$(119)$$

as specified in (89) (cf. (10)), or they are of the form (112) if we are interested in dependencies modelled by $\mu = \mathbb{P}_{(X,Z)}$ for a random covariate Z in \mathbb{R}^k . To enable parallel learning of the arguments f in (118), we may also wish to observe some marginal samples of $\mu_{\mathcal{Y}}$, which appear as \mathfrak{Y} in (89).

We then have the following result on the approximation quality of the structures (107) and (115).

Theorem 5.9. Consider the estimators $\hat{\mathscr{R}}^{\nu}_{\theta} : \mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k) \to \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}^k)$ ($\nu = I, II$) given by

$$\widehat{\mathscr{R}}^{\mathrm{I}}_{(\underline{\lambda},n,N,m)}[f] \coloneqq \widehat{\mathscr{R}}^{\mathrm{I}}_{\underline{\lambda},f}[n,N,m] \qquad and \qquad \widehat{\mathscr{R}}^{\mathrm{II}}_{(\underline{\lambda},M)}[f] \coloneqq \widehat{\mathscr{R}}^{\mathrm{II}}_{\underline{\lambda},f(Y)}[M], \tag{120}$$

as defined in (107) and (115). Then for each $f \in \mathcal{L}^2(\mathbb{P}_Y; \mathbb{R}^k)$, any accuracy $\varepsilon > 0$ and any $\delta > 0$,

$$\sup_{\theta \in \mathcal{A}_{f}^{\nu}(\varepsilon,\delta)} \mathbb{P}\Big(\big\| \mu_{\cdot}(f) - \hat{\mathscr{R}}_{\theta}^{\nu}[f] \big\|_{L^{2}(\mathbb{P}_{X};\mathbb{R}^{k})} \ge \varepsilon \Big) \le \delta$$
(121)

for both $\nu = I, II$. In particular, for each $\varepsilon > 0$ and any given confidence level $q \in [0, 1)$, we have

$$\inf_{\theta \in \Theta_{f}^{\nu}(\varepsilon,q)} \mathbb{P}\Big(\left| \mathbb{E} \big[f(Y) \, \big| \, X \big] - \hat{\mathscr{R}}_{\theta}^{\nu}[f](X) \big| < \varepsilon \Big) \ge q$$
(122)

for both $\nu = I, II$. The inequalities (121) hold for the parameter regimes

$$\begin{aligned} \mathcal{A}_{f}^{\mathrm{I}}(\varepsilon,\delta) &\coloneqq \left\{ \left((\lambda_{1},\lambda_{2}), n, N, m \right) \in \mathbb{R}_{>0}^{2} \times \mathbb{N}^{3} \mid \lambda_{1} \leq \lambda_{f}^{\mathrm{I}}(\varepsilon), \ n \geq n_{f,\lambda_{1}}(\varepsilon,\delta), \ m \geq m_{f,\lambda_{1}}(\varepsilon;\mathfrak{Y}), \\ \lambda_{2} \leq \lambda_{\lambda_{1},m}^{\mathrm{II}}(\varepsilon;\mathfrak{Y}), \ N \geq N_{\lambda_{1},\lambda_{2},m}(\varepsilon,\delta;\mathfrak{Y}) \right\} \\ and \qquad \mathcal{A}_{f}^{\mathrm{II}}(\varepsilon,\delta) \coloneqq \left\{ (\lambda, M) \in \mathbb{R}_{>0} \times \mathbb{N} \mid \lambda \leq \hat{\vartheta}_{\psi_{X,Z}}(\varepsilon/2) \ and \ \tilde{\beta}_{M}(\delta,\lambda) \leq \varepsilon/2 \right\}, \end{aligned}$$

respectively, defined for any $(\tilde{\beta}_m)$ as in Cor. 5.8 (iii) and in the notation of (104), (105), (106). The inequalities (122) hold for the parameter regimes

$$\Theta^{\mathrm{I}}_f(\varepsilon,q) \coloneqq \bigcup_{(\tilde{\varepsilon},\tilde{\delta}) \in I(\varepsilon,q)} \mathcal{A}^{\mathrm{I}}_f(\tilde{\varepsilon},\tilde{\delta}) \qquad and \qquad \Theta^{\mathrm{II}}_f(\varepsilon,q) \coloneqq \bigcup_{(\tilde{\varepsilon},\tilde{\delta}) \in I(\varepsilon,q)} \mathcal{A}^{\mathrm{II}}_f(\tilde{\varepsilon},\tilde{\delta}),$$

respectively, which are defined over $I(\varepsilon,q) \coloneqq \{(\tilde{\varepsilon}, \tilde{\delta}) \in \mathbb{R}^2_{>0} \mid \tilde{\varepsilon}(1-\tilde{\delta}) + \tilde{\delta}\varepsilon^2 \leq (1-q)\varepsilon^2\}.$

Proof. The consistency guarantees (121) and (122) rely on a combination of the results from Sections 5.1 to 5.4. A detailed proof of Theorem 5.9 is provided in Appendix B.3.6. \Box

Remark 5.10 (Time-Discretized Non-iid Data). For simplicity, Theorem 5.9 is based on the standard assumption that the empirical data (119) and (112), from which the estimators (120) are computed, are iid samples of the process pairs (X, Y) and (X, Z), respectively. Yet in practice one usually has access to time-discretized samples of these processes only, often taken over non-equally spaced time grids, and oftentimes only statistically dependent (time-discretized) sample paths of these pairs are available rather than many independent realisations. The above consistency theorem can be extended to hold under these non-standard assumptions, in line with the approaches established in [49, Section 8]. A full elaboration of according modifications is left for future research in order to not exceed the scope of this paper.

6 Numerical Experiments

Numerical experiments and practical example applications will be included in the forthcoming arXiv version of this paper.

Appendix A Notation

_

Below are some of the symbols and terminology that we use throughout the main paper.

Symbol	Meaning	Page
\mathcal{X},\mathcal{Y}	factor spaces of the relation space $\mathcal{X} \times \mathcal{Y}$, usually assumed to be Polish spaces.	4
$\mathcal{B}(\mathcal{Z})$	Borel σ -algebra of a given topological space \mathcal{Z} .	5
$\mathcal{M}_1(\mathcal{Z})$	$:= \{ v : \mathcal{B}(\mathcal{Z}) \to [0,1] v \text{ measure} \}; \text{ cone of all Borel probability measures on } \mathcal{Z}.$	5
$\mathfrak{c}_{\xi}(\mu)$	the conditional dependence of a measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ wrt. the source law $\xi \in \mathcal{M}_1(\mathcal{X})$; often identified with the regression operator (14), denoted by the same symbol.	6
$\mu_{+}(f)$: $\mathcal{X} \ni x \mapsto \int_{\mathcal{Y}} f \mathrm{d} \mu_x$; (conditional) expectation of $f \in \mathcal{L}^2(\xi)$ wrt. the ξ -	6
	disintegration $\mu_{\cdot} = (\mu_x)_{x \in \mathcal{X}}$ of a measure $\mu \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y}).$	
δ_p	$: 2^{\mathbb{Z}} \ni A \mapsto \mathbb{1}_A(p) \in [0,1];$ the Dirac measure at a given point $p \in \mathbb{Z}$.	6
$\mathcal{L}^p(\upsilon; E)$	$\coloneqq \{f: \mathcal{Z} \to E f \text{ is measurable}, \ \int_{\mathcal{Z}} \ f\ _E^p \mathrm{d} v < \infty\}, \ p \in [1,\infty) \cup \{0\}.$	7
$\mathbb{1}_A$	$z \mapsto \mathbb{1}_A(z) \coloneqq \{1 \text{ if } z \in A; 0 \text{ else; indicator function over a set } A.$	8
${\cal C}_d$	$:= \{x : [0,1] \to \mathbb{R}^d \mid x \text{ is continuous}\}; \text{ space of continuous paths in } \mathbb{R}^d, \text{ usually} \\ \text{endowed with the uniform norm } \ x\ _{\infty} := \sup_{t \in [0,1]} x_t \text{ (unless otherwise stated)}.$	10
${\mathcal C}^1_d$	the space of all absolutely continuous paths in \mathbb{R}^d .	10
·	$: \mathbb{R}^n \to \mathbb{R}_+$, Euclidean norm on \mathbb{R}^n ; in partic., $ \cdot $ is the absolute value on $\mathbb{R} = \mathbb{R}^1$.	10
$\alpha \wedge \beta$:= min{ α, β }; the minimum of $\alpha, \beta \in \mathbb{R}$.	12
$\delta_{a,b}$:= {1 if $a = b; 0$ else; the Kronecker delta over two elements a, b .	12
[k]	$:= \{1, \dots, k\}, \text{ and } [k]_0 := [k] \cup \{0\} \ (k \in \mathbb{N}).$	12
$[d]^{\star}$	$\coloneqq \bigcup_{m>0} [d]^{\times m} = \{\emptyset, 1, 2, \dots, 11, 12, 112, \dots\}; \text{ the set of all multi-indices with }$	13
	entries in $[d]$, where $[d]^0 \coloneqq \{\emptyset\}$ and $i_1 \cdots i_m \equiv (i_1, \ldots, i_m)$.	
sig	the time-augmented signature transform.	14
\mathcal{H}_d	codomain of the signature on d -dimensional paths, see (42) for the definition.	15
$\overline{\mathbb{R}}$	$:= \mathbb{R} \cup \{-\infty, \infty\}$; the affine closure of the real numbers.	14
$\underline{\mathfrak{sig}}_{\Lambda}, \underline{\mathfrak{sig}}_{\Xi}$	bounded signature transforms, for feature normalisations Λ and $\Xi.$	16
$L^p_X(\mathcal{H})$	space of all <i>p</i> -integrable, \mathcal{H} -valued X-measurable functions, see (57).	19
$(e_i)_{i \in [k]}$	the standard basis of \mathbb{R}^k (i.e.: $e_i = (\delta_{1i}, \cdots, \delta_{ki}) \in \mathbb{R}^k$, for each $i \in [k]$).	21
$\operatorname{supp}(v)$	the support of a measure $v \in \mathcal{M}_1(\mathcal{Z})$, i.e. the smallest closed subset $C \subseteq \mathcal{Z}$ with $v(C) = 1$ (see e.g. [23, Lemma 1.19]).	23
$\varphi _{\tilde{A}}$	the restriction of a map $\varphi: A \to B$ to a subdomain $\tilde{A} \subseteq A$.	24
$\langle u, v \rangle_2$	$:= u_1 v_1 + \ldots + u_d v_d$; the dot product of two vectors $u = (u_i), v = (v_i)$ in \mathbb{R}^d .	51

Notation A.1. Unless specified otherwise, we assume that both \mathcal{X} and \mathcal{Y} are Polish spaces. For a given topological space \mathcal{Z} , we denote by $\mathcal{B}(\mathcal{Z})$ its Borel σ -algebra and by $\mathcal{M}_1(\mathcal{Z}) \coloneqq \{v : \mathcal{B}(\mathcal{Z}) \to [0,1] | v \text{ measure}\}$ the cone of all Borel probability measures on \mathcal{Z} . Given a measure $v \in \mathcal{M}_1(\mathcal{Z})$, we write $\operatorname{supp}(v)$ for its support, that is the smallest closed subset $C \subseteq \mathcal{Z}$ with v(C) = 1 (e.g. [23, Lemma 1.19]). The Cartesian product $\mathcal{X} \times \mathcal{Y}$ is endowed with the product topology (thus, $\mathcal{B}(\mathcal{X} \times \mathcal{Y}) = \mathcal{B}(\mathcal{X}) \otimes \mathcal{B}(\mathcal{Y})$), the canonical coordinate projections $\mathcal{X} \times \mathcal{Y} \to \mathcal{X}$ and $\mathcal{X} \times \mathcal{Y} \to \mathcal{Y}$ are denoted by $\hat{\pi}_1$ and $\hat{\pi}_2$ respectively, and for a given $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ we write $\mu_{\mathcal{X}} \coloneqq (\hat{\pi}_1)_* \mu \equiv \mu \circ \hat{\pi}_1^{-1}$ and $\mu_{\mathcal{Y}} \coloneqq (\hat{\pi}_2)_* \mu$ for its marginals. Given $\xi \in \mathcal{M}_1(\mathcal{X})$, we write $\mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y}) \coloneqq \{\tilde{\mu} \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}) \mid \tilde{\mu}_{\mathcal{X}} = \xi\}$ for the set of all measures on $\mathcal{X} \times \mathcal{Y}$ with \mathcal{X} -marginal ξ . The law of some \mathcal{Z} -valued (Borel) random variable Z over a probability space $(\Omega, \mathscr{F}, \mathbb{P})$ is denoted by $\mathbb{P}_Z \coloneqq Z_* \mathbb{P} \equiv \mathbb{P} \circ Z^{-1} \in \mathcal{M}_1(\mathcal{Z})$.

References

- R. Abergel, C. Louchet, L. Moisan, and T. Zeng. Total Variation Restoration of Images Corrupted by Poisson Noise with Iterated Conditional Expectations. Springer, 2015.
- [2] P. Artzner, F. Delbaen, J.-M. Eber, and D. Heath. Coherent Measures of Risk. Math. Finance, 9.3:203–228, 1999.
- [3] A. Bain and D. Crisan. Fundamentals of Stochastic Filtering. Stochastic Modelling and Applied Probability 3, Springer, 2009.
- [4] C. Bayer, L. Pelizzari, and J. Schoenmakers. Primal and dual optimal stopping with signatures. arXiv preprint arXiv:2312.03444, 2023.
- [5] A. Borovykh, S. Bohte, and C.W. Oosterlee. Conditional Time Series Forecasting with Convolutional Neural Networks. arXiv preprint arXiv:1703.04691, 2017.
- [6] H.-P. Breuer, E.-M. Laine, J. Piilo, and B. Vacchini. Colloquium: Non-Markovian Dynamics in Open Quantum Systems. *Reviews of Modern Physics*, 88.2:021002, 2016.
- [7] C. Carmeli, E. De Vito, A. Toigo, and V. Umanitá. Vector Valued Reproducing Kernel Hilbert Spaces and Universality. *Analysis and Applications*, 8.01:19–61, 2010.
- [8] T. Cass, T. Lyons, and X. Xu. Weighted Signature Kernels. Ann. Appl. Probab., 34(1A):585–626, 2024.
- [9] K.-T. Chen. Integration of paths—a faithful representation of paths by non-commutative formal power series. Trans. Amer. Math. Soc., 89:395–407, 1958.
- [10] P. Cheridito and B. Gersey. Computation of Conditional Expectations with Guarantees. J. Sci. Comput., 95.12:1–30, 2023.
- [11] I. Chevyrev and T. Lyons. Characteristic functions of measures on geometric rough paths. Ann. Probab., 44.6:4049–4082, 2016.
- [12] I. Chevyrev and H. Oberhauser. Signature moments to characterize laws of stochastic processes. J. Mach. Learn. Res., 23.176:1–42, 2022.
- [13] Samuel N Cohen et al. Nowcasting with signature methods. arXiv preprint arXiv:2305.10256, 2023.

- [14] D. L. Cohn. *Measure Theory*. Second Edition. Birkhäuser Advanced Texts, 2013.
- [15] J. B. Conway. A Course in Functional Analysis. Second Edition. Graduate Texts in Mathematics 96, Springer, 1997.
- [16] F. Cucker and S. Smale. On the Mathematical Foundations of Learning. Bull. Amer. Math. Soc., 39.1:1–49, 2002.
- [17] J. Elstrodt. Maß- und Integrationstheorie. Achte, erweiterte und aktualisierte Ausgabe, Springer Spektrum, 2018.
- [18] P. K. Friz and N. B. Victoir. Multidimensional Stochastic Processes as Rough Paths: Theory and Applications. Cambridge Studies in Advanced Mathematics 120, Cambridge University Press, 2010.
- [19] R. Giles. A Generalization of the Strict Topology. Trans. Amer. Math. Soc., 161:467–474, 1971.
- [20] B. Hambly and T. Lyons. Uniqueness for the signature of a path of bounded variation and the reduced path group. Ann. of Math., 171.1:109–167, 2010.
- [21] T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second Edition, Springer Series in Statistics, 2009.
- [22] S. Janson. Gaussian Hilbert Spaces. Number 129. Cambridge University Press, 1997.
- [23] O. Kallenberg. Foundations of Modern Probability. Springer, 2021.
- [24] I. Karatzas and S. Shreve. Brownian Motion and Stochastic Calculus, volume 113. Springer Science & Business Media, 1998.
- [25] A Kechris. Classical Descriptive Set Theory. Graduate Texts in Mathematics 156, Springer, 1995.
- [26] P. Kidger and T. Lyons. Signatory: differentiable computations of the signature and logsignature transforms, on both CPU and GPU. arXiv: 2001.00706 and https://github.com/ patrick-kidger/signatory; published at ICLR 2021, 2020.
- [27] J.P. Klein and M.L. Moeschberger. Survival Analysis: Techniques for Censored and Truncated Data. Springer, 2003.
- [28] K. Kuratowski. *Topology*. Volume 1. New edition, revised and augmented. Academic Press, New York 1966.
- [29] D. Lee and H. Oberhauser. The Signature Kernel. arXiv preprint arXiv:2305.04625, 2023.
- [30] F. Legoll and T. Lelievre. Effective Dynamics using Conditional Expectations. Nonlinearity, 23.9:2131, 2010.
- [31] D. Levin, T. Lyons, and H. Ni. Learning from the past, predicting the statistics for the future, learning an evolving system. arXiv preprint arXiv:1309.0260, 2013.
- [32] A. Lewis. Option Valuation Under Stochastic Volatility. Finance Press, 2000.
- [33] Z. Li, D. Meunier, M. Mollenhauer, and A. Gretton. Optimal Rates for Regularized Conditional Mean Embedding Learning. Advances in Neural Information Processing Systems, 35:4433– 4445, 2022.

- [34] S. Liao, H. Ni, M. Sabate-Vidales, L. Szpruch, M. Wiese, and B. Xiao. Sig-Wasserstein GANs for Conditional Time Series Generation. *Math. Finance*, 2023.
- [35] T. J. Lyons. Differential Equations Driven by Rough Signals. Revista Matemática Iberoamericana, 14.2:215–310, 1998.
- [36] T. J. Lyons, M. Caruana, and T. Lévy. Differential Equations Driven by Rough Paths. Springer, 2007.
- [37] H. G. Matthies, E. Zander, B. V. Rosić, and A. Litvinenko. Parameter Estimation via Conditional Expectation: a Bayesian Inversion. Advanced Modeling and Simulation in Engineering Sciences, 3.1:1–21, 2016.
- [38] M.J. Neely, E. Modiano, and C.-P. Li. Fairness and Optimal Stochastic Control for Heterogeneous Networks. *IEEE/ACM Transactions On Networking*, 16.2:396–409, 2008.
- [39] M. Ottaviani and P.N. Sørensen. The Strategy of Professional Forecasting. Journal of Financial Economics, 81.2:441–466, 2006.
- [40] H. Owhadi and C. Scovel. Separability of Reproducing Kernel Spaces. Proc. Amer. Math. Soc., 145.5:2131–2138, 2017.
- [41] J. Park and K. Muandet. A Measure-Theoretic Approach to Kernel Conditional Mean Embeddings. Advances in neural information processing systems, 33:21247–21259, 2020.
- [42] J. Park and K. Muandet. Regularised Least-Squares Regression With Infinite-Dimensional Output Space. arXiv preprint arXiv:2010.10973, 2020.
- [43] R. Penrose. On Best Approximate Solutions of Linear Matrix Equations. In Math. Proc. Cambridge Philos. Soc., volume 52.1, pages 17–19. Cambridge University Press, 1956.
- [44] J. M. Robins, S. D. Mark, and W. K. Newey. Estimating Exposure Effects by Modelling the Expectation of Exposure Conditional on Confounders. *Biometrics*, pages 479–495, 1992.
- [45] Peter M Robinson. Nonparametric Estimators for Time Series. J. Time Series Anal., 4.3:185– 207, 1983.
- [46] L. C. G. Rogers and D. Williams. Diffusions, Markov Processes, and Martingales. Volume 1, Cambridge University Press, 2000.
- [47] C. Salvi, T. Cass, J. Foster, T. Lyons, and W. Yang. The Signature Kernel is the Solution of a Goursat PDE. SIAM Journal on Mathematics of Data Science, 3.3:873–899, 2021.
- [48] F. S. Scalora. Abstract Martingale Convergence Theorems. Pacific J. Math., 11.4:347–374, 1961.
- [49] A. Schell and H. Oberhauser. Nonlinear Independent Component Analysis for Discrete-Time and Continuous-Time Signals. (Reference is made to arXiv:2102.02876). Ann. Statist., 51.2:487–518, 2023.
- [50] M.J. Schervish. Theory of Statistics. Springer Series in Statistics, 1995.
- [51] B. Schölkopf, R. Herbrich, and A. J. Smola. A Generalized Representer Theorem. In Computational Learning Theory, pages 416–426. Lecture Notes in Computer Science. Vol. 2111. Springer, 2001.

- [52] L. Schwartz. Sous-Espaces Hilbertiens d'Espaces Vectoriels Topologiques et Noyaux Associés (Noyaux Reproduisants). J. Analyse Math., 13:115–256, 1964.
- [53] M. Shanahan. Talking About Large Language Models. arXiv preprint arXiv:2212.03551, 2022.
- [54] S.E. Shreve. Stochastic Calculus for Finance II: Continuous-Time Models, volume 11. Springer, 2004.
- [55] R. H Shumway and D. S. Stoffer. An Approach to Time Series Smoothing and Forecasting Using the EM Algorithm. J. Time Series Anal., 3.4:253–264, 1982.
- [56] M. Sipser. Introduction to the Theory of Computation. (3 ed.). Boston, MA: Cengage Learning, 2013.
- [57] T.O. To and K.W. Yip. A generalized Jensen's inequality. Pacific J. Math., 58.1:255–259, 1975.
- [58] J. M. A. M. Van Neerven. Stochastic Evolution Equations.
- [59] S. Wolfram. What Is ChatGPT Doing... and Why Does It Work? Stephen Wolfram, 2023.
- [60] J. Yeh. Martingales and Stochastic Analysis. Series on Multivariate Analysis 1, World Scientific Publishing, 1995.
- [61] S. L. Zeger and B. Qaqish. Markov Regression Models for Time Series: A Quasi-Likelihood Approach. *Biometrics*, pages 1019–1031, 1988.
- [62] W. X. Zhao et al. A Survey of Large Language Models. arXiv preprint arXiv:2303.18223, 2023.

Appendix B Further Proofs and Technical Remarks

B.1 Ad Section 2

B.1.1 Remarks on Product-Space Measures and Operator Parametrizations

Remark B.1. For each (Borel) probability measure $\mu \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ with \mathcal{X} -marginal $\xi \in \mathcal{M}_1(\mathcal{X})$, there is a canonical pair (X, Y) of jointly distributed (Borel) rvs $X : \Omega \to \mathcal{X}$ and $Y : \Omega \to \mathcal{Y}$ such that $\mu = \mathbb{P}_{(X,Y)}$ and $\xi = \mathbb{P}_X$. Conversely, for any pair $(\tilde{X}, \tilde{Y}) : \Omega \to \mathcal{X} \times \mathcal{Y}$ of (jointly distributed Borel) random variables, their law $\mathbb{P}_{(\tilde{X}, \tilde{Y})}$ is a Borel probability measure with \mathcal{X} -marginal $\mathbb{P}_{\tilde{X}}$.

Proof. This is trivial: Given $\mu \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$, set $(\tilde{\Omega}, \tilde{\mathscr{F}}, \tilde{\mathbb{P}}) := (\mathcal{X} \times \mathcal{Y}, \mathcal{B}(\mathcal{X} \times \mathcal{Y}), \mathfrak{c}_{\xi}(\mu) \otimes \xi)$, which defines a probability space, and consider the canonical projections $X := \hat{\pi}_1$ and $Y := \hat{\pi}_2$ (thus, $Z := (X, Y) = \operatorname{id}_{\mathcal{X}} \times \operatorname{id}_{\mathcal{Y}})$. (Recall that $\hat{\pi}_1 : \tilde{\Omega} \ni (x, y) \mapsto x \in \mathcal{X}$ and $\hat{\pi}_2 : \tilde{\Omega} \ni (x, y) \mapsto y \in \mathcal{X}$.) Then $X : \tilde{\Omega} \to \mathcal{X}$ and $Y : \tilde{\Omega} \to \mathcal{Y}$ are (Borel) random variables with $\mathbb{P}_X = \xi$ and joint law $\mathbb{P}_Z = \mathbb{P} = \mu$. The converse statement is immediate from the definitions.

Remark B.2 (Injective Parametrization (17) over a Hilbert Space). Recognizing the importance of the injectivity of the parametrization in (semi-)parametric statistical models for the uniqueness, identifiability, and analytical tractability of a model, let us remark that the parametrization (17) of the conditional dependency $\mathfrak{c}_{\xi}(\mu)$ is 'injective' in the following basic sense: Denoting by $(\mathcal{H}_{\kappa}, \|\cdot\|_{\kappa}) \subset \mathcal{L}^2(\mu_{\mathcal{Y}})$ the (q-induced) RKHS with kernel $\kappa : (x, y) \mapsto \langle q(x), q(y) \rangle$, there is a 'parameter' Hilbert space $\tilde{\mathcal{H}}_q$ (that is conveniently explicit and regular for the maps q introduced in Sect. 3.3) such that

$$\mathfrak{c}_{\xi}(\mu) : \overline{\mathcal{H}_{\kappa}}^{L^2} \longrightarrow L^2(\xi) \quad \text{reads} \quad \mathfrak{c}_{\xi}(\mu)|_{\mathcal{H}_{\kappa}} = \psi_q \circ \iota$$
(123)

for an isometric isomorphism $\iota : \mathcal{H}_{\kappa} \to \tilde{\mathcal{H}}_q$ and some q-defined bounded linear operator $\psi_q : \tilde{\mathcal{H}}_q \to L^2(\xi)$, the 'q-parametrisation' of $\mathfrak{c}_{\xi}(\mu)$, where both ι and ψ_q are canonically related to q.

Indeed: Take any $f \in \mathcal{H}_{\kappa}$, whence $f = \langle \ell_f, q(\cdot) \rangle$ for some $\ell_f \in \mathcal{H}_{\mathcal{Y}}$. Writing $[\ell_f] := \{\ell \in \mathcal{H}_{\mathcal{Y}} \mid f = \langle \ell, q(\cdot) \rangle \}$, note that $[\ell_f] = \{\ell \in \mathcal{H}_{\mathcal{Y}} \mid \ell \sim \ell_f\} = \ell_f + F$ under the equivalence relation: $\ell \sim \tilde{\ell} : \iff \langle \ell, q \rangle = \langle \tilde{\ell}, q \rangle \iff (\ell - \tilde{\ell}) \in F := \{\hat{\ell} \in \mathcal{H}_{\mathcal{Y}} \mid \langle \hat{\ell}, q(\cdot) \rangle = 0\} = \ker(\varphi)$, for the bounded linear surjection $\varphi : \mathcal{H}_{\mathcal{Y}} \to \mathcal{H}_{\kappa}, \hat{\ell} \mapsto \langle \hat{\ell}, q(\cdot) \rangle$. Consequently, $\iota : (\mathcal{H}_{\kappa}, \| \cdot \|_{\kappa}) \ni f \mapsto [\ell_f] \in \tilde{\mathcal{H}}_q := (\mathcal{H}_{\mathcal{Y}}/F, \| \cdot \| \|)$, with $\|\|[\ell]\|\| := \inf\{\|\tilde{\ell}\|_{\mathcal{H}_{\mathcal{Y}}} \mid \tilde{\ell} \in [\ell]\}$ the canonical quotient space norm, is a (norm)isometric isomorphism of Hilbert spaces: injectivity is clear, while the definition of the RKHS-norm $\| \cdot \|_{\kappa}$ yields norm-preservation, which in turn (shows that $\||\cdot\|\|$ preserves the parallelogram identity and hence) reveals $\|\|\cdot\|\|$ as an inner product norm; that the quotient \tilde{H} is complete follows since $F (= \ker(\varphi))$ is closed in $\mathcal{H}_{\mathcal{Y}}$. Now the (linear, and bounded due to (16)) evaluation map $\psi_q : \tilde{\mathcal{H}}_q \ni [\ell_f] \mapsto \langle \ell_f, q(\cdot) \rangle \in L^2(\xi)$ is well-defined since clearly $\psi_q(\ell_f) = \psi_q(\ell)$ for each $\ell \in [\ell_f]$, and so the asserted composition (123) holds with Proposition 2.4 (17) and said proposition's final statement.

B.1.2 Proof of Proposition 2.4

Proof. First, it is clear that each measure $\mu_x \in \mathcal{M}_1(\mathcal{Y})$ can be identified with the linear functional $\mu_x : L^2(\mu_{\mathcal{Y}}) \ni f \mapsto \int_{\mathcal{Y}} f \, d\mu_x \in \overline{\mathbb{R}}$ (trivially, since $\{\mathbb{1}_A \mid A \in \mathcal{B}(\mathcal{Y})\} \subset L^2(\mu_{\mathcal{Y}})$). That the functions $\mu_{\cdot}(f) : \mathcal{X} \ni x \mapsto \mu_x(f)$ are in $L^2(\xi)$, for each $f \in L^2(\mu_{\mathcal{Y}})$, follows from Jensen's inequality:

$$\left\|\mu_{\cdot}(f)\right\|_{L^{2}(\xi)}^{2} \leq \int_{\mathcal{X}} \int_{\mathcal{Y}} \left|f(y)\right|^{2} \mu_{x}(\mathrm{d}y)\xi(\mathrm{d}x) = \int_{\mathcal{X}\times\mathcal{Y}} \left|f(y)\right|^{2} \mathrm{d}(\mathfrak{c}_{\xi}(\mu)\otimes\xi) = \left\|f\right\|_{L^{2}(\mu_{\mathcal{Y}})}^{2}$$

The above also shows that the operator $(14) =: T[\mathfrak{c}_{\xi}(\mu)]$ is well-defined and bounded. To prove identification, note first that, since $L^2(\mu_{\mathcal{Y}})$ is separable (e.g. [14, Prop. 3.4.5]; recall that \mathcal{Y} is Polish), there is a dense set $\mathcal{D} := \{g_j \mid j \in J\} \subset L^2(\mu_{\mathcal{Y}})$ with J countable. Now take any $\mu, \tilde{\mu} \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$ with $T[\mathfrak{c}_{\xi}(\mu)] = T[\mathfrak{c}_{\xi}(\tilde{\mu})]$ (thus $L^2(\mu_{\mathcal{Y}}) = L^2(\tilde{\mu}_{\mathcal{Y}})$ in particular). Then for each $j \in J$, we have that $\mu . (g_j) = \tilde{\mu} . (g_j)$ in $L^2(\xi)$ and, thus, $\mu_x(g_j) = \tilde{\mu}_x(g_j)$ for each $x \in \mathcal{X} \setminus \mathcal{N}_j$ with \mathcal{N}_j some ξ -nullset. In particular, for each $x \in \mathcal{X} \setminus \mathcal{N}$ with $\mathcal{N} := \bigcup_{j \in J} \mathcal{N}_j$, we have $\mu_x|_{\mathcal{D}} = \tilde{\mu}_x|_{\mathcal{D}}$ and hence $\mu_x = \tilde{\mu}_x$ (since $\mathcal{D} \subset L^2(\mu_{\mathcal{Y}})$ is dense and each $\mathcal{M}_1(\mathcal{Y}) \ni v : L^2(\mu_{\mathcal{Y}}) \ni f \mapsto \int_{\mathcal{Y}} f \, dv \in \mathbb{R}$ is bounded). Hence, as desired, $\mathfrak{c}_{\xi}(\mu) = \mathfrak{c}_{\xi}(\tilde{\mu})$ in $L^2(\xi)$, since $\xi(\mathcal{N}) = 0$ (as J is countable).

For the remaining identification of (14) with the function (16), take any $\mu, \tilde{\mu} \in \mathcal{M}_1^{\xi}(\mathcal{X} \times \mathcal{Y})$ with $\mu_{\cdot}(q) = \tilde{\mu}_{\cdot}(q)$ in $L^2(\xi; \mathcal{H}_{\mathcal{Y}})$ (i.e. such that $\mu_x(q) = \tilde{\mu}_x(q)$ for each $x \in \mathcal{X} \setminus \hat{\mathcal{N}}$ for some ξ -nullset $\hat{\mathcal{N}}$). Then for each $f \in \mathfrak{H}_q := \{\langle \ell, q(\cdot) \rangle \mid \ell \in \mathcal{H}_{\mathcal{Y}}\}$, say $f = \langle \ell_f, q(\cdot) \rangle$, we have the $L^2(\xi)$ -identities

$$\mathfrak{c}_{\xi}(\mu)(f) = \mu \left(\langle \ell_f, q(\cdot) \rangle \right) = \langle \ell_f, \mu(q) \rangle = \langle \ell_f, \tilde{\mu}(q) \rangle = \mathfrak{c}_{\xi}(\tilde{\mu})(f), \quad \text{i.e.:} \quad \mathfrak{c}_{\xi}(\mu)|_{\mathfrak{H}_q} = \mathfrak{c}_{\xi}(\tilde{\mu})|_{\mathfrak{H}_q} \quad (124)$$

(where the second and the fourth identity holds since Bochner integrals commute with bounded functionals; e.g. [58, Sec. 1.3.1]). But since \mathfrak{H}_q is dense in $L^2(\mu_{\mathcal{V}})$ and the operator (14) is bounded, (124) implies that $\mathfrak{c}_{\xi}(\mu) = \mathfrak{c}_{\xi}(\tilde{\mu})$ (as operators (14) and thus, via the lemma's initial identification, also in the original kernel sense (9)) by uniqueness of extension.

Statement (17) is clear from the second equality in (124) and the assumption $L^2(\mu_{\mathcal{Y}}) = \overline{\mathfrak{H}_q}^{L^2}$. \Box

B.1.3 Proof of Lemma 2.9

Proof. Abbreviating the domain and range of (25) by H_1 and H_2 respectively, note that both $\iota_1 : L^2(\mu_{\mathcal{V}}) \ni f \mapsto f \circ \hat{\pi}_2 \in H_1$ and $\iota_2 : L^2(\xi) \ni g \mapsto g \circ \hat{\pi}_1 \in H_2$ are well-defined Hilbert isometries and onto. Indeed: Well-definedness is clear since $\mu_{\mathcal{V}} = (\hat{\pi}_2)_*\mu$ and $\xi = (\hat{\pi}_1)_*\mu$, and isometry holds since $\langle \iota_1(f), \iota_1(g) \rangle_{H_1} = \int_{\mathcal{X} \times \mathcal{Y}} (fg) \circ \hat{\pi}_2 \, \mathrm{d}\mu = \int_{\mathcal{Y}} fg \, \mathrm{d}[(\hat{\pi}_2)_*\mu] = \langle f, g \rangle_{L^2(\mu_{\mathcal{V}})}$ (and likewise for ι_1); the surjectivity of ι_2 (and likewise that of ι_2) is seen as follows: if $\tilde{f} : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$

is Borel-measurable wrt. $\mathcal{X} \times \mathcal{B}(\mathcal{Y}) \equiv \{\mathcal{X} \times A \mid A \in \mathcal{B}(\mathcal{Y})\}$, then \tilde{f} is constant wrt. its first variable (otherwise, i.e. if $\tilde{f}(x_1, y) \neq \tilde{f}(x_2, y)$ for a $y \in \mathcal{Y}$ and some $x_1, x_2 \in \mathcal{X}$ with $x_1 \neq x_2$, then: $A_{(x_1,y)} \coloneqq \tilde{f}^{-1}(\{\tilde{f}(x_1, y)\}) \cap (\mathcal{X} \cap \{y\}) \in \mathcal{X} \times \mathcal{B}(\mathcal{Y})$ (by measurability of \tilde{f} and since all singletons in \mathcal{Y} are (closed and thus) in $\mathcal{B}(\mathcal{Y})$) and $\hat{\pi}_1(A_{(x_1,y)}) \notin \{\emptyset, \mathcal{X}\}$ — a contradiction), that is: there is $f : \mathcal{Y} \to \mathbb{R}$ such that $\tilde{f}(x, y) = f(y)$ for each $(x, y) \in \mathcal{X} \times \mathcal{Y}$, which is equivalent to $\tilde{f} = f \circ \hat{\pi}_2$.

We now claim that the operator $\mathfrak{c}_{\xi}(\mu)$ from (14) relates to the operator (25) via

$$\mathfrak{c}_{\xi}(\mu) = \iota_2^{-1} \circ P_{\mathcal{Y}|\mathcal{X}} \circ \iota_1.$$
(125)

To verify this, we abbreviate $\tilde{T} \coloneqq \iota_2 \circ \mathfrak{c}_{\xi}(\mu) \circ \iota_1^{-1}$ and need to show that, for any $\tilde{f} \in H_1$,

$$\tilde{f} - \tilde{T}\tilde{f} \in H_2^{\mathsf{T}}$$
, i.e. $\left\langle \tilde{f} - \tilde{T}\tilde{f}, \tilde{g} \right\rangle_{L^2(\mu)} = 0$ for each $\tilde{g} \in H_2$. (126)

(Note here that the $P_{\mathcal{Y}|\mathcal{X}}$ -defining projector P is well-defined since the orthogonal decomposition $L^2(\mu) = H_2 \oplus H_2^{\mathsf{T}}$ exists by the fact (consequent to every L^2 -convergent subsequence admitting an μ -a.e. convergent subsequence) that H_2 is a closed subspace of $L^2(\mu)$.) Proving (126), note that

$$\int_{\mathcal{X}\times\mathcal{Y}} \mathbb{1}_{\tilde{A}} \Delta_{\tilde{f}} \, \mathrm{d}\mu = \int_{A\times\mathcal{Y}} \tilde{f} \, \mathrm{d}\mu - \int_{A} \mu_{x}(f) \, \xi(\mathrm{d}x) = \int_{A\times\mathcal{Y}} \tilde{f} \, \mathrm{d}\mu - \int_{\mathcal{X}} \int_{\mathcal{Y}} f(y) \, \mathbb{1}_{\tilde{A}}(x,y) \, \mu_{x}(\mathrm{d}y) \, \xi(\mathrm{d}x)$$

$$\stackrel{(7)}{=} \int_{A\times\mathcal{Y}} \tilde{f} \, \mathrm{d}\mu - \int_{\mathcal{X}\times\mathcal{Y}} \tilde{f}(x,y) \, \mathbb{1}_{\tilde{A}}(x,y) \, (\mu_{\mathcal{Y}|\mathcal{X}} \otimes \xi)(\mathrm{d}(x,y)) \stackrel{(8)}{=} 0 \qquad \left(f \coloneqq \iota_{1}^{-1}(\tilde{f})\right)$$

for each $\tilde{A} \equiv (A \times \mathcal{Y}) \in \mathcal{B}(\mathcal{X}) \times \mathcal{Y}$ and with $\Delta_{\tilde{f}} \coloneqq \tilde{f} - \tilde{T}\tilde{f}$. This proves (126) and, hence, (125). \Box

B.2 Ad Section 3

B.2.1 Remarks on Stochastic Processes and Weighted Signature-Codomains

Remark B.3 (On Stochastic Processes). Writing $\mathbb{P}_X := \mathbb{P} \circ X^{-1}$ and $X(\omega) \equiv (X_t(\omega))_{t \in [0,1]}$ for each $\omega \in \Omega$, let us note the following general facts on stochastic processes S as in (32).

(i) The rv X is $(\mathscr{F}, \mathcal{B}(\|\cdot\|_{1\text{-var}}))$ -measurable by def. Now $\mathcal{B}(\|\cdot\|_{1\text{-var}}) = \mathcal{B}(\|\cdot\|_{\infty})$ by Lemma 3.2, where $\mathcal{B}(\|\cdot\|_{\infty}) = \sigma(\pi_t \mid t \in [0,1]) =: \mathcal{B}(\mathcal{C}^1_{d_{\mathcal{X}}}) = \mathcal{B}(\mathcal{C}_{d_{\mathcal{X}}}) \cap \mathcal{C}^1_{d_{\mathcal{X}}}$ is the Borel σ -algebra on $(\mathcal{C}^1_d, \|\cdot\|_{\infty})$. (Here, $\pi_t : (x_t)_{t \in [0,1]} \mapsto x_t$ is the t-projection from $\mathcal{C}_{d_{\mathcal{X}}}$ onto $\mathbb{R}^{d_{\mathcal{X}}}$.) Consequently,

$$\mathbb{P}_X \in \mathcal{M}_1(\mathcal{C}^1_{d_{\mathcal{X}}}, \|\cdot\|_{\infty}) \quad \text{and} \quad X_t : \Omega \ni \omega \mapsto X_t(\omega) \in \mathbb{R}^{d_{\mathcal{X}}} \text{ is } (\mathscr{F}, \mathcal{B}(\mathbb{R}^{d_{\mathcal{X}}})) \text{-measurable}$$

for each $t \in [0, 1]$. Hence, we can equivalently define the stochastic process $X : \Omega \to \mathcal{X}$ as an [0, 1]-indexed family $(X_t)_{t \in [0, 1]}$ of (Borel) random vectors $X_t : \Omega \to \mathbb{R}^{d_{\mathcal{X}}}$ such that $t \mapsto X_t(\omega)$ is continuous for each $\omega \in \Omega$, see e.g. [46, Section II.27].

- (ii) In stricter terminology, a stochastic process $X \equiv (X_t(\omega)) : [0,1] \times \Omega \to \mathbb{R}^d$ defined as a $(\mathscr{F}, \mathcal{B}(\mathcal{C}^1_{d_{\mathcal{X}}})$ -measurable map $X : \Omega \to \mathcal{C}_{d_{\mathcal{X}}}$, as we did above, is called *jointly measurable*. If $(\Omega, \mathscr{F}, \mathbb{P})$ is filtered then it can carry stronger measurability notions (such as progressive measurability or predictability, see e.g. [60, Proposition 2.23]), but for our purposes the weak notion of joint measurability will suffice, cf. also Lemma 3.5.
- (iii) It will be no loss of generality for us to assume (if convenient) that in fact

$$X: \Omega \to \mathfrak{D}_X$$
, where $\mathfrak{D}_X \coloneqq \operatorname{supp}(\mathbb{P}_X)$

is the support of X. Indeed: By its definition, the support \mathfrak{D}_X is the smallest closed subset $C \subseteq \mathcal{X}$ for which $\mathbb{P}_X(C) = 1$, see e.g. [23, Lemma 1.19]. Hence $\tilde{\Omega} \coloneqq X^{-1}(\mathfrak{D}_X) \in \mathscr{F}$ is a \mathbb{P} -full set, which implies that X and its \mathfrak{D}_X -valued twins $\tilde{X} \coloneqq \mathbb{1}_{\tilde{\Omega}} \cdot X + \mathbb{1}_{\Omega \setminus \tilde{\Omega}} \cdot x_0 : \Omega \to \mathfrak{D}_X$ (any $x_0 \in \mathfrak{D}_X$ fixed, say $x_0 = 0$ for convenience) are indistinguishable.

(iv) We further assume that the X-induced sub- σ -algebra $\Sigma_X \coloneqq \sigma(X) \subseteq \mathscr{F}$ is \mathbb{P} -complete. Recall that this assumption entails no loss of generality: If $(\Omega, \Sigma_X, \mathbb{P})$ is not complete, we can immediately and 'minimally' complete it as follows. Writing $\mathcal{N}_{\mathbb{P}} \coloneqq \{N \subseteq \Omega \mid \exists A \in \mathscr{F} : N \subseteq A \text{ and } \mathbb{P}(A) = 0\}$ for the system of all subsets of \mathbb{P} -nullsets $[in \mathscr{F}]$, define $\Sigma_X^{\mathbb{P}} \coloneqq \{A \cup N \mid A \in \mathcal{F}_X, N \in \mathcal{N}_{\mathbb{P}}\}$ and $\mathscr{F}^{\mathbb{P}} \coloneqq \{A \cup N \mid A \in \mathscr{F}, N \in \mathcal{N}_{\mathbb{P}}\}$ and $\mathbb{P}: \mathscr{F}^{\mathbb{P}} \to [0,1]$ by $\mathbb{P}(A \cup N) \coloneqq \mathbb{P}(A)$ for all $A, \in \mathscr{F}, N \in \mathcal{N}_{\mathbb{P}}$. Then $\Sigma_X^{\mathbb{P}} \subseteq \mathscr{F}^{\mathbb{P}}$, both $(\Omega, \Sigma_X^{\mathbb{P}}, \mathbb{P})$ and $(\Omega, \mathscr{F}^{\mathbb{P}}, \mathbb{P})$ are complete probability spaces, and each complete extension μ of \mathbb{P} is an extension of \mathbb{P} , e.g. [17, Satz 6.3]. All our objects of interest stay the same when passing to this completion, that is (trivially) $\mathbb{E}[Y \mid \Sigma_X] = \mathbb{E}[Y \mid \Sigma_X^{\mathbb{P}}] \mathbb{P}$ -a.s. and $L^p(\Omega, \Sigma_X, \mathbb{P}; \mathcal{H}) \cong L^p(\Omega, \Sigma_X^{\mathbb{P}}, \mathcal{H})$ (canonically).

Remark B.4 (Alternative Hilbert Codomains). Let $\gamma \equiv (\gamma_m)_{m \geq 0} > 0$ with $0 < \gamma_m \leq \lambda^m$ (all $m \in \mathbb{N}$) for some $\lambda > 0$. Given the gradation (40) of V together with the 'Euclidean identification' of its components $V_m \cong (V_1^{\otimes m}, \langle \cdot, \cdot \rangle_m)$ seen above [right after (34)], with $V_1 \cong (\mathbb{R}^d, \langle \cdot, \cdot \rangle_2)$ and $\langle \cdot, \cdot \rangle_m \equiv \prod_{[m]} \langle \cdot, \cdot \rangle_2$ and $\| \cdot \|_m = \sqrt{\langle \cdot, \cdot \rangle_m^2}$, another natural Hilbert space structure on V is given by

$$\mathcal{H}_{d}^{\gamma} \coloneqq \left\{ \boldsymbol{t} \in V \mid \|\boldsymbol{t}\|_{\gamma} \coloneqq \sqrt{\sum_{m \ge 0} \gamma_{m} \|\pi_{m}(\boldsymbol{t})\|_{m}^{2}} < \infty \right\}$$
(127)

together with the inner product $\langle \boldsymbol{s}, \boldsymbol{t} \rangle_{\gamma} \coloneqq \sum_{m \geq 0} \gamma_m \langle \pi_m(\boldsymbol{s}), \pi_m(\boldsymbol{t}) \rangle_m$. It is clear that $(\mathcal{H}_d^{\gamma}, \langle \cdot, \cdot \rangle)$ is a Hilbert space (as the ℓ_2 -direct sum of the Hilbert spaces $(V_m, \gamma_m \langle \cdot, \cdot \rangle_m), m \geq 0$), see e.g. [15, Prop. I.6.2]), and we denote its topology by τ_{γ} . Clearly $\tau_{\tilde{\gamma}} \subseteq \tau_{\gamma}$ if $\tilde{\gamma} \leq \gamma$, as then $\|\cdot\|_{\tilde{\gamma}} \leq \|\cdot\|_{\gamma}$.

B.2.2 Proof of Lemma 3.2

Proof. The first assertion holds by [18, Propositions 1.31 & 1.32]. In fact, [18, Prop. 1.31] asserts that for $\mathfrak{X} := \mathbb{R}^d \times L^1([0,1]; \mathbb{R}^d)$ and $\mathfrak{Y} := \mathcal{C}_d$, the map $f : \mathfrak{X} \to \mathfrak{Y}$ given by $f(c,v) := c + \int_0^{\cdot} v_s \, ds$ is a Banach space isomorphism (which also proves the norm identity $||x||_{1-\text{var}} = |x_0| + ||\dot{x}||_{L^1}$ on \mathcal{C}_d^1). From this, [25, Theorem 15.1] implies that the image $f(\mathfrak{X}) = \mathcal{C}^1$ is a Borel subset of $(\mathcal{C}^1, ||\cdot||_{\infty})$, i.e. that $\mathcal{C}^1 \in \mathcal{B}(\mathcal{C}_d)$. That $(\mathcal{C}^1, ||\cdot||_{\text{var}})$ is separable and Banach is stated as [18, Corollary 1.35].

For the lemma's second assertion, note first that since $\|\cdot\|_{1-\text{var}} \ge \|\cdot\|_{\infty}$ (which is easy to see), we find that the 1-variation topology on \mathcal{C}^1 is finer than the uniform topology on \mathcal{C}^1 , which of course implies that $\mathcal{B}_{1-\text{var}} \coloneqq \sigma(\mathcal{C}^1, \|\cdot\|_{1-\text{var}}) \supseteq \sigma(\mathcal{C}^1, \|\cdot\|_{\infty}) \eqqcolon \mathcal{B}_{\infty}$. Since the separability of $(\mathcal{C}^1, \|\cdot\|_{1-\text{var}})$ guarantees that the σ -algebra $\mathcal{B}_{1-\text{var}}$ is generated by the closed $\|\cdot\|_{1-\text{var}}$ -balls, the converse inclusion $\mathcal{B}_{1-\text{var}} \subseteq \mathcal{B}_{\infty}$ follows if we can show that

$$B_r^1(x) \coloneqq \{ y \in \mathcal{C}^1 \mid \|y - x\|_{1-\text{var}} \le r \} \in \mathcal{B}_{\infty} \text{ for every } x \in \mathcal{C}^1 \text{ and any } r \ge 0.$$
(128)

To see that this holds, fix any $x \in C^1$ and $r \ge 0$ and recall that, by definition of the 1-variation norm,

$$||z||_{1-\operatorname{var}} = \sup_{\mathcal{I} \in \mathfrak{I}} V_{\mathcal{I}}(z) \text{ with } V_{(t_{\nu})}(z) \coloneqq |z_0| + \sum_{\nu} |z_{t_{\nu+1}} - z_{t_{\nu}}|$$

and where $\mathfrak{I} \coloneqq \{\mathcal{I} = (t_{\nu}) \mid \mathcal{I} \text{ is a (finite) dissection of } [0,1]\}$. Given any $\mathcal{I} \in \mathfrak{I}$ it is clear that the function $Q_{\mathcal{I}} : \mathcal{C}^1 \ni y \mapsto V_{\mathcal{I}}(y-x)$ is continuous wrt. $\|\cdot\|_{\infty}$, whence the level set $C_{\mathcal{I}} \coloneqq \{y \in \mathcal{C}^1 \mid Q_{\mathcal{I}}(y) \leq r\}$ is $\|\cdot\|_{\infty}$ -closed. Combined with this, the immediate identity

$$B_r^1(x) = \bigcap_{\mathcal{I} \in \mathfrak{I}} C_{\mathcal{I}}$$
 implies that $B_r^1(x)$ is closed wrt. $\|\cdot\|_{\infty}$,

which shows that (128) holds as desired.

B.2.3 Proof of Lemma 3.4

Proof. Note that since both $(\mathcal{X}, \|\cdot\|_{1\text{-var}})$ and $(\mathcal{Y}, \|\cdot\|_{1\text{-var}})$ are Polish by Lemma 3.2, so is the product space $(\mathcal{X} \times \mathcal{Y}, \|\cdot\|_{\alpha})$ with $\|(x, y)\|_{\alpha} \coloneqq \max\{\|x\|_{1\text{-var}}, \|y\|_{1\text{-var}}\}$, as the norm $\|\cdot\|_{\alpha}$ induces

the product topology on $(\mathcal{X}, \|\cdot\|_{1-\text{var}}) \times (\mathcal{Y}, \|\cdot\|_{1-\text{var}})$. This implies that \mathcal{Z} itself is Polish, since the norms $\|\cdot\|_{\alpha}$ and $\|\cdot\|_{1-\text{var}}$ are equivalent on $\mathcal{X} \times \mathcal{Y}$. The latter equivalence of norms also gives that

$$\mathcal{B}(\mathcal{Z}, \|\cdot\|_{1-\operatorname{var}}) = \mathcal{B}(\mathcal{Z}, \|\cdot\|_{\alpha}), \tag{129}$$

and since further $\mathcal{B}(\mathcal{Z}, \|\cdot\|_{\alpha}) = \mathcal{B}(\mathcal{X}, \|\cdot\|_{\infty}) \otimes \mathcal{B}(\mathcal{Y}, \|\cdot\|_{\infty}) = \mathcal{B}(\mathcal{Z}, \|\cdot\|_{\beta})$ for the norm $\|z\|_{\beta} := \max\{\|\pi_{\mathcal{X}}(z)\|_{\infty}, \|\pi_{\mathcal{Y}}(z)\|_{\infty}\}$ by Lemma 3.2 (recalling that: (a) the Borel σ -algebra of the product of two (second countable) topological spaces equals the product of their Borel σ -algebras, and (b) the norm $\|\cdot\|_{\beta}$ induces the product topology on $(\mathcal{X}, \|\cdot\|_{\infty}) \times (\mathcal{Y}, \|\cdot\|_{\infty})$) and with the norms $\|\cdot\|_{\infty}$ and $\|\cdot\|_{\beta}$ being equivalent on $\mathcal{X} \times \mathcal{Y}$, we find that $\mathcal{B}(\mathcal{Z}, \|\cdot\|_{1-\operatorname{var}}) = \mathcal{B}(\mathcal{Z}, \|\cdot\|_{\infty})$ as desired.

The claimed characterisation of measurability holds by (129) (upon recalling that $\mathcal{B}(\mathcal{Z}, \|\cdot\|_{\alpha}) = \mathcal{B}(\mathcal{X}, \|\cdot\|_{1-\operatorname{var}}) \otimes \mathcal{B}(\mathcal{Y}, \|\cdot\|_{1-\operatorname{var}}))$ and the fact that a product-space-valued function (here: (X, Y)) is product-measurable iff all of its factor components (here: X and Y) are measurable.

B.2.4 Proof of Lemma 3.12

Proof. Each of the spaces $(V_m, \langle \cdot, \cdot \rangle_m)$ from (40) is Hilbert and separable (with ONB $[d]_m^* := \{w \in [d]^* \mid |w| = m\}$). Thus the space \mathcal{H} is Hilbert and separable — with ONB $[d]^*$ — as the Hilbert direct sum of the family $\{(V_m, \langle \cdot, \cdot \rangle_m) \mid m \in \mathbb{N}_0\}$, see e.g. [15, Proposition I.6.2]. The inclusion $\mathfrak{sig}(\mathcal{C}_d^1) \subset \mathcal{H}$ follows from the factorial decay of the signature coefficients, cf. [11, Corollary 5.5].

The last assertion follows from the usual *p*-variation continuity of sig, see e.g. [11, Corollary 5.5], and the fact that the locally convex topology from [11, Section 2] (defined by the (fundamental) family of semi-norms $\Psi := (||| \cdot |||_{\lambda} | \lambda > 0)$ on *V*, where $||| t |||_{\lambda} := \sum_{m \ge 0} ||\pi_m(t)||_m \cdot \lambda^m$; denote the associated locally *m*-convex topology by τ_{lc}) is finer than the [canonical, i.e. $|| \cdot ||$ -induced] topology on \mathcal{H} (denote this topology by $\tau_{\mathcal{H}}$). To prove the asserted inclusion of topologies: Since τ_{lc} is metrizable, see e.g. [11, Corollary 2.4], the topological space (V, τ_{lc}) is sequential, whence $\tau_{\mathcal{H}} \subseteq \tau_{lc}$ iff every τ_{lc} -convergent sequence in *V* is $\tau_{\mathcal{H}}$ convergent. This clearly holds, however, since for every null-sequence (v_k) in (V, τ_{lc}) there is $k_0 \in \mathbb{N}$ with $\sup_{k \ge k_0} ||v_k|||_{\lambda} < 1$ (for some $\lambda > 1$), whence for $k \ge k_0$ we find that $||v_k||^2 = \sum_{m \ge 0} ||\pi_m(v_k)||_m^2 \le ||v_k|||_{\lambda}$ goes to zero as $k \to \infty$.

B.2.5 Proof for Remark 3.16

Claim: On $\mathcal{H}_d^{\downarrow}$, the locally convex topology τ_{\downarrow} is finer than the (127)-induced subspace topology τ_{γ} .

Proof. Since τ_{\downarrow} is metrizable (Lemma 3.15), the space $(\mathcal{H}_{d}^{\downarrow}, \tau_{\downarrow})$ is sequential, whence $\tau_{\gamma} \subseteq \tau_{\downarrow}$ (iff $\mathrm{id}_{\mathcal{H}_{d}^{\downarrow}} : (\mathcal{H}_{d}^{\downarrow}, \tau_{\downarrow}) \to (\mathcal{H}_{d}^{\downarrow}, \tau_{\gamma}), v \mapsto v$, is continuous) iff every τ_{\downarrow} -convergent sequence in $\mathcal{H}_{d}^{\downarrow}$ is τ_{γ} -convergent. This clearly holds, however, since for every null-sequence (v_{k}) in $(\mathcal{H}_{d}^{\downarrow}, \tau_{\downarrow})$ there is $k_{0} \in \mathbb{N}$ with $\sup_{k \geq k_{0}} |||v_{k}||_{\lambda} < 1$ and hence, for each $k \geq k_{0}$,

$$\|oldsymbol{v}_k\|_{\gamma}^2 = \sum_{m>0} \gamma_m \|\pi_m(oldsymbol{v}_k)\|_m^2 \le \|oldsymbol{v}_k\|_{\lambda} \longrightarrow 0 \quad ext{as} \ k o \infty.$$

As one consequence of the inclusion $\tau_{\gamma} \subseteq \tau_{\downarrow}$, statement (45) also holds for $\mathcal{H}_{d}^{\downarrow}$ replaced by \mathcal{H}_{d}^{γ} . \Box

B.2.6 Proof of Lemma 3.17

Proof. Since the augmentation map $\bar{\iota} = (\theta, \mathrm{id}_{\mathcal{X}}) : \mathcal{C}_d^1 \to \mathcal{C}_d^1$ from (38), with $\theta(x) \coloneqq (t)_{t \in [0,1]}$, is $(\| \cdot \|_{p\text{-var}}, \| \cdot \|_{p\text{-var}})$ -continuous, the continuity (47) follows immediately from (45) and assertion (46).

Let us now prove (46). Recall for this that $\Lambda(t) = \sum_{m\geq 0} \lambda_t^m \pi_m(t)$ by definition, and that, since $(\mathcal{H}_d^{\downarrow}, \tau_{\downarrow})$ is first-countable by Lemma 3.15, the map Λ is continuous if it is sequentially continuous.

Let hence $\mathbf{t}, (\mathbf{t}_k) \subset \mathcal{H}_d^{\downarrow}$ with $\lim_{k \to \infty} \mathbf{t}_k = \mathbf{t}$ in τ_{\downarrow} , i.e. $\lim_{k \to \infty} \||\mathbf{t}_k - \mathbf{t}||_{\tilde{\lambda}} = 0$ for all $\tilde{\lambda} > 0$. Then

$$\|\Lambda(\boldsymbol{t}_{k}) - \Lambda(\boldsymbol{t})\|^{2} \leq 2 \sum_{m \geq 0} \left(\|\lambda_{\boldsymbol{t}_{k}}^{m} \pi_{m}(\boldsymbol{t}_{k}) - \lambda_{\boldsymbol{t}_{k}}^{m} \pi_{m}(\boldsymbol{t})\|_{m}^{2} + \|\lambda_{\boldsymbol{t}_{k}}^{m} \pi_{m}(\boldsymbol{t}) - \lambda_{\boldsymbol{t}}^{m} \pi_{m}(\boldsymbol{t})\|_{m}^{2} \right).$$
(130)

With $\lambda_{\cdot}: t \mapsto \lambda_t \tau_{\downarrow}$ -continuous, we have $\lim_{k \to \infty} |\lambda_{t_k} - \lambda_t| = 0$ and thus $c \coloneqq \max\{\sup_k \lambda_{t_k}, \lambda_t\} < \infty$ ∞ . As noted above, $\lim_{k\to\infty} \||\boldsymbol{t}_k - \boldsymbol{t}\||_c = 0$ and $\||\boldsymbol{t}\||_c < \infty$. The summands $\alpha_{m,k} \coloneqq \|\lambda_{\boldsymbol{t}_k}^m \pi_m(\boldsymbol{t}_k - \boldsymbol{t})\|_m^2$ and $\beta_{m,k} \coloneqq \|(\lambda_{\boldsymbol{t}_k}^m - \lambda_{\boldsymbol{t}}^m)\pi_m(\boldsymbol{t})\|_m^2$ on the right-hand side of (130) compare to

$$\alpha_{m,k} \leq a_{m,k} \coloneqq c^m \|\pi_m(\boldsymbol{q}_k - \boldsymbol{q})\|_m \quad \text{and} \quad \beta_{m,k} \leq b_m \coloneqq 4c^m \|\pi_m(\boldsymbol{q})\|_m, \tag{131}$$

for all $k \in \mathbb{N}$ and each $m \geq m_0$, for some sufficiently large $m_0 \in \mathbb{N}_0$. Hence and from (130) we find

$$\lim_{k \to \infty} \|\Lambda(\boldsymbol{t}_k) - \Lambda(\boldsymbol{t})\|^2 \le 2 \lim_{k \to \infty} \|\boldsymbol{t}_k - \boldsymbol{t}\|_c + \sum_{m \ge 0} \lim_{k \to \infty} \beta_{m,k} = 0,$$
(132)

where interchanging limit and summation for the second summand in (132) is permissible by dominated convergence, which in turn is applicable thanks to the $(\beta_{m,k})$ -domination in (131) and the fact that $\sum_{m>0} |b_m| = 4 |||q|||_c < \infty$. This proves (46), as desired.

B.2.7 Definition of the Strict Topology and Detailed Proof of Proposition 3.19

Let $\mathcal{Z} := (\mathcal{C}^1_d, \|\cdot\|_{1-\text{var}})$, write $C_b(\mathcal{Z})$ for the set of all bounded continuous functions on \mathcal{Z} , and set

$$B_0(\mathcal{Z}) \coloneqq \{ \psi : \mathcal{Z} \to \mathbb{R} \text{ bounded } | \forall \varepsilon > 0 : \exists \mathcal{K} \subset \mathcal{Z} \text{ compact } : \sup_{x \in \mathcal{Z} \setminus \mathcal{K}} |\psi(x)| < \varepsilon \}$$

for the set of all bounded functions on \mathcal{Z} that vanish at infinity.

Definition B.5 (Strict Topology [19]). The strict topology on $C_b(\mathcal{Z})$, denoted by $\tau_{\text{str}}^{\mathcal{Z}}$, is the topology induced by the family of seminorms

$$p_{\psi}(f) \coloneqq \sup_{x \in \mathcal{Z}} |f(x)\psi(x)|, \quad \psi \in B_0(\mathcal{Z}).$$

Note that the strict topology is weaker than the uniform topology on $C_b(\mathcal{Z})$ but stronger than the topology of compact (and thus also pointwise) convergence on $C_b(\mathcal{Z})$, see [19, Theorem 2.4 (i)].

Proof of Proposition 3.19. Let us first note that, indeed,

$$\sum_{w \in [\underline{d}]^*} \underline{\xi}_w^{\Lambda} \cdot w = \Lambda \circ \underline{\mathfrak{sig}}.$$
(133)

Indeed, $\pi_m(\underline{\mathfrak{sig}}_{\Lambda}(x)) = \underline{\lambda}_x^m \sum_{|w|=m} \xi_w(\bar{x}) = \lambda_{\underline{\mathfrak{sig}}(x)}^m [(\sum_{|w|=m} \xi_w) \circ \bar{\iota}](x) = \lambda_{\underline{\mathfrak{sig}}(x)}^m [\pi_m \circ \underline{\mathfrak{sig}}](x) = \lambda_{\underline{\mathfrak{sig}}(x) = \lambda_{\underline{\mathfrak{sig}}(x)}^m [\pi_m \circ \underline{\mathfrak{sig}}](x) = \lambda_{\underline{\mathfrak{sig}}(x) = \lambda_{\underline{\mathfrak{sig}}$ $\pi_m\big(\lambda_{\mathfrak{sig}(x)}^m \cdot \pi_m(\underline{\mathfrak{sig}}(x))\big) = \pi_m\big(\delta_{\lambda_{\mathfrak{sig}(x)}}(\underline{\mathfrak{sig}}(x))\big) = \pi_m\big((\Lambda \circ \underline{\mathfrak{sig}})(x)\big) \text{ for each } (x,m) \in \mathcal{Z} \times \mathbb{N}_0.$ To see $\mathcal{A}_{\Lambda} \subset C(\mathcal{Z})$, note [from (50)] that since each function $\varphi \in \mathcal{A}_{\Lambda}$ can be represented as

$$\varphi = \langle \ell_{\varphi}, \underline{\mathfrak{sig}}_{\Lambda} \rangle \quad \text{for some} \ \ \ell_{\varphi} \in \mathbb{R}[\underline{d}],$$

the desired $\|\cdot\|_{1\text{-var}}$ -continuity of φ follows from (133) and statement (47) of Lemma 3.17. Since for $\mathbf{1} = 1 \cdot \emptyset \in \mathbb{R}[\underline{d}]$ we have $\underline{\xi}_{\emptyset}^{\Lambda} = \langle \mathbf{1}, \underline{\mathfrak{sig}}_{\Lambda} \rangle = \lambda_{\mathbf{1}}^{0} \cdot \langle \emptyset, \underline{\mathfrak{sig}} \rangle \equiv 1$ on \mathcal{Z} , clearly \mathcal{A}_{Λ} is non-vanishing. For \mathcal{A}_{Λ} being point-separating, note that for any $x, y \in \mathcal{Z}$ with $x \neq y$ we have $\underline{\mathfrak{sig}}(x) \neq \underline{\mathfrak{sig}}(y)$ by Lemma 39, and hence also $\underline{\mathfrak{sig}}_{\Lambda}(x) \neq \underline{\mathfrak{sig}}_{\Lambda}(y)$ by (133) and the injectivity of Λ . This implies that there is $w_0 \in [\underline{d}]^*$ such that $\langle w_0, \underline{\mathfrak{sig}}_{\Lambda}(x) \rangle \neq \langle w_0, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle$, whence for $\varphi \coloneqq \underline{\xi}_{w_0}^{\Lambda} \in \mathcal{A}_{\Lambda}$ we find $\varphi(x) \neq \varphi(y)$.

To prove that \mathcal{A}_{Λ} is an algebra, we need to show $\varphi \cdot \psi \in \mathcal{A}_{\Lambda}$ for any two $\varphi, \psi \in \mathcal{A}_{\Lambda}$. And indeed,

$$\varphi \cdot \psi = \sum_{w, \tilde{w} \in [\underline{d}]^*} \langle \ell_{\varphi}, w \rangle \langle \ell_{\psi}, \tilde{w} \rangle \underline{\lambda}^{|\tilde{w}|} \underline{\lambda}^{|\tilde{w}|} \langle w, \underline{\mathfrak{sig}} \rangle \langle \tilde{w}, \underline{\mathfrak{sig}} \rangle$$
(134)

$$= \sum_{w,\tilde{w}\in[d]^*} \underline{\lambda}^{|w|+|\tilde{w}|}_{(\ell_{\varphi},w)} \langle \ell_{\psi}, \tilde{w} \rangle \langle w \sqcup \tilde{w}, \underline{\mathfrak{sig}} \rangle$$
(135)

$$= \sum_{m\geq 0} \underline{\lambda}^{m}_{\cdot, \tilde{w}: |w \sqcup \tilde{w}| = m} \langle \ell_{\varphi}, w \rangle \langle \ell_{\psi}, \tilde{w} \rangle \langle w \sqcup \tilde{w}, \underline{\mathfrak{sig}} \rangle$$
(136)

$$=\sum_{m\geq 0} \underline{\lambda}^{m}_{\cdot} \langle \pi_{m}(\ell_{\varphi\psi}), \underline{\mathfrak{sig}} \rangle$$
(137)

$$= \langle \ell_{\varphi\psi}, \Lambda \circ \underline{\mathfrak{sig}} \rangle, \quad \text{for} \quad \ell_{\varphi\psi} := \sum_{w, \tilde{w} \in [\underline{d}]^*} \langle \ell_{\varphi}, w \rangle \langle \ell_{\psi}, \tilde{w} \rangle \cdot w \sqcup \tilde{w}.$$
(138)

Since both ℓ_{φ} , $\ell_{\psi} \in \mathbb{R}[\underline{d}]$ (and thus $\langle \ell_{\varphi}, w \rangle = 0$ and $\langle \ell_{\psi}, w \rangle = 0$ for almost all $w \in [\underline{d}]^*$), we get that also $\ell_{\varphi\psi} \in \mathbb{R}[\underline{d}]$ and hence $\varphi \cdot \psi \in \mathcal{A}_{\Lambda}$ as desired. A few remarks on this are in order:

While (134) holds simply by linearity of (41), equation (135) involved the character identity

 $\langle \ell_1, \underline{\mathfrak{sig}} \rangle \cdot \langle \ell_2, \underline{\mathfrak{sig}} \rangle = \langle \ell_1 \sqcup \ell_2, \underline{\mathfrak{sig}} \rangle, \quad \text{for any} \ \ \ell_1, \ell_2 \in \mathbb{R}[\underline{d}]$

(see e.g. [36, proof of Thm. 2.15]), where $\[mu: \mathbb{R}[\underline{d}]^{\times 2} \to \mathbb{R}[\underline{d}]\]$ is the so-called *shuffle product*, defined e.g. in [36, eq. (2.5) (p. 35)]. Denoting by $|\ell| \coloneqq \max\{|w| \mid w \in [\underline{d}]^* : w \in \ell\}\]$ the maximal length of any word contained (as a summand) in a given polynomial $\ell \in \mathbb{R}[\underline{d}]\]$, it holds that $|\ell_1 \sqcup \ell_2| =$ $|\ell_1| + |\ell_2|$. This justifies equation (136), where we also used that $[\underline{d}]^* \times [\underline{d}]^* = \bigsqcup_{m \ge 0} \{(w, \tilde{w}) \in$ $[\underline{d}]^* \times [\underline{d}]^* \mid |w \amalg \tilde{w}| = m\}\]$ defines a (disjoint) partition. For equation (137) we used that $\pi_m(\ell_{\varphi\psi}) =$ $\sum_{w \amalg \tilde{w} \mid w \amalg \tilde{w} \mid w \amalg \tilde{w}, w \amalg \tilde{w}$

We thus saw that \mathcal{A}_{Λ} is a subalgebra of $C(\mathcal{Z})$. Now if Λ is in fact an fN of the form (48), then

$$\|\varphi\|_{\infty} \coloneqq \sup_{x \in \mathcal{Z}} \left| \langle \ell_{\varphi}, \underline{\mathfrak{sig}}_{\Lambda}(x) \rangle \right| \le \|\ell_{\varphi}\| \sup \|\Lambda(\underline{\mathfrak{sig}}(\mathcal{Z}))\| \le \|\ell_{\varphi}\| R < \infty$$

by Cauchy-Schwarz, which shows that in this case even $\mathcal{A}_{\Lambda} \subset C_b(\mathcal{Z})$ as claimed.

The asserted denseness of \mathcal{A}_{Λ} in $(C_b(\mathcal{Z}), \tau_{\text{str}}^{\mathcal{Z}})$ is then guaranteed by [19, Theorem 3.1], which generalises the theorem of Stone-Weierstrass to the $\tau_{\text{str}}^{\mathcal{Z}}$ -modulated non-compact setting.

If finally we are in the (unnormalised) special case $\Lambda = \mathrm{id}_{\mathcal{H}_{\downarrow}}$, that is if $\lambda \equiv 1$, then the above arguments show that $\mathcal{A} \coloneqq \mathcal{A}_{\mathrm{id}_{\mathcal{H}_{\downarrow}}}$ is a subalgebra of $C(\mathcal{Z})$, while the bounds (44) yield that, for each $\varphi \in \mathcal{A}$,

$$\|\varphi\|_{\infty;\mathcal{Z}} = \left\| \langle \ell_{\varphi}, \underline{\mathfrak{sig}} \rangle \right\|_{\infty;\mathcal{Z}} \le \|\ell_{\varphi}\| \sup_{x \in \mathcal{Z}} \|\underline{\mathfrak{sig}}(x)\| \le \|\ell_{\varphi}\| \sup_{x \in \mathcal{Z}} \sum_{m \ge 0} \frac{\|\bar{x}\|_{1-\operatorname{var}}^m}{m!} \le \|\ell_{\varphi}\| e^{\kappa_{\mathcal{Z}}+1} < \infty$$

if $\kappa_{\mathcal{Z}} \coloneqq \sup_{x \in \mathcal{Z}} \|x\|_{1-\text{var}}$ is assumed finite, and then $\mathcal{A} \subset C_b(\mathcal{Z})$ as desired.

The following corollary is an immediate consequence of Proposition 3.19.

Corollary B.6. Given $d \in \mathbb{N}$, let \mathcal{Z} be a subset of \mathcal{C}_d^1 . Then for each $f \in C_b(\mathcal{Z})$ we have

$$f = \lim_{k \to \infty} \left\langle \tilde{\ell}_k, \Lambda \circ \underline{\mathfrak{sig}} \right\rangle \ \text{ in } \ \tau_{\mathrm{str}}^{\mathcal{X}}, \quad \text{for some } \ (\tilde{\ell}_k) \subset \mathbb{R}[\underline{d}]$$

B.2.8 Proof of Lemma 4.3

Proof. Let us abbreviate $\varphi := \underline{\operatorname{sig}}_{\Xi}$ and $\phi_w := \underline{\xi}^{\Lambda}_w(X)$ for each $w \in [\underline{d}]^*$. We first prove that $\sigma(X) = \sigma(\varphi(X))$: The inclusion $\overline{\sigma(\varphi(X))} \subseteq \sigma(X)$ is immediate since $\varphi : \mathcal{X} \to \mathcal{H}_{\underline{d}}$ is continuous (and hence Borel-measurable), see Lemma 3.17. For the converse, note that since φ is also an injection, we have that $A_{\varphi} := \varphi(A) \in \mathcal{B}(\mathcal{H}_{\underline{d}}) \cap \varphi(\mathcal{X})$ for any fixed $\|\cdot\|_{1\text{-var}}$ -open set $A \subseteq \mathcal{X}$. Indeed, the set A_{φ} is analytic (as the continuous image of a Borel subset of a Polish space) and so is its complement $A_{\varphi}^c \equiv \varphi(\mathcal{X}) \setminus A_{\varphi} = \varphi(A^c)$, where the last identity holds since φ is injective; this implies that $A_{\varphi} \in \mathcal{B}(\mathcal{H}_{\underline{d}} \cap \varphi(\mathcal{X}))$ by a theorem of Souslin [28, Corollary 3.1 (p. 486)]. Consequently, $X^{-1}(A) = (\varphi(X))^{-1}(A_{\varphi}) \in \sigma(\varphi(X))$ and hence $\sigma(X) \subseteq \sigma(\varphi(X))$, as desired. Let us next prove the second identity in (56).

The inclusion $\sigma(\phi_w \mid w \in [\underline{d}]^*) \subseteq \sigma(\varphi(X))$ is immediate since $\phi_w = \langle w, \varphi(X) \rangle$ for each $w \in [\underline{d}]^*$ (cf. (41)) and each $\langle w, \cdot \rangle : \mathcal{H}_{\underline{d}} \to \mathbb{R}$ is continuous. The converse inclusion holds as, pointwise on Ω ,

$$\|\varphi(X) - \psi_n\| \le \sum_{m=n+1}^{\infty} \|\pi_m(\varphi(X))\|_m \xrightarrow{m \to \infty} 0, \text{ for } \psi_n \coloneqq \sum_{|w| \le n} \phi_w \cdot w$$

(cf. (50) and (47)), i.e. $\varphi(X)$ is the pointwise limit of $\sigma(\phi_w | w)$ -measurable functions and thus $\sigma(\phi_w | w)$ -measurable itself.

B.2.9 Proof of Theorem 4.8

Proof. We begin with the well-known observation that the space $L^2_X(\mathcal{H}_{\mathcal{Y}})$ from (57) is a closed linear subspace of $(L^2(\mathbb{P};\mathcal{H}_{\mathcal{Y}}), \|\cdot\|_{L^2(\mathcal{H}_{\mathcal{Y}})})$, which entails the $\langle\cdot,\cdot\rangle_{L^2(\mathcal{H}_{\mathcal{Y}})}$ -orthogonal decomposition

$$L^{2}(\mathbb{P};\mathcal{H}_{\mathcal{Y}}) = L^{2}_{X}(\mathcal{H}_{\mathcal{Y}}) \oplus L^{2}_{X}(\mathcal{H}_{\mathcal{Y}})^{\perp}.$$

(The closedness of $L^2_X(\mathcal{H}_{\mathcal{Y}})$ follows from the well-known fact (which persists for Hilbert-valued random variables [58, Proposition 2.11]) that L^2 -convergence implies almost sure convergence on a subsequence.) Denoting $E := L^2(\mathbb{P}; \mathcal{H}_{\mathcal{Y}})$ and $G := L^2_X(\mathcal{H}_{\mathcal{Y}})$ for brevity, we then adopt (from the scalar-valued setting, e.g. [58, Section 11.1]) the classical perspective that the (vector-valued) conditional expectation $\mathbb{Y}^{\Lambda}_X := \mathbb{E}[\mathbb{Y}^{\Lambda} \mid X] \in G$ is the orthogonal projection of \mathbb{Y}^{Λ} onto G along G^{\perp} , in symbols:

$$\mathbb{Y}_X^{\Lambda} = \mathbb{P}_G \mathbb{Y}^{\Lambda} \quad \text{for the orthogonal projector } \mathbb{P}_G : E \to E \quad \text{on } G (= \operatorname{im}(\mathbb{P}_G)).$$
(139)

To see that this perspective (139) is true in the present vector-valued setting (54) & (55), denote $\mathbb{Y}_G := \mathbb{P}_G \mathbb{Y}^{\Lambda}$ and notice that then $\Delta := \mathbb{Y}^{\Lambda} - \mathbb{Y}_G \in G^{\perp}$, that is $\langle \Delta, \chi \rangle_{L^2(\mathcal{H}_{\mathcal{Y}})} = 0$ for all $\chi \in G$. In particular,

$$0 = \langle \Delta, w \mathbb{1}_A \rangle_{L^2(\mathcal{H}_{\mathcal{Y}})} = \langle \Delta \mathbb{1}_A, w \rangle_{L^2(\mathcal{H}_{\mathcal{Y}})} = \langle \int_A \mathbb{Y}^\Lambda \, \mathrm{d}\mathbb{P}, w \rangle - \langle \int_A \mathbb{Y}_G \, \mathrm{d}\mathbb{P}, w \rangle \quad \left(A \in \Sigma_X, \ w \in [\tilde{d}]^*\right),$$

where the last identity is due to Bochner integrals commuting with bounded linear functionals, cf. [58, Sec. 1.3.1]; note that each element of E is Bochner-integrable by definition (54) and [58, Prop. 1.16]. Since $[\tilde{d}]^*$ is an orthonormal basis of $\mathcal{H}_{\mathcal{Y}}$, the above implies that: $\int_A \mathbb{Y}^A d\mathbb{P} = \int_A \mathbb{Y}_G d\mathbb{P}$, for all $A \in \Sigma_X$. But the latter property is characteristic also of the vector-valued conditional expectation $\mathbb{E}[\mathbb{Y}^A | \Sigma_X]$, see e.g. [58, Theorem 11.10], which implies that $\mathbb{Y}_G = \mathbb{Y}_X^A$ as claimed in (139).

The above characterisation (139) of \mathbb{Y}_X^{Λ} as the orthogonal projection of \mathbb{Y}^{Λ} onto G implies that

$$\left\| \mathbb{Y}^{\Lambda} - \mathbb{Y}^{\Lambda}_{X} \right\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})} \leq \left\| \mathbb{Y}^{\Lambda} - \mathbb{Z} \right\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})} \quad \text{for all } \mathbb{Z} \in G.$$
(140)

Moreover, the Hilbert projection theorem guarantees that the $\arg \min in (140)$ is unique, so that in fact

$$\Psi_X^{\Lambda} = \underset{\mathbb{Z}\in G}{\operatorname{arg\,min}} \mathbb{E} \| \Psi^{\Lambda} - \mathbb{Z} \|^2.$$
(141)

Now in order to make the variational identity (141) more operational, recall from Prop. 4.5 that

the set
$$\mathfrak{G} \coloneqq \Psi(X) = \left\{ \psi_{\alpha}(X) \mid \alpha \in \mathfrak{L}_X^2 \right\}$$
 from (59) is $\| \cdot \|_{L^2(\mathcal{H}_{\mathcal{Y}})}$ -dense in G . (142)

Let us observe how (142) and (141) imply (65): Abbreviating $\Phi(\mathbb{Z}) \coloneqq \|\mathbb{Y}^{\Lambda} - \mathbb{Z}\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}^{2}$, which defines a function $\Phi : G \to \mathbb{R}_{+}$ that is clearly $\|\cdot\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}$ -continuous and strictly convex, we find that

$$\Phi(\mathbb{Y}_X^{\Lambda}) \stackrel{(141)}{=} \inf_{\mathbb{Z}\in G} \Phi(\mathbb{Z}) \stackrel{(142)}{=} \inf_{\mathbb{Z}\in\mathfrak{G}} \Phi(\mathbb{Z}) = \inf_{\alpha\in\mathfrak{L}_X^2} \Phi(\psi_\alpha(X)) \eqqcolon \gamma.$$

Hence for any minimizing sequence of (66), i.e. any sequence (α_k) in \mathfrak{L}^2_X with $\lim_{k\to\infty} \Phi(\psi_{\alpha_k}(X)) = \gamma$, the functions $(\mathbb{Z}_k) \coloneqq (\psi_{\alpha_k}(X)) \subset G$ are a minimizing sequence for $\inf_{\mathbb{Z}\in G} \Phi(\mathbb{Z})$. Upon recalling that $\|\cdot\|_{L^2(\mathcal{H}_Y)}$ satisfies the parallelogram identity and that G is convex, a quick computation shows

$$\|\mathbb{Z}_n - \mathbb{Z}_m\|_{L^2(\mathcal{H}_{\mathcal{Y}})}^2 \le 2\Phi(\mathbb{Z}_n) + 2\Phi(\mathbb{Z}_m) - 4\Phi(\mathbb{Y}_X^{\Lambda}) \longrightarrow 0 \quad (\text{for } n, m \to \infty), \tag{143}$$

which implies that $(\mathbb{Z}_k)_{k\in\mathbb{N}}$ is Cauchy. Hence, and since G is complete, there is $\mathbb{Z}_{\star} \in G$ such that $\mathbb{Z}_{\star} = \lim_{k\to\infty} \mathbb{Z}_k$ in $\|\cdot\|_{L^2(\mathcal{H}_{\mathcal{V}})}$, whence we have $\Phi(\mathbb{Z}_{\star}) = \lim_{k\to\infty} \Phi(\mathbb{Z}_k) = \gamma$ and thus

$$\mathbb{Z}_{\star} \in \underset{\mathbb{Z}\in G}{\operatorname{arg\,min}} \Phi(\mathbb{Z}), \text{ which, by (141), implies } \mathbb{Z}_{\star} = \mathbb{Y}_X^{\Lambda}$$

(as noted in the lead-up to (140), the above argmin contains exactly one element only). This proves (65). Regarding the final claim on almost sure convergence, note $\|\mathbb{Z}_k - \mathbb{Z}_\star\|_{L^1(\mathcal{H}_{\mathcal{Y}})}^2 \leq \|\mathbb{Z}_k - \mathbb{Z}_\star\|_{L^2(\mathcal{H}_{\mathcal{Y}})}^2 \leq 2(\Phi(\mathbb{Z}_k) - \gamma) =: 2\beta_k$ by (143) [and since $\|\cdot\|_{L^1} \leq \|\cdot\|_{L^2}$], whence if $\sum_{k=0}^{\infty} \sqrt{\beta_k} < \infty$ then, by monotone convergence, $\int_{\Omega} \sum_{k=0}^{\infty} \|\mathbb{Z}_k - \mathbb{Z}_\star\| d\mathbb{P} = \sum_{k=0}^{\infty} \|\mathbb{Z}_k - \mathbb{Z}_\star\|_{L^1(\mathcal{H}_{\mathcal{Y}})} < \infty$, thus $\sum_{k=0}^{\infty} \|\mathbb{Z}_k - \mathbb{Z}_\star\| < \infty$ P-a.s. and hence $\lim_{k\to\infty} \|\mathbb{Z}_k - \mathbb{Z}_\star\| = 0$ a.s., as claimed. \Box

B.2.10 Proof of Corollary 4.11

Proof. Let $f \in C(\mathfrak{D}_Y)$. As $\mathfrak{D}_Y \coloneqq \operatorname{supp}(\mathbb{P}_Y)$ is assumed compact, there is $(\ell_f^{(\nu)})_{\nu} \subset \mathbb{R}[\tilde{d}]$ such that

$$\lim_{\nu \to \infty} \left\| f - \langle \ell_f^{(\nu)}, \underline{\mathfrak{sig}}_{\Lambda} \rangle \right\|_{\infty;\mathfrak{D}_Y} = 0$$
(144)

by Corollary B.6, so minimizing sequences for the optimisation problem (76) exist.

Since $\mathbb{P}(Y \in \mathfrak{D}_Y) = 1$, or even $Y(\Omega) \subseteq \mathfrak{D}_Y$ without loss of generality by Remark B.3 (iii),

$$C \coloneqq \sup_{\mu \in \mathbb{N}} \sup_{\omega \in \Omega} \left| \varsigma_{\mu}(Y(\omega)) \right| < \infty \quad \text{for} \quad \varsigma_{\nu}(y) \coloneqq f(y) - \langle \ell_{f}^{(\nu)}, \underline{\mathfrak{sig}}_{\Lambda}(y) \rangle.$$

Indeed, given $\epsilon > 0$ there is $\nu_0 \in \mathbb{N}$ with $\sup_{\nu \ge \nu_0, \omega} |\varsigma_{\nu}(Y(\omega))| \le \sup_{\nu \ge \nu_0} \|\varsigma_{\nu}\|_{\infty;\mathfrak{D}_Y} \le \epsilon$ by (144), and so $C \le \max_{\nu \le \nu_0} \|\varsigma_{\nu}\|_{\infty;\mathfrak{D}_Y} + \epsilon \le 1 + \|\Lambda\|_{\infty;\mathcal{H}} \cdot \max_{\nu \le \nu_0} \|\ell_f^{(\nu)}\| + \epsilon < \infty$. Thus and with (144),

$$\left| \mathbb{E} \left[f(Y) \, \big| \, X \right] - \mathbb{E} \left[\left\langle \ell_f^{(\nu)}, \underline{\mathfrak{sig}}_{\Lambda}(Y) \right\rangle \, \big| \, X \right] \right| = \left| \mathbb{E} \left[\varsigma_{\nu}(Y) \, \big| \, X \right] \right| \stackrel{\nu \to \infty}{\longrightarrow} 0 \quad \mathbb{P}\text{-a.s}$$

via the conditional dominated convergence theorem, which proves (α) in (77). The convergence (β) follows as in (73), while the corollary's last assertion is immediate from (77) and (65).

B.3 Ad Section 5

B.3.1 Proof of Proposition 5.1

Proof. That \mathfrak{H}_{Ξ} is a vRKHS with reproducing kernel K is shown in [7, Example 3.2 (iii)]. This reference also implies that \mathfrak{H}_{Ξ} is separable: Denoting by $\varphi := \mathfrak{sig}_{\Xi} : \mathcal{X} \to \mathcal{H}_{\mathcal{Y}}$ the feature map of \mathcal{H}_{κ} , we see (e.g. from [40, Lemma 2.1]) that the RKHS \mathcal{H}_{κ} is separable due to the separability of $\varphi(\mathcal{X})$, where the latter holds by the separability of \mathcal{X} (Lemma 3.2) and the fact that φ is continuous (Lemma 3.17). The asserted separability of \mathfrak{H}_{Ξ} now follows from the unitary equivalence $\mathfrak{H}_{\Xi} \cong \bigoplus_{w \in [d]^*} \mathcal{H}_{\kappa}$ (given in [7, Example 3.2 (iii)]) and the fact that the countable direct sum of separable Hilbert spaces is separable.

For a proof of (81), note first that \mathfrak{H}_{Ξ} can be canonically injected into $L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})$ via [42, Lemma 2.1 (i)] by the fact that $\sup_{x \in \mathcal{X}} ||K(x, x)||_{\text{op}} = |\kappa(x, x)| \leq \sup_{t \in \mathcal{H}_{\mathcal{Y}}} ||\Lambda(t)||_{\mathcal{H}_{\mathcal{Y}}}^2 =: R_{\Lambda}^2 < \infty$. (The Borel-measurability of each $f \in \mathfrak{H}_{\Xi}$ holds since $(w \mid w \in [\tilde{d}]^*)$ is a countable ONB of $\mathcal{H}_{\mathcal{Y}}$ and the feature map \mathfrak{sig}_{Ξ} of \mathcal{H}_{κ} is continuous.) Next, take any $g \in L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})$ and $\varepsilon > 0$. Then

$$\sum_{m=0}^{\infty} \beta_m(g) = \|g\|_{L^2(\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}})}^2 < \infty \quad \text{for} \quad \beta_m(g) \coloneqq \sum_{|w|=m} \int_{\mathcal{X}} |\langle g(x), w \rangle|^2 \mathbb{P}_X(\mathrm{d}x)$$

and hence $\sum_{m=m_0+1}^{\infty} \beta_m(g) \leq \varepsilon^2/2$ for some $m_0 \in \mathbb{N}$. As done in (62), we can for all $w \in [\tilde{d}]^*$ find

$$\alpha_w \in \mathbb{R}[d] \quad \text{such that} \quad \left\| \langle g, w \rangle - \langle \alpha_w, \underline{\mathfrak{sig}}_{\Xi} \rangle \right\|_{L^2(\mathbb{P}_X)}^2 = \left\| \langle g(X), w \rangle - \varphi_{\alpha_w} \right\|_{L^2(\mathbb{P})}^2 \leq \varepsilon^2 (2\tilde{d})^{-|w|} / 4$$

by the density (59) (resp. (61)). Hence for $f_{\varepsilon} \coloneqq \sum_{w=0}^{m_0} \sum_{|w|=m} \langle \alpha_w, \underline{\mathfrak{sig}}_{\Xi} \rangle w \in \mathfrak{H}_{\Xi}$ we have that

$$\left\|g - f_{\varepsilon}\right\|_{L^{2}(\mathbb{P}_{X};\mathcal{H}_{\mathcal{Y}})}^{2} = \sum_{|w| \le m_{0}} \left\|\langle g - f_{\varepsilon}, w \rangle\right\|_{L^{2}(\mathbb{P}_{X})}^{2} + \sum_{m=m_{0}+1}^{\infty} \beta_{m}(g) \le \varepsilon^{2} \left(\sum_{m=0}^{\infty} 2^{-m} + 2\right)/4 = \varepsilon^{2}$$

(cf. (63)). Since $g \in L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})$ and $\varepsilon > 0$ have been arbitrary, statement (81) follows.

The proposition's subsequent assertion is an immediate corollary to the above argument.

For the remaining statements on (82), note first that clearly $\Psi_X^* \subseteq L^2(\mathcal{X}, \mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})$ (the arguments are the same as for the inclusion $\Psi_X \subseteq L^2(\mathbb{P}_X; \mathcal{H}_{\mathcal{Y}})$) and $\Psi_X \subseteq \Psi_X^*$. Further, for any $f \in \mathfrak{H}_{\Xi}$ we have

$$\sum_{w \in [\tilde{d}]^*} \mathbb{E}\left[\langle f(X), w \rangle^2\right] = \mathbb{E}\left[\|f(X)\|_{\mathcal{H}_{\mathcal{Y}}}^2\right] = \|f\|_{L^2(\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}})}^2 \le R_{\Lambda}^2 \|f\|_{\mathfrak{H}_{\Xi}}^2 < \infty$$
(145)

by [42, Lemma 2.1 (i)], which shows that the inclusion $\mathfrak{H}_{\Xi} \subseteq L^2(\mathcal{P}_X; \mathcal{H}_{\mathcal{Y}})$ is continuous. Upon recalling that $\mathcal{H}_{\kappa} = \{ \langle u, \underline{\mathfrak{sig}}_{\Xi} \rangle \mid u \in \mathcal{H}_{\mathcal{X}} \}$ (which is the feature-map representation of \mathcal{H}_{κ}), the estimate (145) also shows $\mathfrak{H}_{\Xi} \subseteq \Psi_{\mathcal{X}}^{\star}$. That all inclusions are dense has been shown above. \Box

B.3.2 Proof of Lemma 5.2

Proof. Since the functions $\tilde{\phi}_{f,\lambda}$ and $\phi_{\mu|\lambda}^{(m)}$ are strictly convex and coercive, they both admit a unique minimizer over their respective domains (proving the assertions on $\Upsilon^*_{\mu,\lambda,m}$), see e.g. the proof of [42, Lemma 2.4] for the full argument. In particular, $\tilde{\phi}_{f,\lambda}(h_{\star}) = \min_{h \in \mathcal{H}_{\tilde{\kappa}}} \tilde{\phi}_{f,\lambda}(h) =: \gamma_{\star}$ for some (unique) $h_{\star} \in \mathcal{H}_{\tilde{\kappa}}$. But since $\mathcal{H}_{\tilde{\kappa}}$ is an RKHS with feature map $\underline{\mathfrak{sig}}_{\Lambda}$, we have that (a) $\min_{\ell \in \mathcal{H}_{\mathcal{Y}}} \phi_{\mu_{\mathcal{Y}}|\lambda}^{f}(\ell) = \gamma_{\star}$ (cf. the definition (83)), and (b) there is $\ell_{\star} \in \mathcal{H}_{\mathcal{Y}}$ such that $h_{\star} = \langle \ell_{\star}, \underline{\mathfrak{sig}}_{\Lambda} \rangle$. Hence and since $\phi_{f,\lambda}(h_{\star}) = \phi_{\mu_{\mathcal{Y}}|\lambda}^{f}(\ell_{\star})$, a minimizer $\ell_{\mu_{\mathcal{Y}},\lambda}^{f}$ as in (86) exists.

Suppose now that $\mu_{\mathcal{Y}}$ is finitely supported, i.e. a convex combination of the form $\mu_{\mathcal{Y}} = \sum_{i \in I} p_i \delta_{y_i}$ for some $\{y_i \mid i \in I\} \subset \mathcal{Y}$ with I finite (this representation is valid since the path space \mathcal{Y} is Polish and hence Radon), and note that the associated function (83) then reads

$$\mathcal{H}_{\mathcal{Y}} \ni \ell \longmapsto \phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\ell) = \sum_{i \in I} p_{i} \left| f(y_{i}) - \langle \ell, q_{i} \rangle \right) \right|^{2} + \lambda \|q_{\ell}^{*}\|_{\tilde{\kappa}}^{2} \quad \text{with} \quad q_{i} \coloneqq \underline{\mathfrak{sig}}_{\Lambda}(y_{i})$$
(146)

and $q := \underline{\mathfrak{sig}}_{\Lambda} : \mathcal{Y} \to \mathcal{H}_{\mathcal{Y}}$ and $q_{\ell}^* := \langle \ell, q \rangle \in \mathcal{H}_{\tilde{\kappa}}$. Let further $q_i^* := \langle q_i, q \rangle \in \mathcal{H}_{\tilde{\kappa}}$ for each $i \in I$. The derivation of (88) follows the idea behind the classical representer theorem [51], for which we consider the (closed) subspaces $V_{\mu_{\mathcal{Y}}} = \operatorname{span}\{q_i \mid i \in I\} \subset \mathcal{H}_{\mathcal{Y}}$ and $V_{\mu_{\mathcal{Y}}}^* := \operatorname{span}\{q_i^* \mid i \in I\} \subset \mathcal{H}_{\tilde{\kappa}}$. Then $\mathcal{H}_{\mathcal{Y}} = V_{\mu_{\mathcal{Y}}} \oplus V_{\mu_{\mathcal{Y}}}^{\perp}$ and $\mathcal{H}_{\tilde{\kappa}} = V_{\mu_{\mathcal{Y}}}^* \oplus (V_{\mu_{\mathcal{Y}}}^*)^{\perp}$, so that for any fixed minimizer $\hat{\ell} \in \mathcal{H}_{\mathcal{Y}}$ of (146) we have $\hat{\ell} = \hat{\ell}_1 + \hat{\ell}_2$ and $q_{\hat{\ell}}^* = \hat{q}_1^* + \hat{q}_2^*$ for some pairs $(\hat{\ell}_1, \hat{\ell}_2) \in V_{\mu_{\mathcal{Y}}} \times V_{\mu_{\mathcal{Y}}}^{\perp}$ and $(\hat{q}_1^*, \hat{q}_2^*) \in V_{\mu_{\mathcal{Y}}}^* \times (V_{\mu_{\mathcal{Y}}}^*)^{\perp}$ which are both unique with these properties; said uniqueness implies that $\hat{q}_{\nu}^* = \langle \hat{\ell}_{\nu}, q \rangle$ for $\nu = 1, 2$ (note for this that $\langle \hat{\ell}_2, q \rangle \in (V_{\mu_{\mathcal{Y}}}^*)^{\perp}$ since $\langle q_i^*, \langle \hat{\ell}_2, q \rangle \rangle_{\tilde{\kappa}} = \langle \tilde{\kappa}(y_i, \cdot), \langle \hat{\ell}_2, q \rangle \rangle_{\tilde{\kappa}} = \langle \hat{\ell}_2, q(y_i) \rangle = 0$ for each $i \in I$, by the reproducing property). Clearly $\|\hat{\ell}\|_{\mathcal{H}_{\mathcal{Y}}} \geq \|\hat{\ell}_1\|_{\mathcal{H}_{\mathcal{Y}}}$, and by orthogonality,

$$\langle \hat{\ell}, q_i \rangle = \langle \hat{\ell}_1, q_i \rangle \quad \text{for each } i \in I \quad \text{and} \quad \left\| q_{\hat{\ell}}^* \right\|_{\tilde{\kappa}}^2 = \left\| q_{\hat{\ell}_1}^* \right\|_{\tilde{\kappa}}^2 + \left\| \hat{q}_2^* \right\|_{\tilde{\kappa}}^2 \ge \left\| \langle \hat{\ell}_1, q \rangle \right\|_{\tilde{\kappa}}^2. \tag{147}$$

Combining (146) and (147) implies $\phi_{\mu_{\mathcal{Y}}|\lambda}^{f}(\hat{\ell}) \geq \phi_{\mu_{\mathcal{Y}}|\lambda}^{f}(\hat{\ell}_{1})$ and hence, as $\hat{\ell}$ is a minimizer of $\phi_{\mu_{\mathcal{Y}}|\lambda}^{f}$, that

$$\phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\hat{\ell}_{1}) = \min_{\ell \in \mathcal{H}_{\mathcal{Y}}} \phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\ell) = \min_{\ell \in V_{\mu_{\mathcal{Y}}}} \phi^{f}_{\mu_{\mathcal{Y}}|\lambda}(\ell).$$

Since the above minimizer $\hat{\ell}$ of (146) was arbitrary, the claims surrounding (88) follow for $\ell_{\mu_{\mathcal{Y}},\lambda}^{f,\star} \coloneqq \hat{\ell}_1$ if we can show that the function (146) has only one minimizer over $V_{\mu_{\mathcal{Y}}}$, i.e. that

$$\left| \underset{\ell \in V_{\mu_{\mathcal{Y}}}}{\arg\min} \phi_{\mu_{\mathcal{Y}}|\lambda}^{f}(\ell) \right| = 1.$$
(148)

To this end, suppose $I = \{1, \ldots, |I|\}$ wlog and $p_i \equiv |I|^{-1}$ as in (87), and let $\tilde{\ell} \in V_{\mu_{\mathcal{Y}}}$ be any minimizer of $\phi^f_{\mu_{\mathcal{Y}}|\lambda}$. Then $\tilde{\ell} = \sum_{j \in I} \tilde{\alpha}_j q_j$ for some $\tilde{\alpha}^{(\tilde{\ell})} \equiv (\tilde{\alpha}_j) \in \mathbb{R}^{|I|}$ and $q^*_{\tilde{\ell}} \equiv \langle \tilde{\ell}, q \rangle = \sum_{j \in I} \tilde{\alpha}_j q^*_j \in V^*_{\mu_{\mathcal{Y}}}$ minimizes the function $\tilde{\phi}_{f,\lambda}$ from the lemma's last line. Thus, said $\tilde{\alpha}^{(\tilde{\ell})}$ is a global minimizer of

$$\mathbb{R}^{|I|} \ni (\alpha_j) \longmapsto \tilde{F}(\alpha) \coloneqq \tilde{\phi}_{f,\lambda}(\sum_{j \in I} \alpha_j q_j^*) = \frac{1}{|I|} \sum_{i=1}^{|I|} \left| f(y_i) - \sum_{j=1}^{|I|} \alpha_j \langle q_i, q_j \rangle \right|^2 + \lambda \sum_{i,j=1}^{|I|} \alpha_i \alpha_j \langle q_i, q_j \rangle$$
$$= \frac{1}{|I|} \left[\left\| b - \mathbf{A} \alpha \right\|_2^2 + |I| \lambda \alpha^{\mathsf{T}} \mathbf{A} \alpha \right]$$
(149)

for $\|\cdot\|_2$ the Euclidean norm on $\mathbb{R}^{|I| \times |I|}$ and with the $\mu_{\mathcal{Y}}$ -dependent coefficients

$$\boldsymbol{A} \coloneqq \left(\langle q_i, q_j \rangle \right)_{i,j=1}^{|I|} \quad \text{and} \quad \boldsymbol{b} \coloneqq \left(f(y_i) \right)_{i=1}^{|I|}.$$
(150)

After expanding and rearranging the terms in (149), we obtain for $F \coloneqq |I|\tilde{F}$ that

$$F(\alpha) = \alpha^{\mathsf{T}} (\mathbf{A}^{\mathsf{T}} \mathbf{A} + |I| \lambda \mathbf{A}) \alpha - 2\alpha^{\mathsf{T}} \mathbf{A} b \quad \text{and} \quad \nabla_{\alpha} F = 2(\mathbf{A}^{\mathsf{T}} \mathbf{A} + |I| \lambda \mathbf{A}) \alpha - 2\mathbf{A} b.$$

Since the function $F : \mathbb{R}^{|I|} \to \mathbb{R}$ is convex, all of its local minimizers are global minimizers, and so

$$\underset{\alpha \in \mathbb{R}^{|I|}}{\operatorname{arg\,min}} F(\alpha) = \{ \alpha \in \mathbb{R}^{|I|} \mid \nabla_{\alpha} F = 0 \} = \{ \alpha \in \mathbb{R}^{|I|} \mid (\mathbf{A}^{\mathsf{T}} \mathbf{A} + |I| \lambda \mathbf{A}) \alpha = \mathbf{A}b \}.$$
(151)

Now since $\tilde{\alpha}^{(\tilde{\ell})} \in \arg\min_{\alpha \in \mathbb{R}^{|I|}} F(\alpha)$, the claim (148) thus follows if \boldsymbol{A} (and thus $\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A} + |I|\lambda\boldsymbol{A})$ is invertible. But since \boldsymbol{A} is a Gramian matrix by its definition (150), we know that \boldsymbol{A} is invertible iff

the vectors
$$q_1, \ldots, q_{|I|} \in \mathcal{H}_{\mathcal{Y}}$$
 are linearly independent. (152)

But under the lemma's premise, namely that the points y_i , $i \in I$, be pairwise distinct, the required independence (152) follows by (Theorem 3.9 and the injectivity of any (48) and) [36, Cor. 2.16].

As to the lemma's final assertion, note that for the (unique) minimizer $h_{\star} = \arg \min_{h \in \mathcal{H}_{\mathcal{K}}} \phi_{f,\lambda}(h)$ we saw that $\mathcal{L}_{h_{\star}} := \{\ell \in \mathcal{H}_{\mathcal{Y}} \mid \langle \ell, q \rangle = h_{\star}\} = \{\ell \in \mathcal{H}_{\mathcal{Y}} \mid \tilde{\phi}_{f,\lambda}(\ell) = \gamma_{\star}\}$ and hence

$$\|h_{\star}\|_{\tilde{\kappa}} \stackrel{\text{def}}{=} \inf\{\|\ell\|_{\mathcal{H}_{\mathcal{Y}}} \,|\, \ell \in \mathcal{L}_{h_{\star}}\} \stackrel{(88)}{=} \|\ell_{\mu_{\mathcal{Y}},\lambda}^{f,\star}\|_{\mathcal{H}_{\mathcal{Y}}}$$

as claimed.

B.3.3 Proof of Proposition 5.3

Proof. The equivalence between (the first inclusion in) (92) and (93) has been shown in the proof of Lemma 5.2, see the argumentation around (151). The proposition's remaining assertions follow by a similar adaptation of the representer theorem [51]. Indeed: By the definitions of $\phi_{\mu|\lambda}^{(m)}$ and $\hat{\mu}_{3}$, we have

$$\mathfrak{H}_{\Xi}^{[m]} \ni g \equiv (g_w)_{w \in [\tilde{d}]_{\leq m}^*} \longmapsto \phi_{\hat{\mu}_{\mathfrak{Z}}|\lambda}^{(m)}(g) = \frac{1}{N} \sum_{j=1}^N \left\| \pi_{[m]} \big(\tilde{q}(Y^{(j)}) \big) - g(X^{(j)}) \right\|_{\mathcal{H}_{\mathcal{Y}}}^2 + \lambda \|g\|_{\mathfrak{H}_{\Xi}}^2 = \sum_{k=1}^{m_{\tilde{d}}} \phi_k(g_w)$$
(153)

for each $m \in \overline{\mathbb{N}}$, where $\tilde{q} \coloneqq \underline{\mathfrak{sig}}_{\Lambda}$ and the scalar functions $(\phi_k)_{k \in \mathbb{N}} \subset \mathcal{H}_{\kappa}$ are defined as

$$\phi_k(h) \coloneqq \frac{1}{N} \sum_{j=1}^N \left| \langle \tilde{q}(Y^{(j)}), \eta^{-1}(k) \rangle - h(X^{(j)}) \right|^2 + \lambda ||h||_{\kappa}^2.$$

Note that the second identity in (153) is due to (the definition of $\|\cdot\|_{\mathcal{H}_{\mathcal{Y}}}$ in (42) and) Proposition 5.1, whose first assertion provides the unitary isomorphism $\mathfrak{H}_{\Xi} \ni g \mapsto (\langle g, w \rangle \mid w \in [\tilde{d}]^*) \in \bigoplus_{w \in [\tilde{d}]^*} \mathcal{H}_{\kappa}$ (so that $\|(g_w)\|_{\mathfrak{H}_{\Xi}}^2 = \sum_{w \in [\tilde{d}]^*} \langle g_w, g_w \rangle_{\kappa}^2 = \sum_{w \in [\tilde{d}]^*} \|g_w\|_{\kappa}^2$), cf. also [7, Example 3.2 (iii)]. Therefore,

$$\arg\min_{g\in\mathfrak{H}_{\Xi}^{[m]}}\phi_{\mu_{\mathfrak{Z}}|\lambda}^{(m)}(g) = \arg\min_{(g_w)\in(\mathcal{H}_{\kappa})^{\times m_{\tilde{d}}}} \sum_{k=1}^{m_{\tilde{d}}}\phi_k(g_{\eta^{-1}(k)}) = \left\{\sum_{k=1}^{m_{\tilde{d}}}g_k^*\cdot\eta^{-1}(k) \mid g_k^*\in \argmin_{h\in\mathcal{H}_{\kappa}}\phi_k(h)\right\}.$$
 (154)

Now from the classical representer theorem [51] (cf. also the proof of Lemma 5.2) we know that

$$g_k^* \in \underset{h \in \mathcal{H}_{\kappa}}{\operatorname{arg\,min}} \phi_k(h) \quad \text{only if} \quad g_k^* \in \operatorname{span}\left(\kappa(X^{(j)}, \cdot) \mid j \in [N]\right),$$

which together with (154) (and the definition of κ , see Prop. 5.1) proves (91). In fact, as detailed in the proof of Lemma 5.2, for instance, we have that $g_k^* = \arg \min_{h \in \mathcal{H}_{\kappa}} \phi_k(h)$ if and only if

$$g_k^* = \sum_{j=1}^N \alpha_j^{(k)} \langle q_j, q \rangle \quad \text{for} \quad \tilde{\alpha}_{(k)} \equiv \left(\alpha_j^{(k)}\right) \in \mathbb{R}^N \quad \text{such that} \quad \boldsymbol{C}_{\mathfrak{Z},\lambda} \tilde{\alpha}_{(k)} = \boldsymbol{A}_{\mathfrak{Z}} \boldsymbol{b}_{\mathfrak{Z}}^{(k)}, \tag{155}$$

where $q := \underline{\mathfrak{sig}}_{\Xi}$ and $q_j := q(X^{(j)})$ and $C_{\mathfrak{Z},\lambda} := (A_\mathfrak{Z}^\top A_\mathfrak{Z} + N\lambda A_\mathfrak{Z})$ for $A_\mathfrak{Z}$ as in (96) and $b_\mathfrak{Z}^{(k)} := (\langle \tilde{q}(Y^{(j)}), \eta^{-1}(k) \rangle)_{j \in [N]}$; cf. (151). Now by the definitions (90), (86), (97) and uniqueness (Lem. 5.2),

$$\hat{\Upsilon}_{N,\lambda,m}^{\text{lex}} = \sum_{k=1}^{m_{\tilde{d}}} g_k^* \cdot e_k \stackrel{(155)}{=} \sum_{k=1}^{m_{\tilde{d}}} \left[\sum_{j=1}^N \alpha_j^{(k)} q_j^* \right] \cdot e_k = \left(\tilde{\alpha}_{(1)}^\mathsf{T} \cdot \underline{q}^* \middle| \cdots \middle| \tilde{\alpha}_{(m_{\tilde{d}})}^\mathsf{T} \cdot \underline{q}^* \right)^\mathsf{T} = \tilde{A}_*^\mathsf{T} \cdot \underline{q}^*$$

for the matrix $\tilde{A}_* \equiv \left(\tilde{\alpha}_{(1)} | \cdots | \tilde{\alpha}_{(m_{\tilde{d}})}\right) \in \mathbb{R}^{N \times m_{\tilde{d}}}$ and the function $\underline{q}^* \coloneqq \sum_{j=1}^N q_j^* \cdot \tilde{e}_j : \mathcal{X} \to \mathbb{R}^N$, where $(e_k)_{k \in [m_{\tilde{d}}]}$ and $(e_j)_{j \in [N]}$ are the standard bases of $\mathbb{R}^{m_{\tilde{d}}}$ and \mathbb{R}^N , respectively; this completes the proof of (92). Noting that $B_3 = (b_3^{(1)} | \cdots | b_3^{(m_{\tilde{d}})})$, equation (94) is clear from (155).

B.3.4 Proof of Lemma 5.5

Proof. The statements (i) and (iii) both follow from essentially the same simple observations: By the fact that \mathcal{H}_{κ} is dense in $L^2(\mathbb{P}_Y)$ (Proposition 5.1) and from the definition (84), we for each $\epsilon > 0$ with $\vartheta_f(\epsilon) \neq \infty$ can find $\ell_{\epsilon} \in \mathcal{H}_{\mathcal{Y}}$ such that $\phi_{\mathbb{P}_Y}^f(\ell_{\epsilon}) \leq \epsilon^2/2$ and $\|\langle \ell_{\epsilon}, \underline{\mathfrak{sig}}_{\Lambda} \rangle\|_{\kappa}^2 \leq \epsilon^2/2\vartheta_f(\epsilon)$. Hence for each $\lambda > 0$, we have by definition of ℓ_{λ}^f and the property $\phi_{\mathbb{P}_Y}^f(\ell_{\epsilon}) \leq \epsilon^2/2$ that

$$r_{f}^{I}(\lambda)^{2} = \phi_{\mathbb{P}_{Y}}^{f}(\ell_{\lambda}^{f}) \leq \phi_{\mathbb{P}_{Y}|\lambda}^{f}(\ell_{\lambda}^{f}) = \phi_{\mathbb{P}_{Y}|\lambda}^{f}(\ell_{\lambda}^{f}) - \phi_{\mathbb{P}_{Y}}^{f}(\ell_{\epsilon}) + \phi_{\mathbb{P}_{Y}}^{f}(\ell_{\epsilon}) \\ \leq \phi_{\mathbb{P}_{Y}|\lambda}^{f}(\ell_{\epsilon}) - \phi_{\mathbb{P}_{Y}}^{f}(\ell_{\epsilon}) + \phi_{\mathbb{P}_{Y}}^{f}(\ell_{\epsilon}) \leq \lambda \left\| \langle \ell_{\epsilon}, \underline{\mathfrak{sig}}_{\Lambda} \rangle \right\|_{\kappa}^{2} + \frac{\epsilon^{2}}{2},$$

$$(156)$$

which (by the above choice of ℓ_{ϵ}) implies (i) for the case $\vartheta_f(\epsilon) \neq \infty$. If $\vartheta_f(\epsilon) = \infty$, then for each $\lambda > 0$ there will be some $\ell_{\epsilon}^{\lambda} \in \mathcal{H}_{\mathcal{Y}}$ with $\phi_{\mathbb{P}_{Y}}^{f}(\ell_{\epsilon}^{\lambda}) \leq \epsilon^{2}/2$ and $\|\langle \ell_{\epsilon}, \underline{\mathfrak{sig}}_{\Lambda} \rangle\|_{\kappa}^{2} \leq \epsilon^{2}/2\lambda$, which, by the fact that the inequalities (156) are valid for all $(\lambda, \ell_{\epsilon}) \in \mathbb{R}_{>0} \times (\phi_{\mathbb{P}_{Y}}^{f})^{-1}(0, \epsilon^{2}/2]$, implies (i) also for the case $\vartheta_f(\epsilon) = \infty$. Statement (ii) follows analogously to (i) upon noting that, by statement (ii),

$$R_{m}^{\mathrm{I}}(\lambda)^{2} = \phi_{\mathbb{P}_{(X,Y)}}^{(m)} \left(\pi_{[m]}(\Upsilon_{\lambda}^{*}) \right) - \phi_{\mathbb{P}_{(X,Y)}}^{(m)} \left(\pi_{[m]}(\varphi_{\star}) \right)$$

$$\leq \phi_{\mathbb{P}_{(X,Y)}|\lambda}^{(m)}(g) - \phi_{\mathbb{P}_{(X,Y)}}^{(m)}(g) + \Delta_{g} \leq \lambda \|g\|_{\mathfrak{H}_{\Xi}}^{2} + \Delta_{g}$$
(157)

for the difference $\Delta_g \coloneqq \phi_{\mathbb{P}_{(X,Y)}}^{(m)}(g) - \phi_{\mathbb{P}_{(X,Y)}}^{(m)}(\pi_{[m]}(\varphi_\star))$ and any $\lambda > 0$ and any $g \in L^2(\mathbb{P}_X; \mathcal{H}_{\mathcal{Y}}^{[m]})$; then exactly as shown for (i) above, the density result [following (81)] of Proposition 5.1 allows us to, for any given $\epsilon > 0$ and any $0 < \lambda \leq \tilde{\vartheta}_{\pi_{[m]}(\varphi_\star)}(\epsilon)$, find a $g \in L^2(\mathbb{P}_X; \mathcal{H}_{\mathcal{Y}}^{[m]})$ such that $\lambda ||g||_{\mathfrak{H}_{\Xi}}^2$ and Δ_g are both bounded by $\epsilon^2/2$ (the case $\tilde{\vartheta}_{\pi_{[m]}(\varphi_\star)}(\epsilon) = \infty$ is handled as before), implying (iii).

For the proof of (157), that is of statement (ii), abbreviate $\varphi_{\star}^{[m]} \coloneqq \pi_{[m]}(\varphi_{\star})$ and note that with

$$\varphi_{\star}^{[m]}(X) = \pi_{[m]} \left(\mathbb{E}[\mathbb{Y}^{\Lambda} \mid X] \right) = \mathbb{E}[\pi_{[m]}(\mathbb{Y}^{\Lambda}) \mid X] = \varphi_{m}^{\star}(X) \quad \text{for} \quad \varphi_{m}^{\star} \coloneqq \underset{g \in L^{2}(\mathbb{P}_{X};\mathcal{H}_{\mathcal{Y}}^{[m]})}{\operatorname{small{smaller}}} \phi_{\mathbb{P}_{(X,Y)}}^{(m)}(g)$$

(where the second of the above equalities holds by the commuting property (75) and the linearity of conditional expectations, and the third equality holds by the variational characterisation (cf. (141)) of conditional expectations) we have $\varphi_{\star}^{[m]} = \varphi_{\star}^{\star}$ in $L^2(\mathbb{P}_X; \mathcal{H}_{\mathcal{V}}^{[m]})$ and hence, as claimed, obtain that

$$R_m^{\mathrm{I}}(\lambda)^2 = \mathbb{E}\Big[\|\pi_{[m]}(\Upsilon^*_{\lambda}) - \varphi^*_m(X)\|_{\mathcal{H}_{\mathcal{Y}}}^2 \Big] = \phi^{(m)}_{\mathbb{P}_{(X,Y)}} \big(\pi_{[m]}(\Upsilon^*_{\lambda})\big) - \phi^{(m)}_{\mathbb{P}_{(X,Y)}} \big(\varphi^{[m]}_{\star}\big)$$

for each $\lambda > 0$, by the classical risk decomposition from [16, Proposition 1].

As to statement (iv), denote $\tilde{q}_l := \underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)})$ and note that, by the definitions (98) and (42),

$$\left\|\pi_{\nu}(\hat{\ell}_{n,\lambda,f}^{\star})\right\|_{\nu} \leq \sum_{l=1}^{n} \left|\alpha_{l}^{\star}\right| \|\pi_{\nu}(\tilde{q}_{l})\|_{\nu} \leq \|\alpha^{\star}\|_{2}\gamma_{\nu} \quad \text{for} \quad \gamma_{\nu} \coloneqq \sqrt{\sum_{l=1}^{n} \|\pi_{\nu}(\tilde{q}_{l})\|_{\nu}^{2}} \qquad \left(\nu \in \mathbb{N}\right).$$

Hence by the definitions of $r_{f,\lambda,n}^{\text{III}}(m)$ and $\pi_{[m]}$ and (again) with definition (42), for each $m \in \mathbb{N}_0$,

$$r_{f,\lambda,n}^{\mathrm{III}}(m) = \sqrt{\sum_{\nu=m+1}^{\infty} \left\| \pi_{\nu}(\hat{\ell}_{n,\lambda,f}^{\star}) \right\|_{\nu}^{2}} \le \|\alpha^{\star}\|_{2}\tilde{\beta}_{m} \eqqcolon \beta_{m}^{(\mathrm{cut})} \quad \text{for} \quad \tilde{\beta}_{m} \coloneqq \sqrt{\sum_{\nu=m+1}^{\infty} \gamma_{\nu}^{2}}.$$
 (158)

Now since $\tilde{q}_1, \ldots, \tilde{q}_n \in \mathcal{H}_{\mathcal{Y}}$ and thus $\sum_{\nu=0}^{\infty} \gamma_{\nu}^2 = \sum_{l=1}^n \|\tilde{q}_l\|_{\mathcal{H}_{\mathcal{Y}}} < \infty$, we find that $(\beta_m^{(\text{cut})})_{m=1}^{\infty}$ is a monotone decreasing null sequence which, due to $\alpha^* = (\mathbf{A}_{\mathfrak{Y}}^\top \mathbf{A}_{\mathfrak{Y}} + n\lambda \mathbf{A}_{\mathfrak{Y}})^+ \mathbf{A}_{\mathfrak{Y}} b_{\mathfrak{Y},f}$, depends on $(\mathfrak{Y}, \Lambda, \lambda, f)$ only.

For (v), let $\delta > 0$ and observe that by [42, Proposition 3.2], (102) then holds for the sequences

$$\beta_n^{(\text{out})} \coloneqq \frac{c_\Lambda}{\lambda\sqrt{n\delta}} \quad \text{and} \quad \beta_N^{(\text{in})} \coloneqq \sqrt{\frac{c_{\Xi}c_\Lambda}{\lambda^2 N\delta}} \quad (n, N \in \mathbb{N})$$

featuring the constants $c_Z := \sup_{t \in \mathcal{H}_{\xi}} ||Z(t)||^2_{\mathcal{H}_{\xi}} < \infty$ for $(Z, \xi) \in \{(\Lambda, \mathcal{Y}), (\Xi, \mathcal{X})\}$ (recall Definition 3.18 for their finiteness); recall that the general applicability of [42] is due to Proposition 5.1.

B.3.5 Algorithms I and II are Instances of Definition 2.6

Remark B.7. An approximation of the regression operator (14), the estimator (107) can be written

$$\hat{T}_{N}^{\mathrm{I}}(f,\mathfrak{Y};\mathfrak{Z},\underline{\lambda},n,m) \coloneqq \hat{\mathscr{R}}_{\underline{\lambda},f}^{\mathrm{I}}[n,N,m] = \left\langle \hat{\ell}_{\lambda_{1};\mathfrak{Y}}^{f}, \, \hat{\Xi}_{\lambda_{2},n,m;\mathfrak{Z}}(\cdot) \right\rangle_{2,k} \tag{159}$$

—that is: in the form of Definition 2.6, eq. (18), and (24)—for the bilinear map

$$\langle \cdot, \cdot \rangle_{2,k} : \ell^2(\mathbb{N}_0)^{\times k} \times \ell^2(\mathbb{N}) \longrightarrow \mathbb{R}^k, \ \left((\mathfrak{q}_i), \mathfrak{q}\right) \mapsto \left(\langle \mathfrak{q}_i, \mathfrak{q} \rangle_{\ell^2}\right)_{i=1,\dots,k}$$

and the constituent component representations

$$\hat{\boldsymbol{\ell}}_{\lambda_1;\mathfrak{Y}}^f \coloneqq \left(\sum_{l=1}^n \left[\boldsymbol{\alpha}_{\lambda_1}^f \right]_{li} \eta(\underline{\mathfrak{sig}}_{\Lambda}(Y^{(l)})) \right)_{i=1,\dots,k} \quad \text{and} \quad \hat{\boldsymbol{\Xi}}_{\lambda_2,n,m;\mathfrak{Z}} \coloneqq \hat{\Upsilon}_{N,\lambda_2,m}^{\text{lex}} : \mathcal{X} \to \mathbb{R}^{m_{\tilde{d}}},$$

where $\eta : \mathcal{H}_{\mathcal{Y}} \to \ell^2(\mathbb{N}_0), \ \left(\boldsymbol{t}_w\right)_{w \in [\tilde{d}]^*} \mapsto \left(\boldsymbol{t}_{\eta^{-1}(j)}\right)_{j \in \mathbb{N}_0}$ extends the shortlex ordering used in (97). Note that the estimator (159) is linear in both $\hat{\ell}^f_{\lambda_1;\mathfrak{Y}} \in \ell^2(\mathbb{N}_0)^{\times k}$ and $f \in \mathcal{L}^2(\mathbb{P}_Y)$.

Remark B.8. An approximation of the regression operator (14), the estimator (115) can be written

$$\hat{T}_{M}^{\mathrm{II}}(f;\mathfrak{W}_{f},\lambda) \coloneqq \hat{\mathscr{R}}_{\lambda,f(Y)}^{\mathrm{II}}[M] = \left\langle \hat{\varrho}_{\lambda;\mathfrak{W}_{f}}, \underline{\mathfrak{sig}}_{\Xi}(\cdot) \right\rangle_{\mathcal{H}_{\mathcal{X}},k}$$
(160)

for the data $\mathfrak{W}_f \coloneqq \{(X^{(j)}, f(Y^{(j)})) \mid j \in [M]\}$, which are iid copies of (X, f(Y)), the bilinear map

$$\langle \cdot, \cdot \rangle_{\mathcal{H}_{\mathcal{X}}, k} : \mathcal{H}_{\mathcal{X}}^{\times k} \times \mathcal{H}_{\mathcal{X}} \longrightarrow \mathbb{R}^{k}, \ \left((\mathfrak{q}_{i}), \mathfrak{q}\right) \mapsto \left(\langle \mathfrak{q}_{i}, \mathfrak{q} \rangle_{\mathcal{H}_{\mathcal{X}}}\right)_{i=1, \dots, k}$$

and the constituent \mathfrak{W}_{f} -processing data representation, with $\hat{A}_{*} = (\hat{\alpha}_{ij}^{(\lambda)})$ solving $(114)|_{Z=f(Y)}$,

$$\hat{\boldsymbol{\varrho}}_{\boldsymbol{\lambda};\mathfrak{W}_{f}} \coloneqq \left(\sum_{j=1}^{M} \left[\hat{A}_{*} \right]_{ij} \underbrace{\mathfrak{sig}}_{\Xi}(\boldsymbol{X}^{(j)}) \right)_{i=1,\ldots,k}$$

Note that while the estimator (160) is also linear¹¹ in $f \in \mathcal{L}^2(\mathbb{P}_Y)$ and $\hat{\boldsymbol{\varrho}}_{\lambda;\mathfrak{W}_f} \in \mathcal{H}_{\mathcal{X}}^{\times k}$, unlike the more modular estimator (160) it does not us a factorisation into two separate and possibly independent types of data (i.e., \mathfrak{Z} and $\mathfrak{Y}_f := \{(Y_l, f(Y_l)) \mid Y_l \in \mathfrak{Y}\})$ but instead operates directly on the combined dataset \mathfrak{W}_f , where f(Y) is observed exclusively in association with X.

 $[\]overline{\hat{T}_{M}^{\mathrm{II}} \text{ In the sense that: } \hat{T}_{M}^{\mathrm{II}}(f_{1}+cf_{2};\mathfrak{W}_{f_{1}+cf_{2}},\lambda)} = \hat{T}_{M}^{\mathrm{II}}(f_{1};\mathfrak{W}_{f_{1}},\lambda) + c\hat{T}_{M}^{\mathrm{II}}(f_{2};\mathfrak{W}_{f_{2}},\lambda) \text{ for any } f_{1},f_{2} \in \mathcal{L}^{2}(\mathbb{P}_{Y}), c \in \mathbb{R}.$

B.3.6 Proof of Theorem 5.9

Proof. Denote $\psi := \mu_{\cdot}(f)$, so that $\psi \equiv (\psi_1, \cdots, \psi_k) \in L^2(\mathbb{P}_X; \mathbb{R}^k)$. Showing $(121)|_{\nu=I}$ first, fix any $(\varepsilon, \delta) \in \mathbb{R}^2_{>0}$ and take an arbitrary $(\underline{\lambda}, n, N, m) \in \mathcal{A}^{\mathrm{I}}_f(\varepsilon, \delta)$. Recalling from (159) that

$$\hat{\mathscr{R}}^{\mathrm{I}}_{\underline{\lambda},f}[n,N,m] = \sum_{i=1}^{k} \left\langle \hat{\ell}^{i}_{\lambda_{1},n,m}, \hat{\Upsilon}^{*}_{N,\lambda_{2},m} \right\rangle e_{i}$$
(161)

with $\hat{\Upsilon}^*_{N,\lambda_2,m}$ as in (90) and with $\hat{\ell}^i_{\lambda_1,n,m} \coloneqq \pi_{[m]}(\hat{\ell}^*_{n,\lambda_1,f_i})$ for $\hat{\ell}^*_{n,\lambda_1,f_i}$ as in (98), we find that

$$\mathbb{P}(\|\psi - \hat{\mathscr{R}}_{f,\underline{\lambda}}^{\mathrm{I}}[n,N,m]\| \ge \varepsilon) \le \mathbb{P}(E_1 \ge 3\varepsilon/5) + \mathbb{P}(E_2 \ge 2\varepsilon/5)$$
(162)

(from (161) and by the triangle inequality) with $\|\cdot\| \coloneqq \|\cdot\|_{L^2(\mathbb{P}_X;\mathbb{R}^k)}$ and for the differences

$$E_1 \coloneqq \left\| \psi - \hat{\psi}_{\lambda_1, n, m} \right\| \quad \text{and} \quad E_2 \coloneqq \left\| \sum_{i=1}^k \langle \hat{\ell}^i_{\lambda_1, n, m}, \varphi_\star - \hat{\Upsilon}^*_{N, \lambda_2, m} \rangle e_i \right\|$$

where $\hat{\psi}_{\lambda_1,n,m} \coloneqq \sum_{i=1}^k \langle \hat{\ell}^i_{\lambda_1,n,m}, \varphi_\star \rangle e_i$ and for φ_\star as in (101). Using (75) and the inequality in (74),

$$E_{1} \leq \sum_{i=1}^{k} \tilde{E}_{1,i} \quad \text{for} \quad \tilde{E}_{1,i} \coloneqq \left\| f_{i} - \left\langle \hat{\ell}_{\lambda_{1},n,m}^{i}, \underline{\mathfrak{sig}}_{\Lambda} \right\rangle \right\|_{L^{2}(\mathbb{P}_{Y})}$$

Recalling the definitions in (99) and (100), we obtain, for each $i \in [k]$, that

$$\tilde{E}_{1,i} \leq r_{f_i}^{\mathrm{I}}(\lambda_1) + r_{f_i,\lambda_1}^{\mathrm{II}}(n) + r_{f_i,\lambda_1,n}^{\mathrm{III}}(m) \leq \frac{\varepsilon}{5k} + r_{f_i,\lambda_1}^{\mathrm{II}}(n) + \frac{\varepsilon}{5k},$$

where we used the error bound max $\{r_{f_i}^{\mathrm{I}}(\lambda_1), r_{f_i,\lambda_1,n}^{\mathrm{III}}(m)\} \leq \varepsilon/(5k)$ which is due to Lemma 5.5 (i) and (iv) in conjunction with the $\mathcal{B}_f(\varepsilon, \delta)$ -defining constraints defined in (104) and (105). Hence,

$$\mathbb{P}(E_1 \ge 3\varepsilon/5) \le \mathbb{P}\left(\sum_{i=1}^k r_{f_i,\lambda_1}^{\mathrm{II}}(n) \ge \varepsilon/5\right) \le \sum_{i=1}^k \mathbb{P}\left(r_{f_i,\lambda_1}^{\mathrm{II}}(n) \ge \varepsilon/(5k)\right) \le \frac{\delta}{2}, \quad (163)$$

where the last inequality is due to Lemma 5.5 (v) via the constraint $n \ge n_{f,\lambda_1}(\varepsilon, \delta)$.

For a bound on $\mathbb{P}(E_2 \ge 2\varepsilon/5)$, note that, by Cauchy-Schwarz and the triangle inequality,

$$E_2 \leq c_{\lambda_1,m;\mathfrak{Y}} \|\pi_{[m]}(\varphi_{\star}) - \hat{\Upsilon}^*_{N,\lambda_2,m}\|_{L^2(\mathbb{P}_X;\mathcal{H}_{\mathcal{Y}})} \leq c_{\lambda_1,m;\mathfrak{Y}} \left(R^{\mathrm{I}}_m(\lambda_2) + R^{\mathrm{II}}_{m,\lambda_2}(N)\right)$$
(164)

with $c_{\lambda_1,m;\mathfrak{Y}}$ as defined in (106) and for the errors R_m^{I} and $R_{m,\lambda}^{\mathrm{II}}$ introduced in (101). Lemma 5.5 (iii) and the bound $\lambda_2 \leq \lambda_{\lambda_1,m}^{\mathrm{II}}(\varepsilon;\mathfrak{Y})$ (cf. (106)) ensure $c_{\lambda_1,m;\mathfrak{Y}}R_m^{\mathrm{I}}(\lambda_2) \leq \varepsilon/5$. Consequently,

$$\mathbb{P}(E_2 \ge 2\varepsilon/5) \le \mathbb{P}\left(R_{m,\lambda_2}^{\mathrm{II}}(N) \ge \varepsilon/(5c_{\lambda_1,m};\mathfrak{Y})\right) \le \frac{\delta}{2}$$
(165)

follows by combining (164) and Lemma 5.5 (v), using that $N \ge N_{\lambda_1,\lambda_2,m}(\varepsilon,\delta;\mathfrak{Y})$ (cf. (106)).

Applying the inequalities (163) and (165) to (162) proves the case $\nu = I$ in (121), as desired.

The proof of $(121)|_{\nu=II}$ is very similar but more concise, using Corollary 5.8 applied to $Z \coloneqq f(Y)$ instead of Lemma 5.5: For any fixed $(\lambda, N) \in \mathcal{A}_f(\varepsilon, \delta)$, we know that (since $\psi = \psi_{X,Z}$ by uniqueness)

$$\mathbb{P}(\left\|\psi - \hat{\Upsilon}^*_{f(Y),\lambda;N}\right\| \ge \varepsilon) \le \mathbb{P}(\rho_Z^{\mathrm{I}}(\lambda) \ge \varepsilon/2) + \mathbb{P}(\rho_{Z,\lambda}^{\mathrm{II}}(N) \ge \varepsilon/2) \le 0 + \delta$$

by (the triangle inequality and) Corollary 5.8 (ii) and (iii); this proves the case $\nu = \text{II}$ in (121).

Proving (122), recall that $\hat{\mathscr{R}}^{\mathrm{I}}_{\theta} \coloneqq \hat{\mathscr{R}}^{\mathrm{I}}_{f,\underline{\lambda}}[n,N,m]$ for $\theta = (\underline{\lambda}, n, N, m)$, and $\hat{\mathscr{R}}^{\mathrm{II}}_{\theta} \coloneqq \hat{\mathscr{R}}^{\mathrm{II}}_{\lambda,Z}[M]$ for $\theta = (\lambda, N)$. Then for both $\nu = \mathrm{I}, \mathrm{II}$, the object $\hat{\mathscr{R}}^{\nu}_{\theta}$ is an $L^2(\mathbb{P}_X; \mathbb{R}^k)$ -valued random variable:

$$\hat{\mathscr{R}}^{\nu}_{\theta} = \hat{\mathscr{R}}_{\theta}(\,\cdot\,;\tilde{Z}_{\nu})\,:\,\Omega\longrightarrow L^{2}(\mathbb{P}_{X};\mathbb{R}^{k}),\,\,\omega\mapsto\hat{\mathscr{R}}_{\theta}\big(\,\cdot\,;\tilde{Z}_{\nu}(\omega)\big),$$

for $\tilde{Z}_{\rm I}$ resp. $\tilde{Z}_{\rm II}$ the random vectors from which the data in (89) resp. (112) is sampled, that is

$$\tilde{Z}_{\mathrm{I}} \coloneqq \left(Y^{(1)}, \cdots, Y^{(n)}, (X^{(1)}, Y^{(1)}), \cdots, (X^{(N)}, Y^{(N)})\right) : \Omega \longrightarrow \tilde{\mathcal{Z}}^{\mathrm{I}} \coloneqq \mathcal{Y}^{n} \times \left(\mathcal{X} \times \mathcal{Y}\right)^{N}$$

and $\tilde{Z}_{\mathrm{II}} \coloneqq \left((X^{(1)}, Z_{(1)}), \cdots, (X^{(M)}, Z_{(M)})\right) : \Omega \longrightarrow \tilde{\mathcal{Z}}^{\mathrm{II}} \coloneqq \left(\mathcal{X} \times \mathbb{R}^{k}\right)^{M}.$

Define $A_{\theta,\varepsilon}^{\nu} \coloneqq \left\{ |\mathbb{E}[f(Y) \mid X] - \hat{\mathscr{R}}_{\theta}^{\nu}(X)| \ge \varepsilon \right\} \equiv \left\{ \omega \in \Omega \mid |\mathbb{E}[f(Y) \mid X](\omega) - \hat{\mathscr{R}}_{\theta}^{\nu}(X(\omega); \tilde{Z}_{\nu}(\omega))| \ge \varepsilon \right\}, A_{\theta,\varepsilon}^{\times|\nu} \coloneqq \left\{ (x,z) \in \mathcal{X} \times \tilde{Z}^{\nu} \mid |\psi(x) - \hat{\mathscr{R}}_{\theta}^{\nu}(x;z)| \ge \varepsilon \right\} \text{ and } \Delta_{\theta|\nu}^{z} \coloneqq \left| \psi - \hat{\mathscr{R}}_{\theta}^{\nu}(\cdot;z) \right| \text{ for any fixed } z \in \tilde{Z}^{\nu}.$ With the z-sections $A_{\theta,\varepsilon}^{z|\nu} \coloneqq \left\{ x \in \mathcal{X} \mid (x,z) \in A_{\theta,\varepsilon}^{\times} \right\} (z \in \tilde{Z}^{\nu})$ for $\nu = \mathrm{I}, \mathrm{II}$, we have by that

$$\mathbb{P}(A^{\nu}_{\theta,\varepsilon}) = \int_{\mathcal{X}\times\tilde{\mathcal{Z}}_{\nu}} \mathbb{1}_{A^{\times|\nu}_{\theta,\varepsilon}}(x,z) \,\mathbb{P}_{(X,\tilde{Z}_{\nu})}(\mathrm{d}x,\mathrm{d}z) = \int_{\tilde{\mathcal{Z}}_{\nu}} \mathbb{P}_{X}(A^{z|\nu}_{\theta,\varepsilon}) \,\mathbb{P}_{\tilde{Z}_{\nu}}(\mathrm{d}z),\tag{166}$$

which holds by Fubini and the (89)- resp. (112)-underlying assumption that the random variables X and \tilde{Z}_{I} and \tilde{Z}_{II} are statistically independent. Now for any $(\tilde{\varepsilon}, \tilde{\delta}) \in \mathbb{R}^{2}_{>0}$, the inequalities (121) yield exceptional events $C^{\nu}_{\tilde{\varepsilon},\tilde{\delta}} \in \mathcal{B}(\tilde{Z}^{\nu})$ of measure $\mathbb{P}_{\tilde{Z}_{\nu}}(C^{\nu}_{\tilde{\varepsilon},\tilde{\delta}}) \leq \tilde{\delta}$ such that, for each $\theta \in \mathcal{A}^{\nu}_{f}(\tilde{\varepsilon}, \tilde{\delta})$,

$$\mathbb{P}_{X}(A_{\theta,\varepsilon}^{z|\nu}) \leq \varepsilon^{-2} \left\| \mathbb{1}_{A_{\theta,\varepsilon}^{z|\nu}} \cdot \Delta_{\theta|\nu}^{z} \right\|_{L^{2}(\mathbb{P}_{X};\mathbb{R}^{k})}^{2} \leq \varepsilon^{-2} \|\Delta_{\theta|\nu}^{z}\|_{L^{2}(\mathbb{P}_{X};\mathbb{R}^{k})}^{2} < \frac{\tilde{\varepsilon}}{\varepsilon^{2}}, \quad \text{for all} \quad z \in \left(\mathcal{C}_{\tilde{\varepsilon},\tilde{\delta}}\right)^{c}; \quad (167)$$

this holds for both $\nu = I, II$. The combination of (167) with (166) thus yields that

$$\mathbb{P}(A_{\theta,\varepsilon}^{\nu}) \leq \int_{\left(\mathcal{C}_{\varepsilon,\tilde{\delta}}\right)^{c}} \mathbb{P}_{X}(A_{\theta,\varepsilon}^{z|\nu}) \mathbb{P}_{\tilde{Z}_{\nu}}(\mathrm{d}z) + \tilde{\delta} < \frac{\tilde{\varepsilon}}{\varepsilon^{2}}(1-\tilde{\delta}) + \tilde{\delta} \leq (1-q) \eqqcolon q' \qquad (\nu = \mathrm{I},\mathrm{II})$$

for each $\theta \in \Theta_f^{\nu}(\varepsilon, \delta) \equiv \bigcup_{(\tilde{\varepsilon}, \tilde{\delta}) : \tilde{\varepsilon}(1-\tilde{\delta})+\tilde{\delta}\varepsilon^2 \leq q'\varepsilon^2} \mathcal{A}_f^{\nu}(\tilde{\varepsilon}, \tilde{\delta})$. Since $\mathbb{P}\left(\left|\mathbb{E}[f(Y) \mid X] - \hat{\mathscr{R}}_{\theta}^{\nu}(X)\right| < \varepsilon\right) = 1 - \mathbb{P}(A_{\theta,\varepsilon}^{\nu})$ for $\nu = \mathrm{I}, \mathrm{II}$, the assertion (122) follows.

B.4 'Stabilized' Estimation of Conditional Expectations (Over L²-Balls)

To allow for a statistically consistent sample-based approximation of the optimisation scheme (66), we adopt the proof of Theorem 4.8 for a 'stabilized' version of (65).

Lemma B.9. Given any $R \ge 0$, consider the set defined by

$$\mathfrak{C}_R \coloneqq \Big\{ \mathbb{Z} \in L^2_X(\mathcal{H}_{\mathcal{Y}}) \ \Big| \ \|\mathbb{Z}\|_{\mathcal{L}^{\infty}} \coloneqq \operatorname{ess\,sup}_{\omega \in \Omega} \|\mathbb{Z}(\omega)\| \le R \Big\}.$$

Then the following hold.

- (i) The set \mathfrak{C}_R is a convex and complete subset of $L^2_X(\mathcal{H}_{\mathcal{Y}})$.
- (ii) Denoting $\Psi_R := \{ \psi \cdot \mathbb{1}_{\{ \|\psi\| \le R\}} \mid \psi \in \Psi \}$ (cf. (58)) for $R \ge 0$, then each element of \mathfrak{C}_R is an $\| \cdot \|_{L^2(\mathcal{H}_{\mathcal{Y}})}$ -limit point of the set $\Psi_{2R}(X) := \{ \tilde{\psi}(X) \mid \tilde{\psi} \in \Psi_{2R} \} \subset L^2_X(\mathcal{H}_{\mathcal{Y}}).$
- (iii) For any $R \ge 0$ and with \mathfrak{D}_X the spatial support of X (see Rem. B.3 (iii)), we have that

$$\mathfrak{C}_R = \{\phi(X) \mid \phi \in \mathcal{B}_R\} \quad \text{with} \quad \mathcal{B}_R \coloneqq \{\phi : \mathfrak{D}_X \to \mathcal{H}_{\mathcal{Y}} \mid \phi \ \text{Borel-measurable} : \|\phi\|_{\infty;\mathfrak{D}_X} \le R\},$$

which holds as an identity in $2^{L^2(\mathcal{H}_{\mathcal{Y}})}$, that is up to element-wise inequality on a \mathbb{P} -nullset.

(iv) We have that $\mathbb{E}[\mathbb{Y}^{\Lambda}|X] \in \mathfrak{C}_{R_{\Lambda}}$ for $R_{\Lambda} \coloneqq \sup\{\|\Lambda(t)\| \mid t \in \mathcal{H}_{\mathcal{Y}}\} < \infty$.

Proof. Fix any $R \ge 0$. The convexity of \mathfrak{C}_R is clear since $\|\cdot\|_{\mathcal{L}^{\infty}}$ is homogeneous and subadditive. To see that \mathfrak{C}_R is complete, let $(\mathbb{Z}_k) \subset \mathfrak{C}_R$ be a Cauchy sequence in $(\mathfrak{C}_R, \|\cdot\|_{L^2(\mathcal{H}_{\mathcal{Y}})})$. There is then $\mathbb{Z} \in L^2_X(\mathcal{H}_{\mathcal{Y}})$ with $\lim_{k\to\infty} \|\mathbb{Z} - \mathbb{Z}_k\|_{L^2(\mathcal{H}_{\mathcal{Y}})} = 0$, by the fact that $L^2_X(\mathcal{H}_{\mathcal{Y}})$ is complete. To see that in fact $\mathbb{Z} \in \mathfrak{C}_R$, note that for each $\varepsilon > 0$ and with $A_{\varepsilon} := \{\|\mathbb{Z}\| - R \ge \varepsilon\}$,

$$\mathbb{P}(A_{\varepsilon}) = \frac{1}{\varepsilon^2} \int_{A_{\varepsilon}} \varepsilon^2 \, \mathrm{d}\mathbb{P} \le \frac{1}{\varepsilon^2} \int_{A_{\varepsilon}} \|\mathbb{Z} - \mathbb{Z}_k\|^2 \, \mathrm{d}\mathbb{P} \le \frac{1}{\varepsilon^2} \|\mathbb{Z} - \mathbb{Z}_k\|_{L^2(\mathcal{H}_{\mathcal{Y}})}^2 \quad \text{for all} \quad k \in \mathbb{N},$$

hence $\mathbb{P}(A_{\varepsilon}) = 0$ and thus, as $\varepsilon > 0$ was arbitrary, $\mathbb{P}(||Z|| > R) \le \mathbb{P}(\bigcup_{n \in \mathbb{N}} A_{1/n}) \le \sum_{n \in \mathbb{N}} \mathbb{P}(A_{1/n}) = 0$. This implies that $||\mathbb{Z}||_{\mathcal{L}^{\infty}} \le R$ and hence $\mathbb{Z} \in \mathfrak{C}_R$, concluding the proof of (i).

For a proof of statement (ii), choose an arbitrary $\mathbb{Z}_{\star} \in \mathfrak{C}_{R}$ and any $\varepsilon > 0$. Since by Proposition 4.5 there is $\psi \in \Psi$ with $\|\mathbb{Z}_{\star} - \psi(X)\|_{L^{2}(\mathcal{H}_{\mathcal{V}})} \leq \varepsilon$, the choice $\tilde{\psi} := \psi \cdot \mathbb{1}_{\{\|\psi\| \leq 2R\}}$ gives that

$$\begin{aligned} \left\| \mathbb{Z}_{\star} - \tilde{\psi}(X) \right\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}^{2} &= \int_{\mathcal{A}_{R}} \left\| \mathbb{Z}_{\star} - \psi(X) \right\|^{2} \mathrm{d}\mathbb{P} + \int_{\mathcal{A}_{R}^{c}} \left\| \mathbb{Z}_{\star} \right\|^{2} \mathrm{d}\mathbb{P} \\ &\leq \int_{\mathcal{A}_{R}} \left\| \mathbb{Z}_{\star} - \psi(X) \right\|^{2} \mathrm{d}\mathbb{P} + \int_{\mathcal{A}_{R}^{c}} \left\| \mathbb{Z}_{\star} - \psi(X) \right\|^{2} \mathrm{d}\mathbb{P} \\ &= \left\| \mathbb{Z}_{\star} - \psi(X) \right\|_{L^{2}(\mathcal{H}_{\mathcal{Y}})}^{2} \leq \varepsilon^{2}, \quad \text{where} \quad \mathcal{A}_{R} \coloneqq \{ \| \psi(X) \| \leq 2R \}, \end{aligned}$$

since $\mathcal{A}_R^c \subseteq \left\{ \omega \in \Omega \mid \|\mathbb{Z}_{\star}(\omega)\| \leq R = 2R - R < \left\| \|\psi(X(\omega))\| - \|\mathbb{Z}_{\star}(\omega)\| \right\| \leq \|\psi(X(\omega)) - \mathbb{Z}_{\star}(\omega)\| \right\} \cup \mathcal{N}_{\star}$ for the \mathbb{P} -nullset $\mathcal{N}_{\star} \coloneqq \{\|\mathbb{Z}_{\star}\| > R\}$. This proves (ii), as desired.

As to (iii), note that clearly $\operatorname{ev}_X(\mathcal{B}_R) \subseteq \mathfrak{C}_R$, so let us fix any $\mathbb{Z} \in \mathfrak{C}_R$ for the converse inclusion. Then \mathbb{Z} is $(\Sigma_X, \mathcal{B}(\mathcal{H}_Y))$ -measurable by definition of \mathfrak{C}_R , and hence, by Doob-Dynkin, $\mathbb{Z} = \tilde{\phi}_0(X)$ for some $(\mathcal{B}(\mathfrak{D}_X), \mathcal{B}(\mathcal{H}_Y))$ -measurable function $\tilde{\phi}_0 : \mathfrak{D}_X \to \mathcal{H}_Y$; see e.g. [23, Lemma 1.14]. (Note that this $\tilde{\phi}_0$ is unique \mathbb{P}_X -a.e.: if $\mathbb{Z} = \tilde{\phi}_0(X) = \tilde{\phi}_1(X)$ then $\Delta := (\tilde{\phi}_0 - \tilde{\phi}_1) \circ X = 0$ and thus $1 = \mathbb{P}(\Delta^{-1}\{0\}) = \mathbb{P}_X(\tilde{\phi}_0 = \tilde{\phi}_1)$.) Denoting $B_R := \{\mathbf{t} \in \mathcal{H}_Y \mid \|\mathbf{t}\| \leq R\} \in \mathcal{B}(\mathcal{H}_Y)$ and $\hat{\Omega} :=$ $\mathbb{Z}^{-1}(B_R) \in \mathscr{F}$ and $\tilde{\mathbb{Z}} := \mathbb{1}_{\hat{\Omega}} \cdot \mathbb{Z}$, note that $(\mathbb{Z} = \tilde{\mathbb{Z}} \mathbb{P}$ -a.s. by definition of \mathfrak{C} and thus) $\mathbb{Z} = \tilde{\mathbb{Z}}$ in $L^2(\mathbb{P}; \mathcal{H}_Y)$. Hence, and since $\tilde{\mathbb{Z}} = \phi(X)$ for the map $\phi \equiv \phi(x) := \mathbb{1}_{\tilde{\phi}_0^{-1}(B_R)}(x) \cdot \tilde{\phi}_0(x)$ from \mathcal{B}_R , we have the identity $\mathbb{Z} = \phi(X)$ (in $L^2(\mathbb{P}; \mathcal{H}_Y)$). This proves (iii), as claimed.

Statement (iv) follows via the conditional Jensen's inequality, as stated in [57], which gives that

$$\left\| \mathbb{E} \left[\mathbb{Y}^{\Lambda} \, \big| \, X \right] \right\|_{\mathcal{L}^{\infty}} \leq \operatorname{ess\,sup}_{\omega \in \Omega} \mathbb{E} \left[\left\| \mathbb{Y}^{\Lambda} \right\| \, \big| \, X \right] (\omega) \leq R_{\Lambda}$$

(the last inequality holds by (48) and the monotonicity of the operator P_G from (139)); [57, $f = \|\cdot\|$] is applicable since \mathbb{Y}^{Λ} is Bochner-integrable by [its inclusion in (54) and] [58, Proposition 1.16]. \Box

Given $R \ge 0$ and with notation (58), denote $\psi_{\alpha|R} \coloneqq \psi_{\alpha} \cdot \mathbb{1}_{\{\|\psi_{\alpha}\| \le 2R\}}$ for $\alpha \in \mathfrak{L}^2_X$.

Proposition B.10. In the setting of Theorem 4.8, and with $\psi_{\alpha|R} \coloneqq \psi_{\alpha} \cdot \mathbb{1}_{\{||\psi_{\alpha}|| \leq 2R\}}$ for $\alpha \in \mathfrak{L}^2_X$,

$$\mathbb{E}\left[\mathbb{Y}^{\Lambda} \mid X\right] = \underset{\mathbb{Z} \in \mathfrak{C}_{R_{\Lambda}}}{\operatorname{arg\,min}} \Phi(\mathbb{Z}) = \underset{k \to \infty}{\lim} \psi_{\alpha_{k} \mid R_{\Lambda}}(X) \quad in \ L^{2}_{X}(\mathcal{H}_{\mathcal{Y}})$$
(168)

for any minimizing sequence $(\alpha_k) \subset \mathfrak{L}^2_X$ of the (semi-infinite) linear least squares problem

$$\inf_{\alpha \in \mathfrak{L}^2_X} \mathbb{E} \Big[\| \mathbb{Y}^{\Lambda} - \psi_{\alpha | R_{\Lambda}}(X) \|^2 \Big].$$
(169)

The convergence in (168) holds \mathbb{P} -a.s. if (α_k) is such that $\sum_{k=0}^{\infty} (\Phi(\psi_{\alpha_k|R_{\Lambda}}(X)) - \gamma_{R_{\Lambda}})^{1/2} < \infty$ for $\gamma_{R_{\Lambda}} := \inf_{\alpha \in \mathfrak{L}^2_{X}} \Phi(\psi_{\alpha|R_{\Lambda}}(X)).$

Proof. The above are straightforward consequences of Lemma B.9 and the proof of Theorem 4.8. Indeed: Adopting the notation of said proof, we see that

$$\mathbb{Y}_{X}^{\Lambda} = \underset{\mathbb{Z} \in \mathfrak{C}_{R_{\Lambda}}}{\operatorname{arg\,min}} \mathbb{E}\left[\left\|\mathbb{Y}^{\Lambda} - \mathbb{Z}\right\|^{2}\right] = \underset{\mathbb{Z} \in \mathfrak{C}_{2R_{\Lambda}}}{\operatorname{arg\,min}} \Phi(\mathbb{Z})$$
(170)

from the characterisation (141) and the inclusions $\mathfrak{C}_{2R_{\Lambda}} \subset G$ and $\mathbb{Y}_{X}^{\Lambda} \in \mathfrak{C}_{R_{\Lambda}}$ (Lemma B.9 (iv)).

The convergence (168) then follows by combination of (170) and Lemma B.9 (ii), in complete analogy to how (65) was obtained by combination of (141) and (142): Noting that

$$\Phi(\mathbb{Y}_X^{\Lambda}) \stackrel{(170)}{=} \min_{\mathbb{Z} \in \mathfrak{C}_{R_{\Lambda}}} \Phi(\mathbb{Z}) \stackrel{\text{Lem. B.9 (ii)}}{=} \inf\{\Phi(\mathbb{Z}) \mid \mathbb{Z} \in \Psi_{2R_{\Lambda}}(X)\} = \inf_{\alpha \in \mathfrak{L}_X^2} \Phi(\psi_{\alpha|R_{\Lambda}}(X)) = \gamma_{R_{\Lambda}},$$

we see that, for any minimizing sequence (α_k) of $\{\Phi(\psi_{\alpha|R_{\Lambda}}(X)) \mid \alpha \in \mathfrak{L}^2_X\}$, the sequence $(\mathbb{Z}_k) := (\psi_{\alpha_k|R_{\Lambda}}(X)) \subset \mathfrak{C}_{2R_{\Lambda}}$ is minimizing for $\min\{\Phi(\mathbb{Z}) \mid \mathbb{Z} \in \mathfrak{C}_{2R_{\Lambda}}\}$; hence and with Lemma B.9 (i), both the convergence in (168) and the proposition's final almost-sure assertion follow exactly as in the proof of Theorem 4.8.