

---

# Distribution-Free Distribution Regression

---

Barnabás Póczos<sup>1</sup>

Alessandro Rinaldo<sup>2</sup>

Aarti Singh<sup>3</sup>

Larry Wasserman<sup>4</sup>

Carnegie Mellon University  
Machine Learning Department<sup>1,3,4</sup>  
Department of Statistics<sup>2,4</sup>  
Pittsburgh, PA, USA, 15213

## Abstract

‘Distribution regression’ refers to the situation where a response  $Y$  depends on a covariate  $P$  where  $P$  is a probability distribution. The model is  $Y = f(P) + \mu$  where  $f$  is an unknown regression function and  $\mu$  is a random error. Typically, we do not observe  $P$  directly, but rather, we observe a sample from  $P$ . In this paper we develop theory and methods for distribution-free versions of distribution regression. This means that we do not make strong distributional assumptions about the error term  $\mu$  and covariate  $P$ . We prove that when the effective dimension is small enough (as measured by the doubling dimension), then the excess prediction risk converges to zero with a polynomial rate.

## 1 Introduction

In a standard regression model, we need to predict a real-valued response  $Y$  from a vector-valued covariate (or feature)  $X \in \mathbb{R}^d$ . Recently, there has been interest in extensions of standard regression from finite dimensional Euclidean spaces to other domains. For example, in functional regression (Ferraty and Vieu [2006]) the covariate is a function instead of a finite dimensional vector.

In this paper, we study *distribution regression* where the covariate is a probability distribution  $P$ . This differs from functional regression in two important ways. First,  $P$  is a probability measure on  $\mathbb{R}^k$  rather than a one-dimensional function. Second, and more importantly, we do not observe the covariate  $P$  directly.

---

Appearing in Proceedings of the 16<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2013, Scottsdale, AZ, USA. Volume 31 of JMLR: W&CP 31. Copyright 2013 by the authors.

Rather, we observe a sample from  $P$ , which means that we have a regression model with measurement error (Carroll et al. [2006], Fan and Truong [1993]).

A practical example where this framework can be useful is as follows. Suppose that we need to classify patients in a hospital and diagnose whether they are healthy or suffer from a disease. Traditional machine learning based approaches would make a couple of medical tests, and using the results of these measurements they would form a feature vector for each person and then apply a standard classifier to predict the class labels of the feature vectors. Suppose we have  $m$  patients, and these feature vectors are denoted by  $X_i \in \mathbb{R}^d$ ,  $1 \leq i \leq m$ . Our goal is to predict the class label  $Y \in \{\text{‘healthy’}, \text{‘diseased’}\}$  for a person. The problem with this approach is that our heart rate, blood pressure, chemical concentrations in blood, and many other medical conditions in our body are always changing, and therefore if we repeat these measurements a couple of times, then each time we might get different measurements and different feature vectors for the same person. For the  $i$ th person, let the set of these measurements be denoted by  $\mathcal{X}_i = \{X_{i,1}, \dots, X_{i,n_i}\}$ , where  $X_{i,n_i} \in \mathbb{R}^d$  indicates that we repeated the medical tests  $n_i$  times. Interestingly, traditional feature vector based machine learning algorithms cannot handle well such simple problems. They might construct a new feature vector as the average of the measurements ( $\tilde{X}_i \doteq \frac{1}{n_i} \sum_{j=1}^{n_i} X_{i,j}$ ), but then they lose information. If they want to keep all the measurements in a feature vector, then they cannot just simply stack the feature vectors of each person to a larger vector, because then each of these vectors could have different sizes ( $dn_i$ ). In contrast to the approaches, in our framework we simply say that each person is represented by an unknown distribution  $P_i$ , and those feature vectors are samples from these distributions  $X_{i,j} \sim P_i$  for  $j = 1, \dots, n_i$ . Our goal is to classify these unknown  $P_i$  distributions.

The formal definition of the problem is as follows.

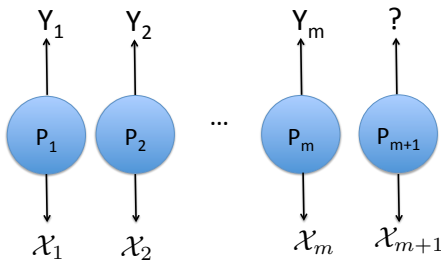


Figure 1: Illustration of the model - distributions  $P_1, \dots, P_m, P_{m+1}$  are unobserved, only the  $\mathcal{X}_1, \dots, \mathcal{X}_m, \mathcal{X}_{m+1}$  sample sets are observable.

We consider a regression problem with variables  $(P_1, Y_1), \dots, (P_m, Y_m)$  where  $Y_i \in \mathbb{R}$  and each  $P_i$  is a probability distribution on a compact subset  $\mathcal{K} \subset \mathbb{R}^k$ . We assume that

$$Y_i = f(P_i) + \mu_i, \quad i = 1, \dots, m,$$

for some functional  $f$ , where  $\mu_i$  is a noise variable with mean 0. We do not observe  $P_i$  directly; rather we observe a sample

$$X_{i1}, \dots, X_{in_i} \stackrel{i.i.d.}{\sim} P_i. \quad (1)$$

Thus the observed data are

$$(\mathcal{X}_1, Y_1), \dots, (\mathcal{X}_m, Y_m) \quad (2)$$

where  $\mathcal{X}_i = \{X_{i1}, \dots, X_{in_i}\}$ . Our goal is to predict a new  $Y_{m+1}$  from a new batch  $\mathcal{X}_{m+1}$  drawn from a new distribution  $P_{m+1}$ . This model is illustrated in Figure 1.

We model the unobservable probability distributions  $P_1, \dots, P_m$  as follows. Let  $\mathbb{D}$  denote the set of all distributions on  $\mathcal{K}$  that have a density with respect to the Lebesgue measure. We assume that the distributions  $P_i$  are an i.i.d. sample from a measure  $\mathcal{P}$  on  $\mathbb{D}$ , that is<sup>1</sup>,

$$P_1, \dots, P_m, P_{m+1} \stackrel{i.i.d.}{\sim} \mathcal{P}.$$

Note that  $f: \mathbb{D} \rightarrow \mathbb{R}$ . If  $Q(\cdot|P)$  denotes the law of  $Y$  given  $P$ , then the joint distribution of  $(Y, P)$  is given by

$$\mathbb{P}(Y \in A, P \in B) = Q(Y \in A|P \in B)\mathcal{P}(P \in B)$$

*Our main result* is a theorem where we prove that when the effective dimension of  $\mathcal{P}$  measured by the

<sup>1</sup>There are some subtle technical difficulties with the definitions of measurability. Using outer expectations these issues can be resolved. In this paper, however, we ignore this question.

doubling dimension is small enough, then the estimator is consistent and the prediction risk converges to zero with a polynomial rate. Our results are *distribution free*, similar to the functional regression case Ferraty and Vieu [2006], in the sense that we do not make any strong distributional assumptions.

**Outline.** In Section 2 we discuss related work. We propose a specific estimator for distribution regression in Section 3. We call this *kernel-kernel estimator* since it makes use of kernels in two different ways. In Section 4 we derive an upper bound on the risk of the estimator. The proofs can be found in Section 5. In Section 6 we analyze the risk bound in terms of the doubling dimension, which is a measure of the intrinsic dimension of the space. We present numerical illustrations in Section 7. Finally, we give some concluding remarks in Section 8. The details of the proofs can be found in the Supplementary material [Póczos et al., 2013].

## 2 Related work

Our framework is related to functional data analysis, which is a new and steadily improving field of statistics. For comprehensive reviews and references, see Ramsay and Silverman [2005], Ferraty and Vieu [2006].

A popular approach to do machine learning, such as classification and regression, on the domain of distributions is to embed the distribution to a Hilbert space, introduce kernels between the distributions, and then use a traditional kernel machine to solve the learning problem. There are both parametric and nonparametric methods proposed in the literature.

Parametric methods, (e.g. Jebara et al. [2004], Moreno et al. [2004], Jaakkola and Haussler [1998]), usually fit a parametric family (e.g. Gaussians distributions or exponential family) to the densities, and using the fitted parameters they estimate the inner products between the distributions. The problem with parametric approaches, however, is that when the true densities do not belong to the assumed parametric families, then this method introduces some unavoidable bias during the estimation of the inner products between the densities.

A couple of nonparametric approaches exist as well. Since our covariates are represented by finite sets, reproducing kernel Hilbert space (RKHS) based set kernels can be used in these learning problems. Smola et al. [2007] proposed to embed the distributions to an RKHS using the mean map kernels. In this framework, the role of universal kernels have been studied by Christmann and Steinwart [2010]. Recently, the representer theorem has also been generalized for the space of probability distributions [Muandet et al., 2012].

Kondor and Jebara [2003] introduced Bhattacharyya’s measure of affinity between finite-dimensional Gaussians in a Hilbert space. In contrast to the previous approaches, Póczos et al. [2012], Póczos et al. [2011] used nonparametric Rényi divergence estimators to solve machine learning problems on the set of distributions.

Although, there are a few algorithms designed for regression on distributions, we know very little about their theoretical properties. To the best of our knowledge, even the simplest, fundamental questions have not been studied yet. For example, we do not know how many training distributions ( $m$ ) and how many samples ( $n_i, i = 1, \dots, m$ ) we need to achieve a target prediction error. Our paper is providing an answer to this question.

### 3 The Kernel-Kernel Estimator

In this section we define an estimator  $\hat{f}$  for the unknown function  $f$ . Let  $\hat{P}_i$  denote an estimator of  $P_i$  based on  $\mathcal{X}_i$ , and let  $\mathcal{X}$  be a sample from a new distribution  $P = P_{m+1}$ . Accordingly, we denote with  $\hat{P}$  an estimator of  $P$  based on  $\mathcal{X}$ . Our predictor for  $Y_{m+1}$  is then  $\hat{Y}_{m+1} = \hat{f}(\hat{P}_{m+1})$ .

Given a bandwidth  $h > 0$  and a kernel function  $K$  (whose properties will be specified later), we define

$$\begin{aligned} \hat{f}(\hat{P}) &= \hat{f}(\hat{P}; \hat{P}_1, \dots, \hat{P}_m) \\ &= \begin{cases} \frac{\sum_i Y_i K\left(\frac{D(\hat{P}_i, \hat{P})}{h}\right)}{\sum_i K\left(\frac{D(\hat{P}_i, \hat{P})}{h}\right)} & \text{if } \sum_i K\left(\frac{D(\hat{P}_i, \hat{P})}{h}\right) > 0 \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

To complete the definition, we need to specify  $\hat{P}_i$ ,  $\hat{P}$  and  $D$ . We will estimate  $P_i$  — or, more precisely, the density  $p_i$  of  $P_i$  — with a kernel density estimator

$$\hat{p}_i(x) = \frac{1}{n_i} \sum_{j=1}^{n_i} \frac{1}{b_i^k} B\left(\frac{\|x - X_{ij}\|_2}{b_i}\right) \quad (3)$$

where  $B$  is an appropriate kernel function (see, e.g. Tsybakov [2010]) with bandwidth  $b_i > 0$ . Here  $\|x\|_2$  denotes the Euclidean norm of  $x \in \mathbb{R}^k$ . Accordingly,  $\hat{P}_i$  is defined by

$$\hat{P}_i(A) = \int_A \hat{p}_i(u) du,$$

for all Borel measurable subsets of  $\mathbb{R}^k$ . For any two probabilities in  $\mathbb{D}$ , we take  $D(P, Q)$  to be the  $L_1$  distance of their densities:  $D(P, Q) = \|p - q\| = \int |p(x) - q(x)| dx$ . Hence,

$$\hat{f}(\hat{P}) = \hat{f}(\hat{P}; \hat{P}_1, \dots, \hat{P}_m) = \frac{\sum_{i=1}^m Y_i K\left(\frac{\|\hat{p} - \hat{p}_i\|}{h}\right)}{\sum_{i=1}^m K\left(\frac{\|\hat{p} - \hat{p}_i\|}{h}\right)} \quad (4)$$

which we call the ‘kernel-kernel estimator’ since it makes use of two kernels,  $B$  and  $K$ .

For simplicity,  $n$  will denote the size of the sample  $\mathcal{X}$ , and  $b$  will be the bandwidth in the estimator of  $\hat{p}$ .

In what follows we will make the following assumptions on  $f, K, \mathcal{P}, \mu_i$ , and  $Y_i$ .

#### Assumptions

- (A1) *Hölder continuous functional.* The unknown functional  $f$  belongs to the class  $\mathcal{M} = \mathcal{M}(L, \beta, D)$  of Hölder continuous functionals on  $\mathbb{D}$ :

$$\mathcal{M} = \left\{ f : |f(P_i) - f(P_j)| \leq L D(P_i, P_j)^\beta \right\},$$

for some  $L > 0$  and  $0 < \beta \leq 1$ , where  $D$  is the above specified  $L_1$  metric on  $\mathbb{D}$ . In the  $\beta = 1$  special case this means that  $f$  is Lipschitz continuous.

- (A2) *Asymmetric boxed and Lipschitz kernel.* The kernel  $K$  satisfies the following properties:  $K : [0, \infty] \rightarrow \mathbb{R}$  is non-negative and Lipschitz continuous with Lipschitz constant  $L_K$ . In addition, there exist constants  $0 < \underline{K} < 1$  and  $0 < r < R < \infty$  such that, for all  $x > 0$ , it holds that

$$\underline{K} I_{\{x \in \mathcal{B}(0, r)\}} \leq K(x) \leq I_{\{x \in \mathcal{B}(0, R)\}}.$$

- (A3) *Hölder class of distributions.* The distribution  $\mathcal{P}$  is supported on the set of distributions with densities that are 1-smooth Hölder functions, as defined in Tsybakov [2010], Rigollet and Vert [2009] for example.

- (A4) *Bounded regression.* We will assume that  $\sup_{P \in \mathcal{P}} |f(P)| < f_{\max}$  for some  $f_{\max} > 0$ . Also,  $\mu_i$  has mean 0 and  $\mathbb{P}(|Y_i| \leq B_Y) = 1$  for some  $B_Y < \infty$ .

- (A5) *Lower bound on  $\min_{1 \leq i \leq m+1} n_i$ .* Let  $n = \min_{1 \leq i \leq m+1} n_i$ . We assume that  $n^{\frac{k}{2+k}} \geq 3 \ln m$ .

- (A6) *Requirements on regression kernel bandwidth  $h$ .* Assume that  $C_* n^{-\frac{1}{2+k}} \leq rh/4$  where  $C_*$  is defined in (9), and  $h \leq H$  where  $H > 0$  is a constant.

- (A7) *Requirement on density kernel bandwidths  $\{b_i\}_{i=1}^m$ .* Assume the bandwidths  $b_i = b := n^{-\frac{1}{k+2}}$ .

### 4 Upper Bound on Risk

We are concerned with upper bounding the risk

$$R(m, n) = \mathbb{E} \left[ |\hat{f}(\hat{P}; \hat{P}_1, \dots, \hat{P}_m) - f(P)| \right],$$

where the expectation is with respect to the joint distribution of the sample  $(\mathcal{X}_1, Y_1), \dots, (\mathcal{X}_m, Y_m)$ , the new covariate  $P = P_{m+1}$  and the new observation  $\mathcal{X}_{m+1}$ . Note that the absolute prediction risk is  $\mathbb{E}|\widehat{Y} - Y| \leq R(m, n) + c$ , where  $c = \mathbb{E}(|\mu|)$  is a constant. So bounding the prediction risk is equivalent to bounding  $R(m, n)$ , which we call the excess prediction risk. In what follows,  $C, c_1, c_2, \dots$  represent constants whose value can be different in different expressions.

Let  $\mathcal{B}(P, h) = \{\tilde{P} \in \mathbb{D} : D(\tilde{P}, P) \leq h\}$  denote the  $L_1$  ball of distributions around  $P$  with radius  $h$ . We will see that the risk depends on the size of the class of probabilities  $\mathbb{D}$ . In particular, the risk depends on the *small ball probability*

$$\Phi_P(h) := \mathcal{P}(\mathcal{B}(P, h)),$$

where  $P$  is a fixed distribution and  $\Phi_P(h)$  is a function of  $P$ .

Our first result, Theorem 1, provides a general upper bound on the risk. In our second result (Section 6) we show that when the effective dimension measured by the doubling dimension is small, then the risk converges to zero. We also derive an upper bound on the rate of convergence.

**Theorem 1** *Suppose that the assumptions (A1)-(A7) stated above hold. Then*

$$\begin{aligned} R(m, n) &\leq \frac{1}{h} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right] C_1 n^{-\frac{1}{2+\kappa}} + C_2 h^\beta \\ &+ C_3 \sqrt{\frac{1}{m}} \sqrt{\mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right]} + \frac{C_4}{m} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right] \\ &+ (m+1) e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}}, \end{aligned}$$

where the constants  $C_i$ 's are specified in the proof.

## 5 Proof of Theorem 1

In this Section we prove our main result, Theorem 1. The main idea of the proof is to use the triangle inequality to write

$$\begin{aligned} R(m, n) &= \mathbb{E}|\widehat{f}(\widehat{P}; \widehat{P}_1, \dots, \widehat{P}_m) - f(P)| \\ &\leq \mathbb{E}|\widehat{f}(\widehat{P}; \widehat{P}_1, \dots, \widehat{P}_m) - \widehat{f}(P; P_1, \dots, P_m)| \quad (5) \\ &+ \mathbb{E}|\widehat{f}(P; P_1, \dots, P_m) - f(P)|. \quad (6) \end{aligned}$$

In Sections 5.2 and 5.3 we will derive upper bounds for (5) and (6), respectively. Section 5.1 contains a series of technical results needed in our proofs.

Throughout, we let  $\widehat{K}_i = K\left(\frac{D(\widehat{P}_i, \widehat{P})}{h}\right)$ ,  $K_i = K\left(\frac{D(P_i, P)}{h}\right)$  and  $\epsilon_i = K_i - \widehat{K}_i$ , for  $i = 1, \dots, m$ . Note that, for ease of readability, we have omitted the dependence on  $h$ .

## 5.1 Technical Results

### 5.1.1 $L_1$ Risk of Density Estimators

In this section we bound  $\mathbb{E}[D(P, \widehat{P})|P] = \mathbb{E}[\int |p - \widehat{p}| |P]$ , the  $L_1$  risk of the density estimator  $\widehat{p}$  of  $p$ , uniformly over all  $P$  in  $\mathbb{D}$ . To this end, suppose that  $n_i \geq n$  for all  $i = 1, 2, \dots, m+1$ , and let  $b_i = b = n^{-\frac{1}{\kappa+2}}$ . In this case, the following lemma provides upper bound on the  $L_1$  risk of the density estimator. Its proof can be found in the supplementary material.

**Lemma 2**

$$\begin{aligned} \mathbb{E}[D(\widehat{P}_i, P_i)|P_i] &\leq \bar{C} n^{-\frac{1}{2+\kappa}}, \quad (7) \\ \mathbb{E}[D(\widehat{P}_i, P_i)] &\leq \bar{C} n^{-\frac{1}{2+\kappa}}, \end{aligned}$$

where

$$\bar{C} = c_0(c_1 + c_2), \quad (8)$$

with  $c_0, c_1$  and  $c_2$  constants specified in the proof.

Next, we show that the terms  $D(\widehat{P}_i, P_i)$  are uniformly bounded by a term of order  $O(h)$ , with high probability.

**Lemma 3** *With probably no smaller than  $1 - (m+1)e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}}$ ,  $D(\widehat{P}_i, P_i) < \frac{rh}{4}$  for all  $i = 1, \dots, m+1$ .*

Notice that by Assumption (A5),  $1 - (m+1)e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}} \rightarrow 1$ .

**Proof.** From McDiarmid's inequality, for any  $\epsilon > 0$  we have that

$$\mathbb{P}(|\|\widehat{p}_i - p_i\|_1 - \mathbb{E}\|\widehat{p}_i - p_i\|_1| > \epsilon) \leq e^{-\epsilon^2/2}$$

(see, for example, section 2.4 of Devroye and Lugosi [2001]). Thus,

$$\mathbb{P}(\|\widehat{p}_i - p_i\|_1 > \mathbb{E}\|\widehat{p}_i - p_i\|_1 + n^{-\frac{1}{2+\kappa}}) \leq e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}},$$

since  $nn^{-\frac{2}{2+\kappa}} = n^{\frac{\kappa}{2+\kappa}}$ . This implies that

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq i \leq m+1} \|\widehat{p}_i - p_i\|_1 > \mathbb{E}\|\widehat{p}_i - p_i\|_1 + n^{-\frac{1}{2+\kappa}}\right) \\ \leq (m+1) e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}} \rightarrow 0, \end{aligned}$$

by assumption (A5). Therefore,

$$\begin{aligned} 1 - (m+1) e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}} \\ \leq \mathbb{P}\left(\max_{1 \leq i \leq m+1} \|\widehat{p}_i - p_i\|_1 \leq \mathbb{E}\|\widehat{p}_i - p_i\|_1 + n^{-\frac{1}{2+\kappa}}\right) \\ \leq \mathbb{P}\left(\max_{1 \leq i \leq m+1} \|\widehat{p}_i - p_i\|_1 \leq (1 + c_0(c_1 + c_2)) n^{-\frac{1}{2+\kappa}}\right). \end{aligned}$$

This implies that with

$$C_* = (1 + c_0(c_1 + c_2)) \quad (9)$$

and using assumption (A6), we have that

$$D(\widehat{P}_i, P_i) \leq C_* n^{-\frac{1}{\kappa+2}} \leq \frac{rh}{4} \quad \text{for all } i \quad (10)$$

on an event  $\Omega_{m,n}$ , where  $\mathbb{P}(\Omega_{m,n}^c) \leq (m+1)e^{-\frac{1}{2}n^{\frac{\kappa}{2+\kappa}}}$ . Here  $\Omega_{m,n}^c$  denotes the complement of  $\Omega_{m,n}$ .  $\square$

### 5.1.2 Other Lemmata

Throughout this section we will make use of the constant  $\bar{C}$ , defined in (8). In what follows, we will need a few lemmas that we list below. Their proofs can be found in the supplementary material.

The following lemma provides an upper bound on  $\mathbb{P}(\sum_{i=1}^m K_i = 0)$  with the help of small ball probabilities.

#### Lemma 4

$$\mathbb{P}\left(\sum_{i=1}^m K_i = 0\right) \leq \mathbb{P}\left(\sum_{i=1}^m K_i < \underline{K}\right) = \frac{1}{em} \mathbb{E}\left[\frac{1}{\Phi_P(rh)}\right].$$

We will also need the following lemma.

#### Lemma 5

$$\mathbb{E}\left[\frac{1}{\sum_i K_i} I_{\{\sum_i K_i \geq \underline{K}\}}\right] \leq \frac{1+1/\underline{K}}{m\underline{K}} \mathbb{E}\left[\frac{1}{\Phi_P(rh)}\right].$$

The following lemma provides an upper bound on  $|\epsilon_i|$ .

**Lemma 6** *Assume that the kernel function  $K$  is Lipschitz continuous with Lipschitz constant  $L_K$ . We have that*

$$|\epsilon_i| \leq \frac{L_K}{h} (D(P, \widehat{P}) + D(P_i, \widehat{P}_i)).$$

By definition,  $|\epsilon_i| = |K_i - \widehat{K}_i| = |K(\frac{D(P, P_i)}{h}) - K(\frac{D(\widehat{P}, \widehat{P}_i)}{h})|$ , which is a deterministic function of random variables  $P, P_i, \widehat{P}$ , and  $\widehat{P}_i$ . We will denote this deterministic relationship as  $\epsilon_i = \epsilon_i(P, \widehat{P}, P_i, \widehat{P}_i)$ . The following lemma shows that for any  $\kappa > 0$ ,

$$\mathbb{P}\left(\sum_i |\epsilon_i(P, \widehat{P}, P_i, \widehat{P}_i)| < \kappa | \{P_i\}_{i=1}^m, P\right)$$

can be lower bounded by a non-trivial quantity that does not depend on  $P$  and  $\{P_i\}_{i=1}^m$ .

**Lemma 7** *For any  $\kappa > 0$  we have that*

$$\mathbb{P}\left(\sum_i |\epsilon_i(P, \widehat{P}, P_i, \widehat{P}_i)| < \kappa | \{P_i\}_{i=1}^m, P\right) \geq \eta,$$

where  $\eta = \eta(\kappa, n, m) = 1 - \frac{2L_K m \bar{C}}{h\kappa} n^{-\frac{1}{2+\kappa}}$ .

The following lemma provides an upper bound on the expected value of  $\sum_{i=1}^m |\epsilon_i|$ .

#### Lemma 8

$$\mathbb{E}\left[\sum_{i=1}^m |\epsilon_i| \middle| P, \{P_i\}_{i=1}^m\right] \leq \frac{2L_K \bar{C} m}{h} n^{-\frac{1}{2+\kappa}}.$$

The next lemma shows that  $\mathbb{P}\left(\sum_{i=1}^m \widehat{K}_i < \underline{K}\right)$  can be upper bounded by a small quantity as well. We assume that  $n_i = n$  and  $b_i = b$  for all  $i$ . Define

$$\zeta = \zeta(n, m) = \frac{1}{em} \mathbb{E}\left(\frac{1}{\Phi_P\left(\frac{rh}{2}\right)}\right) + (m+1)e^{-\frac{1}{2}n^{\frac{\kappa}{2+\kappa}}}.$$

#### Lemma 9

$$\mathbb{P}\left(\sum_{i=1}^m \widehat{K}_i = 0\right) \leq \mathbb{P}\left(\sum_{i=1}^m \widehat{K}_i < \underline{K}\right) \leq \zeta.$$

### 5.2 Upper bound on Equation 5

Let  $\Delta \widehat{f} = |\widehat{f}(\widehat{P}; \widehat{P}_1, \dots, \widehat{P}_m) - \widehat{f}(P; P_1, \dots, P_m)|$ . Our goal is to provide an upper bound on  $\mathbb{E}[\Delta \widehat{f}]$ .

Introduce the following events:  $E_0 = \{\sum_i K_i = 0\}$ ,  $E_1 = \{0 < \sum_i K_i < \underline{K}\}$ ,  $E_2 = \{\underline{K} \leq \sum_i K_i\}$ . Similarly,  $\widehat{E}_0 = \{\sum_i \widehat{K}_i = 0\}$ ,  $\widehat{E}_1 = \{0 < \sum_i \widehat{K}_i < \underline{K}\}$ ,  $\widehat{E}_2 = \{\underline{K} \leq \sum_i \widehat{K}_i\}$ . Obviously,  $\mathbb{E}[\Delta \widehat{f}] = \sum_{k=0}^2 \sum_{l=0}^2 \mathbb{E}[\Delta \widehat{f} I_{E_k} I_{\widehat{E}_l}]$ .

Based on the sign of  $\sum_i K_i$  and  $\sum_i \widehat{K}_i$ , there are four different cases. (i) If  $\sum_i K_i > 0$  and  $\sum_i \widehat{K}_i > 0$ , then  $\Delta \widehat{f} = \left| \frac{\sum_i Y_i \widehat{K}_i}{\sum_i \widehat{K}_i} - \frac{\sum_i Y_i K_i}{\sum_i K_i} \right|$ . (ii) If  $\sum_i K_i > 0$  and  $\sum_i \widehat{K}_i = 0$ , then  $\Delta \widehat{f} = \left| \frac{\sum_i Y_i K_i}{\sum_i K_i} \right|$ . (iii) If  $\sum_i K_i = 0$  and  $\sum_i \widehat{K}_i > 0$ , then  $\Delta \widehat{f} = \left| \frac{\sum_i Y_i \widehat{K}_i}{\sum_i \widehat{K}_i} \right|$ , and finally (iv) if  $\sum_i K_i = 0$  and  $\sum_i \widehat{K}_i = 0$ , then  $\Delta \widehat{f} = 0$ . From this it immediately follows that  $\mathbb{E}[\Delta \widehat{f} I_{E_0} I_{\widehat{E}_0}] = 0$ .

When  $\sum_i K_i > 0$ ,  $\left| \frac{\sum_i Y_i K_i}{\sum_i K_i} \right| \leq B_Y$ . Therefore,

$$\begin{aligned} & \mathbb{E}\left[\left|\frac{\sum_i Y_i K_i}{\sum_i K_i}\right| I_{\widehat{E}_0} (I_{E_1} + I_{E_2})\right] \\ & \leq B_Y \mathbb{E}\left[I_{\{\sum_i K_i > 0 \wedge \sum_i \widehat{K}_i = 0\}}\right] \\ & = B_Y \mathbb{P}\left(\sum_i K_i > 0, \sum_i \widehat{K}_i = 0\right) \\ & \leq B_Y \mathbb{P}\left(\sum_{i=1}^m \widehat{K}_i = 0\right) \leq B_Y \zeta(n, m). \end{aligned}$$

Similarly,

$$\mathbb{E}\left[\left|\frac{\sum_i Y_i \widehat{K}_i}{\sum_i \widehat{K}_i}\right| I_{E_0} (I_{\widehat{E}_1} + I_{\widehat{E}_2})\right] \leq \frac{B_Y}{em} \int \frac{dP(P)}{\Phi_P(rh)}.$$

It is also easy to see that

$$\begin{aligned}
 & \mathbb{E} \left[ \Delta \hat{f} I_{E_1} (I_{\hat{E}_1} + I_{\hat{E}_2}) \right] \\
 & \leq \mathbb{E} \left[ \left( \left| \sum_i \frac{Y_i K_i}{\sum_i K_i} \right| + \left| \sum_i \frac{Y_i \hat{K}_i}{\sum_i \hat{K}_i} \right| \right) I_{E_1} (I_{\hat{E}_1} + I_{\hat{E}_2}) \right] \\
 & \leq \mathbb{E} \left[ 2B_Y I_{E_1} (I_{\hat{E}_1} + I_{\hat{E}_2}) \right] \leq 2B_Y \mathbb{E} \left[ I_{E_1} \right] \\
 & = 2B_Y \mathbb{P}(0 < \sum_{i=1}^m K_i < \underline{K}) \leq \frac{2B_Y}{em} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)}.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \mathbb{E} \left[ \Delta \hat{f} I_{\hat{E}_1} (I_{E_1} + I_{E_2}) \right] & \leq 2B_Y \mathbb{P}(0 < \sum_{i=1}^m \hat{K}_i < \underline{K}) \\
 & \leq 2B_Y \zeta(n, m).
 \end{aligned}$$

All that left is to upper bound  $\mathbb{E} \left[ \Delta \hat{f} I_{E_2} I_{\hat{E}_2} \right]$ . The next lemma provides an upper bound for this.

**Lemma 10**

$$\mathbb{E} \left[ \Delta \hat{f} I_{E_2} I_{\hat{E}_2} \right] \leq C_1 \frac{1}{h} \mathbb{E} \left[ \frac{1}{\Phi_P(rh)} \right] n^{-\frac{1}{2+\kappa}}.$$

The proof can be found in the supplementary material.

Finally, putting the pieces together we obtain the following theorem.

**Theorem 11**

$$\begin{aligned}
 & \mathbb{E} |\hat{f}(\hat{P}; \hat{P}_1, \dots, \hat{P}_m) - \hat{f}(P; P_1, \dots, P_m)| \\
 & \leq C_1 \frac{1}{h} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right] n^{-\frac{1}{2+\kappa}} + C_2 \frac{1}{m} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right] \\
 & \quad + (m+1) e^{-\frac{1}{2} n^{\frac{\kappa}{2+\kappa}}}.
 \end{aligned}$$

The proof can be found in the supplementary material.

### 5.3 Upper bound on Equation 6

In this section we show that under the above specified conditions  $\mathbb{E} |\hat{f}(P; P_1, \dots, P_m) - f(P)|$  can be upper bounded by

$$C_1(h^\beta) + C_2 \left( \sqrt{\mathbb{E} \left[ \frac{1}{m\Phi_P(rh/2)} \right]} \right) + \frac{C_3}{m} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right],$$

where the expectation is with respect to the random probability measure  $P$  in  $\mathcal{P}$ .

We have to bound  $\mathbb{E} |\hat{f}(P; P_1, \dots, P_m) - f(P)|$ . Note that  $Y_i = f(P_i) + \mu_i$ , and

$$\mathbb{E} |\hat{f}(P; P_1, \dots, P_m) - f(P)|$$

$$\begin{aligned}
 & = \mathbb{E} \left| \frac{\sum_i Y_i K_i}{\sum_i K_i} I_{\{\sum_i K_i > 0\}} - f(P) \right| \\
 & = \mathbb{E} \left| \frac{\sum_i (f(P_i) + \mu_i) K_i}{\sum_i K_i} I_{\{\sum_i K_i > 0\}} - f(P) \right| \\
 & \leq \mathbb{E} \left[ \left| \frac{\sum_i (f(P_i) - f(P)) K_i}{\sum_i K_i} + \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i > 0\}} \right] \\
 & \quad + \mathbb{E} [|f(P)| I_{\{\sum_i K_i = 0\}}] \\
 & \leq \mathbb{E} \left[ \frac{\sum_i |f(P_i) - f(P)| K_i}{\sum_i K_i} I_{\{\sum_i K_i > 0\}} \right] \\
 & \quad + \mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i > 0\}} \right] + f_{\max} \mathbb{P}(\sum_i K_i = 0).
 \end{aligned}$$

We will bound each of the three terms next. For the first term, since  $f$  is Hölder- $\beta$  we have

$$\begin{aligned}
 & \mathbb{E} \left[ \frac{\sum_i |f(P_i) - f(P)| K_i}{\sum_i K_i} I_{\{\sum_i K_i > 0\}} \right] \\
 & \leq \mathbb{E} \left[ \frac{\sum_i L D(P_i, P)^\beta K_i}{\sum_i K_i} I_{\{\sum_i K_i > 0\}} \right] \leq L (hR)^\beta,
 \end{aligned}$$

where in the last step we used the fact that

$$D(P_i, P)^\beta K_i = D(P_i, P)^\beta K \left( \frac{D(P_i, P)}{h} \right) \leq (hR)^\beta K_i,$$

since  $\text{supp}(K) \subseteq B(0, R)$ .

We now bound the second term.

$$\begin{aligned}
 & \mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i > 0\}} \right] \\
 & = \mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i \geq \underline{K}\}} + \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\underline{K} > \sum_i K_i > 0\}} \right] \\
 & \leq \mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i \geq \underline{K}\}} \right] + B_Y \mathbb{P}(\underline{K} > \sum_i K_i) \\
 & \leq \mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i \geq \underline{K}\}} \right] + \frac{B_Y}{em} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)}.
 \end{aligned}$$

(A4) implies that  $\mathbb{P}(|\mu_i| \leq B_Y) = 1$ , i.e.  $B_Y$  is a bound on the noise. The last step follows from Lemma 4. For the first term in the above expression, we use the following lemma. Its proof can be found in the supplementary material.

**Lemma 12**

$$\mathbb{E} \left[ \left| \frac{\sum_i \mu_i K_i}{\sum_i K_i} \right| I_{\{\sum_i K_i \geq \underline{K}\}} \right] \leq B_Y \sqrt{\frac{1 + 1/\underline{K}}{m\underline{K}}} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)}.$$

Finally, we bound the third term using Lemma 4:

$$f_{\max} \mathbb{P}(\sum_i K_i = 0) \leq \frac{f_{\max}}{em} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)}.$$

Putting everything together, we have

$$\mathbb{E} |\hat{f}(P; P_1, \dots, P_m) - f(P)|$$

$$\begin{aligned}
 &\leq L(hR)^\beta + B_Y \sqrt{\frac{1+1/K}{mK}} \int \frac{d\mathcal{P}(P)}{\Phi_P(Rh)} \\
 &\quad + \frac{B_Y}{em} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)} + \frac{f_{\max}}{em} \int \frac{d\mathcal{P}(P)}{\Phi_P(rh)} \\
 &\leq C_1 h^\beta + C_2 \sqrt{\frac{1}{m} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right]} + \frac{C_3}{m} \mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right].
 \end{aligned}$$

Note that  $\Phi_P(rh/2) \leq \Phi_P(rh) \leq \Phi_P(Rh)$ .

## 6 Doubling Dimension

The upper bound on the risk in Theorem 1 depends on the quantity  $\mathbb{E} \left[ \frac{1}{\Phi_P(rh/2)} \right]$ . In future work, we will show that, without further assumptions, this quantity can be quite large which leads to very slow rates of convergence. This is because the covering number of the class  $\mathcal{H}_k(1)$  is huge. For this paper, we concentrate on the more optimistic case where the support of  $\mathcal{P}$  has small effective dimension.

One way to measure effective dimension is to use the doubling dimension. Following Kpotufe [2011], we say that  $\mathcal{P}$  is a doubling measure with effective dimension  $d$  if, for every  $r > 0$  and  $0 < \epsilon < 1$ ,

$$\frac{\mathcal{P}(\mathcal{B}(s, r))}{\mathcal{P}(\mathcal{B}(s, \epsilon r))} \leq \left(\frac{c}{\epsilon}\right)^d, \quad \forall s. \quad (11)$$

If  $d$  denotes the doubling dimension of measure  $\mathcal{P}$ , then the  $\sqrt{\mathbb{E}[1/(m\Phi_P(rh/2))]}$  term in Theorem 1 can be upper bounded as follows:

$$\begin{aligned}
 \sqrt{\mathbb{E} \left[ \frac{1}{m\Phi_P(rh/2)} \right]} &= \sqrt{\mathbb{E} \left[ \frac{1}{m} \frac{\Phi_P(1)}{\Phi_P(rh/2)} \frac{1}{\Phi_P(1)} \right]} \\
 &\leq \sqrt{\frac{1}{m} C(rh/2)^{-d} \mathbb{E} \left[ \frac{1}{\Phi_P(1)} \right]} \leq \frac{C}{\sqrt{mh^d}}.
 \end{aligned}$$

Note also that when  $mh^d \geq 1$ , then  $\frac{1}{mh^d} \leq \frac{1}{\sqrt{mh^d}}$ . In this case, as a corollary of Theorem 1 and Assumptions (A5)-(A6), we now have that

$$R(m, n) \leq \frac{C_1}{h^{d+1}n^{1/(k+2)}} + C_2 h^\beta + C_3 \sqrt{\frac{1}{mh^d}}, \quad (12)$$

for appropriate constants  $C_1$ ,  $C_2$  and  $C_3$ .

To derive the rates for the risk, we consider two separate cases, depending on whether the third term in the right hand side of (12) dominates the first term or not.

Thus first assume that

$$\sqrt{\frac{1}{mh^d}} = \Omega \left( \frac{C_1}{h^{d+1}n^{1/(k+2)}} \right), \quad (13)$$

so that the risk becomes, asymptotically,  $O \left( h^\beta + \sqrt{\frac{1}{mh^d}} \right)$ . The optimal choice for  $h$  is then  $\Theta \left( m^{-1/(2\beta+d)} \right)$ , yielding a rate for the risk

$$R(m, n) = O \left( m^{-\beta/(2\beta+d)} \right).$$

Notice that this choice of  $h$  ensures that our assumption (A6) is met, since in this case (13) implies that

$$n = \Omega \left( m^{\frac{\beta+d+1}{2\beta+d}(k+2)} \right),$$

from which we obtain that

$$h = \Theta \left( m^{-\frac{1}{2\beta+d}} \right) = \Omega \left( n^{-\frac{1}{(k+2)(\beta+d+1)}} \right) = \Omega \left( n^{-\frac{1}{k+2}} \right).$$

This rate is reasonable because if the number of samples per distribution  $n$  is large compared to the number  $m$  of distributions, then the learning rate is limited by the number of distributions  $m$  and is in fact precisely the same as the rate of learning a standard  $\beta$ -Hölder smooth regression function in  $d$  dimensions. That is, the effect of not knowing the distributions  $P_1, \dots, P_m$  exactly and only having a finite sample from the distributions is negligible.

For the second case, suppose that

$$\sqrt{\frac{1}{mh^d}} = O \left( \frac{1}{h^{d+1}n^{1/(k+2)}} \right). \quad (14)$$

Then,  $R(m, n) = O \left( \frac{1}{h^{d+1}n^{1/(k+2)}} + h^\beta \right)$ , which implies that the optimal choice for  $h$  is  $h = \Theta \left( n^{-\frac{1}{(k+2)(\beta+d+1)}} \right)$ , giving the rate

$$R(m, n) = O \left( n^{-\frac{\beta}{(k+2)(\beta+d+1)}} \right).$$

Just like before, this choice of  $h$  does not violate assumption (A6) since

$$h = \Theta \left( n^{-\frac{1}{(k+2)(\beta+d+1)}} \right) = \Omega \left( n^{-\frac{1}{k+2}} \right).$$

Notice that, (14) also implies that

$$m = \Omega \left( n^{\frac{2\beta+d}{(k+2)(\beta+d+1)}} \right).$$

In this case, the rate is limited by the number of samples per distribution  $n$ , as expected. Notice that the rate gets worse as the dimensionality of each distribution  $k$  grows and as the smoothness  $\beta$  of the regression function deteriorates.

**Remark.** If there is no additive noise, i.e.  $\mu_i = 0$ , similar calculations yield that  $R(m, n) = O \left( m^{-\frac{1}{\beta+d}} \right)$  when  $n = \Omega \left( m^{\frac{\beta+d+1}{(\beta+d)(k+2)}} \right)$ , and  $R(m, n) = O \left( n^{-\frac{\beta}{(k+2)(\beta+d+1)}} \right)$  otherwise. While the rates seem reasonable, establishing optimality of the rates by demonstrating matching lower bounds is an open question that we plan to investigate in future work.

## 7 Numerical Illustrations

The following experiments serve as a proof of concept to demonstrate the applicability of the distribution regression estimator in Section 3. In these experiments, we used triangle kernels ( $k(x) = 1 - |x|$  if  $-1 \leq x \leq 1$ , and 0 otherwise). We set all the  $n, n_1, \dots, n_m$  set sizes and  $b, b_1, \dots, b_m$  bandwidths to the same values, which will be specified below. In the first experiment, we generated 325 sample sets from  $Beta(a, 3)$  distributions where  $a$  was varied between  $[3, 20]$  randomly. We constructed  $m = 250$  sample sets for training, 25 for validation, and 50 for testing. Each sample set contained  $n = 500$   $Beta(a, 3)$  distributed i.i.d. points. Our task in this experiment was to learn the skewness of  $Beta(a, b)$  distributions,  $f = \frac{2(b-a)\sqrt{a+b+1}}{(a+b+2)\sqrt{ab}}$ . We considered the noiseless case, i.e.  $\mu$  was set to zero. Our estimator of course is not aware of that the sample sets are coming from beta distributions, and it does not know the skewness function values in the test sets either; its values are available only in the training and validation sets.

To find appropriate bandwidths  $b$  and  $h$ , we sampled 100 i.i.d. randomly and uniformly distributed values in  $[0, 1]$ , evaluated the MSE performance of the distribution regression estimator on the validation test using these bandwidths parameters, and then chose the bandwidth parameters that lead to the best values on the validation test. To estimate the  $L_2$  distances between  $\hat{p}_i$  and  $p$ , we calculated their estimated values in 4096 points on a uniformly distributed grid between the min and max values in the sample sets, and then estimated the integral  $\int (p(x) - \hat{p}_i(x))^2 d(x)$  with the rectangle method for numerical integration. Figure 2(a) displays the predicted values for the 50 test sample sets, and we also show the true values of the skewness functions. As we can see the true and the estimated values are very close to each other.

In the next experiment, our task was to learn the entropy of Gaussian distributions. We chose a  $2 \times 2$  covariance matrix  $\Sigma = AA^T$ , where  $A \in \mathbb{R}^{2 \times 2}$ , and  $A_{ij}$  was randomly selected from the uniform distribution  $U[0, 1]$ . Just as in the previous experiments we constructed 325 sample sets from  $\{\mathcal{N}(0, R(\alpha_i)\Sigma^{1/2})\}_{i=1}^{325}$ . Where  $R(\alpha_i)$  is a 2d rotation matrix with rotation angle  $\alpha_i = i\pi/325$ . From each  $\mathcal{N}(0, R(\alpha_i)\Sigma^{1/2})$  distribution we sampled 500 2-dimensional i.i.d. points. Similarly to the previous experiment, 250 points was used for training, 25 for selecting appropriate bandwidth parameters, and 50 for training. Our goal was to learn the entropy of the first marginal distribution:  $f = \frac{1}{2} \ln(2\pi e\sigma^2)$ , where  $\sigma^2 = M_{1,1}$  and  $M = R(\alpha_i)\Sigma R^T(\alpha_i) \in \mathbb{R}^{2 \times 2}$ .  $\mu$  was zero in this experiment as well. Figure 2(b) displays the learned en-

tropies of the 50 test sample sets. The true and the estimated values are close to each other in this experiment as well.

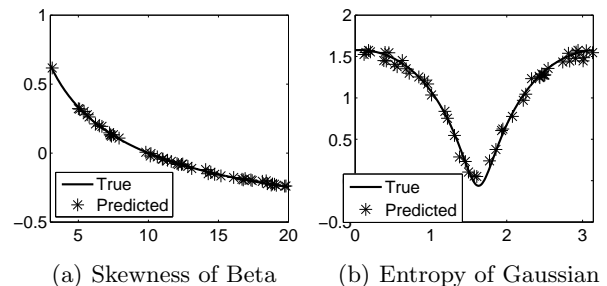


Figure 2: (a) Learned skewness of  $Beta(a, 3)$  distribution. Axis  $x$ : parameter  $a$  in  $[3, 20]$ . Axis  $y$ : skewness of  $Beta(a, 3)$ . (b) Learned entropy of a 1d marginal distribution of a rotated 2d Gaussian distribution. Axes  $x$ : rotation angle in  $[0, \pi]$ . Axis  $y$ : entropy.

## 8 Discussion and Conclusion

We have presented an estimator for distribution regression which is distribution-free in the sense that the estimator makes no strong distributional assumptions on the error variables. We derived upper bounds on the risk of the estimator and, in particular, we analyzed the case with a finite doubling dimension.

We note that our rates are faster than the logarithmic rates that are sometimes obtained in measurement error nonparametric regression models as in Fan and Truong [1993]. The reason is that the logarithmic rates occur when the measurement error is Gaussian. Our measurement error corresponds to  $\|\hat{p}_i - p_i\|$  which is not Gaussian for finite  $n_i$  and which decreases when  $n_i$  increases. In the standard measurement error model, the error is  $O(1)$  and is not decreasing.

In future work, we will prove lower bounds which show that, without further assumptions (such as assumptions about the doubling dimension), the rates can be very slow. We will also verify if the rates in the doubling dimension setting are tight or not. Also, we plan to investigate other estimators such as  $k$ -nn estimators and RKHS estimators.

### Acknowledgements

This research is supported in part by NSF under grants IIS-1247658, IIS-1250350, and DMS 1149677.

### References

R.J. Carroll, D. Ruppert, L.A. Stefanski, and C.M. Crainiceanu. *Measurement error in nonlinear mod-*



- els: a modern perspective*, volume 105. Chapman & Hall/CRC, 2006.
- A. Christmann and I. Steinwart. Universal kernels on non-standard input spaces. In *NIPS*, pages 406–414, 2010.
- L. Devroye and G. Lugosi. *Combinatorial methods in density estimation*. Springer, 2001.
- J. Fan and Y.K. Truong. Nonparametric regression with errors in variables. *The Annals of Statistics*, pages 1900–1925, 1993.
- F. Ferraty and P. Vieu. *Nonparametric Functional Data Analysis: Theory and Practice*. Springer Verlag, 2006.
- L. Györfi, M. Kohler, A. Krzyzak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer, New-york, 2002.
- T. Jaakkola and D. Haussler. Exploiting generative models in discriminative classifiers. In *NIPS*, pages 487–493. MIT Press, 1998.
- T. Jebara, R. Kondor, A. Howard, K. Bennett, and N. Cesa-bianchi. Probability product kernels. *JMLR*, 5:819–844, 2004.
- R. Kondor and T. Jebara. A kernel between sets of vectors. In *ICML*, 2003.
- S. Kpotufe. k-nn regression adapts to local intrinsic dimension. *arXiv preprint arXiv:1110.4300*, 2011.
- P. Moreno, P. Ho, and N. Vasconcelos. A Kullback-Leibler divergence based kernel for SVM classification in multimedia applications. In *NIPS*, 2004.
- K. Muandet, B. Schölkopf, K. Fukumizu, and F. D'Uzso. Learning from distributions via support measure machines. *arXiv.org*, stat.ML, February 2012.
- B. Póczos, L. Xiong, and J. Schneider. Nonparametric divergence estimation with applications to machine learning on distributions. In *UAI*, 2011.
- B. Póczos, L. Xiong, D. Sutherland, and J. Schneider. Nonparametric kernel estimators for image classification. In *Computer Vision and Pattern Recognition*, 2012.
- B. Póczos, A. Rinaldo, A. Singh, and L. Wasserman. Distribution-free distribution regression, 2013. <http://arxiv.org/abs/1302.0082>.
- J.O. Ramsay and B.W Silverman. *Functional data analysis*. Springer, New York, 2nd edition, 2005.
- P. Rigollet and R. Vert. Optimal rates for plug-in estimators of density level sets. *Bernoulli*, 15(4): 1154–1178, 2009.
- A. Smola, A. Gretton, L. Song, and B. Schölkopf. A Hilbert space embedding for distributions. In *ALT*, 2007.
- A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2010.