# Nonlinear Distribution Regression for Remote Sensing Applications

Jose E. Adsuara, Adrián Pérez-Suay, Jordi Muñoz-Marí, Anna Mateo-Sanchis,
Maria Piles, Gustau Camps-Valls

arXiv:2012.06377v1 [cs.LG] 7 Dec 2020

*Abstract*—**In many remote sensing applications one wants to estimate variables or parameters of interest from observations. When the target variable is available at a resolution that matches the remote sensing observations, standard algorithms such as neural networks, random forests or Gaussian processes are readily available to relate the two. However, we often encounter situations where the target variable is only available at the group level, i.e. collectively associated to a number of remotely sensed observations. This problem setting is known in statistics and machine learning as *multiple instance learning* or *distribution regression*. This paper introduces a nonlinear (kernel-based) method for distribution regression that solves the previous problems without making any assumption on the statistics of the grouped data. The presented formulation considers distribution embeddings in reproducing kernel Hilbert spaces, and performs standard least squares regression with the empirical means therein. A flexible version to deal with multisource data of different dimensionality and sample sizes is also presented and evaluated. It allows working with the native spatial resolution of each sensor, avoiding the need of match-up procedures. Noting the large computational cost of the approach, we introduce an efficient version via random Fourier features to cope with millions of points and groups. Real experiments involve SMAP Vegetation Optical Depth data for the estimation of crop production in the US Corn Belt, and MODIS and MISR reflectances for the estimation of Aerosol Optical Depth. An exhaustive empirical evaluation of the method is done against naive (linear and nonlinear) approaches based on input-space means, as well as previously presented methods for multiple-instance learning. We provide source code of our methods in http://isp.uv.es/code/dr.html.**

*Index Terms*—**Kernel methods, distribution regression, crop yield estimation, Aerosol Optical Depth (AOD), Moderate Resolution Imaging Spectro-Radiometer (MODIS), Soil Moisture Active Passive (SMAP), Vegetation Optical Depth (VOD).**

## I. INTRODUCTION

**E**STIMATING variables and bio-geophysical parameters of interest from observations is a central problem in Earth observation [1]–[3]. From a statistical standpoint, the problem reduces to regression and function approximation, for which either physical, statistical or hybrid inversion techniques can be used [3], [4]. In recent years, the amount of available data has allowed tackling complicated remote sensing regression problems learning. When input-output pairs are given, a plethora of algorithms can be readily used, such as neural networks [5], random forests [6] or kernel methods in general [7] and Gaussian processes in particular [8], [9], just to name a few. However, we often encounter problems where the target variable or parameter of interest is only available at the group level, not at the sample level. For example, very often one aims to estimate a bio-geophysical parameter, climate variable or ecological indicator from a set of observations, not just a single one, because of the spatial resolution or the acquisition level (e.g. [4]). Other common problems consider predicting a variable reported on field-based inventories or surveys at the county, region or state level from a set of observations (e.g. crop yield [10] or forest biomass [11]). The inverse situation, where an observation covers multiple samples, is also challenging. This is typically found in satellite retrieval validation, where coarse-scale estimates are compared to short-term intensive field campaign measurements, or to long-term -dense and sparse- station networks of point-scale ground based observations [12]–[14].

Several strategies exist to address this input-output mismatch problem by either 1) *output expansion*, that is replicating the group label for all observations in the group, or alternatively by 2) *input summary* of all group feature vectors. The first approach makes the problem even more ill-posed, as all inputs in a group are associated with the same target variable. We should note here that, while this is actually a valid approach for classification, in regression problems labels have a semantic ordered meaning so 'quantizing' the output space typically fails. When it comes to the second approach, it is customary to summarize all data in each group with its empirical mean or with a set of centroids (one per group) computed by clustering the group datasets. A more detailed description of these strategies are given below in Section II.

While the previous approaches are quite convenient and intuitive, they introduce strong assumptions about the data distributions. First, selecting a good summarizing criterion for all points in a group (e.g. all pixels within a county) is far from trivial. The empirical mean assumes that only the first moment captures all variability inside a group and it is enough to distinguish between group target variables (e.g. crop yield), while clustering assumes a particular data distribution (e.g. via using the Euclidean distance) is valid for all groups. Obviously such assumptions do not necessarily hold in most real case studies. Following with the example of crop yield prediction, adopting the empirical mean strategy could lead to the same (or deemed similar) input average feature vector (spectral signature) for counties having completely different types of crops and areas

Image Processing Laboratory (IPL)
Universitat de València, Catedrático A. Escardino - 46980 Paterna, València (Spain). E-mail: gustau.camps@uv.es

planted. Secondly, by summarizing each set with an average input feature vector one implicitly assumes that all bags or grouped-pixels have the same relevance, despite the number and variability of observations in each one. It may happen, for example, that a smaller county (hence lower number of input spectra available) yields a similar average statistic as a larger county. This leads to systematically biased model estimations, to eventually skewed conclusions about accuracy and fit, and to potentially wrong uncertainty propagation analysis.

The distinct goal of *distribution regression* (DR) is that of exploiting all input data available without explicitly summarizing the groups or replicating target variables. Hence, one is eventually interested in regressing the variable against a distribution of input feature vectors, not just one summarizing vector. Such a setup is ideally suited for the retrieval and validation of EO (Earth Observation) parameters characterized by high spatial variability that are measured in-situ at the point-scale. This is for instance the case of aerosols [15], precipitation [12], greenhouse gases [16], soil moisture [13], land surface temperature [17], ocean salinity [14], among others. DR can also prove very useful for applications where remotely sensed observations are used to predict variables reported in surveys or inventories at a regional scale, especially in regions that are particularly heterogeneous in size and/or composition. This includes for instance prediction of crop yield, carbon stocks, vector-borne diseases, or insect plagues (e.g. desert locust). Other applications of DR may include detection of anomalies and changes of interesting phenomena in datasets, as not single observations but distributions of observations are exploited. In general, distribution regression is an ideally suited framework to tackle problems that need predicting a scalar value from a distribution.

In recent years, the problem of mapping distributions to point estimates has received a broad interest in statistics and machine learning, known under different names such as *area-to-point kriging* (ATPK) [18], *multiple instance learning* (MIL) for classification [19]–[23] and regression [15], [24], [25], yet the field has been recently formalized under the field of *distribution regression* (DR) [26]–[30]. We place our proposal in this later field. This paper follows the principles in [29] and introduces a nonlinear kernel-based DR method for remote sensing applications. The formulation considers distribution embeddings in reproducing kernel Hilbert spaces, and performs regression with the empirical kernel means therein. By virtue of the kernel trick, the explicit mean map embedding is not necessary to solve the problem or apply the final model to new input test data sets [7], [31].

This paper presents a distribution regression framework specifically designed and adapted to the field of remote sensing. It introduces three main novelties. First, a version capable of dealing with multi-source data of different dimensionality and sample sizes is presented and evaluated. It allows working with the native spatial resolution of each sensor, something not allowed by any of the methods mentioned in the previous paragraph, avoiding the need of match-up procedures in which the information of the sensor with higher resolution is lost. Second, noting the large computational cost of the approach and the increasing amount of remote sensing data at high spatio-temporal resolutions available, we also introduce an efficient version via random Fourier features (RFF) that scales well to millions of points. Third, an extensive evaluation of the algorithm is provided to illustrate the applicability of the presented methodology using Soil Moisture Active Passive (SMAP) Vegetation Optical Depth for the prediction of crop yield over the US Corn Belt, and using Multi-angle Imaging Spectro-Radiometer (MISR) and Moderate Resolution Imaging Spectro-Radiometer (MODIS) reflectances from the TERRA satellite for the estimation of Aerosol Optical Depth. The method performs better than common approaches based on input-space means and clustering as well as previously presented methods for multiple-instance learning applied to similar datasets.

The rest of the paper is organized as follows. Section II introduces the basic elements for nonlinear (kernel-based) distribution regression: kernel function and the mean map embedding, as well as related methods such as the maximum mean discrepancy to motivate the DR proposal. Then, in Section III, we introduce our proposed kernel version for distribution regression, a multi-source version to deal with features of different dimensionality and sample sizes, and a fast version for computational efficiency based on random Fourier features. Experiments are detailed in Section IV. We conclude the paper in Section V with some remarks and outline of the future work and opportunities.

## II. KERNEL DISTRIBUTION REGRESSION

### A. Notation

Let us start by fixing the notation adopted in this paper. In distribution regression problems, we are given some sets of observations each of them with a corresponding output target variable to be estimated. Notationally, the training dataset $\mathcal{D}$ is formed by a collection of $B$ bags (or sets) $\mathcal{D} = \{(\mathbf{X}_b \in \mathbb{R}^{n_b \times d}, y_b \in \mathbb{R}) | b = 1, \ldots, B\}$. A training set from a particular group/bag $b$ is formed by $n_b$ examples, and is here denoted as $\mathbf{X}_b = [\mathbf{x}_1, \ldots, \mathbf{x}_{n_b}]^\top \in \mathbb{R}^{n_b \times d}$, where $\mathbf{x}_i \in \mathbb{R}^{d \times 1}$. For later convenience, let us denote all the available data collectively grouped in matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, where $n = \sum_{b=1}^{B} n_b$, and $\mathbf{y} = [y_1, \ldots, y_B]^\top \in \mathbb{R}^{B \times 1}$. This setting hampers the direct application of regression algorithms because not just a single input point $\mathbf{x}_b$ but a set of points $\mathbf{X}_b$ is available for model development, while latter for prediction we may have test points or sets from each bag denoted with a star superscript, $\mathbf{x}_b^* \in \mathbb{R}^{d \times 1}$ or $\mathbf{X}_b^* \in \mathbb{R}^{m_b \times d}$.

The problem in DR reduces to finding a function $f$ that learns the mapping from $\mathbf{x}$ to $y$, that is

$$y = f(\mathbf{x}) + e_i, \quad e_i \sim \mathcal{N}(0, \sigma_n^2).$$

Such a model, however, poses a challenging situation for standard regression since $\mathcal{D}$ contains many-to-one data. As discussed before, two main approaches are typically followed: 1) *output expansion*, that is replicating the label $y_b$ for all points in bag $b$; or 2) *input summary* most notably with the empirical average $\bar{\mathbf{x}}_b = \frac{1}{n_b} \sum_i \mathbf{x}_i$, or a set of centroids $\mathbf{c}_b$, $b = 1, \ldots, B$. The distinct goal of *distribution regression* is to exploit the rich structure in $\mathcal{D}$ by performing regression with

the group distributions directly. Statistically, this boils down to exploit all higher order statistical relationships between the groups, not just the first or second order moments. In this paper we present a method that embeds the bag distribution in a Hilbert space and performs linear regression therein. Some tools from the theory of reproducing kernels and functional analysis are needed, which are reviewed in what follows.

### B. Kernels functions and the mean map embedding

In this work we will rely on the theory of kernel methods to tackle the distribution regression problem. Let us summarize briefly the main needed concepts: kernel function, feature map, representer theorem, mean map embedding and the kernel ridge regression as the method used for estimation. For more comprehensive explanations the reader is addressed to the books [7], [31]–[33]. A detailed treatment of kernel mean embedding of distributions can be found in [34].

*1) Kernel functions and feature maps:* Kernel methods rely on the notion of similarity between samples in a higher (possibly infinite dimensional) Hilbert space. Be a set of empirical data $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \mathcal{X}$, where the points are defined in a $d$-dimensional input space, $\mathbf{x} = [x^1, \ldots, x^d]^\top \in \mathbb{R}^d$. Kernel methods assume the existence of a (dot product) Hilbert space $\mathcal{H}$, where samples are mapped into by means of a feature map $\phi : \mathcal{X} \to \mathcal{H}, \mathbf{x} \mapsto \phi(\mathbf{x})$. The mapping function can be defined explicitly (if some prior knowledge about the problem is available) or implicitly which is often the case in kernel methods. The similarity between the elements in $\mathcal{H}$ can now be measured using its associated dot product $\langle \cdot, \cdot \rangle_\mathcal{H}$. Here, we define a function that computes that similarity kernel, $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that $(\mathbf{x}, \mathbf{x}') \mapsto k(\mathbf{x}, \mathbf{x}')$. This function, often simply called kernel, is required to satisfy the Mercer's Theorem [31]:

$$k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_\mathcal{H}. \tag{1}$$

The mapping $\phi$ is its *feature map*, the space $\mathcal{H}$ is the reproducing Hilbert *feature space*, and $k$ is the reproducing kernel function since it *reproduces* dot products in $\mathcal{H}$ without even mapping the data explicitly therein. The point in kernel methods theory is that one relies on the implicit definition of the kernel function, and hence no feature mapping is explicitly defined. In all our experiments we used the radial basis function (RBF) kernel function, which actually accounts for all higher order monomials. Let us assume the kernel $k(x, x') = (\phi(x), \phi(x')) = \exp(-\gamma(x - x')^2)$, where for simplicity we define $\gamma > 0$. Then, the explicit feature map $\phi(x)$ is infinite dimensional, and can be expressed as

$$\phi(x) = \exp(-\gamma x^2) \left[ 1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \ldots \right]$$

and hence all higher-order relations between $x$ and $x'$ through the consideration of all monomials in the dot product reproduced by the kernel. Figure 1(d) shows a concrete illustrative example of detection of differences between distributions by using an explicit map $\phi(x) = x^2$. In kernel methods, however, the main advantage is that the mapping does not need to be explicitly designed.

*Definition* 1. *Reproducing kernel Hilbert spaces (rkHs) [35]. A Hilbert space $\mathcal{H}$ is said to be a rkHs if: (1) The elements of $\mathcal{H}$ are complex or real valued functions $f(\cdot)$ defined on any set of elements $\mathbf{x}$; And (2) for every element $\mathbf{x}$, $f(\cdot)$ is bounded.*

The name of these spaces come from the so-called *reproducing property*. Indeed, in an rkHs $\mathcal{H}$, there exists a function $k(\cdot, \cdot)$ such that

$$f(\mathbf{x}) = \langle f(\cdot), k(\cdot, \mathbf{x}) \rangle, \quad f \in \mathcal{H} \tag{2}$$

by virtue of the Riesz Representation Theorem [36]. A large class of algorithms has originated from regularization schemes in rkHs. The representer theorem gives us the general form of the solution to the common loss formed by a cost (loss, energy) term and a regularization term.

*Theorem 1. (Representer Theorem) [37] Let $\Omega : [0, \infty) \to \mathbb{R}$ be a strictly monotonic increasing function; let $V : \mathbb{R} \times \mathbb{R} \to \mathbb{R} \cup \{\infty\}$ be an arbitrary loss function; and let $\mathcal{H}$ be a rkHs with reproducing kernel $k$. Then:*

$$f^* = \min_{f \in \mathcal{H}} \left\{ V \left( (f(\mathbf{x}_1), y_1), \ldots, (f(\mathbf{x}_n), y_n) \right) + \Omega(\|f\|_\mathcal{H}^2) \right\} \tag{3}$$

*admits a space of functions (representation) $f$ defined as*

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i), \quad \alpha_i \in \mathbb{R}, \quad \boldsymbol{\alpha} \in \mathbb{R}^{n \times 1} \tag{4}$$

*This is, the function that minimizes the (regularized) optimization functional in Eq. (3) is a linear combination of dot products between data mapped into $\mathcal{H}$.*

*2) Mean map embeddings:* We frame the problem in the theory of mean map embeddings of distributions [7], [34], [38]. Let $\mathcal{B}_\mathcal{X}$ be the set of all probability distributions, then the kernel mean map $\boldsymbol{\mu}$ is defined as

$$\boldsymbol{\mu} : \mathcal{B}_\mathcal{X} \to \mathcal{H}, \quad \mathbb{P} \to \int_\mathcal{X} k(\cdot, \mathbf{x}) d\mathbb{P}(\mathbf{x}) \in \mathcal{H}.$$

Assuming that $k(\cdot, \mathbf{x})$ is bounded for any $\mathbf{x} \in \mathcal{X}$, we can show that for any $\mathbb{P}$, letting $\boldsymbol{\mu}_\mathbb{P} = \boldsymbol{\mu}(\mathbb{P})$, the $\mathbb{E}_P[f] = \langle \boldsymbol{\mu}_\mathbb{P}, f \rangle_\mathcal{H}$, for all $f \in \mathcal{H}$. It is important to note that, as for rkHs, here $\boldsymbol{\mu}$ represents the expectation function on $\mathcal{H}$. Following [34], every probability measure has a unique embedding and the $\boldsymbol{\mu}$ fully determines the corresponding probability measure.

Now the question is how to estimate such mean embeddings from empirical samples. Following our notation before for one particular bag, $\mathbf{X}_b$, drawn i.i.d. from a particular $\mathbb{P}_b$, the empirical mean estimator of $\boldsymbol{\mu}_b$ is given by:

$$\widehat{\boldsymbol{\mu}}_b = \boldsymbol{\mu}_{\mathbb{P}_b} = \int k(\cdot, \mathbf{x}) \widehat{\mathbb{P}}(d\mathbf{x}) \approx \frac{1}{n_b} \sum_{i=1}^{n_b} k(\cdot, \mathbf{x}_i). \tag{5}$$

This is an empirical mean map estimator whose dot product can be computed via kernels:

$$\langle \widehat{\boldsymbol{\mu}}_{\mathbb{P}_b}, \widehat{\boldsymbol{\mu}}_{\mathbb{P}_{b'}} \rangle_\mathcal{H} = \frac{1}{n_b n_{b'}} \sum_{i=1}^{n_b} \sum_{j=1}^{n_{b'}} k(\mathbf{x}_i^b, \mathbf{x}_j^{b'}). \tag{6}$$

Actually one can compute also distances among mean embeddings, which results in a particularly useful kernel

algorithm for hypothesis testing and domain adaptation named Maximum Mean Discrepancy (MMD) [7], [34], which has been previously used in remote sensing for feature selection [39], classification and domain adaptation [40], [41]. In both hypothesis testing and distribution regression we are ultimately concerned about comparing distributions $\mathbb{P}_b$ and $\mathbb{P}_{b'}$. MMD reduces to estimating the distance between the two sample means in a reproducing kernel Hilbert space $\mathcal{H}$ where data are embedded

$$\mathrm{MMD}(\mathbb{P}_b, \mathbb{P}_{b'}) := \|\boldsymbol{\mu}_{\mathbb{P}_b} - \boldsymbol{\mu}_{\mathbb{P}_{b'}}\|_{\mathcal{H}}^2,$$

which can be computed with kernels exploiting the dot product in Eq. (5). Interestingly, MMD tends asymptotically to zero when the two distributions $\mathbb{P}_x$ and $\mathbb{P}_y$ are the same, which allows us to assess differences in distributions of possibly different nature from empirical samples. Figure 1 illustrates the ability of MMD and the mean map embeddings to detect such differences hidden in higher order statistics, and motivates the use of kernel embeddings for distribution regression.

*3) Kernel ridge regression:* Finally, let us now formally define the distribution regression task. For this we perform standard least squares regression using the mean embedded data in Hilbert spaces. As we will see, the solution leads to that of the kernel ridge regression (KRR) algorithm [32] working with mean map embeddings. In our setting we want to minimize a classical regularized functional composed of two terms: the least square errors of the approximation (of the mean embedding) and a regularizer over the class of functions to be learned in Hilbert space $f \in \mathcal{H}$:

$$f^* = \arg\min_{f \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \|y_i - f(\boldsymbol{\mu}_i)\|^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\},$$

where $\lambda > 0$ is the regularization term. The ridge regression objective function has an analytical solution for a test set given a set of training examples:

$$\hat{f}_{\boldsymbol{\mu}_t} = \mathbf{k}(\mathbf{K} + n\lambda\mathbf{I})^{-1}\mathbf{y},$$

where $\boldsymbol{\mu}_t$ is the mean embedding of the test set $\mathbf{X}_t$, $\mathbf{k} = [k(\boldsymbol{\mu}_1, \boldsymbol{\mu}_t), \ldots, k(\boldsymbol{\mu}_n, \boldsymbol{\mu}_t)]^\top \in \mathbb{R}^{n \times 1}$, $\mathbf{K} = [k(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j)] \in \mathbb{R}^{n \times n}$ and $\mathbf{y} = [y_1, \ldots, y_n]^\top$ collectively gathers all outputs.

## III. PROPOSED KERNEL DISTRIBUTION REGRESSION

### A. Formulation

Let us first define a feature map that takes input samples into a Hilbert space, $\boldsymbol{\phi} : \mathbf{x} \in \mathbb{R}^{H \times 1} \to \boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^{H \times 1}$. All mapped data can be collectively grouped in a data matrix in Hilbert space denoted as $\mathcal{H}$: $\boldsymbol{\Phi}_b \in \mathbb{R}^{n_b \times H}$. The goal in kernel distribution regression is to perform a linear regression in $\mathcal{H}$. In order to do this we can summarize the bag feature vectors with the mean map embedding of samples in bag $b$, which is denoted here as $\boldsymbol{\mu}_b = \frac{1}{n_b} \sum_{i=1}^{n_b} \boldsymbol{\phi}(\mathbf{x}_i^b) \in \mathcal{H} \in \mathbb{R}^{H \times 1}$. Now, the collection of all mean embeddings of data (in the same $\mathcal{H}$) is expressed as $\mathbf{M} = [\boldsymbol{\mu}_1|\cdots|\boldsymbol{\mu}_B]^\top \in \mathbb{R}^{B \times H}$.

At this point we aim to perform regression on mean embeddings. We define a linear regression model

$$\hat{y}_b = \boldsymbol{\mu}_b^\top \mathbf{w}, \quad b = 1, \ldots, B,$$

where model weights $\mathbf{w} \in \mathbb{R}^{H \times 1}$ are shared across all bag models. For the sake of convenience, we can define the predictive model in matrix form as $\hat{\mathbf{y}} = \mathbf{M}\mathbf{w}$. The least squares solution corresponds to the normal (Wiener-Hopf) equation, $\mathbf{w} = \mathbf{M}^\dagger \mathbf{y} = (\mathbf{M}^\top \mathbf{M} + \lambda \mathbf{I})^{-1} \mathbf{M}^\top \mathbf{y}$. Nevertheless, noting the high dimensionalty of matrix $\mathbf{M}$, and hence the covariance matrix $\mathbf{M}^\top \mathbf{M}$, the problem cannot be explicitly solved. In order to address this, we follow the standard procedure in kernel methods by which we define a kernel function reproducing a dot product in $\mathcal{H}$, $k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}(\mathbf{x}_i)^\top \boldsymbol{\phi}(\mathbf{x}_j)$ and a representer theorem for model weights, that is $\mathbf{w} = \sum_{b=1}^{B} \alpha_b \boldsymbol{\mu}_b = \mathbf{M}^\top \boldsymbol{\alpha}$, where $\boldsymbol{\alpha} = [\alpha_1, \ldots, \alpha_B] \in \mathbb{R}^{B \times 1}$. Now, the dual solution is given by $\boldsymbol{\alpha} = (\widehat{\mathbf{K}} + \lambda \mathbf{I})^{-1} \mathbf{y}$ with $\widehat{\mathbf{K}} = \mathbf{M}\mathbf{M}^\top \in \mathbb{R}^{B \times B}$, which can be readily recognized as the kernel ridge regression working on a kernel with entries:

$$[\widetilde{\mathbf{K}}]_{b,b'} = \boldsymbol{\mu}_b^\top \boldsymbol{\mu}_{b'} = \frac{1}{n_b n_{b'}} \sum_{i=1}^{n_b} \sum_{j=1}^{n_{b'}} \boldsymbol{\phi}(\mathbf{x}_i^b)^\top \boldsymbol{\phi}(\mathbf{x}_j^{b'})$$

$$= \frac{1}{n_b n_{b'}} \sum_{i=1}^{n_b} \sum_{j=1}^{n_{b'}} k(\mathbf{x}_i^b, \mathbf{x}_j^{b'}) = \frac{1}{n_b n_{b'}} \mathbf{1}_{n_b}^\top \mathbf{K}_{bb'} \mathbf{1}_{n_{b'}},$$

where the matrix $\mathbf{K}_{bb'} \in \mathbb{R}^{n_b \times n_{b'}}$. Therefore, we have an analytic solution of the problem. Now the question is how one can generate predictions for particular samples. Let us define a test sample matrix from a particular bag denoted as $\mathbf{X}_b^* = [\mathbf{x}_1, \ldots, \mathbf{x}_{m_b}]^\top \in \mathbb{R}^{m_b \times d}$, which has a mean embedding $\boldsymbol{\mu}_b^* = \frac{1}{m_b} \sum_{l=1}^{m_b} \boldsymbol{\phi}(\mathbf{x}_l^b) \in \mathcal{H} \in \mathbb{R}^{H \times 1}$. The predictions for the test data can be explicitly derived as follows:

$$\hat{y}_b^* = \boldsymbol{\mu}_b^{*\top} \mathbf{w} = \frac{1}{m_b} \sum_{l=1}^{m_b} \boldsymbol{\phi}^\top(\mathbf{x}_l^b) \mathbf{M}^\top \boldsymbol{\alpha} \tag{7}$$

$$= \frac{1}{m_b} \sum_{l=1}^{m_b} \boldsymbol{\phi}^\top(\mathbf{x}_l^b) [\boldsymbol{\mu}_1|\cdots|\boldsymbol{\mu}_B]\boldsymbol{\alpha}$$

$$= \frac{1}{m_b} \sum_{l=1}^{m_b} \boldsymbol{\phi}^\top(\mathbf{x}_l^b) \left[ \frac{1}{n_1} \sum_{i=1}^{n_1} \boldsymbol{\phi}(\mathbf{x}_i^1) \bigg| \cdots \bigg| \frac{1}{n_B} \sum_{i=1}^{n_B} \boldsymbol{\phi}(\mathbf{x}_i^B) \right] \boldsymbol{\alpha}$$

$$= \left[ \frac{1}{m_b n_1} \sum_{l=1}^{m_b} \sum_{i=1}^{n_1} k(\mathbf{x}_l^b, \mathbf{x}_i^1), \cdots, \frac{1}{m_b n_B} \sum_{l=1}^{m_b} \sum_{i=1}^{n_B} k(\mathbf{x}_l^b, \mathbf{x}_i^B) \right] \boldsymbol{\alpha}$$

$$= \frac{1}{m_b} \sum_{b'=1}^{B} \frac{1}{n_{b'}} \alpha_{b'} \sum_{l=1}^{m_b} \sum_{i=1}^{n_b} k(\mathbf{x}_l^b, \mathbf{x}_i^{b'}) = \frac{1}{m_b n} \mathbf{1}_{m_b}^\top \mathbf{K}_{bb'} \mathbf{1}_{n_{b'}} \boldsymbol{\alpha},$$

where $\mathbf{K}_{bb'} \in \mathbb{R}^{m_b \times n_{b'}}$ which is computed easily given a valid (Mercer) kernel function $k$.

### B. Multisource distribution regression

Assume that now we aim to combine/exploit multi-modal (multi-source) information defining each bag. The sets may have different numbers of both features and sizes, e.g. we aim to combine different spatial, spectral or temporal resolutions. We will illustrate this case in the experimental section by combining at a bag level two distinct sensors for aerosol parameter estimation.

Notationally, now we have access to different matrices $\mathbf{X}_f^b \in \mathbb{R}^{n_b^f \times f}$, $f = 1, \ldots, F$. We propose a multimodal kernel distribution method by embedding each dataset into a mean and exploiting the direct sum of Hilbert spaces in the mean embedding space. To do this, we define $F$ Hilbert
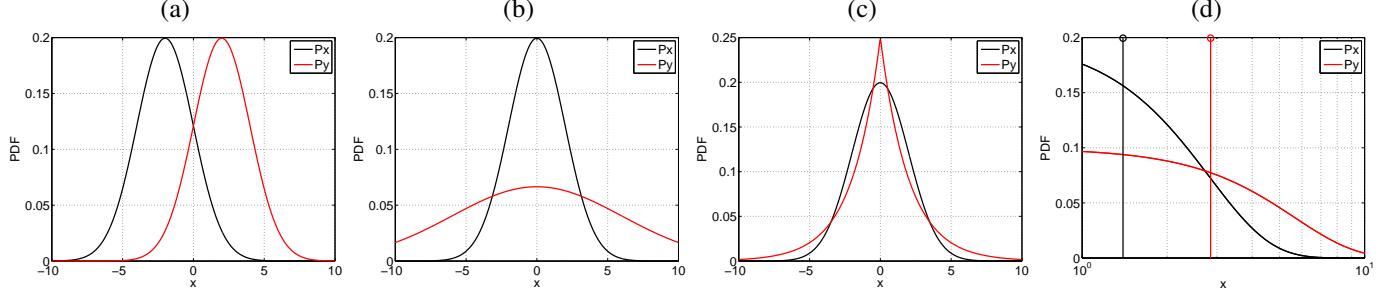
Fig. 1: The two-sample problem reduces to detecting whether two distributions $\mathbb{P}_x$ and $\mathbb{P}_y$ are different or not. Summarizing group distributions for regression using a statistic is common practice, but it may fail because completely different groups can be indistinguishable and make regression strongly ill-posed. (a) Whenever we have two Gaussians with different means one can assess statistical differences by means of a standard $t$-test of different means (the distance between the empirical means is $d_{xy} := \|\mu_x - \mu_y\| = 3.91$) and summarizing the distributions with the empirical mean is good enough; (b) when we have two Gaussians with the same mean but different variance using only the first moment is useless ($d_{xy} = 0.02$), yet one could resort to the second moment -the variance- which is what the $z$-test does; (c) however, we may find situations with the same first and second order yet different in nature (in this case, Gaussian and Laplace distributions with the same mean and variance), which hampers again discrimination and the definition of the summary statistic for regression; (d) the latter can be easily addressed by mapping the random variables to higher order features for discrimination (e.g. for the Gaussian case, second order features of the form $x^2$ suffice for discrimination, $d_{xy} = 13.26$). This simple example motivates the use of kernel mean embeddings for distribution regression and hypothesis testing that are able to estimate all higher order moments without even mapping the data explicitly but resorting to kernel functions only.

spaces $\mathcal{H}_f$, $f = 1, \ldots, F$, and the direct sum of all of them, $\mathcal{H} = \bigoplus_{f=1}^{F} \mathcal{H}_f$. Now the key for multimodal distribution regression is to perform a linear regression in $\mathcal{H}$. We can summarize the bag feature vectors with a set of mean map embeddings of samples in bag $b$, which is denoted here as $\boldsymbol{\mu}_b^f = \frac{1}{n_b^f} \sum_{i=1}^{n_b^f} \boldsymbol{\phi}(\mathbf{x}_i^{b,f}) \in \mathcal{H} \in \mathbb{R}^{H_f \times 1}$. The collection of all mean embeddings in the same $\mathcal{H}$ is defined as

$$\boldsymbol{\mu}_b = [\boldsymbol{\mu}_b^1, \ldots, \boldsymbol{\mu}_b^F] \in \mathcal{H},$$

and then we can use the same formula as before to define the mean map embedding, $\mathbf{M} = [\boldsymbol{\mu}_1 | \cdots | \boldsymbol{\mu}_B]^\top \in \mathbb{R}^{B \times H}$. We now need to compute the multimodal kernel matrix as follows:

$$[\widetilde{\mathbf{K}}]_{b,b'} = \boldsymbol{\mu}_b^\top \boldsymbol{\mu}_{b'} = \sum_{f=1}^{F} \frac{1}{n_b^f n_{b'}^f} \mathbf{1}_{n_b}^{f\top} \mathbf{K}_{bb'} \mathbf{1}_{n_{b'}^f},$$

where the matrix $\mathbf{K}_{bb'} \in \mathbb{R}^{n_b \times n_{b'}}$ is again the same size. Figure 2 illustrates the two DR approaches treated in this paper.

### C. Randomized distribution regression for scalability

The main computational bottleneck of the distribution regression methods presented is the inversion of the $\widetilde{\mathbf{K}}$ matrix of size $B \times B$. This is a cheap operation for standard problems that involve less than a few thousand groups. However, note that each entry $(b, b')$ in matrix $\widetilde{\mathbf{K}}$ involves computing and averaging many kernel matrices and thus each entry scales in time $\mathcal{O}(n_b n_b' d)$ which reduces to $\mathcal{O}(n^2 d)$ if we assume that all bags have the same number of elements $n$ (i.e. $n_b = n_b' = n$). This leads to a total cost of $\mathcal{O}(B^2 n^2 d + B^3)$ where the cubic order in the number of bags $B$ is due to the matrix inversion. This is simply not affordable even for moderate-size problems with a few thousand points per bag. In this section

we introduce a kernel approximation with random Fourier features [42] that alleviates the problem.

An outstanding result in the machine learning literature makes use of a classical definition in harmonic analysis to improve approximation and scalability of kernel methods [42]. The Bochner's theorem states that a continuous kernel $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x} - \mathbf{x}')$ on $\mathbb{R}^d$ is positive definite (p.d.) if and only if $k$ is the Fourier transform of a non-negative measure. If a shift-invariant kernel $k$ is properly scaled, its Fourier transform $p(\mathbf{w})$ is a proper probability distribution. This property is used to approximate kernel functions and matrices with linear projections on a number of $D$ random features, as follows:

$$k(\mathbf{x}, \mathbf{x}') = \int_{\mathbb{R}^d} p(\mathbf{w}) e^{-\mathrm{i}\mathbf{w}^\top (\mathbf{x} - \mathbf{x}')} \mathrm{d}\mathbf{w} \approx \sum_{i=1}^{D} \frac{1}{D} e^{-\mathrm{i}\mathbf{w}_i^\top \mathbf{x}} e^{\mathrm{i}\mathbf{w}_i^\top \mathbf{x}'}$$

where $p(\mathbf{w})$ is set to be the inverse Fourier transform of $k$, $\mathrm{i} = \sqrt{-1}$, and $\mathbf{w}_i \in \mathbb{R}^d$ is randomly sampled from a data-independent distribution $p(\mathbf{w})$ [43]. Note that we can define a $D$-dimensional *randomized* feature map $\mathbf{z}(\mathbf{x}) : \mathbb{R}^d \to \mathbb{C}^D$, which can be explicitly constructed as $\mathbf{z}(\mathbf{x}) := [\exp(\mathrm{i}\mathbf{w}_1^\top \mathbf{x}), \ldots, \exp(\mathrm{i}\mathbf{w}_D^\top \mathbf{x})]^\top$. Other definitions are possible: one could for instance expand the exponentials in pairs $[\cos(\mathbf{w}_i^\top \mathbf{x}), \sin(\mathbf{w}_i^\top \mathbf{x})]$, this increases the mapped data dimensionality to $\mathbb{R}^{2D}$, while approximating exponentials by $[\cos(\mathbf{w}_i^\top \mathbf{x} + b_i)]$, where $b_i \sim \mathcal{U}(0, 2\pi)$, is more efficient (still mapping to $\mathbb{R}^D$) but has proved less accurate [44]. In our experiments we used projections onto the $[\sin, \cos]$ pairs to keep operations in the real domain.

In matrix notation, given $n$ data points, the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ can be approximated with the explicitly mapped data, $\mathbf{Z} = [\mathbf{z}_1 \cdots \mathbf{z}_n]^\top \in \mathbb{R}^{n \times D}$, and will be denoted as $\hat{\mathbf{K}} \approx \mathbf{Z}\mathbf{Z}^\top$. This property can be used to approximate any shift-
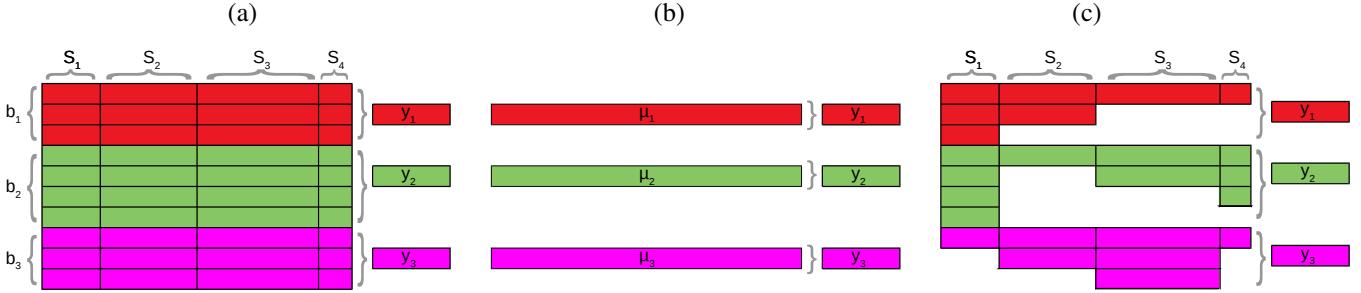
Fig. 2: Distribution regression approaches exemplified. The general DR problem setting is illustrated in (a), cf. §III-A. Consider $B = 3$ bags (different colors: red, green and purple) with different number of samples per bag ($n_1 = 3$, $n_2 = 4$, $n_3 = 3$) and three corresponding target labels, $y_b$, $b = 1, 2, 3$. Columns are different sources of information, in this case there are four sensors acquiring data $S_i$, $i = 1, 2, 3, 4$. Note that the input feature space for each sensor $S_i$ does not necessarily have the same dimensionality (e.g. spectral or temporal resolution). Standard practice summarizes the distributions $\mathbb{P}_b$ with the mean vector $\mu_b$, $b = 1, 2, 3$, where $\mu_b = [\mu_b^{S_1}, \mu_b^{S_2}, \mu_b^{S_3}, \mu_b^{S_4}]$ (b), and then proceed with standard regression methods. This can be done in Hilbert spaces too but with the advantage of considering all moments of the distributions, not just the first one, and accounting for the relations among bag samples (cf. §II-B). In (c) we show the case of multi-source distribution regression (MDR) in which some features are missing for particular bags and samples, which is often encountered when different sensors are combined. This case is addressed in this paper too (cf. §III-B).

invariant kernel. For instance, the RBF kernel, which is the one used in our experiments, can be approximated using $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \sigma^{-2}\mathbf{I})$, $1 \le i \le D$. It is also important to notice that the approximation of $k$ with random Fourier features converges in $\ell_2$-norm error with $\mathcal{O}(D^{-1/2})$ when using an appropriate random parameter sampling distribution [45].

For the case of DR, in principle one should sample $B$ times (one per bag), hence obtaining $B$ sets of vectors $\mathbf{w}^b$ and the associated explicit maps $\mathbf{Z}_b$, $b = 1, \dots, B$. Such approach would be interesting if data from different bags have different dimensions. In practice, however, we will use the same random bases for all bags, thus we only sample once and obtain a single matrix $\mathbf{W} \in \mathbb{R}^{d \times D}$. The randomized DR (RDR) readily reduces to solve the least squares regularized regression problem with the (explicit) means of projected bag data:

$$\mathbf{w} = (\mathbf{Z}^\top \mathbf{Z} + \lambda \mathbf{I}_{2D})^{-1} \mathbf{Z}^\top \mathbf{y},$$

where $\mathbf{w} \in \mathbb{R}^{2D \times 1}$ and $\mathbf{Z} \in \mathbb{R}^{B \times 2D}$ contains the explicit mean over random Fourier feature projections per bag, that is, each row of $\mathbf{Z}$ contains $\mathbf{z}_b = \frac{1}{n_b} \sum_{i=1}^{n_b} \mathbf{z}(\mathbf{x}_i^b)$. Given a test bag dataset $\mathbf{X}_* \in \mathbb{R}^{t \times d}$, one has to explicitly map it onto the same bases, $\mathbf{Z}_* = [\mathbf{z}(\mathbf{x}_1)|\cdots|\mathbf{z}(\mathbf{x}_t)]^\top \in \mathbb{R}^{t \times 2D}$, compute the explicit mean map, $\mathbf{m}_* = \frac{1}{t} \sum_{j=1}^{t} \mathbf{z}(\mathbf{x}_j) \in \mathbb{R}^{2D \times 1}$, and apply the linear prediction model on the explicit mean:

$$\hat{y}_* = \mathbf{m}_* \mathbf{w} \in \mathbb{R}.$$

By considering a number of samples $n \gg D$ against a moderate number of bags $B$, the associated cost by using the random features approximation now reduces to $\mathcal{O}(B^2 D^2 n)$. This is much more convenient than the kernel version that scales quadratically in $n$ as $\mathcal{O}(B^2 n^2 d)$. This situation where the number of samples is dominating the problem is the typicall scenario in real domains.

Finally, and from the point of view of the implementation, despite having guaranteed asymptotic convergence to the RBF

kernel due to the Bochner's theorem, the selection of the parameter $D$ is crucial and may act a regularizer. We cross-validated the value of $D$ to guarantee that we really captured the intrinsic dimensionality of the problem: in some cases, small $D$ values (hence not approximating the RBF kernels) led to improved performance. This effect has been previously noted in [46].

## IV. EXPERIMENTAL RESULTS

This section shows results for several methods and in different problems of distribution regression found in remote sensing. As baseline standard approaches, we consider a least squares linear regression model (LR) and the nonlinear (kernel) counterpart, the kernel ridge regression (KR) method, both working on the means of each bag as input feature vectors. These methods will be compared with our proposals: KDR (kernel distribution regression), its randomized approximation (RDR), and in the last experiment we will add the multi-source DR (MDR) along with stacked-feature approaches. We provide empirical comparisons with standard scores like the root-mean-square-error (RMSE), coefficient of determination ($R^2$), in cross-validation and test sets to assess accuracy and robustness. We provide source code of our methods in http://isp.uv.es/code/dr.html.

### A. Experimental setup and model development

The proposed methods are evaluated in three different scenarios. Firstly, SMAP Vegetation Optical Depth (VOD) data is related to crop production data from the 2015 US agricultural survey (total yield and yield per crop type). Secondly, MISR and MODIS reflectances are used to estimate Aerosol Optical Depth (AOD) (data sets in [25]). The AOD estimation is solved using each sensor data independently, or using multi-sensor combined information through the proposed MDR formulation and a stacked-feature approach.

Evaluation of the algorithms is done as follows. We reserve a percentage of the elements of each bag for test. With the remaining elements, we perform a $k$-fold cross-validation with $k = 5$ also at a bag level. This is done by randomly splitting the data into five subsets: one subset is reserved for validation and the others for training of a predictor. After this procedure, we apply the best model found to the test data. Finally, all this process is repeated ten times, and the average over all test results is computed. Both cross-validation and test errors are reported. We use as evaluation criteria the standard mean error (ME) to account for bias, the root-mean-square-error (RMSE) to assess accuracy, and the coefficient of determination or explained variance ($R^2$) to account for the goodness-of-fit.

### B. Estimation of crop yield

In this experiment, we use satellite-based retrievals of vegetation optical depth (VOD) from SMAP [47] for crop-yield prediction. The vegetation optical depth (VOD) is a measure of the attenuation of soil microwave emissions when they pass through the vegetation canopy; it is sensitive to the amount of living biomass as well as to the amount of water stress experienced by the vegetation [48]. SMAP VOD has been shown to carry information about crop growth and yield in a variety of agro-ecosystems [49], [50].

The data set used in this study was originally presented in [50]. The study area are the extensive croplands of North Dakota, South Dakota, Nebraska, Minnesota, Iowa, Illinois, Indiana and Ohio, within the US Corn Belt (see Fig 3(a)). County-based information on area planted and yield per crop type was obtained from the US Department of Agriculture (USDA-NASS). Yields for a variety of crops are reported, namely corn, soybeans, wheat, oats, beans, barley, peas, canola, flaxseed, sorghum and lentils. A single yield datum was obtained for each county as a weighted average of the reported yield and area planted for each crop, after converting the units of each crop to $kg/m^2$ (see Table S1 in [50]). The study period is one year, starting in April 1, 2015. This corresponds to the first year of SMAP data and the 2015 crop season in the US corn belt. Only SMAP observations over agricultural pixels are considered, following the MODIS IGBP land cover classification. In [50], the yield-VOD relationship was explored using principal components regression with VOD seasonal metrics, and the first principal component allowed explaining 66% of variability. In this experiment, we apply the DR method to explore the yield-VOD relationship at the county scale. Each county is treated as a *bag* for the DR method, containing a number of remotely sensed observations (16 SMAP pixels on average). No information on crop season is used, and the whole VOD time series are used as input. There is a total of 385 counties with yield and satellite data. A 66% of these counties (bags) are used to train/validate and the remaining 33% are used to test.

We show two approaches to the problem: prediction of total yield, and prediction of yield per crop type. For the latter, the three main crops in the region are predicted: corn, soybean and wheat. All the 363, 361 and 204 counties reporting corn, soybean and wheat yield, independently of their relative impor-

TABLE I: Results for Crop-Yield Estimation Using VOD.

| Algorithms | Total crop yield | | |
|---|---|---|---|
| | ME×1000 | RMSE ×100 | $R^2$ |
| LR | 1.19 ± 7.36 | 9.67 ± 0.74 | 0.80 ± 0.02 |
| KR | 2.22 ± 10.77 | 9.34 ± 0.73 | 0.81 ± 0.02 |
| RDR | 2.27 ± 10.98 | 9.37 ± 0.71 | 0.81 ± 0.02 |
| KDR | 2.27 ± 10.95 | 9.35 ± 0.71 | 0.81 ± 0.02 |
| | **Corn yield** | | |
| | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -1.20 ± 5.89 | 7.54 ± 0.50 | 0.85 ± 0.02 |
| KR | 1.68 ± 8.52 | 6.54 ± 0.72 | 0.88 ± 0.02 |
| RDR | 1.52 ± 7.90 | 6.57 ± 0.70 | 0.88 ± 0.02 |
| KDR | 1.59 ± 7.88 | 6.47 ± 0.74 | 0.89 ± 0.02 |
| | **Soybean yield** | | |
| | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -1.99 ± 1.85 | 2.45 ± 0.13 | 0.85 ± 0.03 |
| KR | -0.70 ± 2.92 | 2.47 ± 0.21 | 0.85 ± 0.04 |
| RDR | -0.90 ± 2.58 | 2.44 ± 0.23 | 0.85 ± 0.04 |
| KDR | -0.64 ± 2.43 | 2.40 ± 0.21 | 0.86 ± 0.03 |
| | **Wheat yield** | | |
| | ME×1000 | RMSE ×100 | $R^2$ |
| LR | 2.72 ± 6.65 | 5.46 ± 0.48 | 0.64 ± 0.08 |
| KR | 2.42 ± 8.47 | 5.07 ± 0.38 | 0.69 ± 0.05 |
| RDR | 3.31 ± 7.64 | 5.15 ± 0.43 | 0.68 ± 0.05 |
| KDR | 2.91 ± 7.31 | 5.10 ± 0.40 | 0.69 ± 0.05 |

tance at the county level, were included in the corresponding crop-specific experiments.

Table I shows the crop-yield predictions for the baseline and proposed regression approaches. Notably, these results outperform those obtained in previous literature for corn-soy croplands (see [50] and references therein), even with the simple models (LR or KRR). This can possibly be due to the fact that we are using the whole growing-cycle information contained in the time series to develop the models, and should be a matter of future dedicated studies. Relevant to this work, consistently best results are obtained with the DR models (either RDR or KDR) in all experiments: they are able to explain 81% of total yield, 89% of corn yield, 86% of soybean yield and up to 69% of wheat yield. Everything worked as expected: the higher the capacity of the model, the better results we obtained (understanding this capacity as the treatment of the non-linearity and the grouping by bags). In particular, the advantage of non-linear approaches is more evident in the more complex prediction of independent crop yields in mixed-crops, than in the prediction of total yield, where the contributions of all crops are accounted for.

Results of the best regression model between VOD and official corn yields at county level are illustrated in Fig. 3. Except in some specific counties, it can be seen that the corn predictions are reasonably good, with relative errors below 3%. We anticipate that the proposed DR approaches will be particularly useful for regional crop forecasting in areas covering different agro-climatic conditions and fragmented agricultural landscapes (e.g. Europe), where scale effects need to be properly addressed for adequate analysis and predictions [10].

(a) Study area



(b) Official corn yield $(kg/m^2)$



(c) KDR relative error (%)
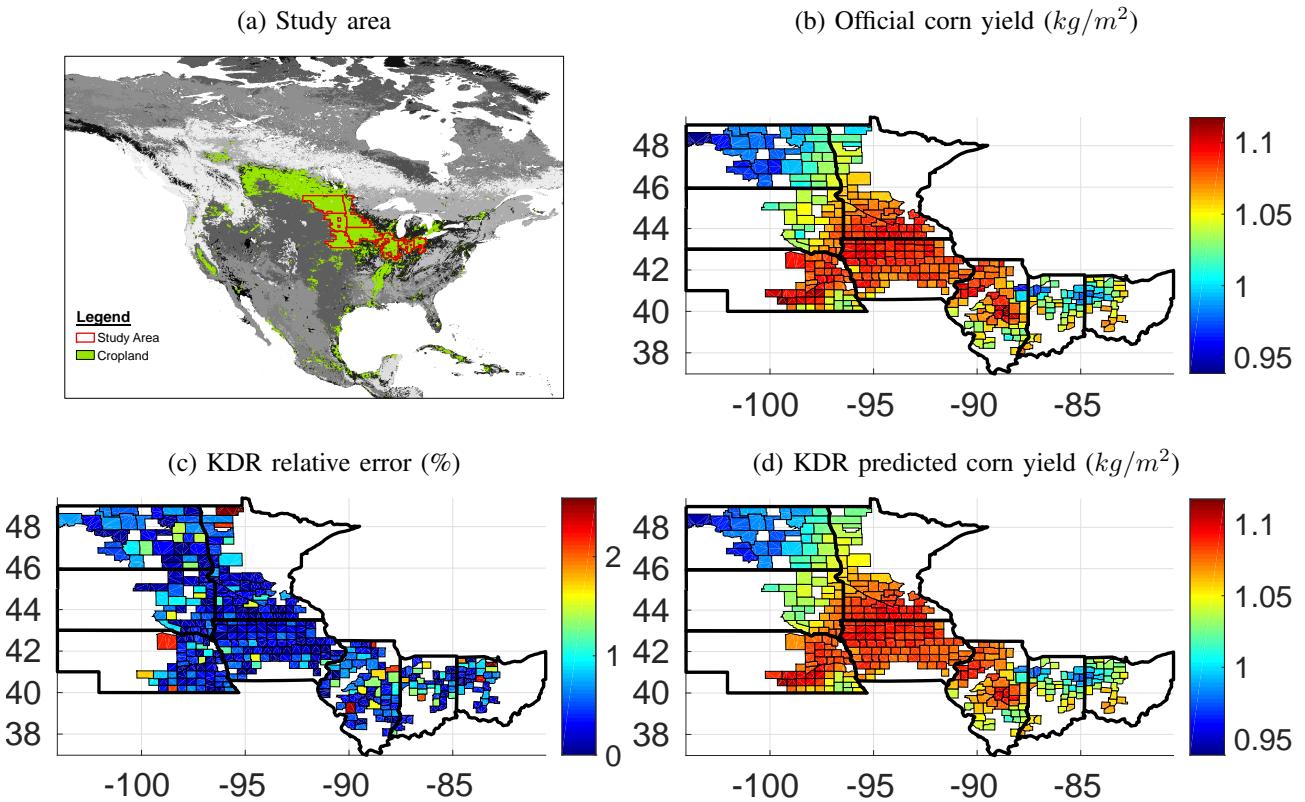


(d) KDR predicted corn yield $(kg/m^2)$



Fig. 3: (a) Study Area including the 8 states and cropland mask following the MODIS IGBP land cover classification. (b) Map of corn yield for year 2015 from USDA-NASS survey $(kg/m^2)$. (c) KDR relative error prediction per county (%). (d) KDR predicted corn yield per county $(kg/m^2)$

## C. Experiment 2: Multisensor aerosol optical depth estimation

This second application consists on the prediction of the Aerosol Optical Depth (AOD) from remotely-sensed data. Aerosol data consist of *in situ* data obtained from the Aerosol Robotic Network [51] (AERONET), and remote sensing data are in this case measurements from the Multi-angle Imaging SpectroRadiometer (MISR), and from the Moderate Resolution Imaging Spectroradiometer (MODIS) sensors on-board the TERRA satellite. Aerosol estimation is one of the biggest challenges of current climate research. Aerosols both reflect and absorb incoming solar radiation, and their presence directly affects the Earth's radiation budget. Three different datasets of study are used here and originally presented in [25]:

- The AOD-MISR1 data set is a collection of 800 bags collected at 35 AERONET ground sites within the continental U.S. between 2001 and 2004 from the MISR sensor. Each bag consists of 100 instances, representing randomly selected pixels within 20-km radius around a given AERONET site, which leads to a total number of 76040 instance samples. The instance attributes are 12 reflectances from the three middle MISR cameras as well as four solar and view zenith angles. The bag target value is the AOD measured by the AERONET instrument within 30 min of the satellite overpass.
- The AOD-MISR2 data set has the same properties as AOD-MISR1. The only difference is that the 100 in-

stances in each bag are sampled only from the cloud-free pixels. It consists on a total of 800 bags and 76040 instance samples. Since cloudy pixels are known to be noisy and lead to reduced retrieval quality, this data set is expected to provide better predictions.

- The AOD-MODIS data set is a collection of 1364 bags collected at 45 AERONET sites within the continental U.S. between 2002 and 2004 from the MODIS sensor. Each bag consists of 100 instances. The instance attributes are seven MODIS reflectance bands and five solar and view zenith angles, and the bag label was the corresponding AERONET AOD measurement.

TABLE II: Cross-Validation Results on the Three Remote Sensing Data Sets (See Text for Details) for AOD Prediction.

| Algorithms | MISR1 | MISR2 | MODIS |
|---|---|---|---|
| | RMSE ×100 | RMSE ×100 | RMSE ×100 |
| LR | 21.59 ± 0.48 | 21.01 ± 0.38 | 21.65 ± 0.37 |
| KR | 9.07 ± 0.46 | 7.70 ± 0.76 | 11.64 ± 0.63 |
| RDR | 7.90 ± 0.50 | 7.21 ± 0.65 | 9.69 ± 0.39 |
| Best in [25] | 7.5 ± 0.1 | 7.3 ± 0.1 | 9.5 ± 0.1 |
| Best in [29] | 7.90 ± 1.63 | | |

The way to proceed in applications of this type consists of predicting the current year AOD from the previous ones. Although the current data are not labeled per year, they correspond to four years for MISR1 and MISR2 (three years

for MODIS). We reserved the equivalent of one year of data, 25% of them (33% for MODIS) for testing, and the remaining 75% (66% for MODIS) of data for train/validation. It is worth to note here that these sets will be different each of the 10 times we repeat the experiment. Reported results are the average of these 10 trials. We normalize the data.

We show two approaches to the problem. First, we will apply the proposed DR methods using information from each sensor separately (unisensor approach), and compare results with the ones reported in [25], as well as to baseline regression approaches. Second, we will use the multisource DR for combining information from different sensors, specifically MISR2 and MODIS (multi-sensor approach). As this is only possible if we have land information for both sensors, the first thing we do is to build a new dataset retaining only those bags having a common land measure. Originally, we had 800 bags for MISR2 and 1364 for MODIS. After the combination, we retain only 289 bags. At this point, we want to emphasize that the combination of sensors with multiple spatial resolutions is not easy in a normal setting, and quite challenging in DR settings too. The naive approach is to adjust the two data meshes and concatenate them (feature-stacking approach), so one can use any standard regression model on this feature vector. We can do this straightforwardly because each feature is normalized. While this is practical, we are losing the information of the sensor with higher resolution. The alternative approach proposed here is to use the direct sum of kernels in our MDR model to be able to treat each input independently even if they have different resolutions and number of data points (composite approach). We will also apply our methods to the sensors separately, using only the common 289 bags, to have the single-source performance when using this particular subset of bags. We show the results for the two approaches in what follows.

*1) Single source approach:* Tables II and III show the results. Table II shows results in the training/validation set obtained by cross-validation, whereas Table III shows results in the unseen test set. For cross-validation we only show the RMSE, that is the criterion used for the optimization of the parameters of the models. We also show the best results obtained by Wang et al. [25] and Szabó et al. [29]. The results obtained by our methods, LR, KR and KDR, are comparable to previous proposals, although KDR does not improve previous results. It can also be seen that the composition of kernels proposed by Szabó in [29] obtains similar results.

For the test results in Table III, which are evaluated in test data not seen when training the models, we also show the mean error (ME) and the squared Pearson correlation coefficient (R$^2$). The proposed models generalize well, and, in fact, in this most difficult case, the two proposed methods perform better than the best methods presented so far.

*2) Multi-source approach:* The results are shown in Table IV. The first rows present the results of the sensors MISR2 and MODIS separately, and the last ones the results of the two possible combinations of both sensors: stacked and composite.

On the one hand, stacking inputs makes conclusions more articulated, as results depend on the models used at a great extent. If the models used are the simplest ones, e.g LR

and KR, improvements are not noticeable. On the contrary, when using DR methods, either RDR and KDR, slightly better results are obtained. However, the improvements with respect the use of the sensors separately is insignificant. On the other hand, using the composite approach, the proposed MDR always improves results with respect to all previous cases, i.e., the combination of sources is better than its individual use, and the kernel additive approach is better than the stacking one.

## V. Conclusions

This paper presented several methods based on kernels for distribution regression, and illustrated their performance in several scenarios of interest in remote sensing: 1) estimation of total crop and crop-specific yield from SMAP Vegetation Optical Depth, 2) estimation of Aerosol Optical Depth from MISR and MODIS reflectances, separately; and 3) combination of MISR and MODIS reflectances for the estimation of Aerosol Optical Depth. Our proposed methods successfully outperforms naive approaches based on input-space means and clustering, and is competitive with previously presented methods for multiple-instance learning.

The methods can be used in many other problems in geoscience, remote sensing and environmental sciences that require predicting a scalar from a set of vectors. We foresee many potential applications of interest beyond crop yield prediction or Aerosol Optical Depth from remote sensing data: for example predicting poverty, carbon stocks, plagues and diseases, population density, wealth or urbanization at county, region or country levels from multi-temporal and multi-sensor remote sensing data, are some possible examples. Currently, we are extending the method and application to work jointly with Sentinel-1 and 2 data for crop yield prediction in Europe. Methodologically, our agenda is tied to extend the framework to deal with missing data and features.

## References

[1] S. Liang, *Advances in Land Remote Sensing: System, Modeling, Inversion and Applications*. Germany: Springer Verlag, 2008.

[2] C. D. Rodgers, *Inverse Methods for Atmospheric Sounding: Theory and Practice*. World Scientific Publishing Co. Ltd., 2000.

[3] G. Camps-Valls, D. Tuia, L. Gómez-Chova, and J. Malo, Eds., *Remote Sensing Image Processing*. Morgan & Claypool, Sept 2011.

[4] J. Verrelst, L. Alonso, G. Camps-Valls, J. Delegido, and J. Moreno, "Retrieval of vegetation biophysical parameters using Gaussian process techniques," *IEEE Trans. Geosc. Rem. Sens.*, vol. 50, no. 5 PART 2, pp. 1832–1843, 2012.

[5] F. Ratle, G. Camps-Valls, and J. Weston, "Semisupervised neural networks for efficient hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 48, no. 5, pp. 2271–2282, 2010, cited By 59.

[6] G. Tramontana, K. Ichii, G. Camps-Valls, E. Tomelleri, and D. Papale, "Uncertainty analysis of gross primary production upscaling using random forests, remote sensing and eddy covariance data," *Remote Sensing of Environment*, vol. 168, pp. 360–373, 2015, cited By 1.

[7] J. Rojo-Álvarez, M. Martínez-Ramón, J. Muñoz-Marí, and G. Camps-Valls, *Digital Signal Processing with Kernel Methods*. UK: Wiley & Sons, Apr 2017.

[8] G. Camps-Valls, J. Verrelst, J. Muñoz-Marí, V. Laparra, F. Mateo-Jimenez, and J. Gomez-Dans, "A survey on gaussian processes for earth observation data analysis: A comprehensive investigation," *IEEE Geoscience and Remote Sensing Magazine*, no. 6, June 2016.

[9] A. Mateo-Sanchis, J. Muñoz-Marí, A. Pérez-Suay, and G. Camps-Valls, "Warped gaussian processes in remote sensing parameter estimation and causal inference," *IEEE Geoscience and Remote Sensing Letters*, pp. 1–5, 2018.

TABLE III: Test Results on the Three Remote Sensing Data Sets (See Text for Details) for AOD Prediction.

| Algorithms | MISR1 | | | MISR2 | | | MODIS | | |
|---|---|---|---|---|---|---|---|---|---|
| | ME×1000 | RMSE ×100 | $R^2$ | ME×1000 | RMSE ×100 | $R^2$ | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -6.14 ± 8.02 | 11.38 ± 0.90 | 0.64 ± 0.05 | -6.85 ± 6.28 | 9.90 ± 0.85 | 0.73 ± 0.05 | -1.64 ± 5.13 | 13.32 ± 0.92 | 0.54 ± 0.05 |
| KR | -2.33 ± 8.45 | 9.63 ± 0.65 | 0.74 ± 0.03 | -1.73 ± 6.50 | 8.87 ± 1.46 | 0.78 ± 0.06 | -0.75 ± 7.51 | 12.15 ± 1.16 | 0.62 ± 0.06 |
| RDR | -4.28 ± 7.40 | 8.30 ± 0.85 | 0.81 ± 0.03 | -1.71 ± 6.37 | 8.16 ± 1.26 | 0.81 ± 0.05 | -2.08 ± 5.61 | 10.09 ± 0.78 | 0.73 ± 0.06 |
| KDR | -4.34 ± 7.43 | 8.20 ± 0.83 | 0.81 ± 0.03 | -1.72 ± 6.31 | 8.18 ± 1.31 | 0.81 ± 0.05 | -2.42 ± 5.82 | 10.27 ± 1.46 | 0.73 ± 0.08 |

TABLE IV: Test Results on the Combined Dataset with 289 Common Bags for MISR2 and MODIS (See Text for Details) for AOD Prediction.

| Algorithms | MISR2 | | |
|---|---|---|---|
| | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -10.98 ± 5.14 | 8.13 ± 1.55 | 0.79 ± 0.05 |
| KR | 1.12 ± 8.03 | 6.00 ± 0.68 | 0.87 ± 0.04 |
| RDR | 0.86 ± 8.56 | 5.86 ± 0.77 | 0.88 ± 0.04 |
| KDR | 0.77 ± 8.55 | 5.85 ± 0.74 | 0.88 ± 0.04 |
| **MODIS** | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -7.33 ± 6.62 | 10.81 ± 1.78 | 0.61 ± 0.08 |
| KR | 3.29 ± 11.83 | 10.65 ± 1.71 | 0.64 ± 0.09 |
| RDR | 0.73 ± 9.14 | 8.74 ± 0.99 | 0.74 ± 0.06 |
| KDR | 0.88 ± 9.38 | 8.74 ± 0.96 | 0.74 ± 0.06 |
| **MISR2+MODIS (stack)** | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -6.76 ± 6.95 | 7.88 ± 0.80 | 0.79 ± 0.04 |
| KR | 2.74 ± 6.41 | 6.06 ± 0.66 | 0.87 ± 0.04 |
| RDR | 2.21 ± 6.35 | 5.85 ± 0.65 | 0.87 ± 0.04 |
| KDR | 2.41 ± 6.39 | 5.83 ± 0.64 | 0.87 ± 0.04 |
| **MISR2+MODIS (compo)** | ME×1000 | RMSE ×100 | $R^2$ |
| LR | -1.11 ± 9.18 | 7.17 ± 0.98 | 0.83 ± 0.03 |
| KR | 0.44 ± 8.45 | 6.07 ± 0.80 | 0.87 ± 0.04 |
| RDR | 0.72 ± 7.44 | 5.72 ± 0.61 | 0.89 ± 0.03 |
| KDR | 0.70 ± 7.28 | 5.69 ± 0.62 | 0.89 ± 0.03 |

[10] R. López-Lozano, G. Duveiller, L. Seguini, M. Meroni, S. García-Condado, J. Hooker, O. Leo, and B. Baruth, "Towards regional grain yield forecasting with 1km-resolution eo biophysical products: Strengths and limitations at pan-european level," *Agricultural and Forest Meteorology*, vol. 206, pp. 12 – 32, 2015.

[11] G. Galidaki, D. Zianis, I. Gitas, K. Radoglou, V. Karathanassi, M. Tsakiri–Strati, I. Woodhouse, and G. Mallinis, "Vegetation biomass estimation with remote sensing: focus on forest and other wooded land over the mediterranean ecosystem," *International Journal of Remote Sensing*, vol. 38, no. 7, pp. 1940–1966, 2017.

[12] P. A. Kucera, E. E. Ebert, F. J. Turk, V. Levizzani, D. Kirschbaum, F. J. Tapiador, A. Loew, and M. Borsche, "Precipitation from space: Advancing earth system science," *Bulletin of the American Meteorological Society*, vol. 94, no. 3, pp. 365–375, 2013.

[13] W. T. Crow, A. A. Berg, M. H. Cosh, A. Loew, B. P. Mohanty, R. Panciera, P. de Rosnay, D. Ryu, and J. P. Walker, "Upscaling sparse ground-based soil moisture observations for the validation of coarse-resolution satellite soil moisture products," *Reviews of Geophysics*, vol. 50, no. 2, 2012.

[14] W. Tang, S. H. Yueh, A. G. Fore, and A. Hayashi, "Validation of aquarius sea surface salinity with in situ measurements from argo floats and moored buoys," *Journal of Geophysical Research: Oceans*, vol. 119, no. 9, pp. 6171–6189, 2014.

[15] Z. Wang, "Aerosol Optical Depth Prediction from Satellite Observations by Multiple Instance Regression," *SIAM International Conference on Data Mining, SIAM*, pp. 165–176, 2008.

[16] N. Charkovska, J. Horabik-Pyzel, R. Bun, O. Danylo, Z. Nahorski, M. Jonas, and X. Xiangyang, "High-resolution spatial distribution and associated uncertainties of greenhouse gas emissions from the agricultural sector," *Mitigation and Adaptation Strategies for Global Change*, Jan 2018.

[17] G. Moser, M. D. Martino, and S. B. Serpico, "Estimation of air surface temperature from remote sensing images and pixelwise modeling of the estimation uncertainty through support vector machines," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 8, no. 1, pp. 332–349, Jan 2015.

[18] P. Goovaerts, "Combining areal and point data in geostatistical interpolation: Applications to soil science and medical geography." *Mathematical geosciences*, vol. 42 5, pp. 535–554, 2010.

[19] J. Bolton and P. Gader, "Application of multiple-instance learning for hyperspectral image analysis," *IEEE Geoscience and Remote Sensing Letters*, vol. 8, no. 5, pp. 889–893, Sept 2011.

[20] A. Manandhar, P. A. Torrione, L. M. Collins, and K. D. Morton, "Multiple-instance hidden markov model for gpr-based landmine detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 4, pp. 1737–1745, April 2015.

[21] C. Jiao and A. Zare, "Functions of multiple instances for learning target signatures," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 8, pp. 4670–4686, Aug 2015.

[22] S. E. Yuksel, J. Bolton, and P. Gader, "Multiple-instance hidden markov models with applications to landmine detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 12, pp. 6766–6775, Dec 2015.

[23] X. Liu, L. Jiao, J. Zhao, J. Zhao, D. Zhang, F. Liu, S. Yang, and X. Tang, "Deep multiple instance learning-based spatial-spectral classification for pan and ms imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 1, pp. 461–473, Jan 2018.

[24] K. L. Wagstaff and T. Lane, "Salience Assignment for Multiple-Instance Regression," *ICML '2007 Workshop on Constrained Optimization and Structured Output Spaces*, 2007.

[25] Z. Wang, L. Lan, and S. Vucetic, "Mixture model for multiple instance regression and applications in remote sensing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 6, pp. 2226–2237, 2012.

[26] B. Poczos, A. Singh, A. Rinaldo, and L. Wasserman, "Distribution-free distribution regression," in *Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, C. M. Carvalho and P. Ravikumar, Eds., vol. 31. Scottsdale, Arizona, USA: PMLR, 29 Apr–01 May 2013, pp. 507–515.

[27] J. Oliva, B. Poczos, and J. Schneider, "Distribution to distribution regression," in *Proceedings of the 30th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, S. Dasgupta and D. McAllester, Eds., vol. 28, no. 3. Atlanta, Georgia, USA: PMLR, 17–19 Jun 2013, pp. 1049–1057.

[28] J. Oliva, W. Neiswanger, B. Poczos, J. Schneider, and E. Xing, "Fast Distribution To Real Regression," in *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, S. Kaski and J. Corander, Eds., vol. 33. Reykjavik, Iceland: PMLR, 22–25 Apr 2014, pp. 706–714.

[29] Z. Szabó, B. K. Sriperumbudur, B. Póczos, and A. Gretton, "Learning theory for distribution regression," *Journal of Machine Learning Research*, vol. 17, no. 152, pp. 1–40, 2016.

[30] H. C. L. Law, D. J. Sutherland, D. Sejdinovic, and S. Flaxman, "Bayesian Approaches to Distribution Regression," *ArXiv e-prints*, 2017.

[31] B. Schölkopf and A. Smola, *Learning with Kernels – Support Vector Machines, Regularization, Optimization and Beyond*. MIT Press Series, 2002.

[32] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. New York, NY, USA: Cambridge University Press, 2004.

[33] G. Camps-Valls and L. Bruzzone, Eds., *Kernel methods for Remote Sensing Data Analysis*. UK: Wiley & Sons, Dec. 2009.

[34] K. Muandet, K. Fukumizu, B. Sriperumbudur, and B. Schölkopf, *Kernel Mean Embedding of Distributions: A Review and Beyond*, 2016.

[35] N. Aronszajn, "Theory of reproducing kernels," *Trans. Amer. Math. Soc.*, vol. 68, no. 3, pp. 337–404, May 1950.

[36] F. Riesz and B. S. Nagy, *Functional Analysis*. Frederick Ungar Publishing, 1955.

[37] G. Kimeldorf and G. Wahba, "Some results on tchebycheffian spline functions," *Journal of Mathematical Analysis and Applications*, vol. 33, no. 1, pp. 82–95, 1971.

[38] Z. Harchaoui, F. Bach, O. Cappe, and E. Moulines, "Kernel-based methods for hypothesis testing: A unified view," *IEEE Signal Processing Magazine*, vol. 30, no. 4, pp. 87–97, July 2013.

[39] G. Camps-Valls, J. Mooij, and B. Schölkopf, "Remote sensing feature selection by kernel dependence estimation," *IEEE Geoscience and Remote Sensing Letters*, vol. 7, no. 3, pp. 587–591, July 2010.

[40] G. Matasci, M. Volpi, M. Kanevski, L. Bruzzone, and D. Tuia, "Semisupervised transfer component analysis for domain adaptation in remote sensing image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 7, pp. 3550–3564, July 2015.

[41] C. Persello and L. Bruzzone, "Kernel-based domain-invariant feature selection in hyperspectral images for transfer learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 5, pp. 2615–2626, May 2016.

[42] A. Rahimi and B. Recht, "Random features for large-scale kernel machines," in *Proceedings of the 20th International Conference on Neural Information Processing Systems*, ser. NIPS'07. USA: Curran Associates Inc., 2007, pp. 1177–1184.

[43] A. Rahimi and B. Recht, "Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning," in *Advances in Neural Information Processing Systems 21*, D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, Eds., 2009, pp. 1313–1320.

[44] J. Sutherland and J. Schneider, "On the error of random fourier features," in *UAI*, 2015, pp. 862–871.

[45] L. K. Jones, "Annals of statistics," *A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training*, vol. 20, pp. 608–613, 1992.

[46] P. Morales-Álvarez, A. Pérez-Suay, R. Molina, and G. Camps-Valls, "Remote sensing image classification with large-scale gaussian processes," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 2, pp. 1103–1114, Feb 2018.

[47] A. G. Konings, M. Piles, N. Das, and D. Entekhabi, "L-band vegetation optical depth and effective scattering albedo estimation from smap," *Remote Sensing of Environment*, vol. 198, pp. 460 – 470, 2017.

[48] T. Jackson and T. Schmugge, "Vegetation effects on the microwave emission of soils," *Remote Sensing of Environment*, vol. 36, no. 3, pp. 203 – 212, 1991.

[49] M. Piles, G. Camps-Valls, D. Chaparro, D. Entekhabi, A. G. Konings, and T. Jagdhuber, "Remote sensing of vegetation dynamics in agro-ecosystems using smap vegetation optical depth and optical vegetation indices," in *IGARSS17*, July 2017, pp. 4346–4349.

[50] D. Chaparro, M. Piles, M. Vall-llossera, A. Camps, A. G. Konings, and D. Entekhabi, "L-band vegetation optical depth seasonal metrics for crop yield assessment," *Remote Sensing of Environment*, vol. 212, pp. 249–259, June 2018.

[51] B. Holben, T. Eck, I. Slutsker, D. Tanré, J. Buis, A. Setzer, E. Vermote, J. Reagan, Y. Kaufman, T. Nakajima, F. Lavenu, I. Jankowiak, and A. Smirnov, "Aeronet—a federated instrument network and data archive for aerosol characterization," *Remote Sensing of Environment*, vol. 66, no. 1, pp. 1 – 16, 1998.

**Adrián Pérez-Suay** obtained his B.Sc. degree in Mathematics (2007), Master degree in Advanced Computing and Intelligent Systems (2010) and the Ph.D. degree in Computational Mathematics and Computer Science (2015), all from the Universitat de València. He is assistant professor in the Department of Mathematics in the Universitat de València. He is currently a Postdoctoral Researcher at the Image Processing Laboratory working on dependence estimation, kernel methods and causal inference for remote sensing data analysis.

**Jordi Muñoz-Marí** was born in Valencia, Spain in 1970, and received a B.Sc. degree in Physics (1993), a B.Sc. degree in Electronics Engineering (1996), and a Ph.D. degree in Electronics Engineering (2003) from the Universitat de València (UV). At present he is an Associate Professor in the Electronics Engineering Department at UV, where he teaches machine learning, big data, digital signal processing and electronics. He is a research member of the Image Processing Laboratory (https://isp.uv.es/). His research interests are tied to machine learning, statistical methods and digital signal processing applied to data analysis in general and remote sensing in particular. He is a skilled programmer in different computer languages such as C/C++, Java, Python and others. He is author and co-author of many journal papers, book chapters, and conference papers. Please visit https://www.uv.es/jordi/ for more information.

**Anna Mateo-Sanchis** received the B.Sc. degree in cartography and geodesy engineering from Universitat Politècnica de València (UPV), Spain. In 2017, she joined the Image Processing Laboratory, Universitat de València (UV). She is currently working toward the Ph.D. degree at Universitat de València, working on multivariate/multioutput learning methods for regression problems applied to remote sensing. Her main research interests include image processing and machine learning algorithms for Earth Observation global monitoring.

**Maria Piles** (S'05-M'11-SM'19) received the M.Sc. degree (2005) in telecommunication engineering from Universitat Politècnica de València, Spain, and the Ph.D. degree (2010) in signal theory and communications, from the Universitat Politècnica de Catalunya (UPC), Spain, mastering in remote sensing.

In 2010, she was Research Fellow at University of Melbourne, Australia. From 2011 to 2015, she was Research Scientist at UPC and Affiliated Scientist at Massachusetts Institute of Technology, Cambridge. In 2016, she joined the Institute of Marine Sciences, CSIC, as a Research Scientist. Since 2017, she is with the Image Processing Laboratory, Universitat de València, as a Ramón y Cajal Senior Researcher. She has wide experience in the retrieval of the water content in soils and vegetation from low-frequency microwaves, and has been actively involved within the scientific activities of the ESA's SMOS and the NASA's SMAP missions. She is a member of the CIMR Mission Advisory Group. Her research interests include microwave remote sensing, estimation of soil moisture and vegetation bio-geophysical parameters, and development of multi-sensor techniques for enhanced retrievals with focus on agriculture, forestry, wildfire prediction, extreme detection and climate studies. She is currently serving as president of the Spanish chapter of the IEEE Geoscience and Remote Sensing Society.

**Jose E. Adsuara** received the M.Sc. degree in computer science and engineering from Universitat Politècnica de Catalunya (UPC), Spain. He spent a few years working in industry and education. After this, he also received the M.Sc. degree in mathematics and the Ph.D. degree in advanced physics, astrophysics, both from Universitat de València (UV), Spain. He is currently an assistant professor in the Department of Computer Science at UV. He is also with the Image Processing Laboratory (IPL), Universitat de València, since 2018 as a postdoctoral researcher working on machine learning for remote sensing.

**Gustau Camps-Valls** (M'04−SM'07−F'18) received the Ph.D. degree in physics from the Universitat de València, Valencia, Spain, in 2002. He is currently a Full Professor of electrical engineering and a Coordinator of the Image and Signal Processing Group, Image Processing Laboratory, with the Universitat de València. He is involved in the development of machine learning algorithms for geoscience and remote sensing data analysis. He has authored 200 journal papers, more than 200 conference papers, and 20 international book chapters. He holds a Hirsch's index, h = 60 (source: Google Scholar), entered the ISI list of Highly Cited Researchers in 2011, and Thomson Reuters ScienceWatch identified one of his papers on Kernel-based analysis of hyperspectral images as a Fast Moving Front research.

In 2015, he was a recipient of the Prestigious European Research Council (ERC) Consolidator Grant on Statistical Learning for Earth Observation Data Analysis. He is/has been the Associate Editor of the IEEE TRANSACTIONS ON SIGNAL PROCESSING, the IEEE GEOSCIENCE AND REMOTE SENSING LETTERS, the IEEE SIGNAL PROCESSING LETTERS, and the Invited Guest Editor for IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING in 2012 and IEEE *Geoscience and Remote Sensing Magazine* in 2015. He serves as the Editor for the books *Kernel Methods engineering*, *Signal and Image Processing* (IGI, 2007), *Kernel Methods for Remote Sensing Data Analysis* (Wiley & Sons, 2009), *Remote Sensing Image Processing* (MC, 2011), and *Digital Signal Processing with Kernel Methods* (Wiley & Sons, 2018).