

An SVM approach for forecasting the evolution of an unknown ergodic dynamical system from observations with unknown noise

Ingo Steinwart and Marian Anghel
Information Sciences, CCS-3
Los Alamos National Laboratory
{ingo,manghel}@lanl.gov

March 28, 2007

Abstract

We consider the problem of forecasting the next (observable) state of an unknown ergodic dynamical system from a noisy observation of the present state. Our main result shows that support vector machines (SVMs) using Gaussian RBF kernels can learn the best forecaster from a sequence of noisy observations if *a*) the unknown observational noise processes is bounded and has a summable α -mixing rate and *b*) the unknown ergodic dynamical system is defined by a Lipschitz continuous function on some compact subset of \mathbb{R}^d and has a summable decay of correlations for Lipschitz continuous functions. In order to prove this result we first establish a general learning theorem for SVMs and all stochastic processes that satisfy a mixing notion that is substantially weaker than α -mixing.

1 Introduction

Let us assume that we have an ergodic dynamical system described by the sequence $(F^n)_{n \geq 0}$ of iterates of an (essentially) unknown map $F : M \rightarrow M$, where $M \subset \mathbb{R}^d$ is compact and the corresponding ergodic measure μ is assumed to be unique. Furthermore, assume that all observations \tilde{x} of this dynamical system are corrupted by some stationary, \mathbb{R}^d -valued, additive noise process $\mathcal{E} = (\varepsilon_n)_{n \geq 0}$ whose distribution ν we assume to be independent of the state, but otherwise *unknown*, too. In other words *all possible* observations of the system at time $n \geq 0$ are of the form

$$\tilde{x}_n = F^n(x_0) + \varepsilon_n, \tag{1}$$

where x_0 is a true but unknown state at time 0. Now, given an observation of the system at some arbitrary time our goal is to forecast the next *observable* state¹, i.e., given $x + \varepsilon$ we want to forecast $F(x) + \varepsilon'$, where ε and ε' are the observational errors for x and its successor $F(x)$. Of course, if we know neither F nor ν then this task is impossible, and hence we assume that we have a finite sequence $T = (\tilde{x}_0, \dots, \tilde{x}_{n-1})$ of noisy observations from a trajectory of the dynamical system, i.e., all \tilde{x}_i , $i = 0, \dots, n-1$, are given by (1) for a conjoint initial state x_0 . Now, informally speaking our goal is to use T to build a forecaster $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ whose average forecasting performance on noisy

¹We will see later that under some circumstances this is equivalent to forecasting the next *true* state. For the moment, however, we deal with observable states since under our above assumptions these are the only ones we have access to.

observations is as small as possible. In order to render this more precisely we need a loss function $L : \mathbb{R}^d \rightarrow [0, \infty)$ such that

$$L(F(x) + \varepsilon' - f(x + \varepsilon))$$

gives a value for the discrepancy between the forecast $f(x + \varepsilon)$ and the observed next state $F(x) + \varepsilon'$. In the following, we always assume implicitly, that small values of $L(F(x) + \varepsilon' - f(x + \varepsilon))$ correspond to small values of $\|F(x) + \varepsilon' - f(x + \varepsilon)\|_2$, where $\|\cdot\|_2$ denotes the Euclidean distance in \mathbb{R}^d . Now, by the stationarity of \mathcal{E} the average forecasting performance is given by the L -risk

$$\mathcal{R}_{L,P}(f) := \int \int L(F(x) + \varepsilon_1 - f(x + \varepsilon_0)) \nu(d\varepsilon) \mu(dx), \quad (2)$$

where $\varepsilon = (\varepsilon_i)_{i \geq 0}$ and $P := \nu \otimes \mu$. Obviously, the smaller the risk the better the forecaster is, and hence we ideally would like to have a forecaster $f_L^* : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that attains the minimal L -risk

$$\mathcal{R}_{L,P}^* := \inf \{ \mathcal{R}_{L,P}(f) \mid f : \mathbb{R}^d \rightarrow \mathbb{R}^d \text{ measurable} \}. \quad (3)$$

Now assume that we have a method \mathcal{L} that assigns to every training set $T \in (\mathbb{R}^d)^n$ a forecaster f_T . Then the method \mathcal{L} achieves our goal asymptotically, if it is *consistent* in the sense of

$$\lim_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_T) = \mathcal{R}_{L,P}^*, \quad (4)$$

where the limit is in probability P . To our best knowledge this forecasting problem has not been considered in the literature. Moreover, even the observational noise model itself has only been considered sporadically, though it clearly “captures important features of many experimental situations”, [26]. One of the reasons for this lack of consideration may simply be the fact that unlike dynamical noise which typically leads to a Markov chain, the observational noise does not change the deterministic character and the long range dependence of the system, and hence the observational noise model cannot be treated by more traditional time series techniques. Moreover, most of the existing work on the observational noise model deals with the question of denoising [23, 17, 34, 22, 24, 25, 26]. In particular, [24, 25, 26] provide both positive and negative results on the existence of consistent denoising procedures. Finally, for the least squares loss [31] presents methods for finding f_L^* for systems of the form $Z_{i+1} := F(Z_i) + \varepsilon_{i+1}$, $i \geq 0$, where (F^i) is a dynamical system and (ε_i) is some additive centered i.i.d. dynamical noise. In particular, consistency of two histogram-based methods is established if *a)* $F : M \rightarrow M$ is continuous and (ε_i) is bounded, or *b)* F is bounded and ε_i is absolutely continuous, respectively. Note that the first case shows that in the absence of dynamical and observational noise there is a method which can identify the map F whenever it is continuous but otherwise unknown. However, it is unclear how to extend the methods of [31] to deal with observational noise.

Variants of the forecasting problem for general stationary ergodic processes (Z_i) have been extensively studied in the literature. One often considered variant is *static autoregression* (see [21, p. 569ff] and the references therein) where the goal is to find sequences $\hat{f}_m(Z_{-1}, \dots, Z_{-m})$ of estimators that converge almost surely to the conditional expectation $\mathbb{E}(Z_0 | Z_{-1}, \dots, Z_{-\infty})$, which is the least squares optimal one-step-ahead forecaster using an infinite past of observations. However, even if we consider forecasters using a longer history of observations the goal of static autoregression cannot be compared to our concept of consistency because different notions of convergence are considered. Indeed, in static autoregression the goal is to find a near optimal prediction for \tilde{x}_0 using the previously observed $\tilde{x}_{-1}, \dots, \tilde{x}_{-m}$ of the *same* trajectory, whereas our goal is to use the observations to build a predictor which predicts near optimal for *all* future observations. In machine learning terminology static autoregression is thus an “on-line” learning problem whereas

our notion of consistency defines a “batch” learning problem. Learning methods for forecasting goals closer related to our notion of consistency were considered by, e.g., [29, 28] but unfortunately these methods require α - or β -mixing conditions that cannot be satisfied by non-trivial dynamical systems. Finally, a result by Nobel [30] shows that there is no method that is universally consistent for, e.g., classification and regression problems where the data is generated by an arbitrary stationary ergodic process.

If the observational noise process \mathcal{E} is mixing in the ergodic sense then it is not hard to check that the process described by (1) is ergodic and hence it satisfies a strong law of large numbers by Birkhoff’s theorem. Using the recent results in [38] we then see that there exists a support vector machine (see the next section for a description) *depending* on F and \mathcal{E} which is consistent in the sense of (4). However, [38] does not provide an explicit method for finding a consistent SVM even if both F and \mathcal{E} are known. Consequently, it is fair to say that though SVMs do not have principal limitations for the forecasting problem there is currently no theoretically sound way to use them. The goal of this work is to address this issue by showing that certain SVMs are consistent for all Lipschitz continuous F and bounded \mathcal{E} that have a sufficiently fast decay of correlations for Lipschitz continuous functions. In particular, we show that these SVMs are consistent for, e.g., all uniformly smooth expanding or hyperbolic dynamics F and all bounded i.i.d. noise processes \mathcal{E} .

The rest of this work is organized as follows: In Section 2 we recall the definition of support vector machines (SVMs). Then, in Section 3 we present a consistency result for SVMs and general stochastic processes which have a sufficiently fast decay of correlations. This result is then applied to the above forecasting problem in Section 4, where we also briefly review some dynamical systems with sufficiently fast decay of correlations. Possible future extensions of this work are discussed in Section 5. Finally, the proofs of the two main results can be found in Section 6 and Section 7, respectively.

2 Support Vector Machines

The goal of this section is to briefly describe support vector machines which were first introduced by [7, 15] as a method for learning binary classification tasks. Since then they were generalized to other problem domains such as regression and anomaly detection, and nowadays they are considered to be one of the state-of-the-art machine learning methods for these problem domains. For an thorough introduction to SVMs we refer the reader to the books [41, 16, 35].

Let us begin by introducing some notions related to SVMs. To this end let us fix two non-empty closed sets $X \subset \mathbb{R}^d$ and $Y \subset \mathbb{R}$, and a measurable function $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$, which in the following is called loss function². Moreover, let H be the reproducing kernel Hilbert space (RKHS) of a measurable kernel $k : X \times X \rightarrow \mathbb{R}$ (see [1] for such spaces). In addition, for a finite sequence $T \in (X \times Y)^n$ and a function $f : X \rightarrow \mathbb{R}$ we define the empirical L -risk by

$$\mathcal{R}_{L,T}(f) := \frac{1}{n} \sum_{i=0}^{n-1} L(x_i, y_i, f(x_i)).$$

Now, for given such a T and a regularization parameter $\lambda > 0$ support vector machines construct a function $f_{T,\lambda,H} : X \rightarrow \mathbb{R}$ satisfying

$$\lambda \|f_{T,\lambda,H}\|_H^2 + \mathcal{R}_{L,T}(f_{T,\lambda,H}) = \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,T}(f). \quad (5)$$

²Note that this is a more general concept of a loss function than the informal notion of a loss function used in the introduction.

It is well-known that if L is a *convex loss function* in the sense that $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ is convex for all $(x, y) \in X \times Y$, then there exists a unique $f_{T, \lambda, H}$. Moreover, in this case (5) becomes a strictly convex optimization problem which can be solved by, e.g., simple gradient descent algorithms. However, for specific losses including the least squares loss other, more efficient algorithmic approaches are used in practice, see [42], [41], [35], and [40].

Let us now introduce two additional properties of loss functions which will be used in this work.

Definition 2.1 A loss function $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ is called *differentiable* if $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ is differentiable for all $(x, y) \in X \times Y$. In this case the derivative is denoted by $L'(x, y, \cdot)$.

Definition 2.2 A loss function $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ is called *locally Lipschitz continuous* if for all $a \geq 0$ there exists a constant $c_a \geq 0$ such that for all $x \in X$, $y \in Y$ and all $t, t' \in [-a, a]$ we have

$$|L(x, y, t) - L(x, y, t')| \leq c_a |t - t'|.$$

Moreover, the smallest possible constant c_a is denoted by $|L|_{a,1}$.

Finally, L is called *Lipschitz continuous* if we have $|L|_1 := \sup_{a \geq 0} |L|_{a,1} < \infty$.

Let us now summarize some assumptions on the loss function L which we will use frequently.

Assumption L: The loss $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ is convex, differentiable and locally Lipschitz continuous in the above sense, and it also satisfies $L(x, y, 0) \leq 1$ for all $(x, y) \in X \times Y$. Moreover, for the derivative L' there exists a constant $c \in [0, \infty)$ such that for all $(x, y, t), (x', y', t') \in X \times Y \times \mathbb{R}$ we have

$$|L'(x, y, t) - L'(x', y', t')| \leq c \|(x, y, t) - (x', y', t')\|_2 \quad (6)$$

and $|L'(x, y, 0)| \leq c$.

Note that combining the two assumptions on L' yields $|L'(x, y, t)| \leq c(1 + |t|)$ for all $(x, y, t) \in X \times Y \times \mathbb{R}$, and from this it is not hard to conclude that $|L|_{a,1} \leq c(1 + a)$ for all $a > 0$.

Since the Assumption **L** is rather complex let us now illustrate it for two particular classes of loss functions used in many SVM variants.

Example 2.3 A loss $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ of the form $L(x, y, t) = \varphi(yt)$ for a suitable function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ and all $x \in X$, $y \in Y := \{-1, 1\}$ and $t \in \mathbb{R}$, is called *margin-based*. Obviously, L is convex, continuous, (locally) Lipschitz continuous, or differentiable if and only if φ is. In addition, convexity of L implies local Lipschitz continuity of L . Furthermore recall that [6] showed that L is suitable for binary classification tasks if and only if φ is differentiable at 0 with $\varphi'(0) < 0$.

Let us now consider Assumption **L**. Obviously, the first part is satisfied if and only if φ is convex and differentiable, and also satisfies $\varphi(0) \leq 1$. Moreover, the latter can always be ensured by rescaling φ . Furthermore, we have $L'(x, y, t) = y\varphi'(yt)$ and by considering the cases $y = y'$ and $y \neq y'$ separately we see that (6) is satisfied if and only if φ' is Lipschitz continuous and satisfies

$$|\varphi'(t) + \varphi'(t')| \leq c(1 + |t + t'|), \quad t, t' \in \mathbb{R},$$

for a constant $c > 0$. Finally, the condition $|L'(x, y, 0)| = |\varphi'(0)| \leq c$ is always satisfied for sufficiently large c . From these consideration we conclude that the classical SVM losses defined by $\varphi(t) = (1 - t)_+$ and $\varphi(t) = (1 - t)_+^2$, where $(x)_+ := \max\{0, x\}$, do *not* satisfy Assumption **L**, whereas the least square loss and the logistic loss defined $\varphi(t) = (1 - t)^2$ and $\varphi(t) = \ln(1 + \exp(-t))$, respectively, fulfill **L**. \triangleleft

Example 2.4 A loss $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ of the form $L(y, t) = \psi(y - t)$ for a suitable function $\psi : \mathbb{R} \rightarrow \mathbb{R}$ and all $x \in X$, $y \in Y \subset \mathbb{R}$ and $t \in \mathbb{R}$, is called *distance-based*. Recall that distance-based losses such as the least squares loss $\psi(r) = r^2$, Huber's insensitive loss $\psi(r) = \min\{r^2, \max\{1, 2|r| - 1\}\}$, the logistic loss

$\psi(r) = \ln((1+e^r)^2 e^{-r}) - \ln 4$, or the ϵ -insensitive loss $\psi(r) = (|r| - \epsilon)_+$ are usually used for regression. In the following we assume that Y is a compact subset of \mathbb{R} . Then it is easy to see that L is convex, differentiable, or locally Lipschitz continuous if and only if ψ is.

Let us now consider Assumption **L**. Obviously, the first part is satisfied if and only if ψ is convex and differentiable, and also satisfies $\sup_{y \in Y} \psi(y) \leq 1$. Note that the latter can always be ensured by rescaling ψ since the convexity of ψ implies its continuity and Y is assumed to be compact. Furthermore, we have $L'(x, y, t) = -\psi'(y - t)$, and hence we see that (6) is satisfied if and only if ψ' is Lipschitz continuous. Finally, every convex and differentiable function is continuously differentiable and hence we can always ensure $|L'(x, y, 0)| = |\psi'(y)| \leq c$. From these considerations we immediately see that all of the above distance-based losses besides the ϵ -insensitive loss satisfy Assumption **L**. \triangleleft

Finally, in the following we are mainly interested in Gaussian RBF kernels $k_\sigma : X \times X \rightarrow \mathbb{R}$ defined by

$$k_\sigma(x, x') := \exp(-\sigma^2 \|x - x'\|_2^2), \quad x, x' \in X,$$

where $X \subset \mathbb{R}^d$ is a non-empty subset and $\sigma > 0$ is a free parameter called the width. We write $H_\sigma(X)$ for the corresponding RKHSs which are described in some detail in [39]. Finally, for SVMs using a Gaussian kernel we write $f_{T, \lambda, \sigma} := f_{T, \lambda, H_\sigma(X)}$ in order to simplify notations.

3 Consistency of SVMs for a General Class of Stochastic Processes

The goal of this section is to establish consistency of SVMs for a class of stochastic processes satisfying a certain assumption on the decay of certain correlations. This result will then be used to establish consistency of SVMs for the forecasting problem and suitable combinations of dynamical systems F and noise processes \mathcal{E} .

Let us begin with some notations. To this end let us assume that we have a probability space $(\Omega, \mathcal{A}, \mu)$, a measurable space (Z, \mathcal{B}) , and a measurable map $T : \Omega \rightarrow Z$. Then $\sigma(T)$ denotes the smallest σ -algebra on Ω for which T is measurable. Moreover, μ_T denotes the T -image measure of μ , which is defined by $\mu_T(B) := \mu(T^{-1}(B))$, $B \subset Z$ measurable. Recall that a *stochastic process* $\mathcal{Z} := (Z_n)_{n \geq 0}$, i.e., a sequence of measurable maps $Z_n : \Omega \rightarrow Z$, $n \geq 0$, is called *identically distributed* if $\mu_{Z_n} = \mu_{Z_m}$ for all $n, m \geq 0$. In this case we write $P := \mu_{Z_0}$ in the following. Moreover, \mathcal{Z} is called *stationary in the wide sense* if $\mu_{(Z_{i_1+i}, Z_{i_2+i})} = \mu_{(Z_{i_1}, Z_{i_2})}$ for all $i_1, i_2, i \geq 1$, and it is said to be *stationary* if $\mu_{(Z_{i_1+i}, \dots, Z_{i_n+i})} = \mu_{(Z_{i_1}, \dots, Z_{i_n})}$ for all $n, i, i_1, \dots, i_n \geq 1$.

The following definition introduces a notion of correlation for stochastic processes that will be used throughout this work.

Definition 3.1 *Let $(\Omega, \mathcal{A}, \mu)$ be a probability space, (Z, \mathcal{B}) be a measurable space, \mathcal{Z} be a Z -valued, identically distributed process on Ω and $P := \mu_{Z_0}$. Then for $\varphi, \psi \in L_2(P)$ the n -th correlation, $n \geq 0$, is defined by*

$$\text{cor}_{\mathcal{Z}, n}(\psi, \varphi) := \int_{\Omega} \psi(Z_0) \cdot \varphi(Z_n) d\mu - \int_Z \psi dP \cdot \int_Z \varphi dP.$$

Moreover, $(\text{cor}_{\mathcal{Z}, n}(\psi, \varphi))_{n \geq 0}$ is called the sequence of correlations of ψ and φ .

Obviously, if \mathcal{Z} is an i.i.d. process we have $\text{cor}_{\mathcal{Z}, n}(\psi, \varphi) = 0$ for all $\varphi, \psi \in L_2(P)$ and $n \geq 0$, and this remains true if $\psi \circ Z_0$ and $\varphi \circ Z_n$ are uncorrelated. Consequently, if $\lim_{n \rightarrow \infty} \text{cor}_{\mathcal{Z}, n}(\psi, \varphi) = 0$ the corresponding speed of convergence provides information how fast $\psi \circ Z_0$ becomes uncorrelated from $\varphi \circ Z_n$. This idea has been intensively used in the statistical literature in terms of, e.g., the α -mixing coefficients

$$\alpha(\mathcal{Z}, n) := \sup_{\substack{A \in \mathcal{F}_0^0 \\ B \in \mathcal{F}_n^\infty}} |\mu(A \cap B) - \mu(A)\mu(B)|,$$

where \mathcal{F}_i^j is the initial σ -algebra of Z_i, \dots, Z_j . These and related coefficients together with examples including, e.g., certain Markov chains, ARMA processes, and GARCH processes are discussed in some detail in the survey article [10] and the books [8, 20, 11]. Moreover, for processes \mathcal{Z} satisfying $\alpha(\mathcal{Z}, n) \leq cn^{-\alpha}$ for some constant $c > 0$ and all $n \geq 1$ it was recently described in [38] how to find a regularization sequence (λ_n) for which the corresponding SVM is consistent. Unfortunately, however, it is well-known that every non-trivial ergodic dynamical system does *not* satisfy $\lim_{n \rightarrow \infty} \alpha(\mathcal{Z}, n) = 0$, and therefore the result of [38] cannot be used to investigate consistency for the forecasting problem. On the other hand various dynamical systems enjoy a uniform decay rate over smaller sets of functions such as Lipschitz continuous functions (see Section 4 for some examples). This leads to the following definition.

Definition 3.2 *Let $(\Omega, \mathcal{A}, \mu)$ be a probability space, $Z \subset \mathbb{R}^d$ be a compact set, \mathcal{Z} be a Z -valued, identically distributed process on Ω and $P := \mu_{Z_0}$. Moreover, let $(\gamma_i)_{i \geq 0}$ be a strictly positive sequence converging to 0. Then \mathcal{Z} is said to have a decay of correlations of the order (γ_i) if for all $\psi, \varphi \in \text{Lip}(Z)$ there exists a constant $\kappa_{\psi, \varphi} \in [0, \infty)$ such that*

$$|\text{cor}_{\mathcal{Z}, i}(\psi, \varphi)| \leq \kappa_{\psi, \varphi} \gamma_i, \quad i \geq 0, \quad (7)$$

where $\text{Lip}(Z)$ denotes the set of all \mathbb{R} -valued Lipschitz continuous functions defined on Z .

Let us now summarize our assumptions on the process \mathcal{Z} which we will make in the rest of this section.

Assumption Z: The process $\mathcal{Z} = (X_i, Y_i)_{i \geq 0}$ is defined on the probability space $(\Omega, \mathcal{A}, \mu)$ and is $X \times Y$ -valued, where $X \subset \mathbb{R}^d$ and $Y \subset \mathbb{R}$ are compact subsets. Moreover \mathcal{Z} is stationary in the wide sense.

Finally, we will need the following two, mutually exclusive assumptions on the regularization sequence and the kernel width of SVMs.

Assumption S1: For a fixed strictly positive sequence $(\gamma_i)_{i \geq 0}$ converging to 0 and a locally Lipschitz continuous loss L the sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfy $\sigma_n \geq \ln(n+1)$ for all $n \geq 0$,

$$\lim_{n \rightarrow \infty} \frac{\lambda_n \sigma_n^{4d}}{|L|_{\lambda_n^{-1/2}, 1}} = 0, \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{|L|_{\lambda_n^{-1/2}, 1}^3 \sigma_n^4}{n \lambda_n^4} \sum_{i=0}^{n-1} \gamma_i = 0.$$

Assumption S2: For a fixed strictly positive sequence $(\gamma_i)_{i \geq 0}$ converging to 0 and a locally Lipschitz continuous loss L the sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfy $\sigma_n \geq \ln(n+1)$ for all $n \geq 0$, $\lim_{n \rightarrow \infty} \lambda_n \sigma_n^d = 0$,

$$\lim_{n \rightarrow \infty} \frac{\lambda_n \sigma_n^{4d}}{|L|_{\lambda_n^{-1/2}, 1}} = \infty, \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{|L|_{\lambda_n^{-1/2}, 1}^6 \sigma_n^{4+12d}}{n \lambda_n} \sum_{i=0}^{n-1} \gamma_i = 0.$$

In order to illustrate the assumptions **S1** and **S2** let us first assume that L is Lipschitz continuous. Furthermore, we assume $(\gamma_i) \in \ell_1$ as well as $\lambda_n := n^{-\alpha}$ and $\sigma_n := n^\beta$ for $n \geq 1$ and constants $\alpha, \beta > 0$. Then Assumption **S1** is met if $\alpha > 4d\beta$ and $\alpha + \beta < 1/4$, whereas Assumption **S2** is met if $\alpha + (4 + 12d)\beta < 1$ and $d\beta < \alpha < 4d\beta$. Moreover, if we take $\lambda_n := n^{-\alpha}$ and $\sigma_n := (\ln n)^{1+\beta}$ then Assumption **S1** is met if $0 < \alpha < 1/4$.

Let us now consider an arbitrary loss L satisfying Assumption **L**. Then we discussed after Assumption **L** that we have $|L|_{\lambda^{-1/2},1} \leq c(1 + \lambda^{-1/2})$ for a constant $c > 0$ and all $\lambda > 0$. Let us again consider the case $(\gamma_i) \in \ell_1$, $\lambda_n := n^{-\alpha}$, and $\sigma_n := n^\beta$ for $n \geq 1$ and constants $\alpha, \beta > 0$. Then Assumption **S1** is met if $\alpha > \frac{8}{3}d\beta$ and $\frac{11}{8}\alpha + \beta < \frac{1}{4}$, whereas Assumption **S2** is met if $d\beta < \alpha < \frac{8}{3}d\beta$ and $(1 + 3d)\beta + \alpha < \frac{1}{4}$. Moreover, if we take $\lambda_n := n^{-\alpha}$ and $\sigma_n := (\ln n)^{1+\beta}$ then Assumption **S1** is met if $0 < \alpha < 2/11$.

The illustrations above show that both Assumptions **S1** and **S2** consist of two contrary conditions, namely one which implies that λ_n tends to 0 with a sufficiently fast speed and another one which ensures that this speed is not too fast. Roughly speaking the first condition guarantees that the approximation error tends to zero (see Lemma 6.5), but since this simultaneously means that the statistical error becomes larger, the second condition is needed to ensure that the latter error still tends to zero (see the proof of Theorem 3.3). This trade-off between approximation and statistical error is typical for consistent learning algorithms (see the books [19] and [21] for several such examples).

With the help of these assumptions we can now establish the announced consistency of SVMs:

Theorem 3.3 *Let $\mathcal{Z} = (X_i, Y_i)_{i \geq 0}$ be stochastic process satisfying Assumption **Z**. We write $P := \mu_{(X_0, Y_0)}$ and assume that \mathcal{Z} has a decay of correlations of some order (γ_i) . In addition, let $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a loss satisfying Assumption **L**. Then for all sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfying either **S1** or **S2** and all $\epsilon \in (0, 1]$ we have*

$$\lim_{n \rightarrow \infty} \mu \left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega), \lambda_n, \sigma_n}) - \mathcal{R}_{L,P}^*| > \epsilon \right) = 0,$$

where $T(\omega) := ((X_0(\omega), Y_0(\omega)), \dots, (X_{n-1}(\omega), Y_{n-1}(\omega)))$.

Theorem 3.3 in particular applies to stochastic processes which are α -mixing with rate (γ_i) . However, the Assumptions **S1** and **S2** ensuring consistency are substantially stronger than the ones obtained in [38] for such processes. On the other hand, there are quite a few stochastic processes which are not α -mixing but still enjoy a reasonably fast decay of correlations. Since we are mainly interested in the forecasting problem we will delay the discussion of such examples to the next section.

4 Consistency of SVMs for the Forecasting Problem

In this section we present our main result which establishes consistency of SVM for the forecasting problem if the dynamical system enjoys a certain decay of correlations. In addition, we discuss some examples of such dynamical systems.

We begin by first revisiting our informal problem description given in the introduction. To this end let $M \subset \mathbb{R}^d$ be a compact set and $F : M \rightarrow M$ be map such that the dynamical system $\mathcal{D} := (F^i)_{i \geq 0}$ has an ergodic measure μ . Moreover, let $\mathcal{E} = (\varepsilon_i)_{i \geq 0}$ be a \mathbb{R}^d -valued stochastic process which is (stochastically) independent of \mathcal{D} . Then the process that generates the noisy observations (1) is $(F^i + \varepsilon_i)_{i \geq 0}$. In particular, a sequence of observations $(\tilde{x}_0, \dots, \tilde{x}_n)$ generated by this process is of the form (1) for a conjoint initial state. Now recall that given an observation of the system at some arbitrary time our goal is to forecast the next *observable* state. Consequently, we will use the training set

$$T(x, \varepsilon) := ((\tilde{x}_0, \tilde{x}_1), \dots, (\tilde{x}_{n-1}, \tilde{x}_n)) = ((x + \varepsilon_0, F(x) + \varepsilon_1), \dots, (F^{n-1}(x) + \varepsilon_{n-1}, F^n(x) + \varepsilon_n)) \quad (8)$$

whose input/output pairs are consecutive observable states. Now note that a single sample $(F^{i-1}(x) + \varepsilon_{i-1}, F^i(x) + \varepsilon_i)$ depends on the pair $(\varepsilon_i, \varepsilon_{i+1})$ and therefore we have to consider the

process of such pairs. The following assumption summarizes the needed requirements of this process \mathcal{N} .

Assumption N: For the \mathbb{R}^{2d} -valued stochastic process \mathcal{N} there exist a constant $B > 0$ and a probability measure ν on $[-B, B]^{d\mathbb{N}_0}$ such that the coordinate process $\mathcal{E} := (\pi_0 \circ S^i)_{i \geq 0}$ is stationary with respect to ν and satisfies $\mathcal{N} = (\pi_0 \circ S^i, \pi_0 \circ S^{i+1})_{i \geq 0}$, where S denotes the shift operator $(x_i)_{i \geq 0} \mapsto (x_{i+1})_{i \geq 0}$ and π_0 denotes the projection $(x_i)_{i \geq 0} \mapsto x_0$.

Before we state our main result we furthermore note that the input variable $x + \varepsilon$ and the output variable $F(x) + \varepsilon'$ are d -dimensional vectors. Consequently, our notion of a loss function introduced in Section 2 needs a refinement which captures the ideas of the introduction. To this end we state the following assumption.

Assumption LD: For the function $L : \mathbb{R}^d \rightarrow [0, \infty)$ there exists a distance-based loss satisfying Assumption **L** such that its representing function $\psi : \mathbb{R}^d \rightarrow [0, \infty)$ satisfies

$$L(r_1, \dots, r_d) = \psi(r_1) + \dots + \psi(r_d), \quad (r_1, \dots, r_d) \in \mathbb{R}^d. \quad (9)$$

Moreover, ψ has a unique global minimum at 0.

Obviously, if L satisfies Assumption **LD** then L is a loss in the sense of the introduction. Moreover note that the specific form (9) makes it possible to consider the coordinates of the output variable *separately*. Consequently, we will use the forecaster

$$\bar{f}_{T, \lambda, \sigma} := (f_{T_1, \lambda, \sigma}, \dots, f_{T_d, \lambda, \sigma}), \quad (10)$$

where $f_{T_j, \lambda, \sigma}$ is the SVM solution obtained by considering the distance-based loss defined by ψ and the training set $T_j := ((\tilde{x}_0, \pi_j(\tilde{x}_1)), \dots, (\tilde{x}_{n-1}, \pi_j(\tilde{x}_n)))$ which is obtained by projecting the output variable of T onto its j^{th} -coordinate via the projection π_j . In other words we build the forecaster $\bar{f}_{T, \lambda, \sigma}$ by training d different SVMs on the training sets T_1, \dots, T_d .

With the help of these preparations we can now present our main result which establishes consistency for such a forecaster.

Theorem 4.1 *Let $M \subset \mathbb{R}^d$ be a compact set, $F : M \rightarrow M$ be a Lipschitz continuous map such that the dynamical system $\mathcal{D} := (F^i)_{i \geq 0}$ has a unique ergodic measure μ , and \mathcal{N} be a stochastic process satisfying Assumption **N**. Assume that both processes \mathcal{D} and \mathcal{N} have a decay of correlations of the order (γ_i) . Moreover, let $L : \mathbb{R}^d \rightarrow [0, \infty)$ be a function satisfying Assumption **LD**. Then for all sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfying either **S1** or **S2** and all $\epsilon \in (0, 1]$ we have*

$$\lim_{n \rightarrow \infty} \mu \otimes \nu \left((x, \varepsilon) \in M \times [-B, B]^{d\mathbb{N}} : |\mathcal{R}_{L, P}(\bar{f}_{T(x, \varepsilon), \lambda_n, \sigma_n}) - \mathcal{R}_{L, P}^*| > \epsilon \right) = 0,$$

where $T(x, \varepsilon)$ is defined by (8) and the risks are given by (2) and (3).

Note that if \mathcal{E} is an i.i.d. process then \mathcal{N} has a decay of correlations of any order. Moreover, if \mathcal{E} is α -mixing with mixing rate (γ_i) then \mathcal{N} has a decay of correlations of order (γ_i) . Finally, if \mathcal{D} has a decay of correlations $(\gamma_i^{(1)})$ and \mathcal{N} has a decay of correlations $(\gamma_i^{(2)})$ then they obviously have both a decay of correlations (γ_i) , where $\gamma_i := \max\{\gamma_i^{(1)}, \gamma_i^{(2)}\}$. In particular, if the noise has slowly decaying correlations this will slow down learning even though the system itself, which we want to forecast may have a fast decay of correlations.

Obviously, the function $L(r) := \|r\|_2^2$, $r \in \mathbb{R}^d$, satisfies Assumption **LD** since the least squares loss satisfies Assumption **L** as we have discussed in Example 2.4. Let us now additionally assume

that the noise is *pairwise independent*, i.e. ε_i and $\varepsilon_{i'}$ are independent if $i \neq i'$, and centered, i.e. it satisfies $\mathbb{E}_{\varepsilon \sim \nu} \pi_0(\varepsilon) = 0$. For a forecaster $f = (f_1, \dots, f_d) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we then obtain

$$\begin{aligned} \mathcal{R}_{L,P}(f) &= \int \int \sum_{j=1}^d (\pi_j(F(x) + \varepsilon_1) - f_j(x + \varepsilon_0))^2 \nu(d\varepsilon) \mu(dx) \\ &= \int \int \sum_{j=1}^d (\pi_j(F(x)) - f_j(x + \varepsilon_0))^2 \nu(d\varepsilon) \mu(dx) + \int \|\varepsilon_0\|_2^2 \nu(d\varepsilon) \\ &=: \mathcal{R}_{L,\bar{P}}(f) + \int \|\varepsilon_0\|_2^2 \nu(d\varepsilon), \end{aligned}$$

where $\pi_j : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes the j^{th} -coordinate projection. Consequently, a forecaster f which approximately minimizes the L -risk is also an approximate forecaster of the *true* next state in the sense of $\mathcal{R}_{L,\bar{P}}(\cdot)$. Before we combine this observation with Theorem 4.1 let us first rephrase Assumptions **S1** and **S2** for the least squares loss.

Assumption S1-LS: For a fixed strictly positive sequence $(\gamma_i)_{i \geq 0}$ converging to 0 and a locally Lipschitz continuous loss L the sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfy $\sigma_n \geq \ln(n+1)$ for all $n \geq 0$,

$$\lim_{n \rightarrow \infty} \lambda_n \sigma_n^{8d/3} = 0, \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\sigma_n^4}{n \lambda_n^{11/2}} \sum_{i=0}^{n-1} \gamma_i = 0.$$

Assumption S2-LS: For a fixed strictly positive sequence $(\gamma_i)_{i \geq 0}$ converging to 0 and a locally Lipschitz continuous loss L the sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfy $\sigma_n \geq \ln(n+1)$ for all $n \geq 0$, $\lim_{n \rightarrow \infty} \lambda_n \sigma_n^d = 0$,

$$\lim_{n \rightarrow \infty} \lambda_n \sigma_n^{8d/3} = \infty, \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\sigma_n^{4+12d}}{n \lambda_n^4} \sum_{i=0}^{n-1} \gamma_i = 0.$$

Now we can state a result showing that SVMs using a least squares loss can be used to forecast the next *true* state of the dynamical system if the observational noise is sufficiently benign.

Corollary 4.2 *Let $M \subset \mathbb{R}^d$ be a compact set, $F : M \rightarrow M$ be a Lipschitz continuous map such that the dynamical system $\mathcal{D} := (F^i)_{i \geq 0}$ has a unique ergodic measure μ . Moreover, let $\mathcal{E} = (\varepsilon_i)_{i \geq 0}$ be an i.i.d. process of $[-B, B]^d$ -valued and centered random variables. Assume that \mathcal{D} has a decay of correlations of the order (γ_i) . Moreover, let $L : \mathbb{R}^d \rightarrow [0, \infty)$ be defined by $L(r) := \|r\|_2^2$, $r \in \mathbb{R}^d$. Then for all sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfying either **S1-LS** or **S2-LS** and all $\epsilon \in (0, 1]$ we have*

$$\lim_{n \rightarrow \infty} \mu \otimes \nu \left((x, \varepsilon) \in M \times [-B, B]^{d\mathbb{N}} : |\mathcal{R}_{L,\bar{P}}(\bar{f}_{T(x,\varepsilon),\lambda_n,\sigma_n}) - \mathcal{R}_{L,\bar{P}}^*| > \epsilon \right) = 0,$$

where $T(x, \varepsilon)$ is defined by (8).

It is interesting to note that the above corollary does *not* require the noise to be symmetric. Instead it only requires centered noise, i.e. the observations are not systematically biased in a certain direction. The following remark rephrases Theorem 4.1 and its corollary for situations with summable decays of correlations.

Remark 4.3 (Universal consistency for summable correlations) It is important to note that if the sequence (γ_i) bounding the correlation sequence is *summable*, i.e., $\sum \gamma_i < \infty$ then the Assumptions **S1** and **S2**, or **S1-LS** and **S2-LS**, respectively, are *independent* of both the dynamical system and the observational noise process. Consequently, using sequences satisfying these assumptions yields an SVM which is consistent for *all* such pairs of dynamical systems and observational noise processes. In other words, such an SVM can learn the optimal forecaster without knowing specifics of the dynamical systems and the observational noise. To be a bit more specific, let us assume that we use the least squares loss and sequences $\lambda_n := n^{-\alpha}$ and $\sigma_n := (1 + \ln n)^{1+\beta}$, $n \geq 1$, for fixed $0 < \alpha < 2/11$ and $\beta > 0$. Then the corresponding SVM is consistent for all bounded observational noise processes having a summable α -mixing rate and all ergodic dynamical systems on M which are defined by a Lipschitz continuous $F : M \rightarrow M$ and have a summable decay of correlations. Moreover, if the noise process is also i.i.d. and centered then this SVM actually learns to forecast the next *true* state.

Let us finally discuss some examples of classes of dynamical systems enjoying at least a polynomial decay of correlations. Since the existing literature on such systems is vast these examples are only meant to be illustrations for situations where Theorem 4.1 can be applied, and are *not* intended to provide an overview of known results. However, compilations of known results can be found in the survey articles of V. Baladi [3] and S. Luzzatto [27], and the book [2] of V. Baladi.

Example 4.4 (Smooth expanding dynamics) Let M be a compact connected Riemannian manifold and $F : M \rightarrow M$ be $C^{1+\varepsilon}$ for some $\varepsilon > 0$. Furthermore assume that there exists constants $c > 0$ and $\lambda > 1$ such that

$$\max\{\|DF_x^n(v)\| : x \in M, v \in T_x M \text{ with } \|v\| = 1\} \geq c\lambda^n$$

for all $n \geq 0$, where $T_x M$ denotes the tangent space of M at x and DF_x^n denotes the derivative of F^n at x . Then it is a classical result that F possesses a unique ergodic measure which is absolutely continuous with respect to the Riemannian volume. Moreover, it is well-known (see, e.g., [32] and the references mentioned in [27, Theorem 5]) that there exists a $\tau \in (0, 1)$ such that the dynamical system has decay of correlations of the order (τ^i) . Generalizations of this result to piecewise smooth and piecewise (non)-uniformly expanding dynamics are discussed in [3]. Finally, [27, Theorem 10] recalls results (together with references) for non-uniformly expanding dynamics having either exponential or polynomial decay of correlations.

Example 4.5 (Smooth hyperbolic dynamics) Assume that F is a topologically mixing $C^{1+\varepsilon}$ Anosov or Axiom A diffeomorphism. Then it is well-known (see, e.g., [9, 33]) that there exists a $\tau \in (0, 1)$ such that the dynamical system has decay of correlations of the order (τ^i) for some $\tau \in (0, 1)$. Moreover, [3] lists various extensions of this result to, e.g., smooth non-uniformly hyperbolic systems and hyperbolic systems with singularities.

Besides these classical results and their extensions, [3] also compiles a list of “parabolic” or “intermittent” systems which have a polynomial decay of correlations.

5 Discussion

The goal of this work was to show that in principle support vector machines can learn how to predict one-step-ahead noisy observation of a dynamical system without knowing specifics of the dynamical system or the observational noise besides a certain, rather general stochasticity. However, there remain several open questions which can be subject to further research:

More General Losses and Kernels. In the “statistical part” of our analysis, we used an approach which is based on a “stability” argument. However, it is also possible to use a “skeleton” argument based on covering numbers, instead. Utilizing the latter it seems possible to relax the assumptions on the loss L by making stronger assumptions on both (λ_n) and (σ_n) . A particular loss which is interesting in this direction would be the ϵ -insensitive loss

used in classical SVMs for regression. Another possible extension of our work is considering different kernels, such as the kernels that generate Sobolev spaces. In fact, we only focused on Gaussian RBF kernels since these kernels are most commonly used in practice.

Learning Rates. So far we have only shown that the risk of the SVM solution converges to the smallest possible risk, however, for practical considerations the *speed* of this convergence is of great importance, too. The proof we utilized already gives such learning rates if a *quantitative* version of the Approximation Lemma 6.5 is available, which is possible if, e.g., quantitative assumptions on the smoothness of F and the regularity of ν are made. However, since we conjecture that the statistical part of our analysis is not sharp we did not present a corresponding result. In this regard we note that recently [14] established a concentration result for piecewise regular expanding and topologically mixing maps of the interval $[0, 1]$ which is substantially stronger than our elementary Chebyshev inequality of Lemma 6.9. We strongly believe that such a concentration result can be used to substantially sharpen the statistical part of our analysis.

Dynamic Noise Systems. Another extension of the current work is to consider dynamical systems which are perturbed by some noise. Our general consistency result in Theorem 3.3 suggests that such an extension is possible whenever the perturbed system has a decay of correlations. In this regard we note that for certain perturbed systems of expanding maps decays of correlations have already been established in [5], and it would be interesting to check whether they can be used to prove consistency of SVMs.

Longer Past. So far we only used the present observation to forecast the next observation, but it is not hard to check that in almost any system/noise combination the minimal risk $\mathcal{R}_{L,P}^*$ reduces if one allows to use additional past observations. On the other hand it appears that the learning problem becomes harder in this case since we have to approximate a function which lives on a higher dimensional input space, and hence there seems to be a trade-off for finite sample sizes. While investigating this trade-off in more detail seems to be possible with the techniques developed in this work we again assume that the statistical part of our analysis is not sharp enough to obtain a meaningful picture.

Acknowledgment.

The authors gratefully thank V. Baladi for pointing us to the unpublished note [13] of P. Collet.

6 Proof of Theorem 3.3

The goal of this section is to prove Theorem 3.3. Since the proof requires several preliminary results we divided this section into subsections, which provide these prerequisites.

6.1 Some Basics on the Decay of Correlations

The main goal of this section is to establish some *uniform* bounds on the correlation sequence. Let us begin with the following lemma which establishes some basic properties of the correlation sequence.

Lemma 6.1 *Let $(\Omega, \mathcal{A}, \mu)$ be a probability space, (Z, \mathcal{B}) be a measurable space, \mathcal{Z} be a Z -valued, identically distributed process on Ω and $P := \mu_{Z_0}$. Then the sequence of correlations defines a bilinear operator $\text{cor}_{\mathcal{Z}} : L_2(P) \times L_2(P) \rightarrow \ell_\infty$ called the correlation operator.*

Proof: The bi-linearity is obvious and $\text{cor}_{\mathcal{Z}}(\psi, \varphi) \in \ell_\infty$ follows from

$$\left| \int_{\Omega} \psi(Z_0) \cdot \varphi(Z_n) d\mu \right| \leq \|\psi \circ Z_0\|_{L_2(\mu)} \|\varphi \circ Z_n\|_{L_2(\mu)} = \|\psi\|_{L_2(P)} \|\varphi\|_{L_2(P)}.$$

■

The following key theorem which goes back to an unpublished note [13] of P. Collet (see also p. 101 in [4]) can be used to establish continuity of the correlation operator. Before we present this result let us first recall that a Banach space E is said to be continuously embedded into the Banach space F if $E \subset F$ and the natural inclusion map $\text{id} : E \rightarrow F$ is continuous.

Theorem 6.2 *Let $(\Omega, \mathcal{A}, \mu)$ be a probability space, (Z, \mathcal{B}) be a measurable space, \mathcal{Z} be a Z -valued, identically distributed process on Ω and $P := \mu_{Z_0}$. Moreover, let E_1 and E_2 be Banach spaces that are continuously embedded into $L_2(P)$ and let F be a Banach space that is continuously embedded into ℓ_∞ . If for all $\psi \in E_1$ and all $\varphi \in E_2$ the correlation operator satisfies*

$$\text{cor}_{\mathcal{Z}}(\psi, \varphi) \in F$$

then there exists a constant $c \in [0, \infty)$ such that

$$\|\text{cor}_{\mathcal{Z}}(\psi, \varphi)\|_F \leq c \cdot \|\psi\|_{E_1} \|\varphi\|_{E_2}$$

for all $\psi \in E_1$ and all $\varphi \in E_2$. The smallest such constant c is denoted by $\|\text{cor}_{\mathcal{Z}} : E_1 \times E_2 \rightarrow F\|$.

For the sake of completeness the proof of this key result can be found in the Appendix. The most obvious examples of Banach spaces F in the above theorem are the spaces ℓ_p . However, in the literature on dynamical systems results on the sequence of correlations are usually stated in the form

$$|\text{cor}_{\mathcal{Z},n}(\psi, \varphi)| \leq \kappa_{\psi,\varphi} \gamma_n, \quad n \geq 0$$

where (γ_n) is a strictly positive sequence converging to 0 and $\kappa_{\psi,\varphi}$ is a constant depending on ψ and φ . To apply Theorem 6.2 in this situation we obviously need Banach spaces which capture such a behaviour of $\text{cor}_{\mathcal{Z}}(\cdot, \cdot)$. Therefore, let us fix a strictly positive sequence $\gamma := (\gamma_n)_{n \geq 0}$ such that $\lim_{n \rightarrow \infty} \gamma_n = 0$. For a sequence $b := (b_n) \subset \mathbb{R}$ we define

$$\|b\|_{\Lambda(\gamma)} := \sup_{n \geq 0} \frac{|b_n|}{\gamma_n}.$$

Moreover, we write

$$\Lambda(\gamma) := \{(b_n) \subset \mathbb{R} : \|(b_n)\|_{\Lambda(\gamma)} < \infty\}.$$

The following lemma establishes some basic properties of the pair $(\Lambda(\gamma), \|\cdot\|_{\Lambda(\gamma)})$.

Lemma 6.3 *The pair $(\Lambda(\gamma), \|\cdot\|_{\Lambda(\gamma)})$ is a Banach space and we have $\|\text{id} : \Lambda(\gamma) \rightarrow \ell_\infty\| \leq \|\gamma\|_\infty$.*

Proof: The fact that $(\Lambda(\gamma), \|\cdot\|_{\Lambda(\gamma)})$ is a normed space is elementary to prove. Moreover, we have

$$\|b\|_{\Lambda(\gamma)} = \sup_{n \geq 0} \frac{|b_n|}{\gamma_n} \geq \sup_{n \geq 0} \frac{|b_n|}{\|\gamma\|_\infty} = \frac{\|b\|_\infty}{\|\gamma\|_\infty},$$

and hence we find $\|\text{id} : \Lambda(\gamma) \rightarrow \ell_\infty\| \leq \|\gamma\|_\infty$. Finally, let $(b^{(i)})_{i \geq 1}$ be a Cauchy sequence in $\Lambda(\gamma)$. The previous step shows that it is also a Cauchy sequence in ℓ_∞ , and by the completeness of ℓ_∞

there consequently exists a sequence $b := (b_n) \in \ell_\infty$ such that $\lim_{i \rightarrow \infty} \|b^{(i)} - b\|_\infty = 0$. Let us now fix an $\varepsilon > 0$. Then there exists an index $i_0 \geq 0$ such that for all $i, j \geq i_0$ we have $\|b^{(i)} - b^{(j)}\|_{\Lambda(\gamma)} \leq \varepsilon$. Consequently, for fixed $N \geq 0$ we have

$$\sup_{n=0, \dots, N} \frac{|b_n^{(i)} - b_n^{(j)}|}{\gamma_n} \leq \|b^{(i)} - b^{(j)}\|_{\Lambda(\gamma)} \leq \varepsilon,$$

and by taking the limit $j \rightarrow \infty$ we conclude

$$\sup_{n=0, \dots, N} \frac{|b_n^{(i)} - b_n|}{\gamma_n} \leq \varepsilon.$$

However, N was arbitrary and hence we find $\|b^{(i)} - b\|_{\Lambda(\gamma)} \leq \varepsilon$ for all $i \geq i_0$. In other words we have shown that $(b^{(i)})_{i \geq 1}$ converges to b in $\|\cdot\|_{\Lambda(\gamma)}$, i.e., $(\Lambda(\gamma), \|\cdot\|_{\Lambda(\gamma)})$ is complete. \blacksquare

Combing the above lemma with Theorem 6.2 we immediately obtain the following corollary.

Corollary 6.4 *Let $(\Omega, \mathcal{A}, \mu)$ be a probability space, (Z, \mathcal{B}) be a measurable space, \mathcal{Z} be a Z -valued, identically distributed process on Ω and $P := \mu_{\mathcal{Z}_0}$. Moreover, let E_1 and E_2 be Banach spaces that are continuously embedded into $L_2(P)$. In addition, let $\gamma := (\gamma_n)_{n \geq 0}$ be a strictly positive sequence such that $\lim_{n \rightarrow \infty} \gamma_n = 0$. If for all $\psi \in E_1$ and all $\varphi \in E_2$ there exists a constant $\kappa_{\psi, \varphi} \in [0, \infty)$ such that*

$$|\text{cor}_{\mathcal{Z}, n}(\psi, \varphi)| \leq \kappa_{\psi, \varphi} \gamma_n$$

for all $n \geq 0$, then there exists a constant $c \in [0, \infty)$ such that

$$|\text{cor}_{\mathcal{Z}, n}(\psi, \varphi)| \leq c \|\psi\|_{E_1} \cdot \|\varphi\|_{E_2} \cdot \gamma_n$$

for all $\psi \in E_1$, $\varphi \in E_2$ and all $n \geq 0$.

6.2 Some Properties of Gaussian RBF Kernels

In this subsection we establish some properties of Gaussian RBF kernels which will be heavily used in the proof of Theorem 3.3. Let us begin with an approximation result.

Lemma 6.5 *Let $X \subset \mathbb{R}^d$ and $Y \subset \mathbb{R}$ be compact subsets, $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a convex locally Lipschitz continuous loss and P be a probability measure on $X \times Y$ such that $\mathcal{R}_{L, P}(0) < \infty$. Then for all sequences $(\lambda_n) \subset (0, 1]$ and $(\sigma_n) \subset [1, \infty)$ satisfying*

$$\lim_{n \rightarrow \infty} \lambda_n \sigma_n^d = 0 \tag{11}$$

we have

$$\lim_{n \rightarrow \infty} \left(\inf_{f \in H} \lambda_n \|f\|_{H_{\sigma_n}(X)}^2 + \mathcal{R}_{L, P}(f) \right) = \mathcal{R}_{L, P}^*.$$

Proof: For $\sigma > 0$ we write $\mathcal{R}_{L, P, H_\sigma(X)}^* := \inf\{\mathcal{R}_{L, P}(f) : f \in H_\sigma(X)\}$. Since L is locally Lipschitz continuous and $\mathcal{R}_{L, P}(0) < \infty$ the discussion after (4) in [37] shows that it is a P -integrable Nemitski loss in the sense of [37]. Now recall (see [36]) that $H_\sigma(X)$ is universal, i.e. it is dense in $C(X)$, and hence [37, Corollary 1] shows $\mathcal{R}_{L, P, H_\sigma(X)}^* = \mathcal{R}_{L, P}^*$ for all $\sigma > 0$. Let us now fix an $\varepsilon > 0$. The above discussion then shows that there exists an $f_\varepsilon \in H_1(X)$ such that

$$\mathcal{R}_{L, P}(f_\varepsilon) \leq \mathcal{R}_{L, P}^* + \varepsilon.$$

Furthermore (11) implies the existence of an $n_0 \geq 0$ such that

$$\lambda_n \sigma_n^d \leq \varepsilon \|f_\varepsilon\|_{H_1(X)}^{-2}$$

for all $n \geq n_0$. Since $\sigma_n \geq 1$ we also know $f_\varepsilon \in H_{\sigma_n}(X)$ and $\|f_\varepsilon\|_{H_{\sigma_n}(X)}^2 \leq \sigma_n^d \|f_\varepsilon\|_{H_1(X)}^2$ by [39, Corollary 6], and therefore we obtain

$$\mathcal{R}_{L,P}^* \leq \inf_{f \in H} \lambda_n \|f\|_{H_{\sigma_n}(X)}^2 + \mathcal{R}_{L,P}(f) \leq \lambda_n \|f_\varepsilon\|_{H_{\sigma_n}(X)}^2 + \mathcal{R}_{L,P}(f_\varepsilon) \leq \mathcal{R}_{L,P}^* + 2\varepsilon$$

for all $n \geq n_0$. From this we easily deduce the assertion. \blacksquare

Before we establish the next result let us recall that a function $f : X \rightarrow \mathbb{R}$ on a subset $X \subset \mathbb{R}^d$ is called Lipschitz continuous if there exist a constant $c \in [0, \infty)$ such that $|f(x) - f(x')| \leq c \|x - x'\|_2$ for all $x, x' \in X$. In the following the smallest such constant is denoted by $|f|_1$ and the set of all Lipschitz continuous functions is denoted by $\text{Lip}(X)$. Moreover, recall that if X is compact then $\text{Lip}(X)$ together with the norm $\|f\|_{\text{Lip}(X)} := \max\{\|f\|_\infty, |f|_1\}$ forms a Banach space. Moreover, in this case $\text{Lip}(X)$ is also closed under multiplication. Indeed, for $f, g \in \text{Lip}(X)$ and $x, x' \in X$ we have

$$\begin{aligned} |f(x)g(x) - f(x')g(x')| &\leq |f(x)g(x) - f(x)g(x')| + |f(x)g(x') - f(x')g(x')| \\ &\leq \|f\|_\infty \cdot |g|_1 |x - x'| + |f|_1 |x - x'| \cdot \|g\|_\infty \\ &\leq 2\|f\|_{\text{Lip}(X)} \|g\|_{\text{Lip}(X)} |x - x'|, \end{aligned}$$

and hence we find $fg \in \text{Lip}(X)$ with $\|fg\|_{\text{Lip}(X)} \leq 2\|f\|_{\text{Lip}(X)} \|g\|_{\text{Lip}(X)}$.

Our next result shows that every function in $H_\sigma(X)$ is Lipschitz continuous.

Lemma 6.6 *Let $X \subset \mathbb{R}^d$ be a non-empty set and $\sigma > 0$. Then every $f \in H_\sigma(X)$ is Lipschitz continuous with $|f|_1 \leq \sqrt{2}\sigma \|f\|_{H_\sigma(X)}$.*

Proof: Let us write $\Phi : X \rightarrow H_\sigma(X)$ for the canonical feature map defined by $\Phi(x) := k_\sigma(x, \cdot)$. Now recall that Φ satisfies the reproducing property

$$f(x) = \langle \Phi(x), f \rangle, \quad x \in X, f \in H_\sigma(X),$$

and hence in particular $k_\sigma(x', x) = \langle \Phi(x), \Phi(x') \rangle$ for all $x, x' \in X$. Using these equalities together with $1 - e^{-t} \leq t$ for $t \geq 0$ we obtain

$$\begin{aligned} |f(x) - f(x')| &= |\langle \Phi(x) - \Phi(x'), f \rangle| \leq \|f\|_{H_\sigma(X)} \cdot \|\Phi(x) - \Phi(x')\|_{H_\sigma(X)} \\ &= \|f\|_{H_\sigma(X)} \sqrt{\langle \Phi(x), \Phi(x) \rangle + \langle \Phi(x'), \Phi(x') \rangle - 2\langle \Phi(x), \Phi(x') \rangle} \\ &= \|f\|_{H_\sigma(X)} \sqrt{2 - 2\exp(-\sigma^2 \|x - x'\|_2^2)} \\ &\leq \sqrt{2}\sigma \|f\|_{H_\sigma(X)} \|x - x'\|_2, \end{aligned}$$

i.e., we have proved the assertion. \blacksquare

In the following we consider certain orthonormal bases (ONBs) of $H_\sigma(X)$. To this end let us first recall that in [39, Theorem 5] it was shown that $(e_n)_{n \geq 0}$, where $e_n : \mathbb{R} \rightarrow \mathbb{R}$ is defined by

$$e_n(x) := \sqrt{\frac{2^n \sigma^{2n}}{n!}} x^n e^{-\sigma^2 x^2}, \quad x \in \mathbb{R}, \quad (12)$$

forms an ONB of $H_\sigma(\mathbb{R})$. Moreover, it was shown that if $X \subset \mathbb{R}$ has a non-empty interior the restrictions of e_n to X form an ONB of $H_\sigma(X)$. The following lemma establishes upper bounds on $\|e_n\|_\infty$ if X is a closed interval.

Lemma 6.7 Let $\sigma > 0$ and $a > 0$ be fixed real numbers and $(e_n)_{n \geq 0}$ be the ONB of $H_\sigma([-a, a])$, where e_n is defined by the restriction of (12) to $[-a, a]$. Then we have $\|e_n\|_\infty \leq (2\pi n)^{-1/4}$ for all $n \geq 1$ and

$$\|e_n\|_\infty \leq \sqrt{\frac{2^n a^{2n} \sigma^{2n}}{n!}} e^{-a^2 \sigma^2} \quad (13)$$

for all $n \geq 2a^2 \sigma^2$. In addition, for $n \geq 8ea^2 \sigma^2$ we have

$$\left(\sum_{i=n+1}^{\infty} \|e_i\|_\infty^2 \right)^{1/2} \leq \left(\frac{2}{\pi(n+1)} \right)^{1/4} 2^{-(n+1)} e^{-a^2 \sigma^2}, \quad (14)$$

and for $a\sigma \geq 1$ we also have

$$\left(\sum_{i=0}^{\infty} \|e_i\|_\infty^2 \right)^{1/2} \leq \sqrt{6a\sigma}. \quad (15)$$

Proof: Elementary calculus shows

$$e'_n(x) = \sqrt{\frac{2^n \sigma^{2n}}{n!}} x^{n-1} e^{-\sigma^2 x^2} (n - 2\sigma^2 x^2)$$

for all $n \geq 1$ and $x \in \mathbb{R}$. From this we conclude $e'_n(x^*) = 0$ if and only if $x^* = \pm \sqrt{\frac{n}{2\sigma^2}}$. Now $n < 2a^2 \sigma^2$ implies $|x^*| < a$ and then it is not hard to see that $|e_n|$ attains its maximum at these x^* . Consequently, we obtain

$$\|e_n\|_\infty \leq \sqrt{\frac{n^n}{n!}} e^{-n/2} \leq \sqrt{\frac{n^n}{\sqrt{2\pi n n^n e^{-n}}}} e^{-n/2} = (2\pi n)^{-1/4}$$

by Stirling's formula. On the other hand, $n \geq 2a^2 \sigma^2$ implies $|x^*| \geq a$ and in this case it is not hard to see that $|e_n|$ attains its maximum at $\pm a$. From these considerations we conclude (13). For the first part of the lemma it thus remains to show that $\|e_n\|_\infty \leq (2\pi n)^{-1/4}$ for all $n \geq 2a^2 \sigma^2$. To this end we apply Stirling's formula to the already obtained (13). This yields

$$\begin{aligned} \|e_n\|_\infty &\leq \sqrt{\frac{2^n a^{2n} \sigma^{2n}}{n!}} e^{-a^2 \sigma^2} \leq \sqrt{\frac{2^n a^{2n} \sigma^{2n} e^n}{\sqrt{2\pi n n^n e^{-n}}}} e^{-a^2 \sigma^2} = (2\pi n)^{-1/4} \left(\frac{2a^2 \sigma^2 e^{1-2a^2 \sigma^2/n}}{n} \right)^{n/2} \\ &\leq (2\pi n)^{-1/4}, \end{aligned}$$

where in the last step we used that $te^{1-t} \leq 1$ for all $t \geq 0$.

For the proof of (14) we recall that the remainder of the Taylor series of the exponential function satisfies

$$\sum_{i=n+1}^{\infty} \frac{y^i}{i!} \leq 2 \frac{|y|^{n+1}}{(n+1)!}$$

for $|y| \leq 1 + n/2$. Since $n \geq 8ea^2 \sigma^2$ implies $2a^2 \sigma^2 \leq 1 + n/2$ we consequently obtain

$$\begin{aligned} \sum_{i=n+1}^{\infty} \|e_i\|_\infty^2 &\leq \sum_{i=n+1}^{\infty} \frac{2^i a^{2i} \sigma^{2i}}{i!} e^{-2a^2 \sigma^2} \leq \frac{2^{n+2} a^{2(n+1)} \sigma^{2(n+1)}}{(n+1)!} e^{-2a^2 \sigma^2} \\ &\leq 2 \frac{2^{n+1} a^{2(n+1)} \sigma^{2(n+1)} e^{(n+1)}}{\sqrt{2\pi(n+1)} (n+1)^{(n+1)}} e^{-2a^2 \sigma^2} \\ &= \left(\frac{2}{\pi(n+1)} \right)^{1/2} \left(\frac{2ea^2 \sigma^2}{n+1} \right)^{n+1} e^{-2a^2 \sigma^2} \\ &\leq \left(\frac{2}{\pi(n+1)} \right)^{1/2} 4^{-(n+1)} e^{-2a^2 \sigma^2}. \end{aligned}$$

From this we easily deduce (14). Finally, for the proof of (15) we observe that

$$\begin{aligned}
\sum_{i=0}^{\lceil 8ea^2\sigma^2 \rceil} \|e_i\|_\infty^2 &\leq 1 + (2\pi)^{-1/2} + \sum_{i=2}^{\lceil 8ea^2\sigma^2 \rceil} (2\pi i)^{-1/2} \\
&\leq 1 + (2\pi)^{-1/2} + (2\pi)^{-1/2} \int_1^{8ea^2\sigma^2+1} x^{-1/2} dx \\
&\leq 1 + (2\pi)^{-1/2} + (e/\pi)^{-1/2} 4a\sigma \\
&\leq 3/2 + 4a\sigma.
\end{aligned}$$

Combining this estimate with (14) we then obtain

$$\begin{aligned}
\sum_{i=0}^{\infty} \|e_i\|_\infty^2 &= \sum_{i=0}^{\lceil 8ea^2\sigma^2 \rceil} \|e_i\|_\infty^2 + \sum_{i=\lceil 8ea^2\sigma^2 \rceil+1}^{\infty} \|e_i\|_\infty^2 \\
&\leq 3/2 + 4a\sigma + \left(\frac{2}{\pi(\lceil 8ea^2\sigma^2 \rceil + 1)} \right)^{1/2} 4^{-(\lceil 8ea^2\sigma^2 \rceil+1)} e^{-2a^2\sigma^2} \\
&\leq 3/2 + 4a\sigma + \left(\frac{1}{8e\pi a^2\sigma^2} \right)^{1/2} 4^{-8ea^2\sigma^2} e^{-2a^2\sigma^2} \\
&\leq 2 + 4a\sigma,
\end{aligned}$$

and from the latter we easily obtain (15). ■

Our next goal is to generalize the above result to the multi-dimensional case. To this end recall that the tensor product $f \otimes g : X \times X \rightarrow \mathbb{R}$ of two functions $f, g : X \rightarrow \mathbb{R}$ is defined by $f \otimes g(x, x') := f(x)g(x')$, $x, x' \in X$. Obviously, for bounded functions we have $\|f \otimes g\|_\infty = \|f\|_\infty \|g\|_\infty$.

In the following we will deal with multi-indexes $\eta = (n_1, \dots, n_d) \in \mathbb{N}_0^d$ and use the notation $\eta \geq n$ if $n \in \mathbb{N}_0$ and $n_i \geq n$ for all $i = 1, \dots, d$. Moreover, for $\eta = (n_1, \dots, n_d) \in \mathbb{N}_0^d$ we write

$$e_\eta := e_{n_1} \otimes \dots \otimes e_{n_d},$$

where e_{n_i} is defined by (12). Then [39, Theorem 5] shows that $(e_\eta)_{\eta \in \mathbb{N}_0^d}$ is an ONB of $H_\sigma(\mathbb{R}^d)$ and the restrictions of the members of this ONB to $[-a, a]^d$ form an ONB of $H_\sigma([-a, a]^d)$. The following lemma generalizes the estimates of Lemma 6.7 to this multi-dimensional ONB.

Corollary 6.8 *For $\sigma > 0$ and $a > 0$ satisfying $a\sigma \geq 1$, and $d \in \mathbb{N}$ let $(e_\eta)_{\eta \in \mathbb{N}_0^d}$ be the restriction of the above ONB to $[-a, a]^d$. Then for $n \geq 8ea^2\sigma^2$ we have*

$$\left(\sum_{\substack{\eta \in \mathbb{N}_0^d \\ \exists i: \eta_i > n}} \|e_\eta\|_\infty^2 \right)^{1/2} \leq \sqrt{d} e^{-a^2\sigma^2} (6a\sigma)^{(d-1)/2} \left(\frac{2}{\pi(n+1)} \right)^{1/4} 2^{-(n+1)}.$$

Proof: Using $\|e_{i_1} \otimes \dots \otimes e_{i_d}\|_\infty = \|e_{i_1}\|_\infty \dots \|e_{i_d}\|_\infty$ we obtain

$$\begin{aligned}
\sum_{\substack{\eta \in \mathbb{N}_0^d \\ \exists i: \eta_i > n}} \|e_\eta\|_\infty^2 &\leq d \sum_{i_1=n+1}^{\infty} \sum_{i_2=0}^{\infty} \dots \sum_{i_d=0}^{\infty} \prod_{j=1}^d \|e_{i_j}\|_\infty^2 = d \left(\sum_{i=n+1}^{\infty} \|e_i\|_\infty^2 \right) \left(\sum_{i=0}^{\infty} \|e_i\|_\infty^2 \right)^{d-1} \\
&\leq d \left(\frac{2}{\pi(n+1)} \right)^{1/2} 2^{-2(n+1)} e^{-2a^2\sigma^2} (6a\sigma)^{d-1}
\end{aligned}$$

by Lemma 6.7. From this we immediately obtain the assertion. ■

6.3 Some concentration inequalities in RKHSs

In this subsection we will establish a concentration inequality for RKHS-valued functions and for processes which have a certain decay of correlations. This concentration result will then be the key ingredient in the statistical analysis of the proof of Theorem 3.3.

Let us begin by showing a simple concentration result in the form of Chebyshev's inequality for \mathbb{R} -valued functions and processes with a decay of correlations.

Lemma 6.9 *Let $\mathcal{Z} = (Z_i)_{i \geq 0}$ be a Z -valued process on $(\Omega, \mathcal{A}, \mu)$ that is stationary in the wide sense. Then for $P := \mu_{Z_0}$, $f \in L_2(P)$, $n \geq 1$, and $\delta > 0$ we have*

$$\mu\left(\omega \in \Omega : \left| \frac{1}{n} \sum_{i=0}^{n-1} f \circ Z_i(\omega) - \mathbb{E}_P f \right| \geq \delta\right) \leq \frac{2}{n\delta^2} \sum_{i=0}^{n-1} \text{cor}_{\mathcal{Z},i}(f, f). \quad (16)$$

Proof: By Markov's inequality we immediately obtain

$$\mu\left(\omega \in \Omega : \left| \frac{1}{n} \sum_{i=0}^{n-1} f \circ Z_i(\omega) - \mathbb{E}_P f \right| \geq \delta\right) \leq \frac{1}{n^2\delta^2} \mathbb{E}_\mu \left(\sum_{i=0}^{n-1} f \circ Z_i - \mathbb{E}_P f \right)^2.$$

Moreover, a simple calculation shows

$$\begin{aligned} \mathbb{E}_\mu \left(\sum_{i=0}^{n-1} f \circ Z_i - \mathbb{E}_P f \right)^2 &= \sum_{i=0}^{n-1} (\mathbb{E}_\mu (f \circ Z_i)^2 - (\mathbb{E}_P f)^2) + 2 \sum_{i=0}^{n-1} \sum_{j=0}^{i-1} (\mathbb{E}_\mu f \circ Z_i f \circ Z_j - \mathbb{E}_P f \mathbb{E}_P f) \\ &\leq 2n \sum_{i=0}^{n-1} \text{cor}_{\mathcal{Z},i}(f, f). \end{aligned}$$

Combining both estimates immediately yields the assertion. ■

We will see in Subsection 6.4 that the proof of Theorem 3.3 heavily relies on the estimate

$$\|f_{P,\lambda,H} - f_{T,\lambda,H}\|_H \leq \frac{1}{\lambda} \|\mathbb{E}_P h_\lambda \Phi - \mathbb{E}_T h_\lambda \Phi\|_H,$$

where $f_{P,\lambda,H}$ is the SVM solution (see Theorem 6.13 for an exact definition) one obtains by replacing the empirical risk $\mathcal{R}_{L,T}(\cdot)$ with the true risk $\mathcal{R}_{L,P}(\cdot)$ in (5), h_λ is a function independent of the training set T , $\Phi : X \rightarrow H$ is the canonical feature map $x \mapsto k(x, \cdot)$ of a kernel k , and \mathbb{E}_T denotes the expectation with respect to the empirical measure defined by T , i.e. $\mathbb{E}_T g := n^{-1} \sum_{i=1}^n g(x_i, y_i)$. Consequently, our next goal is to estimate terms of the form $\|\mathbb{E}_P h_\lambda \Phi - \mathbb{E}_T h_\lambda \Phi\|_H$. To this end we begin with the following lemma which, roughly speaking, will be used to reduce RKHS-valued functions to \mathbb{R} -valued functions.

Lemma 6.10 *Let H be the separable RKHS of a bounded measurable kernel $k : X \times X \rightarrow \mathbb{R}$, let $\Phi : X \rightarrow H$ be the corresponding canonical feature map and $(e_i)_{i \geq 0}$ be an ONB of H . Moreover, let Y be another measurable space, P and Q be probability measures on $X \times Y$, and $h \in L_1(P) \cap L_1(Q)$. Then for all $n \geq 0$ we have*

$$\|\mathbb{E}_P h \Phi - \mathbb{E}_Q h \Phi\|_H \leq \left(\sum_{i=0}^n |\mathbb{E}_P h e_i - \mathbb{E}_Q h e_i|^2 \right)^{1/2} + \left(\sum_{i=n+1}^{\infty} \|e_i\|_\infty^2 \right)^{1/2} (\mathbb{E}_P |h| + \mathbb{E}_Q |h|).$$

Proof: Let us define $S_n : H \rightarrow \text{span}\{e_0, \dots, e_n\}$ by $\sum_{i \geq 0} \langle f, e_i \rangle e_i \mapsto \sum_{i=0}^n \langle f, e_i \rangle e_i$. Then we have

$$\|S_n \Phi(x) - \Phi(x)\|_H^2 = \left\| \sum_{i=n+1}^{\infty} \langle \Phi(x), e_i \rangle e_i \right\|_H^2 = \sum_{i=n+1}^{\infty} |\langle \Phi(x), e_i \rangle|^2 = \sum_{i=n+1}^{\infty} |e_i(x)|^2$$

by the reproducing property and hence we obtain

$$\begin{aligned} \|\mathbb{E}_P h \Phi - \mathbb{E}_Q h \Phi\|_H &\leq \|\mathbb{E}_P h \Phi - \mathbb{E}_P h S_n \Phi\|_H + \|\mathbb{E}_P h S_n \Phi - \mathbb{E}_Q h S_n \Phi\|_H + \|\mathbb{E}_Q h S_n \Phi - \mathbb{E}_Q h \Phi\|_H \\ &\leq \mathbb{E}_P |h| \|\Phi - S_n \Phi\|_H + \|\mathbb{E}_P h S_n \Phi - \mathbb{E}_Q h S_n \Phi\|_H + \mathbb{E}_Q |h| \|\Phi - S_n \Phi\|_H \\ &\leq \|\mathbb{E}_P h S_n \Phi - \mathbb{E}_Q h S_n \Phi\|_H + \left(\sum_{i=n+1}^{\infty} \|e_i\|_{\infty}^2 \right)^{1/2} (\mathbb{E}_P |h| + \mathbb{E}_Q |h|). \end{aligned}$$

Moreover, using the reproducing property we have $\langle \mathbb{E}_P h \Phi, e_i \rangle = \mathbb{E}_P h e_i$ and $\langle \mathbb{E}_Q h \Phi, e_i \rangle = \mathbb{E}_Q h e_i$, and thus we conclude

$$\begin{aligned} \|\mathbb{E}_P h S_n \Phi - \mathbb{E}_Q h S_n \Phi\|_H^2 &= \left\| \sum_{i=0}^n \langle \mathbb{E}_P h \Phi - \mathbb{E}_Q h \Phi, e_i \rangle e_i \right\|_H^2 = \sum_{i=0}^n |\langle \mathbb{E}_P h \Phi - \mathbb{E}_Q h \Phi, e_i \rangle|^2 \\ &= \sum_{i=0}^n |\mathbb{E}_P h e_i - \mathbb{E}_Q h e_i|^2. \end{aligned}$$

Combining this equality with the previous estimate we then obtain the assertion. \blacksquare

Before we can establish the concentration inequality for RKHS-valued functions we finally need the following simple lemma.

Lemma 6.11 *For $d \geq 1$ and $t > 18 d \ln(d)$ we have $t^{-1/4} 2^{-t} \leq t^{-2d}$.*

Proof: Since $2^{-t} = e^{-t \ln 2}$ and $t^{-2d+1/4} = e^{(1/4-2d) \ln t}$ the assertion is equivalent to

$$t \ln 2 + (1/4 - 2d) \ln t \geq 0. \quad (17)$$

Let us first prove the case $d = 1$. Then (17) reduces to the assertion $h(t) := t \ln 2 - \frac{7}{4} \ln t \geq 0$. To establish the latter, note that we have $h'(t) = \ln 2 - \frac{7}{4} t^{-1}$ and hence $h'(t^*) = 0$ holds if and only if $t^* = \frac{7}{4 \ln 2}$. Simple considerations then show that h has its only global minimum at t^* and therefore we have $h(t) \geq h(t^*) \geq \frac{7}{4} - \frac{7}{4} \ln\left(\frac{7}{\ln 16}\right) > 0$.

Let us now consider the case $d \geq 2$. To this end we fix a $t > 18 d \ln(d)$. Then there exists a unique $x > 18$ with $t = x d \ln(d)$, and hence we obtain

$$\begin{aligned} t \ln 2 + (1/4 - 2d) \ln t &= x d \ln(d) \ln 2 + 1/4 \ln(x d \ln(d)) - 2d \ln(x d \ln(d)) \\ &> x d \ln(d) \ln 2 - 2d \ln(x d \ln(d)) \\ &= d(x \ln(d) \ln 2 - 2 \ln x - 2 \ln d - 2 \ln(\ln(d))) \\ &> d\left(x \ln(d) \ln 2 - 2 \frac{\ln d}{\ln 2} \ln x - 2 \ln d - 2 \ln d\right) \\ &= d \ln(d) \left(x \ln 2 - \frac{2}{\ln 2} \ln x - 4\right), \end{aligned}$$

where in the last estimate we used $d \geq 2$. Now it is elementary to check that $x \mapsto x \ln 2 - \frac{2}{\ln 2} \ln x - 4$ is increasing on $[2(\ln 2)^{-2}, \infty)$ and since $18 \ln 2 - \frac{2}{\ln 2} \ln 18 - 4 > 0$ we then obtain (17). \blacksquare

Theorem 6.12 For $\sigma > 0$ and $a > 0$ satisfying $a\sigma \geq 1$, and $d \geq 1$ let $\Phi : [-a, a]^d \rightarrow H_\sigma([-a, a]^d)$ be the canonical feature map of the Gaussian RBF kernel and let $(e_\eta)_{\eta \in \mathbb{N}_0^d}$ be the ONB of $H_\sigma([-a, a]^d)$ which is considered in Corollary 6.8. In addition, let Y be a measurable space and let $\mathcal{Z} = (X_i, Y_i)_{i \geq 0}$ be a $[-a, a]^d \times Y$ -valued process on $(\Omega, \mathcal{A}, \mu)$ that is stationary in the wide sense. Furthermore, let $(\gamma_i)_{i \geq 0}$ be a strictly positive null sequence, $h : [-a, a]^d \times Y \rightarrow \mathbb{R}$ be a bounded measurable function, and $K_h \in [1, \infty)$ be a constant such that

$$\|h\|_\infty \leq K_h \quad (18)$$

and

$$\text{cor}_{\mathcal{Z}, i}(he_\eta, he_\eta) \leq K_h \gamma_i \quad (19)$$

for all $i \geq 0$, $\eta \in \mathbb{N}_0^d$. Then for all $\epsilon > 0$ satisfying both $\epsilon \leq (1 + 8ea^2\sigma^2)^{-2d}$ and $\epsilon \leq (18d \ln d)^{-2d}$, and all $n \geq 1$ we have

$$\mu\left(\omega \in \Omega : \|\mathbb{E}_P h\Phi - \mathbb{E}_{T(\omega)} h\Phi\|_H \leq \epsilon\right) \geq 1 - \frac{2(1 + \frac{1}{8ea^2\sigma^2})^d K_h C_{a\sigma, d, h}^3}{n\epsilon^3} \sum_{i=0}^{n-1} \gamma_i,$$

where $\mathbb{E}_{T(\omega)}$ denotes the expectation operator with respect to the empirical measure associated to $T(\omega) := (Z_0(\omega), \dots, Z_{n-1}(\omega))$, i.e. $\mathbb{E}_{T(\omega)} g := \frac{1}{n} \sum_{i=0}^{n-1} g(X_i(\omega), Y_i(\omega))$, $\mathbb{E}_P h\Phi$ denotes the Bochner integral of $h\Phi$ in H , and

$$C_{a\sigma, d, h} := \left(1 + \frac{1}{8ea^2\sigma^2}\right)^{\frac{d}{2}} + 2\sqrt{d} e^{-a^2\sigma^2} (6a\sigma)^{\frac{d-1}{2}} K_h.$$

Proof: Let us write

$$\delta := \left(\frac{\epsilon}{C_{a\sigma, d, h}}\right)^{5/4}.$$

Using $C_{a\sigma, d, h} \geq (1 + \frac{1}{8ea^2\sigma^2})^{\frac{d}{2}} \geq 1$ and $\epsilon \leq (1 + 8ea^2\sigma^2)^{-2d}$ we then find

$$\delta \leq (1 + 8ea^2\sigma^2)^{-5d/2},$$

and consequently, there exists a natural number $m \geq 8ea^2\sigma^2$ such that $(m+1)^{-5d/2} \leq \delta < m^{-5d/2}$. Moreover, for later use we note that using $C_{a\sigma, d, h} \geq 1$ and $\epsilon \leq (18d \ln d)^{-2d}$ yields

$$\delta^{-\frac{2}{5d}} \geq 18d \ln d.$$

Let us now consider an $\omega \in \Omega$ such that

$$|\mathbb{E}_P h e_\eta - \mathbb{E}_{T(\omega)} h e_\eta| < \delta \quad (20)$$

for all $\eta \in \{0, \dots, m\}^d$. By Lemma 6.10 and Corollary 6.8 we then obtain

$$\begin{aligned} \|\mathbb{E}_P h\Phi - \mathbb{E}_{T(\omega)} h\Phi\|_H &\leq \left(\sum_{\eta \leq m} |\mathbb{E}_P h e_\eta - \mathbb{E}_{T(\omega)} h e_\eta|^2\right)^{1/2} \\ &\quad + \left(\sum_{\substack{\eta \in \mathbb{N}_0^d \\ \exists i: \eta_i > m}} \|e_\eta\|_\infty^2\right)^{1/2} (\mathbb{E}_P |h| + \mathbb{E}_{T(\omega)} |h|) \\ &\leq (m+1)^{d/2} \delta + 2\sqrt{d} e^{-a^2\sigma^2} (6a\sigma)^{(d-1)/2} \left(\frac{2}{\pi(m+1)}\right)^{1/4} 2^{-(m+1)} K_h \\ &\leq \left(1 + \frac{1}{8ea^2\sigma^2}\right)^{d/2} \delta^{4/5} + 2\sqrt{d} e^{-a^2\sigma^2} (6a\sigma)^{(d-1)/2} \delta^{\frac{1}{10d}} 2^{-\delta^{-\frac{2}{5d}}} K_h, \end{aligned}$$

where in the last step we used the inequalities $8ea^2\sigma^2 \leq m < \delta^{-\frac{2}{5d}} \leq m + 1$. Using Lemma 6.11 for $t := \delta^{-\frac{2}{5d}}$ we consequently obtain

$$\|\mathbb{E}_P h\Phi - \mathbb{E}_{T(\omega)} h\Phi\|_H \leq \left(\left(1 + \frac{1}{8ea^2\sigma^2}\right)^{d/2} + 2\sqrt{d}e^{-a^2\sigma^2}(6a\sigma)^{(d-1)/2}K_h \right) \delta^{4/5} = \epsilon.$$

Moreover, by Lemma 6.9 and a simple union bound argument we see that the probability of ω satisfying (20) for all $\eta \in \{0, \dots, m\}^d$ simultaneously is not smaller than

$$1 - \sum_{\eta \in \{0, \dots, m\}^d} \frac{2}{n\delta^2} \sum_{i=0}^{n-1} \text{cor}_{Z,i}(he_\eta, he_\eta).$$

In addition, we have

$$\sum_{\eta \in \{0, \dots, m\}^d} \frac{2}{n\delta^2} \sum_{i=0}^{n-1} \text{cor}_{Z,i}(he_\eta, he_\eta) \leq \sum_{\eta \in \{0, \dots, m\}^d} \frac{2}{n\delta^2} \sum_{i=0}^{n-1} K_h \gamma_i \leq \frac{2(m+1)^d}{n\delta^2} \sum_{i=0}^{n-1} K_h \gamma_i,$$

and since $21 < 8ea^2\sigma^2 \leq m < \delta^{-\frac{2}{5d}}$ we additionally have

$$\begin{aligned} \frac{2(m+1)^d}{n\delta^2} &\leq \frac{2(1 + \frac{1}{8ea^2\sigma^2})^d m^d}{n\delta^2} \leq \frac{2(1 + \frac{1}{8ea^2\sigma^2})^d}{n\delta^{12/5}} \\ &= \frac{2(1 + \frac{1}{8ea^2\sigma^2})^d \left((1 + \frac{1}{8ea^2\sigma^2})^{\frac{d}{2}} + 2\sqrt{d}e^{-a^2\sigma^2}(6a\sigma)^{\frac{d-1}{2}}K_h \right)^3}{n\epsilon^3}. \end{aligned}$$

Combining these estimates we then obtain the assertion. \blacksquare

6.4 Proof of Theorem 3.3

For the proof of Theorem 3.3 we need some final preparations. Let us begin with the following result on the existence and uniqueness of infinite sample SVMs which is a slight extension of similar results established in [18, 12]:

Theorem 6.13 *Let $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a convex, locally Lipschitz continuous loss function satisfying $L(x, y, 0) \leq 1$ for all $(x, y) \in (X \times Y)$, and let P be a distribution on $X \times Y$. Furthermore, let H be a RKHS of a bounded measurable kernel over X . Then for all $\lambda > 0$ there exists exactly one element $f_{P,\lambda,H} \in H$ such that*

$$\lambda \|f_{P,\lambda,H}\|_H^2 + \mathcal{R}_{L,P}(f_{P,\lambda,H}) = \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f). \quad (21)$$

Furthermore, we have $\|f_{P,\lambda,H}\|_H \leq \lambda^{-1/2}$.

Note that the above theorem in particular yields $\|f_{T,\lambda,H}\|_H \leq \lambda^{-1/2}$ by considering the empirical measure associated to a training set $T \in (X \times Y)^n$. The following result which was (essentially) shown in [18, 12] describes the stability of the empirical SVM solutions.

Theorem 6.14 *Let X be a separable metric space, $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a convex, locally Lipschitz continuous loss function satisfying $L(x, y, 0) \leq 1$ for all $(x, y) \in (X \times Y)$, and let P be a distribution on $X \times Y$. Furthermore, let H be the RKHS of continuous kernel $k : X \times X \rightarrow \mathbb{R}$*

satisfying $\|k\|_\infty \leq 1$ and let $\Phi : X \rightarrow H$ be the corresponding canonical feature map. Then for all $\lambda > 0$ the function $h_\lambda : X \times Y \rightarrow \mathbb{R}$ defined by

$$h_\lambda(x, y) := L'(x, y, f_{P,\lambda}(x)), \quad (x, y) \in X \times Y, \quad (22)$$

is bounded and satisfies and

$$\|f_{P,\lambda,H} - f_{T,\lambda,H}\|_H \leq \frac{1}{\lambda} \|\mathbb{E}_P h_\lambda \Phi - \mathbb{E}_T h_\lambda \Phi\|_H \quad (23)$$

for all training sets $T = ((x_0, y_0), \dots, (x_{n-1}, y_{n-1})) \in (X \times Y)^n$, where \mathbb{E}_T denotes the expectation operator with respect to the empirical measure associated to T , i.e. $\mathbb{E}_T g := \frac{1}{n} \sum_{i=0}^{n-1} g(x_i, y_i)$, and $\mathbb{E}_P h_\lambda \Phi$ is a Bochner integral in H .

Proof of Theorem 3.3: Obviously, it suffices to consider sets X of the form $X = [-a, a]^d$ for some $a \geq 1$. For $\sigma > 0$ and $\lambda > 0$ we write $h_{\lambda,\sigma}$ for the function we obtain by Theorem 6.14 for $H := H_\sigma(X)$. By the local Lipschitz continuity of L , $\|k_\sigma\|_\infty \leq 1$, Theorem 6.13, and (23) we then have

$$\begin{aligned} |\mathcal{R}_{L,P}(f_{T,\lambda,\sigma}) - \mathcal{R}_{L,P}(f_{P,\lambda,\sigma})| &\leq |L|_{\lambda^{-1/2},1} \|f_{T,\lambda,\sigma} - f_{P,\lambda,\sigma}\|_\infty \\ &\leq |L|_{\lambda^{-1/2},1} \|f_{T,\lambda,\sigma} - f_{P,\lambda,\sigma}\|_{H_\sigma(X)} \\ &\leq \frac{|L|_{\lambda^{-1/2},1}}{\lambda} \|\mathbb{E}_P h_{\lambda,\sigma} \Phi - \mathbb{E}_T h_{\lambda,\sigma} \Phi\|_H \end{aligned} \quad (24)$$

for all $\sigma > 0$, $\lambda > 0$ and all $T = ((x_0, y_0), \dots, (x_{n-1}, y_{n-1})) \in (X \times Y)^n$. Moreover, using (22), (6), and Lemma 6.6 we have

$$\begin{aligned} |h_{\lambda,\sigma}(x, y) - h_{\lambda,\sigma}(x', y')| &= |L'(x, y, f_{P,\lambda,\sigma}(x)) - L'(x', y', f_{P,\lambda,\sigma}(x'))| \\ &\leq c \cdot \left(|x - x'|^2 + |y - y'|^2 + |f_{P,\lambda,\sigma}(x) - f_{P,\lambda,\sigma}(x')|^2 \right)^{1/2} \\ &\leq c \cdot \left(|x - x'|^2 + |y - y'|^2 + 2\sigma^2 \|f_{P,\lambda,\sigma}\|_{H_\sigma(X)}^2 |x - x'|^2 \right)^{1/2} \\ &\leq c \cdot \left(|x - x'|^2 + |y - y'|^2 + 2\sigma^2 \lambda^{-1} |x - x'|^2 \right)^{1/2} \\ &\leq 2c\sigma\lambda^{-1/2} \|(x, y) - (x', y')\|_2 \end{aligned}$$

for all $\sigma \geq 1$, $\lambda \in (0, 1]$ and all $(x, y), (x', y') \in X \times Y$. Consequently, we find $|h_{\lambda,\sigma}|_1 \leq 2c\sigma\lambda^{-1/2}$. Moreover, by the remark after Assumption **L** we have

$$|h_{\lambda,\sigma}(x, y)| = |L'(x, y, f_{P,\lambda,\sigma}(x))| \leq c(1 + |f_{P,\lambda,\sigma}(x)|) \leq c(1 + \lambda^{-1/2}) \leq 2c\sigma\lambda^{-1/2}$$

for all $\sigma \geq 1$, $\lambda \in (0, 1]$ and all $(x, y) \in X \times Y$. Combining the last two estimates we thus have $\|h_{\lambda,\sigma}\|_{\text{Lip}(X \times Y)} \leq 2c\lambda^{-1/2}\sigma$ for all $\sigma \geq 1$, $\lambda \in (0, 1]$. By Lemma 6.6 we then find

$$\|h_{\lambda,\sigma} e_\eta^{(\sigma)}\|_{\text{Lip}(X \times Y)} \leq 2\|h_{\lambda,\sigma}\|_{\text{Lip}(X \times Y)} \|e_\eta^{(\sigma)}\|_{\text{Lip}(X \times Y)} \leq 4\sqrt{2}c\lambda^{-1/2}\sigma^2,$$

where $e_\eta^{(\sigma)}$ denotes the η -th element, $\eta \in \mathbb{N}_0^d$, of the ONB of $H_\sigma(X)$ considered in Corollary 6.8. Moreover, by Corollary 6.4 we may assume without loss of generality that $\kappa_{\psi,\varphi}$ is of the form

$$\kappa_{\psi,\varphi} = c_Z \|\psi\|_{\text{Lip}(X \times Y)} \|\varphi\|_{\text{Lip}(X \times Y)},$$

where $c_{\mathcal{Z}}$ is a constant only depending on \mathcal{Z} and (γ_i) . Consequently, we obtain

$$|\text{cor}_{\mathcal{Z},i}(h_{\lambda,\sigma}e_{\eta}^{(\sigma)}, h_{\lambda,\sigma}e_{\eta}^{(\sigma)})| \leq 32c_{\mathcal{Z}}c^2\lambda^{-1}\sigma^4$$

for all $\sigma \geq 1$, $\lambda \in (0, 1]$, and $\eta \in \mathbb{N}_0^d$. Since, in addition, we have $\|h_{\lambda,\sigma}\|_{\infty} \leq 2c\sigma\lambda^{-1/2}$ we see that (18) and (19) are satisfied for $K_{h_{\lambda,\sigma}} := \tilde{c}\lambda^{-1}\sigma^4$, where $\sigma \geq 1$, $\lambda \in (0, 1]$, and \tilde{c} is a constant independent of λ and σ . For $n \geq 1$ and $\epsilon > 0$ satisfying both

$$\epsilon \leq (1 + 8ea^2\sigma^2)^{-2d}|L|_{\lambda^{-1/2},1}\lambda^{-1} \quad (25)$$

$\epsilon \leq (18d \ln d)^{-2d}$, Theorem 6.12 together with (24) thus yields

$$\begin{aligned} & \mu\left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega),\lambda,\sigma}) - \mathcal{R}_{L,P}(f_{P,\lambda,\sigma})| > \epsilon\right) \\ & \leq \mu\left(\omega \in \Omega : \|\mathbb{E}_P h\Phi - \mathbb{E}_{T(\omega)} h\Phi\|_H > \frac{\lambda\epsilon}{|L|_{\lambda^{-1/2},1}}\right) \\ & \leq \frac{2\tilde{c}\left(1 + \frac{1}{8ea^2\sigma^2}\right)^{\frac{d}{2}} \tilde{C}_{\lambda,\sigma,d,a}^3 |L|_{\lambda^{-1/2},1}^3 \sigma^4}{\epsilon^3 n \lambda^4} \sum_{i=0}^{n-1} \gamma_i, \end{aligned} \quad (26)$$

where $\tilde{C}_{\lambda,\sigma,d,a} := \left(1 + \frac{1}{8ea^2\sigma^2}\right)^{\frac{d}{2}} + 2c\sqrt{d}e^{-a^2\sigma^2}(6a\sigma)^{\frac{d-1}{2}}\lambda^{-1}\sigma^4$. Now the last condition of Assumption **S1** or **S2** implies that $\lambda_n \geq n^{-1}$ for sufficiently large n . Using the assumption $\sigma_n \geq \ln(n+1)$ contained in both **S1** and **S2** we hence obtain

$$\begin{aligned} \tilde{C}_{\lambda_n,\sigma_n,d,a} &= \left(1 + \frac{1}{8ea^2\sigma_n^2}\right)^{\frac{d}{2}} + 2c\sqrt{d}e^{-a^2\sigma_n^2}(6a\sigma_n)^{\frac{d+7}{2}}\lambda_n^{-1} \\ &\leq \left(1 + \frac{1}{8ea^2(\ln(n+1))^2}\right)^{\frac{d}{2}} + 2c\sqrt{d}n^{-a^2(\ln(n+1))^2}(6a \ln(n+1))^{\frac{d+7}{2}}n \end{aligned}$$

for all sufficiently large n . Since the last expression tends to 1 for $n \rightarrow \infty$ we then obtain $\lim_{n \rightarrow \infty} \tilde{C}_{\lambda_n,\sigma_n,d,a} = 1$. Analogously, we obtain $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{8ea^2\sigma_n^2}\right)^{d/2} = 1$. Let us now consider the case where Assumption **S1** is fulfilled. Then we have

$$(1 + 8ea^2\sigma_n^2)^{-2d}|L|_{\lambda_n^{-1/2},1}\lambda_n^{-1} \rightarrow \infty$$

and hence Condition (25) is satisfied for all sufficiently large n . Moreover, by the remark after Assumption **L** we have $|L|_{\lambda^{-1/2},1} \leq c(1 + \lambda^{-1/2})$ for all $\lambda > 0$ and hence the first assumption of **S1** implies $\lim_{n \rightarrow \infty} \lambda_n \sigma_n^d = 0$. By Lemma 6.5 we thus find

$$\lim_{i \rightarrow \infty} \mathcal{R}_{L,P}(f_{P,\lambda_n,\sigma_n}) = \mathcal{R}_{L,P}^*.$$

Consequently, (26) together with the above proven limit relations $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{8ea^2\sigma_n^2}\right)^{d/2} = 1$ and $\lim_{n \rightarrow \infty} \tilde{C}_{\lambda_n,\sigma_n,d,a} = 1$ shows that for sufficiently large n we have

$$\begin{aligned} & \mu\left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega),\lambda_n,\sigma_n}) - \mathcal{R}_{L,P}^*| > 2\epsilon\right) \\ & \leq \mu\left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega),\lambda_n,\sigma_n}) - \mathcal{R}_{L,P}(f_{P,\lambda_n,\sigma_n})| > \epsilon\right) \\ & \leq 4\tilde{c} \frac{|L|_{\lambda_n^{-1/2},1}^3 \sigma_n^4}{\epsilon^3 n \lambda_n^4} \sum_{i=0}^{n-1} \gamma_i, \end{aligned}$$

and hence we obtain the assertion by the last condition of **S1**.
Conversely, if Assumption **S2** is fulfilled we have

$$\epsilon_n := (1 + 8ea^2\sigma_n^2)^{-2d} |L|_{\lambda_n^{-1/2}, 1} \lambda_n^{-1} \rightarrow 0,$$

and consequently for a fixed $\epsilon > 0$ we have $\epsilon_n \leq \epsilon$ for all sufficiently large n . Analogously to the first case we thus find

$$\begin{aligned} & \mu\left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega),\lambda_n,\sigma_n}) - \mathcal{R}_{L,P}^*| > 2\epsilon\right) \\ & \leq \mu\left(\omega \in \Omega : |\mathcal{R}_{L,P}(f_{T(\omega),\lambda_n,\sigma_n}) - \mathcal{R}_{L,P}(f_{P,\lambda_n,\sigma_n})| > \epsilon_n\right) \\ & \leq 4\tilde{c} \frac{|L|_{\lambda_n^{-1/2}, 1}^3 \sigma_n^4}{\epsilon_n^3 n \lambda_n^4} \sum_{i=0}^{n-1} \gamma_i \\ & \leq 2^{31d+2} c \frac{|L|_{\lambda_n^{-1/2}, 1}^6 \sigma_n^{4+12d}}{n \lambda_n} \sum_{i=0}^{n-1} \gamma_i \end{aligned}$$

for all sufficiently large n , and hence we obtain the assertion by the last condition of **S2**. \blacksquare

7 Proof of Theorem 4.1

For the proof of Theorem 4.1 we need to bound the correlation sequences for stochastic processes which are the sum of a dynamical system and an observational noise process. This is the goal of the following results. We begin with a lemma which computes the correlation of a joint process from the correlations of its components.

Lemma 7.1 *Let $\mathcal{X} = (X_i)_{i \geq 0}$ be an X -valued, identically distributed stochastic process defined on $(\Omega, \mathcal{A}, \mu)$ and $\mathcal{Y} = (Y_i)_{i \geq 0}$ be a Y -valued, identically distributed stochastic process defined on $(\Theta, \mathcal{B}, \nu)$. Then the stochastic process $\mathcal{Z} = (Z_i)_{i \geq 0}$ defined on $(\Omega \times \Theta, \mathcal{A} \otimes \mathcal{B}, \mu \otimes \nu)$ by $Z_i := (X_i, Y_i)$ is identically distributed with $P := (\mu \otimes \nu)_{Z_0} = \mu_{X_0} \otimes \nu_{Y_0}$. Moreover, for $\psi, \varphi \in L_2(P)$ the i -th coordinate of the correlation sequence of \mathcal{Z} is given by*

$$\text{cor}_{\mathcal{Z},i}(\psi, \varphi) = \mathbb{E}_\nu \text{cor}_{\mathcal{X},i}(\psi(\cdot, Y_0), \varphi(\cdot, Y_i)) + \mathbb{E}_\mu \mathbb{E}_\mu \text{cor}_{\mathcal{Y},i}(\psi(X_0, \cdot), \varphi(X'_0, \cdot)),$$

where X'_0 is an independent copy of X_0 .

Proof: The first assertion regarding P is obvious. For the second assertion we fix an independent copy $\mathcal{X}' = (X'_i)_{i \geq 0}$ of \mathcal{X} . Then an easy calculation using the fact that both \mathcal{X} and \mathcal{Y} are identically distributed yields

$$\begin{aligned} \text{cor}_{\mathcal{Z},i}(\psi, \varphi) &= \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X_i, Y_i) - \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \cdot \mathbb{E}_\mu \mathbb{E}_\nu \varphi(X_0, Y_0) \\ &= \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X_i, Y_i) - \mathbb{E}_\mu \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X'_0, Y_i) \\ &\quad + \mathbb{E}_\mu \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X'_0, Y_i) - \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \cdot \mathbb{E}_\mu \mathbb{E}_\nu \varphi(X_0, Y_0) \\ &= \mathbb{E}_\nu (\mathbb{E}_\mu \psi(X_0, Y_0) \varphi(X_i, Y_i) - \mathbb{E}_\mu \mathbb{E}_\mu \psi(X_0, Y_0) \varphi(X'_0, Y_i)) \\ &\quad + \mathbb{E}_\mu \mathbb{E}_\mu \mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X'_0, Y_i) - \mathbb{E}_\mu \mathbb{E}_\mu (\mathbb{E}_\nu \psi(X_0, Y_0) \cdot \mathbb{E}_\nu \varphi(X'_0, Y_0)) \\ &= \mathbb{E}_\nu (\mathbb{E}_\mu \psi(X_0, Y_0) \varphi(X_i, Y_i) - \mathbb{E}_\mu \psi(X_0, Y_0) \cdot \mathbb{E}_\mu \varphi(X_0, Y_i)) \\ &\quad + \mathbb{E}_\mu \mathbb{E}_\mu (\mathbb{E}_\nu \psi(X_0, Y_0) \varphi(X'_0, Y_i) - \mathbb{E}_\nu \psi(X_0, Y_0) \cdot \mathbb{E}_\nu \varphi(X'_0, Y_0)) \\ &= \mathbb{E}_\nu \text{cor}_{\mathcal{X},i}(\psi(\cdot, Y_0), \varphi(\cdot, Y_i)) + \mathbb{E}_\mu \mathbb{E}_\mu \text{cor}_{\mathcal{Y},i}(\psi(X_0, \cdot), \varphi(X'_0, \cdot)), \end{aligned}$$

i.e., we have proved the assertion. \blacksquare

The following elementary lemma establishes the Lipschitz continuity of a certain type of function which is important when considering the process that generates noisy observations of a dynamical system.

Lemma 7.2 *Let $M \subset \mathbb{R}^d$ be a compact subset and $F : M \rightarrow M$ be a Lipschitz continuous map. For $B > 0$ and a fixed $j \in \{1, \dots, d\}$ we write $X := M + [-B, B]^d$ and $Y := \pi_j(X)$, where $\pi_j : \mathbb{R}^d \rightarrow \mathbb{R}$ denotes the j -th coordinate projection. For $h \in \text{Lip}(X \times Y)$ we define the function $\bar{h} : M \times [-B, B]^{2d} \rightarrow \mathbb{R}$ by*

$$\bar{h}(x, \varepsilon_0, \varepsilon_1) := h(x + \varepsilon_0, \pi_j(F(x) + \varepsilon_1)), \quad x \in M, \varepsilon_0, \varepsilon_1 \in [-B, B]^d. \quad (27)$$

Then for all $x \in M$ and $\varepsilon_0, \varepsilon_1 \in [-B, B]^d$ we have

$$\begin{aligned} \|\bar{h}(x, \cdot, \cdot)\|_{\text{Lip}([-B, B]^{2d})} &\leq (1 + \|F\|_{\text{Lip}(M)}) \|h\|_{\text{Lip}(X \times Y)} \\ \|\bar{h}(\cdot, \varepsilon_0, \varepsilon_1)\|_{\text{Lip}(M)} &\leq \|h\|_{\text{Lip}(X \times Y)}. \end{aligned}$$

Proof: For the first assertion we observe that for $(\varepsilon_0, \varepsilon_1), (\varepsilon'_0, \varepsilon'_1) \in [-B, B]^d \times [-B, B]^d$ we obviously have

$$\begin{aligned} &|h(x + \varepsilon_0, \pi_j(F(x) + \varepsilon_1)) - h(x + \varepsilon'_0, \pi_j(F(x) + \varepsilon'_1))| \\ &\leq \|h\|_{\text{Lip}(X \times Y)} (\|\varepsilon_0 - \varepsilon'_0\|_2^2 + |\pi_j(F(x) + \varepsilon_1) - \pi_j(F(x) + \varepsilon'_1)|^2)^{1/2} \\ &\leq \|h\|_{\text{Lip}(X \times Y)} (\|\varepsilon_0 - \varepsilon'_0\|_2^2 + \|\varepsilon_1 - \varepsilon'_1\|_2^2)^{1/2} \\ &= \|h\|_{\text{Lip}(X \times Y)} \|(\varepsilon_0, \varepsilon_1) - (\varepsilon'_0, \varepsilon'_1)\|_2. \end{aligned}$$

Analogously, for $x, x' \in M$ we have

$$\begin{aligned} &|h(x + \varepsilon_0, \pi_j(F(x) + \varepsilon_1)) - h(x' + \varepsilon_0, \pi_j(F(x') + \varepsilon_1))| \\ &\leq \|h\|_{\text{Lip}(X \times Y)} (\|x - x'\|_2^2 + |\pi_j(F(x) + \varepsilon_1) - \pi_j(F(x') + \varepsilon_1)|^2)^{1/2} \\ &= \|h\|_{\text{Lip}(X \times Y)} (\|x - x'\|_2^2 + \|F(x) - F(x')\|_2^2)^{1/2} \\ &\leq \|h\|_{\text{Lip}(X \times Y)} (1 + \|F\|_{\text{Lip}(M)}) \|x - x'\|_2. \end{aligned}$$

From these estimates we easily obtain the assertions. ■

The following theorem bounds the correlation sequence for functions defined by (27). It will be the key to apply Theorem 3.3 in the proof of Theorem 4.1.

Theorem 7.3 *Let $M \subset \mathbb{R}^d$ be a compact subset and $F : M \rightarrow M$ be a Lipschitz continuous map such that the dynamical system $\mathcal{X} := (F^i)_{i \geq 0}$ has an ergodic measure μ . Moreover, let $\gamma = (\gamma_i)_{i \geq 0}$ be a strictly positive null sequence such that*

$$\text{cor}_{\mathcal{X}}(\psi, \varphi) \in \Lambda(\gamma), \quad \psi, \varphi \in \text{Lip}(M). \quad (28)$$

Furthermore, let $\mathcal{E} = (\varepsilon_i)_{i \geq 0}$ be an in the wide sense stationary $[-B, B]^d$ -valued stochastic process on $(\Theta, \mathcal{B}, \nu)$ such that the $[-B, B]^{2d}$ -valued process $\mathcal{Y} = (Y_i)_{i \geq 0}$ on $(\Theta, \mathcal{B}, \nu)$ that is defined by $Y_i(\vartheta) = (\varepsilon_i(\vartheta), \varepsilon_{i+1}(\vartheta))$, $i \geq 0$, $\vartheta \in \Theta$, satisfies

$$\text{cor}_{\mathcal{Y}}(\psi, \varphi) \in \Lambda(\gamma), \quad \psi, \varphi \in \text{Lip}([-B, B]^{2d}). \quad (29)$$

For a fixed $j \in \{1, \dots, d\}$ we write $X := M + [-B, B]^d$ and $Y := \pi_j(X)$. Moreover, we define the process $\bar{Z} = (Z_i)_{i \geq 0}$ on $(\Omega \times \Theta, \mathcal{A} \otimes \mathcal{B}, \mu \otimes \nu)$ by $\bar{Z}_i = (F^i, \varepsilon_i, \varepsilon_{i+1})$, $i \geq 0$. Then for all $\psi, \varphi \in \text{Lip}(X \times Y)$ we have

$$\text{cor}_{\bar{Z}}(\bar{\psi}, \bar{\varphi}) \in \Lambda(\gamma),$$

where $\bar{\psi}$ and $\bar{\varphi}$ are defined by (27).

Proof: Let $c_{\mathcal{X}}$ and $c_{\mathcal{Y}}$ be the constants we obtain by applying (28) and (29) to Theorem 6.2. Moreover, since \mathcal{E} is stationary in the wide sense we observe that \mathcal{Y} is identically distributed. Applying Lemma 7.1 to the processes \mathcal{X} and \mathcal{Y} then yields

$$\begin{aligned} |\text{cor}_{\bar{Z},i}(\bar{\psi}, \bar{\varphi})| &\leq \left| \mathbb{E}_{\nu} \text{cor}_{\mathcal{X},i}(\bar{\psi}(\cdot, Y_0), \bar{\varphi}(\cdot, Y_i)) \right| + \left| \mathbb{E}_{x \sim \mu} \mathbb{E}_{x' \sim \mu} \text{cor}_{\mathcal{Y},i}(\bar{\psi}(F^0(x), \cdot), \bar{\varphi}(F^0(x'), \cdot)) \right| \\ &\leq c_{\mathcal{X}} \mathbb{E}_{\nu} \|\bar{\psi}(\cdot, \varepsilon_0, \varepsilon_1)\|_{\text{Lip}(M)} \|\bar{\varphi}(\cdot, \varepsilon_i, \varepsilon_{i+1})\|_{\text{Lip}(M)} \cdot \gamma_i \\ &\quad + c_{\mathcal{Y}} \mathbb{E}_{x \sim \mu} \mathbb{E}_{x' \sim \mu} \|\bar{\psi}(x, \cdot)\|_{\text{Lip}([-B, B]^{2d})} \|\bar{\varphi}(x', \cdot)\|_{\text{Lip}([-B, B]^{2d})} \cdot \gamma_i \\ &\leq c_{\mathcal{X}} \|\psi\|_{\text{Lip}(X \times Y)} \|\varphi\|_{\text{Lip}(X \times Y)} \gamma_i \\ &\quad + c_{\mathcal{Y}} (1 + \|F\|_{\text{Lip}(M)}) \|\psi\|_{\text{Lip}(X \times Y)} \|\varphi\|_{\text{Lip}(X \times Y)} \cdot \gamma_i, \end{aligned}$$

where in the last step we used Lemma 7.2. ■

Note that using the estimate of Theorem 7.3 in Lemma 6.9 we need that the above process \mathcal{Y} is stationary in the wide sense. Obviously, the latter is satisfied if the process \mathcal{E} is stationary.

Proof of Theorem 4.1: For a fixed $j \in \{1, \dots, d\}$ we write $X := M + [-B, B]^d$ and $Y := \pi_j(X)$. Moreover, we define the $X \times Y$ -valued process $\mathcal{Z} = (X_i, Y_i)_{i \geq 0}$ on $(M \times [-B, B]^{d\mathbb{N}}, \mu \otimes \nu)$ by $X_i := F^i + \pi_0 \circ S^i$ and $Y_i := \pi_j(F^{i+1} + \pi_0 \circ S^{i+1})$, and in addition, we write $P_j := (\mu \otimes \nu)_{(X_0, Y_0)}$. Let us further consider the $M \times [-B, B]^{2d}$ -valued stationary process $\bar{\mathcal{Z}} := (F^i, \pi_0 \circ S^i, \pi_0 \circ S^{i+1})$ which is defined on $(M \times [-B, B]^{d\mathbb{N}}, \mu \otimes \nu)$. For $\psi, \varphi \in \text{Lip}(X \times Y)$ Theorem 7.3 together with our decay of correlations assumptions then shows

$$|\text{cor}_{\bar{\mathcal{Z}},i}(\bar{\psi}, \bar{\varphi})| \leq \kappa_{\psi, \varphi} \gamma_i$$

for all $i \geq 0$, where $\kappa_{\psi, \varphi} \in [0, \infty)$ is a constant independent of i . Moreover, our construction ensures $\text{cor}_{\mathcal{Z},i}(\psi, \varphi) = \text{cor}_{\bar{\mathcal{Z}},i}(\bar{\psi}, \bar{\varphi})$ for all $i \geq 0$ and hence Theorem 3.3 yields

$$\lim_{n \rightarrow \infty} \mu \otimes \nu \left((x, \varepsilon) \in M \times [-B, B]^{d\mathbb{N}} : |\mathcal{R}_{L, P_j}(f_{T_j(x, \varepsilon), \lambda_n, \sigma_n}) - \mathcal{R}_{L, P_j}^*| > \epsilon \right) = 0$$

for all $\epsilon > 0$. Using Assumption **LD** and the definition (10) we then easily obtain the assertion. ■

Appendix: Proof of Theorem 6.2

In the following B_E denotes the closed unit ball of a Banach space E . Recall that a linear operator $S : E \rightarrow F$ acting between two Banach spaces E and F is continuous if and only if it is *bounded*, i.e.,

$$\|S\| := \sup_{x \in B_E} \|Sx\| < \infty.$$

Our first goal is to recall another equivalent condition which in practice is often easier to check. To this end we need the following definition.

Definition 7.4 Let E and F be Banach spaces and $S : E \rightarrow F$ be a linear map. Then S is said to have a closed graph if for all $x \in E$, $y \in F$ and all sequences $(x_n) \subset E$ satisfying $x_n \rightarrow x$ and $Sx_n \rightarrow y$ we have $Sx = y$.

Obviously, every continuous linear operator has a closed graph. The following fundamental theorem from functional analysis shows the converse implication.

Theorem 7.5 (Closed Graph Theorem) Let E and F be Banach spaces and $S : E \rightarrow F$ be a linear map that has a closed graph. Then S is continuous.

Our next goal is to establish an analogous result for bilinear maps. To this end we first recall the principle of uniform boundedness which is also known as Banach-Steinhaus Theorem.

Theorem 7.6 (Principle of Uniform Boundedness) Let E and F be Banach spaces, A be a non-empty set, and $S_\alpha : E \rightarrow F$, $\alpha \in A$, be bounded linear operators. If the family $(S_\alpha)_{\alpha \in A}$ satisfies

$$\sup_{\alpha \in A} \|S_\alpha x\| < \infty$$

for all $x \in E$ then we actually have

$$\sup_{\alpha \in A} \sup_{x \in B_E} \|S_\alpha x\| < \infty.$$

Let us now recall that a map $S : E_1 \times E_2 \rightarrow F$ between Banach spaces E_1 , E_2 , and F is called *bilinear* if the maps $S(x_1, \cdot) : E_2 \rightarrow F$ and $S(\cdot, x_2) : E_1 \rightarrow F$ are linear for all $x_1 \in E_1$ and $x_2 \in E_2$. In order to state a closed graph theorem for bilinear maps we also need a notion which describes a closed graph property for bilinear maps:

Definition 7.7 Let E_1 , E_2 , and F be Banach spaces and $S : E_1 \times E_2 \rightarrow F$ be a bilinear map. Then S is said to have a *partially closed graph* if the maps $S(x_1, \cdot) : E_2 \rightarrow F$ and $S(\cdot, x_2) : E_1 \rightarrow F$ have closed graphs for all $x_1 \in E_1$ and $x_2 \in E_2$.

With these preparations we can now state and prove the announced closed graph theorem for bilinear maps:

Theorem 7.8 Let E_1 , E_2 , and F be Banach spaces and $S : E_1 \times E_2 \rightarrow F$ be a bilinear map that has a partially closed graph. Then there exists a constant $c \in [0, \infty)$ such that

$$\|S(x_1, x_2)\|_F \leq c \|x_1\|_{E_1} \cdot \|x_2\|_{E_2}$$

for all $x_1 \in E_1$ and $x_2 \in E_2$.

Proof: By applying the Closed Graph Theorem we see that the maps $S(x_1, \cdot) : E_2 \rightarrow F$ and $S(\cdot, x_2) : E_1 \rightarrow F$ are bounded linear operators for all $x_1 \in E_1$ and $x_2 \in E_2$. In particular, the boundedness of the operators $S(\cdot, x_2) : E_1 \rightarrow F$ yields

$$\sup_{x_1 \in B_{E_1}} \|S(x_1, x_2)\| < \infty$$

for all $x_2 \in E_2$. Applying the Principle of Uniform Boundedness to the family of bounded operators $(S(x_1, \cdot))_{x_1 \in B_{E_1}}$ thus shows

$$c := \sup_{x_1 \in B_{E_1}} \sup_{x_2 \in B_{E_2}} \|S(x_1, x_2)\| < \infty.$$

Using the bi-linearity of S we then obtain the assertion. ■

With these preparations we can now present the proof of Theorem 6.2.

Proof of Theorem 6.2: Obviously, $\text{cor}_{\mathcal{Z}} : E_1 \times E_2 \rightarrow F$ is a well-defined bilinear operator. In view of Theorem 7.8 it suffices to show that this operator has a partially closed graph. We begin by showing that $\text{cor}_{\mathcal{Z}}(\psi, \cdot) : E_2 \rightarrow F$ has a closed graph for all $\psi \in E_1$. To this end let us fix some $\psi \in E_1$, $\varphi \in E_2$, a sequence $b := (b_n)_{n \geq 0} \in F$ and a sequence $(\varphi_i)_{i \geq 1} \subset E_2$ such that

$$\begin{aligned} \lim_{i \rightarrow \infty} \|\varphi_i - \varphi\|_{E_2} &= 0 \\ \lim_{i \rightarrow \infty} \|\text{cor}_{\mathcal{Z}}(\psi, \varphi_i) - b\|_F &= 0. \end{aligned} \tag{30}$$

Obviously, $\text{cor}_{\mathcal{Z}}(\psi, \cdot) : E_2 \rightarrow F$ has a closed graph if $\text{cor}_{\mathcal{Z}}(\psi, \varphi) = b$. To show this equality we first observe that for fixed $n \geq 0$ and $i \rightarrow \infty$ we have

$$\left| \int_{\mathcal{Z}} \varphi_i dP - \int_{\mathcal{Z}} \varphi dP \right| \leq \|\varphi - \varphi_i\|_{L_1(P)} \leq \|\text{id} : E_2 \rightarrow L_2(P)\| \cdot \|\varphi - \varphi_i\|_{E_2} \rightarrow 0$$

and

$$\begin{aligned} \left| \int_{\Omega} \psi(Z_0) \cdot \varphi(Z_n) d\mu - \int_{\Omega} \psi(Z_0) \cdot \varphi_i(Z_n) d\mu \right| &\leq \|\psi\|_{L_2(P)} \cdot \|\varphi - \varphi_i\|_{L_2(P)} \\ &\leq \|\psi\|_{L_2(P)} \cdot \|\text{id} : E_2 \rightarrow L_2(P)\| \cdot \|\varphi - \varphi_i\|_{E_2} \rightarrow 0. \end{aligned}$$

From this we conclude $\lim_{i \rightarrow \infty} \text{cor}_{\mathcal{Z},n}(\psi, \varphi_i) = \text{cor}_{\mathcal{Z},n}(\psi, \varphi)$ for the n -th coordinate of sequences of correlations. Moreover, F is continuously included in ℓ_{∞} and hence (30) implies $\lim_{i \rightarrow \infty} \text{cor}_{\mathcal{Z},n}(\psi, \varphi_i) = b_n$ for all $n \geq 0$. Combining these considerations yields $\text{cor}_{\mathcal{Z},n}(\psi, \varphi) = b_n$ for all $n \geq 0$, i.e. we have shown that $\text{cor}_{\mathcal{Z}}(\psi, \cdot) : E_2 \rightarrow F$ has a closed graph. Since the fact that all $\text{cor}_{\mathcal{Z}}(\cdot, \varphi) : E_1 \rightarrow F$ have a closed graph can be shown completely analogous the proof is completed. \blacksquare

References

- [1] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68:337–404, 1950.
- [2] V. Baladi. *Positive Transfer Operators and Decay of Correlations*. World Scientific, Singapore, 2000.
- [3] V. Baladi. Decay of correlations. In *1999 AMS Summer Institute on Smooth Ergodic Theory and Applications*, pages 297–325. AMS, 2001.
- [4] V. Baladi, M. Benedicks, and V. Maume-Deschamps. Almost sure rates of mixing for i.i.d. unimodal maps. *Ann. E.N.S.*, 35:77–126, 2002.
- [5] V. Baladi, A. Kondah, and B Schmitt. Random correlations for small perturbations of expanding maps. *Random Comput. Dynam.*, 4:179–204, 1996.
- [6] P.L. Bartlett, M.I. Jordan, and J.D. McAuliffe. Convexity, classification, and risk bounds. *J. Amer. Statist. Assoc.*, 101:138–156, 2006.
- [7] B.E. Boser, I.Guyon, and V. Vapnik. A training algorithm for optimal margin classifiers. In *Computational Learning Theory*, pages 144–152, 1992.

- [8] D. Bosq. *Nonparametric Statistics for Stochastic Processes*. Springer, New York, 2nd edition, 1998.
- [9] R. Bowen. *Equilibrium States and the Ergodic Theory of Anosov Diffeomorphisms*. Springer, Berlin, 1975.
- [10] R.C. Bradley. Basic properties of strong mixing conditions. A survey and some open questions. *Probability Surveys*, 2:107–144, 2005.
- [11] R.C. Bradley. *Introduction to Strong Mixing Conditions*, volume 1-3. Technical Report, Department of Mathematics, Indiana University, Bloomington, Custom Publishing of I.U., Bloomington, 2005.
- [12] A. Christmann and I. Steinwart. Consistency and robustness of kernel based regression. *Bernoulli*, to appear, 2007. http://www.c3.lanl.gov/ml/pubs/2005_regression/paper.pdf.
- [13] P. Collet. A remark about uniform de-correlation prefactors. Technical report, 1999.
- [14] P. Collet, S. Martinez, and B. Schmitt. Exponential inequalities for dynamical measures of expanding maps of the interval. *Probab. Theory Related Fields*, 123:301–322, 2002.
- [15] C. Cortes and V. Vapnik. Support vector networks. *Machine Learning*, 20:273–297, 1995.
- [16] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines*. Cambridge University Press, 2000.
- [17] M. Davies. Noise reduction schemes for chaotic time series. *Physica D. Nonlinear Phenomena*, 79:174–192, 1994.
- [18] E. DeVito, L. Rosasco, A. Caponnetto, M. Piana, and A. Verri. Some properties of regularized kernel methods. *Journal of Machine Learning Research*, 5:1363–1390, 2004.
- [19] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, New York, 1996.
- [20] J. Fan and Q. Yao. *Nonlinear Time Series*. Springer, New York, 2003.
- [21] L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-free Theory of Nonparametric Regression*. Springer, New York, 2002.
- [22] E.J. Kostelich and T. Schreiber. Noise reduction in chaotic time-series data: A survey of common methods. *Physical Review E*, 48:1752–1763, 1993.
- [23] E.J. Kostelich and J.A. Yorke. Noise reduction: finding the simplest dynamical system consistent with the data. *Physica D. Nonlinear Phenomena*, 41:183–196, 1990.
- [24] S.P. Lalley. Beneath the noise, chaos. *Ann. Statist.*, 27:461–479, 1999.
- [25] S.P. Lalley. Removing the noise from chaos plus noise. In *Nonlinear Dynamics and Statistics*, pages 233–244. Birkhäuser, 2001.
- [26] S.P. Lalley and A.B. Nobel. Denoising deterministic time series. *Dynamics of partial differential equations*, 3:259–279, 2006.

- [27] S. Luzzatto. Stochastic-like behaviour in nonuniformly expanding maps. In B. Hasselblatt and A. Katok, editors, *Handbook of Dynamical Systems 1B*, pages 265–326. Elsevier, Amsterdam, 2006.
- [28] R. Meir. Nonparametric time series prediction through adaptive model selection. *Machine learning*, 39:5–34, 2000.
- [29] D.S. Modha and E. Masry. Memory-universal prediction of stationary random processes. *IEEE Transactions on Information Theory*, 44:117–133, 1998.
- [30] A.B. Nobel. Limits to classification and regression estimation from ergodic processes. *Ann. Statist.*, 27:262–273, 1999.
- [31] A.B. Nobel. Consistent estimation of a dynamical map. In *Nonlinear dynamics and statistics*, pages 267–280. Birkhäuser, Boston, 2001.
- [32] D. Ruelle. The thermodynamic formalism for expanding maps. *Comm. Math. Phys.*, 125:239–262, 1989.
- [33] D. Ruelle. *Thermodynamic Formalism*. Cambridge University Press, 2nd edition, 2004.
- [34] T. Sauer. A noise reduction method for signals from nonlinear systems. *Physica D. Nonlinear Phenomena*, 58:193–201, 1992.
- [35] B. Schölkopf and A.J. Smola. *Learning with Kernels*. MIT Press, 2002.
- [36] I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *J. Mach. Learn. Res.*, 2:67–93, 2001.
- [37] I. Steinwart, D. Hush, and C. Scovel. Function classes that approximate the Bayes risk. In *Proceedings of the 19th Annual Conference on Learning Theory, COLT 2006*, pages 79–93. Springer, 2006.
- [38] I. Steinwart, D. Hush, and C. Scovel. Learning from dependent observations. Technical report, Los Alamos National Laboratory, 2006. http://www.c3.lanl.gov/ml/pubs/2006_dependent/paper.pdf.
- [39] I. Steinwart, D. Hush, and C. Scovel. The reproducing kernel Hilbert space of the Gaussian RBF kernel. *IEEE Trans. Inform. Theory*, 52:4635–4643, 2006.
- [40] J.A.K. Suykens, T. Van Gestel, J. De Brabanter, B. De Moor, and J. Vandewalle. *Least Squares Support Vector Machines*. World Scientific, 2002.
- [41] V.N. Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [42] G. Wahba. *Spline Models for Observational Data*. Series in Applied Mathematics, Vol. 59, SIAM, Philadelphia, 1990.