
Two-stage sampled learning theory on distributions

Zoltán Szabó¹
¹Gatsby Unit, UCL

Arthur Gretton¹
²Machine Learning Department, CMU

Barnabás Póczos²

Bharath Sriperumbudur³
³Department of Statistics, PSU

Abstract

We focus on the distribution regression problem: regressing to a real-valued response from a probability distribution. Although there exist a large number of similarity measures between distributions, very little is known about their generalization performance in specific *learning tasks*. Learning problems formulated on distributions have an inherent two-stage sampled difficulty: in practice only samples from sampled distributions are observable, and one has to build an estimate on similarities computed between sets of points. To the best of our knowledge, the only existing method with consistency guarantees for distribution regression requires kernel density estimation as an intermediate step (which suffers from slow convergence issues in high dimensions), and the domain of the distributions to be compact Euclidean. In this paper, we provide theoretical guarantees for a remarkably simple algorithmic alternative to solve the distribution regression problem: embed the distributions to a reproducing kernel Hilbert space, and learn a ridge regressor from the embeddings to the outputs. Our main contribution is to prove the consistency of this technique in the two-stage sampled setting under mild conditions (on separable, topological domains endowed with kernels). For a given total number of observations, we derive convergence rates as an explicit function of the problem difficulty. As a special case, we answer a 15-year-old open question: we establish the consistency of the classical set kernel [Haussler, 1999; Gärtner et al., 2002] in regression, and cover more recent kernels on distributions, including those due to [Christmann and Steinwart, 2010].

1 INTRODUCTION

We address the learning problem of *distribution regression* in the two-stage sampled setting [1]: we regress

Appearing in Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS) 2015, San Diego, CA, USA. JMLR: W&CP volume 38. Copyright 2015 by the authors. The ordering of the second through fourth authors is alphabetical.

from probability measures to real-valued responses, where we only have bags of samples from the probability distributions. Many classical problems in machine learning and statistics can be analysed in this framework. On the machine learning side, multiple instance learning [2, 3, 4] can be thought of in this way, in the case where each instance in a labeled bag is an i.i.d. (independent identically distributed) sample from a distribution. On the statistical side, tasks might include point estimation of statistics on a distribution (e.g., its entropy or a hyperparameter), where a supervised learning method can help in parameter estimation problems without closed form analytical expressions, or if simulation-based results are computationally expensive.

Before reviewing the existing techniques in the literature, let us start with a somewhat informal definition of the distribution regression problem, and an intuitive phrasing of our goal. Let us suppose that our data consist of $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^l$, where x_i is a probability distribution, $y_i \in \mathbb{R}$, and each (x_i, y_i) pair is i.i.d. sampled from a meta distribution \mathcal{M} . However, we do not observe x_i directly; rather, we observe a sample $x_{i,1}, \dots, x_{i,N_i} \stackrel{i.i.d.}{\sim} x_i$. Thus the observed data are $\hat{\mathbf{z}} = \{(\{x_{i,n}\}_{n=1}^{N_i}, y_i)\}_{i=1}^l$. Our goal is to predict a new y_{l+1} from a new batch of samples $x_{l+1,1}, \dots, x_{l+1,N_{l+1}}$ drawn from a new distribution x_{l+1} . For example, in a medical application the i^{th} patient might be identified with a probability distribution (x_i) , which can be periodically accessed, measured by blood tests $(\{x_{i,n}\}_{n=1}^{N_i})$. We are also given some health indicator of the patient (y_i) , which might be inferred from his/her blood measurements. Based on the observations $(\hat{\mathbf{z}})$, one might try to learn the mapping from the set of blood tests to the health indicator; and the hope is that by observing more patients (larger l) and performing a larger number of tests (larger N_i) the estimated mapping $(\hat{f} = \hat{f}(\hat{\mathbf{z}}))$ becomes more “precise”.

The performance of the estimated mapping (\hat{f}) depends on the assumed function class (\mathcal{H}) , the family of \hat{f} candidates. Let $f_{\mathcal{H}}$ denote the best estimator from \mathcal{H} given infinite training samples ($l = \infty, N_i = \infty$), and let $\mathcal{E}[f_{\mathcal{H}}]$ be its prediction error. Our goal is to

obtain upper bounds for the $0 \leq \mathcal{E}[\hat{f}] - \mathcal{E}[f_{\mathcal{X}}]$ quantity which hold with high probability. More precisely, we are aiming at

1. deriving upper bounds on the excess risk, proving consistency: We construct $\mathcal{E}[\hat{f}] - \mathcal{E}[f_{\mathcal{X}}] \leq r(l, N, \lambda)$ bounds, where λ is a regularization parameter converging to zero as we see more samples ($l \rightarrow \infty$, $N = N_i \rightarrow \infty$), and choose the (l, N, λ) triplet appropriately to drive $r(l, N, \lambda)$ and hence $\mathcal{E}[\hat{f}] - \mathcal{E}[f_{\mathcal{X}}]$ to 0.
2. obtaining convergence rates: We establish convergence rates for a general prior family $\mathcal{P}(b, c)$ [5], where b captures the effective input dimension, and larger c means smoother $f_{\mathcal{X}}$. In particular, when $l = N^a$ ($a > 0$), the effective dimension is small (large b), and the total number of samples processed $t = lN = N^{a+1}$ is fixed, one obtains a rate of $1/t^{2/7}$ for a smooth regression function ($c = 2$), $1/t^{1/5}$ in the non-smooth case ($c = 1$).

The motivation for considering the $\mathcal{P}(b, c)$ family is two-fold:

1. it does not assume parametric distributions, still certain complexity terms can be explicitly upper bounded in the family. This property will be exploited in our analysis.
2. (for special input distributions) parameter b can be related to the spectral decay of Gaussian Gram matrices, thus available analysis techniques [6] might give alternative prior characterizations.

Briefly, we focus on the following question:

Can the distribution regression problem be solved consistently under mild conditions?

Despite the large number of available “solutions” and applications of distribution regression dating back to 1999 [7], surprisingly this pretty fundamental question has hardly been touched. In our paper we give affirmative answer to the question by presenting the analysis of a *simple* kernel ridge regression approach [see Eq. (3)] in the two-stage sampled ($\mathcal{M} \rightarrow \mathbf{z} \rightarrow \hat{\mathbf{z}}$) setting.

Review of approaches to learning on distributions: A number of methods have been proposed over the years to compute the similarity of distributions or bags of samples. As a first approach, one could fit a parametric model to the bags, and estimate the similarity of the bags based on the obtained parameters. It is then possible to define learning algorithms on the basis of these similarities, which often take analytical form. Typical examples with explicit formulas include Gaussians, finite mixtures of Gaussians, and distributions from the exponential family (with

known log-normalizer function and zero carrier measure) [8, 9, 10, 11]. A major limitation of these methods, however, is that they apply quite simple parametric assumptions, which may not be sufficient or verifiable in practise.

A heuristic related to the parametric approach is to assume that the training distributions are Gaussians in a reproducing kernel Hilbert space; see for example [10, 12] and references therein. This assumption is algorithmically appealing, as many divergence measures for Gaussians can be computed in closed form using only inner products, making them straightforward to kernelize. A fundamental shortfall of kernelized Gaussian divergences is the lack of their consistency analysis in specific learning algorithms.

A more theoretically grounded approach to learning on distributions has been to define positive definite kernels [13] on the basis of statistical divergence measures on distributions, or by metrics on non-negative numbers; these can then be used in kernel algorithms. This category includes work on semigroup kernels [14], nonextensive information theoretical kernel constructions [15], and kernels based on Hilbertian metrics [16]. For example, in [14] the intuition is as follows: if two measures or sets of points overlap, then their sum is expected to be more concentrated. The value of dispersion can be measured by entropy or inverse generalized variance. In the second type of approach [16], homogeneous Hilbert metrics on the non-negative real line are used to define the similarity of probability distributions. While these techniques guarantee to provide valid kernels on certain restricted domains of measures, the performance of learning algorithms based on finite sample estimates of these kernels remains a challenging open question. One might also plug into learning algorithms (based on similarities of distributions) consistent Rényi and Tsallis divergence estimates [17, 18], but these similarity indices are *not* kernels, and their consistency in specific learning tasks, similarly to the previous works, is open.

To the best of our knowledge, the only prior work addressing the consistency of regression on distributions requires kernel density estimation [1, 19], assumes that the response variable is scalar-valued¹, and the covariates are nonparametric continuous distributions on \mathbb{R}^d . As in our setting, the exact forms of these distributions are unknown; they are available only through finite sample sets. Póczos et al. estimated these distributions through a kernel density estimator (assuming these distributions to have a density) and then constructed a kernel regressor that acts on these kernel

¹[20] considers the case where the responses are also distributions.

density estimates.² Using the classical bias-variance decomposition analysis for kernel regressors, they show the consistency of the constructed kernel regressor, and provide a polynomial upper bound on the rates, assuming the true regressor to be Hölder continuous, and the meta distribution that generates the covariates x_i to have finite doubling dimension [22].³

An alternative paradigm in learning when the inputs are “bags of objects” is to simply treat each input as a finite set, and to define kernel learning algorithms based on set kernels [23] (also called multi-instance kernels or ensemble kernels, and instances of convolution kernels [7]). In this case, the similarity of two sets is measured by the average pairwise point similarities between the sets. From a *theoretical* perspective, very little has been done to establish the consistency of set kernels in learning since their introduction in 1999 [7, 23]: i.e. in what sense (and with what rates) is the learning algorithm consistent, when the number of items per bag, and the number of bags, is allowed to increase?

It is possible, however, to view set kernels in a distribution setting, as they represent valid kernels between (mean) embeddings of empirical probability measures into a reproducing kernel Hilbert space (RKHS) [24]. The *population limits are well-defined* as being dot products between the embeddings of the generating distributions [25], and for characteristic kernels the distance between embeddings defines a *metric* on probability measures [26, 27]. When bounded kernels are used, mean embeddings *exist for all probability measures* [28]. When we consider the distribution regression setting, however, there is no reason to limit ourselves to set kernels. Embeddings of probability measures to RKHS are used by [29] in defining a yet larger class of easily computable kernels on distributions, via operations performed on the embeddings and their distances. Note that the relation between set kernels and kernels on distributions has been applied by [30] for classification on distribution-valued inputs, however consistency was not studied in that work.

Our **contribution** in this paper is to establish the consistency of an algorithmically simple, mean embedding based ridge regression method (described in Section 2) for the distribution regression problem. This result applies both to the basic set kernels of [7, 23], the distri-

²We would like to clarify that the kernels used in their work are classical smoothing kernels (extensively studied in non-parametric statistics [21]) and not the reproducing kernels that appear throughout our paper.

³Using a random kitchen sinks approach, with orthonormal basis projection estimators and RBF kernels [19] proposes a distribution regression algorithm that can computationally handle large scale datasets; as with [1], this approach is based on density estimation in \mathbb{R}^d .

bution kernels of [29], and additional related kernels proposed herein. We provide two-stage sampled excess error bounds, consistency proof and convergence rates in Section 4, and break down the various trade-offs arising in different sample size and problem difficulties. The principal challenge in proving theoretical guarantees arises from the two-stage sampled nature of the inputs. In our analysis, we make use of [5], who provide error bounds for the one-stage sample setup. These results will make our analysis somewhat shorter (but still rather challenging) by giving upper bounds for some of the upcoming objective terms. Even the verification of these conditions requires care (Section 3) since the inputs in the ridge regression are themselves distribution embeddings (i.e., functions in a reproducing kernel Hilbert space).

Due to the differences in the assumptions made and the loss function used, a direct comparison of our theoretical result and that of [1]³ remains an open question, however we make two observations. First, our approach is more general, since we may regress from any probability measure defined on a separable, topological domain endowed with a kernel. Póczos et al.’s work is restricted to compact domains of finite dimensional Euclidean spaces, and requires the distributions to admit probability densities; distributions on strings, time series, graphs, and other structured objects are disallowed. Second, density estimates in high dimensional spaces suffer from slow convergence rates [31, Section 6.5]. Our approach avoids this problem, as it works directly on distribution embeddings, and does not make use of density estimation as an intermediate step.

2 THE DISTRIBUTION REGRESSION PROBLEM

In this section, we define the distribution regression problem, for a general RKHS on distributions. In Section 3, we will provide examples of valid kernels for this RKHS, including set kernels [7, 23], the kernels from [29], and further related kernels. Below, we first introduce some notation and then formally discuss the distribution regression problem.

Notation: Let (\mathcal{X}, τ) be a topological space and let $\mathcal{B}(\mathcal{X}) := \mathcal{B}(\tau)$ be the Borel σ -algebra induced by the topology τ . $\mathcal{M}_1^+(\mathcal{X})$ denotes the set of Borel probability measures on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. The weak topology $(\tau_w = \tau_w(\mathcal{X}, \tau))$ on $\mathcal{M}_1^+(\mathcal{X})$ is defined as the weakest topology such that the $L_h : (\mathcal{M}_1^+(\mathcal{X}), \tau_w) \rightarrow \mathbb{R}$, $L_h(x) = \int_{\mathcal{X}} h(u) dx(u)$ mapping is continuous for all $h \in C_b(\mathcal{X}) = \{(\mathcal{X}, \tau) \rightarrow \mathbb{R} \text{ bounded, continuous functions}\}$. Let $H = H(k)$ be the RKHS [6] with $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ as the reproducing

kernel. Denote by

$$X = \mu(\mathcal{M}_1^+(\mathcal{X})) = \{\mu_x : x \in \mathcal{M}_1^+(\mathcal{X})\} \subseteq H$$

the set of $\mu_x = \int_{\mathcal{X}} k(\cdot, u) dx(u) = \mathbb{E}_{u \sim x}[k(\cdot, u)] \in H$ mean embeddings [24] of the distributions to the space H , and let $Y = \mathbb{R}$. Intuitively, μ_x is the canonical feature map $[k(\cdot, u)]$ averaged according to the probability measure $[dx(u)]$. Let $\mathcal{H} = \mathcal{H}(K)$ be the RKHS of functions with $K : X \times X \rightarrow \mathbb{R}$ as the reproducing kernel. $\mathcal{L}(\mathcal{H})$ is the space of $\mathcal{H} \rightarrow \mathcal{H}$ bounded linear operators, and δ_{μ_a} denotes the evaluation functional at μ_a ($a \in \mathcal{M}_1^+(\mathcal{X})$). For $M \in \mathcal{L}(\mathcal{H})$ the operator norm is defined as $\|M\|_{\mathcal{L}(\mathcal{H})} = \sup_{0 \neq q \in \mathcal{H}} \frac{\|Mq\|_{\mathcal{H}}}{\|q\|_{\mathcal{H}}}$. Given (U_1, \mathcal{S}_1) and (U_2, \mathcal{S}_2) measurable spaces the $\mathcal{S}_1 \otimes \mathcal{S}_2$ product σ -algebra [6, page 480] on the product space $U_1 \times U_2$ is the σ -algebra generated by the cylinder sets $U_1 \times S_2$, $S_1 \times U_2$ ($S_1 \in \mathcal{S}_1$, $S_2 \in \mathcal{S}_2$). $\mathbb{E}[\cdot]$ denotes expectation.

Distribution regression: In the *distribution regression* problem, we are given samples $\hat{\mathbf{z}} = \{(\{x_{i,n}\}_{n=1}^{N_i}, y_i)\}_{i=1}^l$ with $x_{i,1}, \dots, x_{i,N_i} \stackrel{i.i.d.}{\sim} x_i$ where $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^l$ with $x_i \in \mathcal{M}_1^+(\mathcal{X})$ and $y_i \in Y$ drawn i.i.d. from a joint meta distribution \mathcal{M} defined on the measurable space $(\mathcal{M}_1^+(\mathcal{X}) \times \mathbb{R}, \mathcal{B}(\tau_w) \otimes \mathcal{B}(\mathbb{R}))$. Unlike in classical supervised learning problems, the problem at hand involves two levels of randomness, wherein first \mathbf{z} is drawn from \mathcal{M} and then $\hat{\mathbf{z}}$ is generated by sampling points from x_i for all $i = 1, \dots, l$. The **goal** is to learn the relation between the random distribution x and scalar response y based on the observed $\hat{\mathbf{z}}$. For notational simplicity, we will assume that $N = N_i$ ($\forall i$).

As in the classical regression task ($\mathbb{R}^d \rightarrow \mathbb{R}$), distribution regression can be tackled as a kernel ridge regression problem (using squared loss as the discrepancy criterion). The kernel (say $k_{\mathcal{G}}$) is defined on $\mathcal{M}_1^+(\mathcal{X})$, and the regressor is then modelled by an element in the RKHS $\mathcal{G} = \mathcal{G}(k_{\mathcal{G}})$ of functions mapping from $\mathcal{M}_1^+(\mathcal{X})$ to \mathbb{R} . In this paper, we choose $k_{\mathcal{G}}(x, x') = K(\mu_x, \mu_{x'})$ where $x, x' \in \mathcal{M}_1^+(\mathcal{X})$ and so that the function (in \mathcal{G}) to describe the (x, y) random relation is constructed as a composition

$$\mathcal{M}_1^+(\mathcal{X}) \xrightarrow{\mu} X (\subseteq H = H(k)) \xrightarrow{f \in \mathcal{H} = \mathcal{H}(K)} \mathbb{R}.$$

In other words, the distribution $x \in \mathcal{M}_1^+(\mathcal{X})$ is first mapped to $X \subseteq H$ by the mean embedding μ , and the result is composed with f , an element of the RKHS $\mathcal{H} = \mathcal{H}(K)$. Assuming that $f_{\mathcal{G}}$, the minimizer of the expected risk (\mathcal{E}) over \mathcal{G} exists, then a function $f_{\mathcal{H}}$ also exists, and satisfies

$$\begin{aligned} \mathcal{E}[f_{\mathcal{H}}] &= \inf_{f \in \mathcal{H}} \mathcal{E}[f] = \inf_{f \in \mathcal{H}} \mathbb{E}_{(x,y) \sim \mathcal{M}} [f(\mu_x) - y]^2 \\ &= \inf_{g \in \mathcal{G}} \mathbb{E}_{(x,y) \sim \mathcal{M}} [g(x) - y]^2 = \inf_{g \in \mathcal{G}} \mathcal{E}[g] = \mathcal{E}[g_{\mathcal{G}}]. \end{aligned}$$

The classical regularization approach is to optimize

$$f_{\mathbf{z}}^{\lambda} = \arg \min_{f \in \mathcal{H}} \frac{1}{l} \sum_{i=1}^l [f(\mu_{x_i}) - y_i]^2 + \lambda \|f\|_{\mathcal{H}}^2 \quad (1)$$

instead of \mathcal{E} , based on samples \mathbf{z} . Since \mathbf{z} is not accessible, we consider the objective function defined by the observable quantity $\hat{\mathbf{z}}$,

$$f_{\hat{\mathbf{z}}}^{\lambda} = \arg \min_{f \in \mathcal{H}} \frac{1}{l} \sum_{i=1}^l [f(\mu_{\hat{x}_i}) - y_i]^2 + \lambda \|f\|_{\mathcal{H}}^2, \quad (2)$$

where $\hat{x}_i = \frac{1}{N} \sum_{n=1}^N \delta_{x_{i,n}}$ is the empirical distribution determined by $\{x_{i,n}\}_{n=1}^N$. Algorithmically, ridge regression is quite simple [32]: given training samples $\hat{\mathbf{z}}$, the prediction for a new t test distribution is

$$\begin{aligned} (f_{\hat{\mathbf{z}}}^{\lambda} \circ \mu)(t) &= [y_1, \dots, y_l] (\mathbf{K} + \lambda \mathbf{I}_l)^{-1} \mathbf{k} \in \mathbb{R}, \\ \mathbf{K} &= [K(\mu_{\hat{x}_i}, \mu_{\hat{x}_j})] \in \mathbb{R}^{l \times l}, \\ \mathbf{k} &= [K(\mu_{\hat{x}_1}, \mu_t); \dots; K(\mu_{\hat{x}_l}, \mu_t)] \in \mathbb{R}^l. \end{aligned} \quad (3)$$

Remarks:

1. It is important to note that the algorithm has access to the sample points *only via their mean embeddings* $\{\mu_{\hat{x}_i}\}_{i=1}^l$ in Eq. (2).
2. There is a *two-stage sampling difficulty* to tackle: The transition from $f_{\mathcal{H}}$ to $f_{\mathbf{z}}^{\lambda}$ represents the fact that we have only l distribution samples (\mathbf{z}); the transition from $f_{\mathbf{z}}^{\lambda}$ to $f_{\hat{\mathbf{z}}}^{\lambda}$ means that the x_i distributions can be accessed only via samples ($\hat{\mathbf{z}}$).
3. While ridge regression can be performed using the kernel $k_{\mathcal{G}}$, the two-stage sampling makes it difficult to work with arbitrary $k_{\mathcal{G}}$. By contrast, our choice of $k_{\mathcal{G}}(x, x') = K(\mu_x, \mu_{x'})$ enables us to handle the two-stage sampling by estimating μ_x with an empirical estimator and using it in the algorithm as shown above.

The main **goal** of this paper is to analyse the excess risk $\mathcal{E}[f_{\hat{\mathbf{z}}}^{\lambda}] - \mathcal{E}[f_{\mathcal{H}}]$, i.e., the regression performance compared to the best possible estimation from \mathcal{H} , and to establish consistency and rates of convergence as a function of the (l, N, λ) triplet, and of the difficulty of the problem in the sense of [5].

3 ASSUMPTIONS

In this section we detail our assumptions on the (\mathcal{X}, k, K) triplet, and show that regressing with set kernels fit into the studied problem family. Our analysis will rely on existing ridge regression results [5] which focus on problem (1), where only a single-stage sampling is present; hence we have to verify the associated conditions. Though we make use of these results, the analysis still remains rather challenging; the available bounds can moderately shorten our

proof. We must also take particular care in verifying that [5]’s conditions are met, since they must hold for the space of *mean embeddings of the distributions* ($X = \mu(\mathcal{M}_1^+(\mathcal{X}))$), whose properties as a function of \mathcal{X} and H must themselves be established. Our assumptions:

- $\exists f_{\mathcal{H}}$ such that $\mathcal{E}[f_{\mathcal{H}}] = \inf_{f \in \mathcal{H}} \mathcal{E}(f)$.
- (\mathcal{X}, τ) is a separable, topological domain.
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is bounded ($\exists B_k < \infty$ such that $\sup_{u \in \mathcal{X}} k(u, u) \leq B_k$) and continuous.
- $K : X \times X \rightarrow \mathbb{R}$ is bounded, i.e., $\exists B_K < \infty$ such that

$$K(\mu_a, \mu_a) \leq B_K, \quad (\forall \mu_a \in X), \quad (4)$$

and $\Psi(\mu_c) := K(\cdot, \mu_c) : X \rightarrow \mathcal{H}$ is Hölder continuous, i.e., $\exists L > 0, h \in (0, 1]$ such that for $\forall(\mu_a, \mu_b) \in X \times X$

$$\|\Psi(\mu_a) - \Psi(\mu_b)\|_{\mathcal{H}} \leq L \|\mu_a - \mu_b\|_H^h. \quad (5)$$

- y is bounded: $\exists C < \infty$ such that $|y| \leq C$ almost surely.
- $X = \mu(\mathcal{M}_1^+(\mathcal{X})) \in \mathcal{B}(H)$.

Discussion of the assumptions: We give a short insight into the consequences of our assumptions and present some concrete examples.

1. The boundedness and continuity of k imply the measurability of $\mu : (\mathcal{M}_1^+(\mathcal{X}), \mathcal{B}(\tau_w)) \rightarrow (H, \mathcal{B}(H))$, which using the $X \in \mathcal{B}(H)$ condition guarantees that the ρ , the measure induced by \mathcal{M} on $X \times \mathbb{R}$ is well-defined (see the supplementary material).
2. For a linear kernel, $K(\mu_a, \mu_b) = \langle \mu_a, \mu_b \rangle_H$, ($\mu_a, \mu_b \in X$), one can verify (see the supplementary material) that Hölder continuity holds with $L = 1, h = 1$. Also, since $K(\mu_a, \mu_b) \leq B_k$ for any $a, b \in \mathcal{M}_1^+(\mathcal{X})$, we can choose $B_K = B_k$. Evaluating the kernel, K at the $\mu_{\hat{x}_i} = \int_{\mathcal{X}} k(\cdot, u) d\hat{x}_i(u) = \frac{1}{N} \sum_{n=1}^N k(\cdot, x_{i,n})$ empirical embeddings yields the standard set kernel:

$$K(\mu_{\hat{x}_i}, \mu_{\hat{x}_j}) = \frac{1}{N^2} \sum_{n,m=1}^N k(x_{i,n}, x_{j,m}).$$

3. One can also prove (see the supplement) by using the properties of negative/positive definite functions [33] that many K functions on $X \times X$ are kernels and (in case of compact metric \mathcal{X} domains) Hölder continuous.⁴ Some examples are listed in Table 1; these kernels are the natural extensions to distributions of the Gaussian [29], exponential, Cauchy, generalized t-student and inverse multiquadratic kernels.

⁴To guarantee the Hölder property of K -s, we assume the continuity of μ . For example, if \mathcal{X} is a compact metric space and k is universal, then μ metrizes the weak topology τ_w [34, Theorem 23, page 1552], hence μ is continuous. In this case $X = \mu(\mathcal{M}_1^+(\mathcal{X}))$ is compact metric (see the supplement), thus closed and hence $X \in \mathcal{B}(H)$ also holds.

4. $Y = \mathbb{R}$ is a separable Hilbert space hence Polish, and thus the $\rho(y|\mu_a)$ conditional distribution ($y \in \mathbb{R}, \mu_a \in X$) is well-defined; see [6, Lemma A.3.16, page 487].

5. The separability of \mathcal{X} and the continuity of k implies the separability of H [6, Lemma 4.33, page 130]. Also, since $X \subseteq H, X$ is separable; hence so is \mathcal{H} due to the continuity of K .

Verification of [5]’s conditions: Below we prove that [5]’s conditions hold under our assumptions.

1. $Y = \mathbb{R}$ and \mathcal{H} are separable Hilbert spaces – as we have seen.
2. By the bilinearity of $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and the reproducing property of K , the measurability of $(\mu_x, \mu_t) \mapsto \langle K(\cdot, \mu_x)w, K(\cdot, \mu_t)v \rangle_{\mathcal{H}} = wK(\mu_x, \mu_t)v \quad (\forall w, v \in \mathbb{R})$ is equivalent to that of $(\mu_x, \mu_t) \mapsto K(\mu_x, \mu_t)$; the latter follows from the Hölder continuity of Ψ (see the supplement).
3. Due to the boundedness of y , we have $\int_{X \times \mathbb{R}} y^2 d\rho(\mu_x, y) \leq \int_{X \times \mathbb{R}} C^2 d\rho(\mu_x, y) \leq C^2 < \infty$, and $\exists \Sigma > 0, \exists M > 0$ such that

$$\int_{\mathbb{R}} e^{\frac{|y - f_{\mathcal{H}}(\mu_x)|}{M}} - \frac{|y - f_{\mathcal{H}}(\mu_x)|}{M} - 1 d\rho(y|\mu_x) \leq \frac{\Sigma^2}{2M^2} \quad (6)$$

for ρ_X -almost $\mu_x \in X$, where $\rho(\mu_x, y) = \rho(y|\mu_x)\rho_X(\mu_x)$ is factorized into conditional and marginal distributions. (6) is a model of the noise of the output y ; it is satisfied, for example in case of bounded noise [5, page 9]. By the boundedness of y and that of kernel K this property holds: $|y - f_{\mathcal{H}}(\mu_x)| \leq |y| + |f_{\mathcal{H}}(\mu_x)| \leq C + \|f_{\mathcal{H}}\|_H \sqrt{B_K}$, where we used the triangle inequality and Lemma 4.23 (page 124) from [6].

4 ERROR BOUNDS, CONSISTENCY, CONVERGENCE RATE

In this section, we present our main result: we derive high probability upper bound for the excess risk $\mathcal{E}[f_{\frac{\lambda}{2}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}]$ of the mean embedding based ridge regression (MERR) method, see our main theorem. We also illustrate the upper bound for particular classes of prior distributions, resulting in sufficient conditions for convergence and concrete convergence rates (see Consequences 1-2). We first give a high-level sketch of our convergence analysis and the results are stated with their intuitive interpretation. Then an outline of the main proof ideas follows; technical details of the proof steps may be found in the supplement.

At a high level, our convergence analysis takes the following form: Having explicit expressions for $f_{\mathbf{z}}^\lambda, f_{\frac{\lambda}{2}}^\lambda$ [see Eq. (9)-(10)], we will decompose the excess risk

Table 1: Nonlinear kernels on mean embedded distributions: $K = K(\mu_a, \mu_b)$; $\theta > 0$. For the Hölder continuity, we assume that \mathcal{X} is a compact metric space and μ is continuous (the latter is implied e.g., by a universal k).

K_G	K_e	K_C	K_t	K_i
$e^{-\frac{\ \mu_a - \mu_b\ _H^2}{2\theta^2}}$	$e^{-\frac{\ \mu_a - \mu_b\ _H}{2\theta^2}}$	$(1 + \ \mu_a - \mu_b\ _H^2 / \theta^2)^{-1}$	$(1 + \ \mu_a - \mu_b\ _H^\theta)^{-1}$	$(\ \mu_a - \mu_b\ _H^2 + \theta^2)^{-\frac{1}{2}}$
$h = 1$	$h = \frac{1}{2}$	$h = 1$	$h = \frac{\theta}{2} (\theta \leq 2)$	$h = 1$

$\mathcal{E}[f_{\mathcal{Z}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}]$ into five terms:

$$\begin{aligned} \mathcal{E}[f_{\mathcal{Z}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}] &\leq 5[S_{-1} + S_0 + \mathcal{A}(\lambda) + S_1 + S_2], \\ S_{-1} &= \|\sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(g_{\mathbf{z}} - g_{\mathbf{x}})\|_{\mathcal{H}}^2, \\ S_0 &= \|\sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(T_{\mathbf{x}} - T_{\mathbf{x}})f_{\mathcal{Z}}^\lambda\|_{\mathcal{H}}^2, \\ S_1 &= \|\sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(g_{\mathbf{z}} - T_{\mathbf{x}}f_{\mathcal{H}})\|_{\mathcal{H}}^2, \\ S_2 &= \|\sqrt{T}(T_{\mathbf{x}} + \lambda)^{-1}(T - T_{\mathbf{x}})(f^\lambda - f_{\mathcal{H}})\|_{\mathcal{H}}^2, \\ \mathcal{A}(\lambda) &= \|\sqrt{T}(f^\lambda - f_{\mathcal{H}})\|_{\mathcal{H}}^2, \end{aligned}$$

where $f^\lambda = \arg \min_{f \in \mathcal{H}} \mathcal{E}[f] + \lambda \|f\|_{\mathcal{H}}^2$, $T_{\mu_a} = K(\cdot, \mu_a)\delta_{\mu_a}$ [$T_{\mu_a}(f) = K(\cdot, \mu_a)f(\mu_a)$, $\mu_a \in X$],

$$T = \int_X T_{\mu_a} d\rho_X(\mu_a) \in \mathcal{L}(\mathcal{H}), \quad T_{\mu_a} \in \mathcal{L}(\mathcal{H}). \quad (7)$$

- Three of the terms ($S_1, S_2, \mathcal{A}(\lambda)$) will be identical to terms in [5], hence their bounds can be applied.
- The two new terms (S_{-1}, S_0), the result of two-stage sampling, will be upper bounded by making use of the convergence of the empirical mean embeddings, and the Hölder property of K .

These bounds will lead to the following results:

Main theorem (bound on the excess risk). *Let M, Σ and T be as in (6), (7). Let $\Psi(\mu_a) = K(\cdot, \mu_a) : X \rightarrow \mathcal{H}$ be Hölder continuous with constants L, h . Let $l \in \mathbb{N}, N \in \mathbb{N}, \lambda > 0, 0 < \eta < 1, C > 0, \delta > 0, C_\eta = 32 \log^2(6/\eta), |y| \leq C$ (a.s.), $\mathcal{A}(\lambda)$ the residual as above, and define $\mathcal{B}(\lambda) = \|f^\lambda - f_{\mathcal{H}}\|_{\mathcal{H}}^2$ the reconstruction error, $\mathcal{N}(\lambda) = \text{Tr}[(T + \lambda)^{-1}T]$ the effective dimension. Then with probability at least $1 - \eta - e^{-\delta}$*

$$\begin{aligned} \mathcal{E}[f_{\mathcal{Z}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}] &\leq \\ &\leq 5 \left\{ \frac{4L^2C^2 \left(1 + \sqrt{\log(l) + \delta}\right)^{2h} (2B_k)^h}{\lambda N^h} \left[1 + \frac{4(B_K)^2}{\lambda^2} \right] + \right. \\ &\quad \left. \mathcal{A}(\lambda) + C_\eta \left[\frac{B_K^2 \mathcal{B}(\lambda)}{l^2 \lambda} + \frac{B_K \mathcal{A}(\lambda)}{4l \lambda} + \frac{B_K M^2}{l^2 \lambda} + \frac{\Sigma^2 \mathcal{N}(\lambda)}{l} \right] \right\} \\ &\quad \times \left[1 + \frac{4(B_K)^2}{\lambda^2} \right] + R\lambda^c + C_\eta \times \end{aligned}$$

provided that $l \geq 2C_\eta B_K \mathcal{N}(\lambda) / \lambda, \lambda \leq \|T\|_{\mathcal{L}(\mathcal{H})}, N \geq (1 + \sqrt{\log(l) + \delta})^2 2^{\frac{h+6}{h}} B_k (B_K)^{\frac{1}{h}} L^{\frac{2}{h}} / \lambda^{\frac{2}{h}}$.

Below we specialize our bound on the excess risk for a general prior class, which captures the difficulty of the

regression problem as defined in [5]. This $\mathcal{P}(b, c)$ class is described by two parameters b and c : intuitively, larger b means faster decay of the eigenvalues of the covariance operator T [(7)], hence smaller effective input dimension; larger c corresponds to smoother $f_{\mathcal{H}}$. Formally:

Definition of the $\mathcal{P}(b, c)$ class: Let us fix the positive constants $M, \Sigma, R, \alpha, \beta$. Then given $1 < b, c \in [1, 2]$, the $\mathcal{P}(b, c)$ class is the set of probability distributions ρ on $Z = X \times \mathbb{R}$ such that (i) the (μ_x, y) assumption holds with M, Σ in (6), (ii) there is a $g \in \mathcal{H}$ such that $f_{\mathcal{H}} = T^{\frac{c-1}{2}}g$ with $\|g\|_{\mathcal{H}}^2 \leq R$, (iii) in the $T = \sum_{n=1}^N t_n \langle \cdot, e_n \rangle_{\mathcal{H}} e_n$ spectral theorem based decomposition ($\{e_n\}_{n=1}^N$ is a basis of $\ker(T)^\perp$), $N = +\infty$, and the eigenvalues of T satisfy $\alpha \leq n^b t_n \leq \beta \quad (\forall n \geq 1)$.

We can provide a simple example of when the source decay conditions hold, in the event that the distributions are normal with means m_i and identical variance ($x_i = N(m_i, \sigma^2 I)$). When Gaussian kernels (k) are used with linear K , then $K(\mu_{x_i}, \mu_{x_j}) = e^{-c\|m_i - m_j\|^2}$ [30, Table 1, line 2] (Gaussian, with arguments equal to the difference in means). Thus, this Gram matrix will correspond to the Gram matrix using a Gaussian kernel between points m_i . The spectral decay of the Gram matrix will correspond to that of the Gaussian kernel, with points drawn from the meta-distribution over the m_i . Thus, the source conditions are analysed in the same manner as for Gaussian Gram matrices, e.g. see [6] for a discussion of the spectral decay properties.

In the $\mathcal{P}(b, c)$ family, the behaviour of $\mathcal{A}(\lambda), \mathcal{B}(\lambda)$ and $\mathcal{N}(\lambda)$ is known; specializing our theorem we get:⁵

Consequence 1 (Excess risk in the $\mathcal{P}(b, c)$ class).

$$\begin{aligned} \mathcal{E}[f_{\mathcal{Z}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}] &\leq 5 \left\{ \frac{4L^2C^2 \left(1 + \sqrt{\log(l) + \delta}\right)^{2h} (2B_k)^h}{\lambda N^h} \right. \\ &\quad \left. \left[\frac{B_K^2 R \lambda^{c-2}}{l^2} + \frac{B_K R \lambda^{c-1}}{4l} + \frac{B_K M^2}{l^2 \lambda} + \frac{\Sigma^2 \beta b}{(b-1)l \lambda^{\frac{1}{b}}} \right] \right\}. \end{aligned}$$

⁵In what follows, we assume the conditions of the main theorem and $\rho \in \mathcal{P}(b, c)$.

By choosing λ appropriately as a function of l and N , the excess risk $\mathcal{E}[f_{\hat{\mathbf{z}}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}]$ converges to 0, and we can use Consequence 1 to obtain convergence rates: the task reduces to the study of

$$r(l, N, \lambda) = \frac{\log^h(l)}{N^h \lambda^3} + \lambda^c + \frac{1}{l^2 \lambda} + \frac{1}{l \lambda^{\frac{1}{b}}} \rightarrow 0, \quad (8)$$

subject to $l \geq \lambda^{-\frac{1}{b}-1}$.⁶ By matching two terms in (8), solving for λ and plugging the result back to the bound (see the supplementary material), we obtain:

Consequence 2 (Consistency and convergence rate in $\mathcal{P}(b, c)$). *Let $l = N^a$ ($a > 0$). The excess risk can be upper bounded (constant multipliers are discarded) by the quantities given in the last column of Table 2.*

Note: in function r [Eq. (8)] (i) the first term comes from the error of the mean embedding estimation, (ii) the second term corresponds to $\mathcal{A}(\lambda)$, a complexity measure of $f_{\mathcal{H}}$, (iii) the third term is due to the S_1 bound, (iv) the fourth term expresses $\mathcal{N}(\lambda)$, a complexity index of the hypothesis space \mathcal{H} according to the marginal measure ρ_X . As an example, let us take two rows from Table 2:

1. First row: In this case the first and second terms dominate $r(l, N, \lambda)$ in (8); in other words the error is determined by the mean embedding estimation process and the complexity of $f_{\mathcal{H}}$. Let us assume that b is large in the sense that $1/b \approx 0$, $(b+1)/b \approx 1$ (hence, the effective dimension of the input space is small); and assume that K is Lipschitz ($h = 1$). Under these conditions the lower bound for a is approximately $\max(c/(c+3), 1/(c+3)) = c/(c+3) \leq a$ (since $c \geq 1$). Using such an a (i.e., the exponent in $l = N^a$ is not too small), then the convergence rate is $[\log(N)/N]^{\frac{c}{c+3}}$. Thus, for example, if $c = 2$ ($f_{\mathcal{H}} = T^{\frac{c-1}{2}}g$ is smoothed by T from a $g \in \mathcal{H}$), then $a = \frac{2}{2+3} = 0.4$ and the convergence rate is $[\log(N)/N]^{0.4}$; in other words the rate is approximately $1/N^{0.4}$. If c takes its minimal value ($c = 1$; $f_{\mathcal{H}}$ is less smooth), then $a = \frac{1}{1+3} = \frac{1}{4}$ results in an approximate rate of $1/N^{0.25}$. Alternatively, if we keep the total number of samples processed $t = lN = N^{a+1}$ fixed, $r(t) \approx 1/N^a = 1/t^{a/(a+1)} = 1/t^{1-1/(a+1)}$, i.e., the convergence rate becomes larger for smoother regression problems (increasing c).

2. Last row: At this extreme, two terms dominate: the complexity of \mathcal{H} according to ρ_X , and a term from the bound on S_1 . Under this condition, although one can solve the matching criterion for λ , and it is possible to drive the individual terms of r to zero, l cannot

be chosen large enough (within the analysed $l = N^a$ ($a > 0$) scheme) to satisfy the $l \geq \lambda^{-\frac{1}{b}-1}$ constraint; thus convergence fails.

Proof of main theorem: We present the main steps of the proof of our theorem; detailed derivations can be found in the supplementary material. Let us define $\mathbf{x} = \{x_i\}_{i=1}^l$ and $\hat{\mathbf{x}} = \{\{x_{i,n}\}_{n=1}^N\}_{i=1}^l$ as the ‘x-part’ of \mathbf{z} and $\hat{\mathbf{z}}$. One can express $f_{\mathbf{z}}^\lambda$ [5], and similarly $f_{\hat{\mathbf{z}}}^\lambda$ as

$$f_{\mathbf{z}}^\lambda = (T_{\mathbf{x}} + \lambda)^{-1} g_{\mathbf{z}}, \quad T_{\mathbf{x}} = \frac{1}{l} \sum_{i=1}^l T_{\mu_{x_i}}, \quad (9)$$

$$f_{\hat{\mathbf{z}}}^\lambda = (T_{\hat{\mathbf{x}}} + \lambda)^{-1} g_{\hat{\mathbf{z}}}, \quad T_{\hat{\mathbf{x}}} = \frac{1}{l} \sum_{i=1}^l T_{\mu_{\hat{x}_i}}, \quad (10)$$

$$g_{\mathbf{z}} = \frac{1}{l} \sum_{i=1}^l K(\cdot, \mu_{x_i}) y_i, \quad g_{\hat{\mathbf{z}}} = \frac{1}{l} \sum_{i=1}^l K(\cdot, \mu_{\hat{x}_i}) y_i. \quad (11)$$

In Eqs. (9), (10), (11), $T_{\mathbf{x}}, T_{\hat{\mathbf{x}}} : \mathcal{H} \rightarrow \mathcal{H}$, $g_{\mathbf{z}}, g_{\hat{\mathbf{z}}} \in \mathcal{H}$.

• **Decomposition of the excess risk:** We derive the upper bound for the excess risk

$$\mathcal{E}[f_{\hat{\mathbf{z}}}^\lambda] - \mathcal{E}[f_{\mathcal{H}}] \leq 5[S_{-1} + S_0 + \mathcal{A}(\lambda) + S_1 + S_2]. \quad (12)$$

• **It is sufficient to upper bound S_{-1} and S_0 :** [5] has shown that $\forall \eta > 0$ if $l \geq \frac{2C_\eta B_K \mathcal{N}(\lambda)}{\lambda}$ and $\lambda \leq \|T\|_{\mathcal{L}(\mathcal{H})}$, then $\mathbb{P}(\Theta(\lambda, \mathbf{z}) \leq 1/2) \geq 1 - \eta/3$, where

$$\Theta(\lambda, \mathbf{z}) = \|(T - T_{\mathbf{x}})(T + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})}$$

and one can obtain upper bounds on S_1 and S_2 which hold with probability $1 - \eta$. For $\mathcal{A}(\lambda)$ no probabilistic argument was needed.

• **Probabilistic bounds on $\|g_{\hat{\mathbf{z}}} - g_{\mathbf{z}}\|_{\mathcal{H}}^2$, $\|T_{\mathbf{x}} - T_{\hat{\mathbf{x}}}\|_{\mathcal{L}(\mathcal{H})}^2$, $\|\sqrt{T}(T_{\hat{\mathbf{x}}} + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})}^2$, $\|f_{\hat{\mathbf{z}}}^\lambda\|_{\mathcal{H}}^2$:** By using the $\|Mu\|_{\mathcal{H}} \leq \|M\|_{\mathcal{L}(\mathcal{H})} \|u\|_{\mathcal{H}}$ ($M \in \mathcal{L}(\mathcal{H}), u \in \mathcal{H}$) inequality, we bound S_{-1} and S_0 as

$$S_{-1} \leq \|\sqrt{T}(T_{\hat{\mathbf{x}}} + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})}^2 \|g_{\hat{\mathbf{z}}} - g_{\mathbf{z}}\|_{\mathcal{H}}^2,$$

$$S_0 \leq \|\sqrt{T}(T_{\hat{\mathbf{x}}} + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})}^2 \|T_{\mathbf{x}} - T_{\hat{\mathbf{x}}}\|_{\mathcal{L}(\mathcal{H})}^2 \|f_{\mathbf{z}}^\lambda\|_{\mathcal{H}}^2.$$

For the terms on the r.h.s., we can derive the upper bounds [for α see Eq. (13)]:

$$\begin{aligned} \|g_{\hat{\mathbf{z}}} - g_{\mathbf{z}}\|_{\mathcal{H}}^2 &\leq L^2 C^2 \frac{(1 + \sqrt{\alpha})^{2h} (2B_k)^h}{N^h}, \\ \|\sqrt{T}(T_{\hat{\mathbf{x}}} + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})} &\leq \frac{2}{\sqrt{\lambda}}, \\ \|T_{\mathbf{x}} - T_{\hat{\mathbf{x}}}\|_{\mathcal{L}(\mathcal{H})}^2 &\leq \frac{(1 + \sqrt{\alpha})^{2h} 2^{h+2} (B_k)^h B_K L^2}{N^h}, \\ \|f_{\mathbf{z}}^\lambda\|_{\mathcal{H}}^2 &\leq \frac{C^2 B_K}{\lambda^2}. \end{aligned}$$

The bounds hold under the following conditions:

⁶Note that the $N \geq \log(l)/\lambda^{\frac{2}{h}}$ constraint has been discarded; it is implied by the convergence of the first term in r [Eq. (8)] (see the supplementary material).

Table 2: Convergence conditions, convergence rates. Rows from top: 1 – 2, 1 – 3, 1 – 4, 2 – 3, 2 – 4, 3 – 4th terms are matched in $r(l, N, \lambda)$, the upper bound on the excess risk; see Eq. (8). First column: convergence condition. Second column: conditions for the dominance of the matched terms *while* they also converge to zero. Third column: convergence rate of the excess risk.

Convergence condition	Dominance + convergence condition	Convergence rate
$\max\left(\frac{h}{(c+3)\min(2,b)}, \frac{h(b+1)}{(c+3)b}\right) \leq a$	$\max\left(\frac{h(\frac{1}{b}+c)}{c+3}, \frac{h(b+1)}{(c+3)b}\right) \leq a$	$\left[\frac{\log(N)}{N}\right]^{\frac{hc}{c+3}}$
$\max\left(\frac{h}{6}, \frac{h}{2(b+1)}, \frac{h(b+1)}{2(2b+1)}\right) \leq a < \frac{h}{2}$	$\max\left(\frac{h}{6}, \frac{h(b+1)}{2(2b+1)}\right) \leq a < \min\left(\frac{h}{2} - \frac{h}{c+3}, \frac{\frac{h}{2}(\frac{1}{b}-1)}{\frac{1}{b}-2}\right)$	$\frac{1}{N^{3a-\frac{h}{2}} \log^{\frac{h}{2}}(N)}$
$\max\left(\frac{hb}{7b-2}, \frac{h}{3b}, \frac{h(b+1)}{4b}\right) \leq a < h$	$\max\left(\frac{h(b-1)}{4b-2}, \frac{h}{3b}, \frac{h(b+1)}{4b}\right) \leq a < \frac{h(bc+1)}{3b+bc}$	$\frac{1}{N^{a+\frac{a-h}{3b-1}} \log^{\frac{h}{3b-1}}(N)}$
$a < \frac{h(c+1)}{6}, 1 > \frac{2(b+1)}{(c+1)b}$	never	never
$a < \frac{h(bc+1)}{3b}, 1 > \frac{b+1}{bc+1}$	$a < \frac{h(bc+1)}{3b+bc}, 1 > \frac{b+1}{bc+1}$	$\frac{1}{N^{\frac{abc}{bc+1}}}$
never	never	never

1. $\|g_{\mathbf{z}} - g_{\mathbf{z}}\|_{\mathcal{H}}^2$: if the empirical mean embeddings are close to their population counterparts, i.e.,

$$\|\mu_{x_i} - \mu_{\hat{x}_i}\|_H \leq \frac{(1 + \sqrt{\alpha})\sqrt{2B_k}}{\sqrt{N}} \tag{13}$$

for $\forall i = 1, \dots, l$. This event has probability $1 - le^{-\alpha}$ over all l samples by a union bound.

2. $\|T_{\mathbf{x}} - T_{\hat{\mathbf{x}}}\|_{\mathcal{L}(\mathcal{H})}^2$: (13) is assumed.

3. $\|\sqrt{T}(T_{\hat{\mathbf{x}}} + \lambda)^{-1}\|_{\mathcal{L}(\mathcal{H})}^2$: $\frac{(1+\sqrt{\alpha})^2 2^{\frac{h+6}{h}} B_k(B_K)^{\frac{1}{h}} L^{\frac{2}{h}}}{(\lambda)^{\frac{2}{h}}} \leq N$, (13), and $\Theta(\lambda, \mathbf{z}) \leq \frac{1}{2}$.

4. $\|f_{\mathbf{z}}^\lambda\|_{\mathcal{H}}^2$: This upper bound always holds (under the model assumptions).

• **Union bound:** By applying an $\alpha = \log(l) + \delta$ reparameterization, and combining the received upper bounds with [5]’s results for S_1 and S_2 , the theorem follows with a union bound.

Finally, we note that

- existing results were used at two points to simplify our analysis: bounding S_1 , S_2 , $\Theta(\lambda, \mathbf{z})$ [5] and $\|\mu_{x_i} - \mu_{\hat{x}_i}\|_H$ [25].

- although the primary focus of our paper is clearly theoretical, we have provided some illustrative experiments in the supplementary material. These include
 1. a comparison with the only alternative, theoretically justified distribution regression method [1]³ on supervised entropy learning, where our approach gives better performance,
 2. an experiment on aerosol prediction based on satellite images, where we perform as well as recent domain-specific, engineered methods [35] (which themselves beat state-of-the-art multiple instance learning alternatives).

5 CONCLUSION

In this paper we established the learning theory of distribution regression under mild conditions, for probability measures on separable, topological domains endowed with kernels. We analysed an algorithmically simple and parallelizable⁷ ridge regression scheme defined on the embeddings of the input distributions to a RKHS. As a special case of our analysis, we proved the consistency of regression for set kernels [7, 23] in the distribution-to-real regression setting (which was a 15-year-old open problem), and for a recent kernel family [29], which we have expanded upon (Table 1). To keep the presentation simple we focused on the quadratic loss (\mathcal{E}), bounded kernels (k, K), real-valued labels (Y), and mean embedding (μ) based distribution regression with i.i.d. samples ($\{x_{i,n}\}_{n=1}^N$). In future work, we will relax these assumptions, and also consider deriving bounds with approximation error (capturing the richness of class \mathcal{H} in the bounds).⁸ Another exciting open question is whether (i) lower bounds on convergence can be proved, (ii) optimal convergence rates can be derived, (iii) one can obtain error bounds for non-point estimates.

Acknowledgements

This work was supported by the Gatsby Charitable Foundation, and by NSF grants IIS1247658 and IIS1250350. The work was carried out while Bharath K. Sriperumbudur was a research fellow in the Statistical Laboratory, Department of Pure Mathematics and Mathematical Statistics at the University of Cambridge, UK.

⁷Recently, [36] has constructed theoretically sound parallelization algorithms for kernel ridge regression.

⁸The extension to separable Hilbert output spaces and the misspecified case with approximation error are already available [37].

References

- [1] Barnabás Póczos, Alessandro Rinaldo, Aarti Singh, and Larry Wasserman. Distribution-free distribution regression. *AISTATS (JMLR W&CP)*, 31:507–515, 2013.
- [2] Thomas G. Dietterich, Richard H. Lathrop, and Tomás Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial Intelligence*, 89:31–71, 1997.
- [3] Soumya Ray and David Page. Multiple instance regression. In *ICML*, pages 425–432, 2001.
- [4] Daniel R. Dooly, Qi Zhang, Sally A. Goldman, and Robert A. Amar. Multiple-instance learning of real-valued data. *Journal of Machine Learning Research*, 3:651–678, 2002.
- [5] Andrea Caponnetto and Ernesto De Vito. Optimal rates for regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7:331–368, 2007.
- [6] Ingo Steinwart and Andres Christmann. *Support Vector Machines*. Springer, 2008.
- [7] David Haussler. Convolution kernels on discrete structures. Technical report, Department of Computer Science, University of California at Santa Cruz, 1999. (<http://cbse.soe.ucsc.edu/sites/default/files/convolutions.pdf>).
- [8] Fei Wang, Tanveer Syeda-Mahmood, Baba C. Vemuri, David Beymer, and Anand Rangarajan. Closed-form Jensen-Rényi divergence for mixture of Gaussians and applications to group-wise shape registration. *Medical Image Computing and Computer-Assisted Intervention*, 12:648–655, 2009.
- [9] Risi Kondor and Tony Jebara. A kernel between sets of vectors. In *ICML*, pages 361–368, 2003.
- [10] Tony Jebara, Risi Kondor, and Andrew Howard. Probability product kernels. *Journal of Machine Learning Research*, 5:819–844, 2004.
- [11] Frank Nielsen and Richard Nock. A closed-form expression for the Sharma-Mittal entropy of exponential families. *Journal of Physics A: Mathematical and Theoretical*, 45:032003, 2012.
- [12] Shaohua Kevin Zhou and Rama Chellappa. From sample similarity to ensemble similarity: Probabilistic distance measures in reproducing kernel Hilbert space. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28:917–929, 2006.
- [13] Bernhard Schölkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2002.
- [14] Marco Cuturi, Kenji Fukumizu, and Jean-Philippe Vert. Semigroup kernels on measures. *Journal of Machine Learning Research*, 6:11691198, 2005.
- [15] André F. T. Martins, Noah A. Smith, Eric P. Xing, Pedro M. Q. Aguiar, and Mário A. T. Figueiredo. Nonextensive information theoretical kernels on measures. *Journal of Machine Learning Research*, 10:935–975, 2009.
- [16] Matthias Hein and Olivier Bousquet. Hilbertian metrics and positive definite kernels on probability measures. In *AISTATS*, pages 136–143, 2005.
- [17] Barnabás Póczos, Liang Xiong, and Jeff Schneider. Nonparametric divergence estimation with applications to machine learning on distributions. In *UAI*, pages 599–608, 2011.
- [18] Barnabás Póczos, Liang Xiong, Dougal Sutherland, and Jeff Schneider. Support distribution machines. Technical report, Carnegie Mellon University, 2012. (<http://arxiv.org/abs/1202.0302>).
- [19] Junier B. Oliva, Willie Neiswanger, Barnabás Póczos, Jeff Schneider, and Eric Xing. Fast distribution to real regression. *AISTATS (JMLR W&CP)*, 33:706–714, 2014.
- [20] Junier Oliva, Barnabás Póczos, and Jeff Schneider. Distribution to distribution regression. *ICML (JMLR W&CP)*, 28:1049–1057, 2013.
- [21] László Györfi, Michael Kohler, Adam Krzyżak, and Harro Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer, New-york, 2002.
- [22] Samory Kpotufe. k-NN regression adapts to local intrinsic dimension. Technical report, Max Planck Institute for Intelligent Systems, 2011. (<http://arxiv.org/abs/1110.4300>).
- [23] Thomas Gärtner, Peter A. Flach, Adam Kowalczyk, and Alexander Smola. Multi-instance kernels. In *ICML*, pages 179–186, 2002.
- [24] Alain Berlinet and Christine Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer, 2004.

- [25] Yasemin Altun and Alexander Smola. Unifying divergence minimization and statistical inference via convex duality. In *COLT*, pages 139–153, 2006.
- [26] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.
- [27] Bharath K. Sriperumbudur, Kenji Fukumizu, and Gert R. G. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12:2389–2410, 2011.
- [28] Kenji Fukumizu, Francis Bach, and Michael Jordan. Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces. *Journal of Machine Learning Research*, 5:73–99, 2004.
- [29] Andreas Christmann and Ingo Steinwart. Universal kernels on non-standard input spaces. In *NIPS*, pages 406–414, 2010.
- [30] Krikamol Muandet, Kenji Fukumizu, Francesco Dinuzzo, and Bernhard Schölkopf. Learning from distributions via support measure machines. In *NIPS*, pages 10–18, 2012.
- [31] Larry Wasserman. *All of Nonparametric Statistics*. Springer, 2006.
- [32] Felipe Cucker and Steve Smale. On the mathematical foundations of learning. *Bulletin of the American Mathematical Society*, 39:1–49, 2002.
- [33] Christian Berg, Jens Peter Reus Christensen, and Paul Ressel. *Harmonic Analysis on Semigroups*. Springer-Verlag, 1984.
- [34] Bharath Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Gert Lanckriet, and Bernhard Schölkopf. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 11:1517–1561, 2010.
- [35] Zhuang Wang, Liang Lan, and Slobodan Vucetic. Mixture model for multiple instance regression and applications in remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 50:2226–2237, 2012.
- [36] Yuchen Zhang, John C. Duchi, and Martin J. Wainwright. Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. Technical report, University of California, Berkeley, 2014. (<http://arxiv.org/abs/1305.5029>).
- [37] Zoltán Szabó, Bharath Sriperumbudur, Barnabás Póczos, and Arthur Gretton. Learning theory for distribution regression. Technical report, Gatsby Unit, University College London, 2014. (<http://arxiv.org/abs/1411.2066>).
- [38] Michael Reed and Barry Simon. *Methods of Modern Mathematical Physics – Functional Analysis*. Academic Press, 1980.
- [39] Gabriel Nagy. Real analysis (lecture notes): Chapter III: Measure theory, Section 3: Measurable spaces and measurable maps. Technical report, Kansas State University. (<http://www.math.ksu.edu/~nagy/real-an/3-03-measbl.pdf>).
- [40] K.R. Parthasarathy. *Probability Measures on Metric Spaces*. Academic Press, 1967.
- [41] John L. Kelley. *General Topology*. Springer, 1975.
- [42] John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [43] Zoltán Szabó. Information theoretical estimators toolbox. *Journal of Machine Learning Research*, 15:283–287, 2014.