

# Learning Structures in Earth Observation Data with Gaussian Processes

Fernando Mateo, Jordi Muñoz-Marí, Valero Laparra,  
Jochem Verrelst, and Gustau Camps-Valls \*

Image Processing Laboratory, University of Valencia,  
C/ Catedrático José Beltrán 2, 46980 Paterna, Spain  
{fmateo, jordi, lapeva, jverrelst, gcamps}@uv.es  
<http://isp.uv.es/>

**Abstract.** Gaussian Processes (GPs) has experienced tremendous success in geoscience in general and for bio-geophysical parameter retrieval in the last years. GPs constitute a solid Bayesian framework to formulate many function approximation problems consistently. This paper reviews the main theoretical GP developments in the field. We review new algorithms that respect the signal and noise characteristics, that provide feature rankings automatically, and that allow applicability of associated uncertainty intervals to transport GP models in space and time. All these developments are illustrated in the field of geoscience and remote sensing at a local and global scales through a set of illustrative examples.

**Keywords:** Kernel methods, Gaussian Process Regression (GPR), Bio-geophysical parameter estimation.

## 1 Introduction

Spatio-temporally explicit, quantitative retrieval methods of Earth surface and atmosphere characteristics are a requirement in a variety of Earth observation applications. Optical sensors mounted on-board Earth observation (EO) satellites are being endowed with high temporal, spectral and spatial resolutions, and thus enable the retrieval and monitoring of climate and bio-geophysical variables [9, 25]. With the super-spectral Copernicus Sentinel-2 (S2) [10] and the forthcoming Sentinel-3 missions [8], among other planned space missions, an unprecedented data stream for land, ocean and atmosphere monitoring will soon become available to a diverse user community. This vast data streams require enhanced processing techniques. Statistical inference methods play an important role in this area of research. Understanding is more challenging than predicting, and thus statistical models should not only be accurate but also capture plausible physical relations and explain the problem at hand.

---

\* Paper published in *Advanced Analysis and Learning on Temporal Data. AALTD 2015. Lecture Notes in Computer Science*, vol 9785. Springer, Cham. [https://doi.org/10.1007/978-3-319-44412-3\\_6](https://doi.org/10.1007/978-3-319-44412-3_6)

Over the last few decades a wide diversity of bio-geophysical retrieval methods have been developed, but only a few of them made it into operational processing chains. Essentially, we may find two main approaches to the inverse problem of estimating biophysical parameters from spectra: *parametric physically-based models* and *non-parametric statistical models*. Lately, machine learning has attained outstanding results in the estimation of climate variables and related bio-geophysical parameters at local and global scales [7]. For example, current operational vegetation products, like leaf area index (LAI), are typically produced with neural networks [2], Gross Primary Production (GPP) as the largest global CO<sub>2</sub> flux driving several ecosystem functions is estimated using ensembles of random forests and neural networks [3, 15], biomass has been estimated with stepwise multiple regression [24], PCA and piecewise linear regression for sun-induced fluorescence (SIF) estimation [12], support vector regression showed high efficiency in modelling LAI, fractional vegetation cover (fCOVER), evapotranspiration [11, 35], relevance vector machines were successful in ocean chlorophyll estimation [5], and recently, Gaussian Processes (GPs) [21] provided excellent results in vegetation properties estimation [22, 30–32]. The family of Bayesian non-parametrics, and of Gaussian processes in particular [21], have been paid wide attention in the last years in remote sensing data analysis. We will review the main developments in GPs for EO data analysis in this paper.

The remainder of the paper is organized in two main parts: Section II reviews the main notation and theory of GP regression. Section III presents some of the most recent advances of GP models applied to remote sensing data processing. Section IV presents ways to extract knowledge from those GP models. We conclude in Section V with a discussion about the upcoming challenges and research directions.

## 2 Gaussian Process Regression

### 2.1 Gaussian processes: a gentle introduction

Gaussian processes (GPs) are state-of-the-art tools for discriminative machine learning. They can be interpreted as a family of kernel methods with the additional advantage of providing a full conditional statistical description for the predicted variable. Standard regression approximates observations (often referred to as *outputs*)  $\{y_n\}_{n=1}^N$  as the sum of some unknown latent function  $f(\mathbf{x})$  of the inputs  $\{\mathbf{x}_n \in \mathbb{R}^D\}_{n=1}^N$  plus *constant power* (*homoscedastic*) Gaussian noise, i.e.

$$y_n = f(\mathbf{x}_n) + \varepsilon_n, \quad \varepsilon_n \sim \mathcal{N}(0, \sigma^2). \quad (1)$$

GP regression proceeds in a Bayesian, non-parametric way, to fit the observed data. A zero mean<sup>1</sup> GP prior is placed on the latent function  $f(\mathbf{x})$  and a Gaussian prior is used for each latent noise term  $\varepsilon_n$ ,  $f(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, k_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}'))$ , where  $k_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}')$  is a covariance function parametrized by  $\boldsymbol{\theta}$  and  $\sigma^2$  is a hyperparameter that specifies the noise power. Essentially, a GP is a stochastic process whose marginals are distributed as a multivariate Gaussian. In particular, given the priors  $\mathcal{GP}$ , samples drawn from  $f(\mathbf{x})$  at the set of locations  $\{\mathbf{x}_n\}_{n=1}^N$  follow a joint multivariate Gaussian with zero mean and covariance (sometimes referred as to *kernel*) matrix  $\mathbf{K}_{\mathbf{ff}}$  with  $[\mathbf{K}_{\mathbf{ff}}]_{ij} = k_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{x}_j)$ .

If we consider a test location  $\mathbf{x}_*$  with corresponding output  $y_*$ , priors  $\mathcal{GP}$  induce a prior distribution between the observations  $\mathbf{y} \equiv \{y_n\}_{n=1}^N$  and  $y_*$ . Collecting available data in  $\mathcal{D} \equiv \{\mathbf{x}_n, y_n | n = 1, \dots, N\}$ , it is possible to analytically compute the posterior distribution over the unknown output  $y_*$ :

$$p(y_* | \mathbf{x}_*, \mathcal{D}) = \mathcal{N}(y_* | \mu_{\text{GP}*}, \sigma_{\text{GP}*}^2) \quad (2)$$

$$\mu_{\text{GP}*} = \mathbf{k}_{\mathbf{f}*}^{\top} (\mathbf{K}_{\mathbf{ff}} + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{y} = \mathbf{k}_{\mathbf{f}*}^{\top} \boldsymbol{\alpha} \quad (3)$$

$$\sigma_{\text{GP}*}^2 = \sigma^2 + k_{**} - \mathbf{k}_{\mathbf{f}*}^{\top} (\mathbf{K}_{\mathbf{ff}} + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{k}_{\mathbf{f}*}. \quad (4)$$

which is computable in  $\mathcal{O}(n^3)$  time (this cost arises from the inversion of the  $n \times n$  matrix  $(\mathbf{K}_{\mathbf{ff}} + \sigma^2 \mathbf{I})$ , see [21]). In addition to the computational cost, GPs require large memory since in naive implementations one has to store the training kernel matrix, which amounts to  $\mathcal{O}(n^2)$ .

## 2.2 On the model selection

The corresponding hyperparameters  $\{\boldsymbol{\theta}, \sigma_n\}$  are typically selected by Type-II Maximum Likelihood, using the marginal likelihood (also called evidence) of the observations, which is also analytical (explicitly conditioning on  $\boldsymbol{\theta}$  and  $\sigma_n$ ):

$$\log p(\mathbf{y} | \boldsymbol{\theta}, \sigma_n) = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}_{\mathbf{ff}} + \sigma_n^2 \mathbf{I}). \quad (5)$$

When the derivatives of (5) are also analytical, which is often the case, conjugated gradient ascend is typically used for optimization.

## 2.3 On the covariance function

The core of a kernel method like GPs is the appropriate definition of the covariance (or kernel) function. A standard, widely used covariance function is the squared exponential,  $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / (2\sigma^2))$ , which captures sample similarity well in most of the (unstructured) problems, and only one hyperparameter  $\sigma$  needs to be tuned.

<sup>1</sup> It is customary to subtract the sample mean to data  $\{y_n\}_{n=1}^N$ , and then to assume a zero mean model.

**Table 1.** Some kernel functions used in the literature.

Kernel function	Expression
Linear	$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}' + c$
Polynomial	$k(\mathbf{x}, \mathbf{x}') = (\alpha \mathbf{x}^\top \mathbf{x}' + c)^d$
Gaussian	$k(\mathbf{x}, \mathbf{x}') = \exp(-\ \mathbf{x} - \mathbf{x}'\ ^2 / (2\sigma^2))$
Exponential	$k(\mathbf{x}, \mathbf{x}') = \exp(-\ \mathbf{x} - \mathbf{x}'\  / (2\sigma^2))$
Rational Quadratic	$k(\mathbf{x}, \mathbf{x}') = 1 - (\ \mathbf{x} - \mathbf{x}'\ ^2) / (\ \mathbf{x} - \mathbf{x}'\ ^2 + c)$
Multiquadric	$k(\mathbf{x}, \mathbf{x}') = \sqrt{\ \mathbf{x} - \mathbf{x}'\ ^2 + c^2}$
Inv. Multiquad.	$k(\mathbf{x}, \mathbf{x}') = 1 / (\sqrt{\ \mathbf{x} - \mathbf{x}'\ ^2 + \theta^2})$
Power	$k(\mathbf{x}, \mathbf{x}') = -\ \mathbf{x} - \mathbf{x}'\ ^d$
Log	$k(\mathbf{x}, \mathbf{x}') = -\log(\ \mathbf{x} - \mathbf{x}'\ ^d + 1)$

In the context of GPs, kernels with more hyperparameters can be efficiently inferred. This is an opportunity to exploit asymmetries in the feature space by including a parameter per feature, as in the very common anisotropic squared exponential (SE) kernel function:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \nu \exp\left(-\sum_{f=1}^F \frac{(x_i^f - x_j^f)^2}{2\sigma_f^2}\right) + \sigma_n^2 \delta_{ij},$$

where  $\nu$  is a scaling factor,  $\sigma_n$  is the standard deviation of the (estimated) noise, and a  $\sigma_f$  is the length-scale per input features,  $f = 1, \dots, F$ . This is a very flexible covariance function that typically suffices to tackle most of the problems. Table 1 summarizes the most common kernel functions in standard applications with kernel methods.

## 2.4 Gaussian processes exemplified

Let us illustrate the solution of GP regression (GPR) in a toy example. In Fig. 1 we include an illustrative example with 6 training points in the range between  $-2$  and  $+2$ . We firstly depict several random functions drawn from the GP prior and then we include functions drawn from the posterior. We have chosen an isotropic Gaussian kernel and  $\sigma_\nu = 0.1$ . We have plotted the mean function plus/minus two standard deviations (corresponding to a 95% confidence interval). Typically, the hyperparameters are unknown, as well as the mean, covariance and likelihood functions. We assumed a Squared Exponential (SE) covariance function and learned the optimal hyperparameters by minimizing the negative log marginal likelihood (NLML) w.r.t. the hyperparameters. We observe three different regions in the figure. Below  $x = -1.5$ , we do not have samples and the GPR provides the solution given by the prior (zero mean and  $\pm 2$ ). At the center, where most of the data points lie, we have a very accurate view of the latent function with small error bars (close to  $\pm 2\sigma_\nu$ ). For  $x > 0$ , we do not have training samples neither so we have same behaviour. GPs typically provide

**Fig. 1.** Example of a Gaussian process. Left: some functions drawn at random from the GP prior. Right: some random functions drawn from the posterior, i.e. the prior conditioned on 6 noise-free observations indicated in big black dots. The shaded area represents the point-wise mean plus and minus two times the standard deviation for each input value (corresponding to the 95 confidence region). It can be noted that the confidence intervals become large for regions far from the observations. *Note: This is an animated figure that only works in Acrobat reader.*

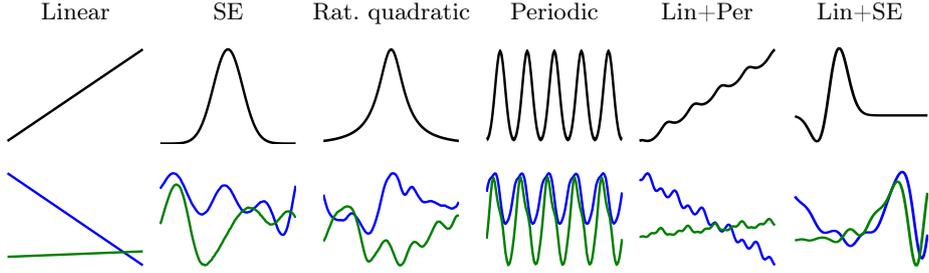
an accurate solution where the data lies and high error bars where we do not have available information and, consequently, we presume that the prediction in that area is not accurate. This is why in regions of the input space without points the confidence intervals are wide resembling the prior distribution.

### 3 Advances in Gaussian Process Regression

In this section, we review some recent advances in GPR especially suited for remote sensing data analysis. We will review the main aspects to design covariance functions that capture non-stationarities and multiscale time relations in EO data, as well as GPs that can learn arbitrary transformations of the observed variable and noise models.

#### 3.1 Structured, non-stationary and multiscale

Commonly used kernels families include the squared exponential (SE), periodic (Per), linear (Lin), and rational quadratic (RQ), cf. Table 1. Illustration of the base kernel and drawings from the GP prior is shown in Fig. 2. These base kernels can be actually combined following simple operations: summation, multiplication or convolution. This way one may build sophisticated covariances from simpler ones. Note that the same essential property of kernel methods apply here: a valid covariance function must be positive semidefinite. In general, the design of the kernel should rely on the information that we have for each estimation problem and should be designed to get the most accurate solution with the least amount of samples.



**Fig. 2.** Base kernels (top) and two random draws from a GP with each respective kernel and combination of kernels (bottom).

In Fig. 2, all the base kernels are one-dimensional. Nevertheless, kernels over multidimensional inputs can be actually constructed by adding and multiplying kernels over individual dimensions: (a) linear, (b) squared exponential (or RBF), (c) rational quadratic, and (d) periodic. See Table 1 for the explicit functional form of each kernel. Some simple kernel combinations are represented in the two last columns of the figure: a linear plus periodic covariances may capture structures that are periodic with trend (e), while a linear plus squared exponential covariances can accommodate structures with increasing variation (f). By summing kernels, we can model the data as a superposition of independent functions, possibly representing different structures in the data. For example, in multitemporal image analysis, one could for instance dedicate a kernel for the time domain (perhaps trying to capture trends and seasonal effects) and another kernel function for the spatial domain (equivalently capturing spatial patterns and auto-correlations). In time series models, sums of kernels can express superposition of different processes, possibly operating at different scales: very often changes in geophysical variables through time occur at different temporal resolutions (hours, days, etc.), and this can be incorporated in the prior covariance with those simple operations. In multiple dimensions, summing kernels gives additive structure over different dimensions, similar to generalized additive models [13]. Alternatively, multiplying kernels allows us to account for interactions between different input dimensions or different notions of similarity. In the following section, we show how to design kernels that incorporate particular time resolutions, trends and periodicities.

### 3.2 Time-based covariance for GPR

Signals to be processed typically show particular characteristics, with time-dependent cycles and trends. One could include time  $t_i$  as an additional feature in the definition of the input samples. This *stacked approach* [4] essentially relies on a covariance function  $k(\mathbf{z}_i, \mathbf{z}_j)$ , where  $\mathbf{z}_i = [t_i, \mathbf{x}_i]^\top$ . The shortcoming is that the time relations are naively left to the nonlinear regression algorithm, and hence no explicit time structure model

is assumed. To cope with this, one can use a linear combination (or composite) of different kernels: one dedicated to capture the different temporal characteristics, and the other to the feature-based relations.

The issue here is how to design kernels capable of dealing with non-stationary processes. A possible approach is to use a *stationary* covariance operating on the variable of interest after being mapped with a nonlinear function engineered to discount such undesired variations. This approach was used in [23] to model *spatial patterns* of solar radiation with GPR. It is also possible to adopt a squared exponential (SE) as stationary covariance acting on the *time* variable mapped to a two-dimensional *periodic space*  $\mathbf{z}(t) = [\cos(t), \sin(t)]^\top$ , as explained in [21],

$$k(t_i, t_j) = \exp\left(-\frac{\|\mathbf{z}(t_i) - \mathbf{z}(t_j)\|^2}{2\sigma_t^2}\right), \quad (6)$$

which gives rise to the following periodic covariance function

$$k(t_i, t_j) = \exp\left(-\frac{2\sin^2[(t_i - t_j)/2]}{\sigma_t^2}\right), \quad (7)$$

where  $\sigma_t$  is a hyper-parameter characterizing the periodic scale and needs to be inferred. It is not clear, though, that the seasonal trend is exactly periodic, so we modify this equation by taking the product with a squared exponential component, to allow a decay away from exact periodicity:

$$k_2(t_i, t_j) = \gamma \exp\left(-\frac{2\sin^2[\pi(t_i - t_j)]}{\sigma_t^2} - \frac{(t_i - t_j)^2}{2\sigma_d^2}\right), \quad (8)$$

where  $\gamma$  gives the magnitude,  $\sigma_t$  the smoothness of the periodic component,  $\sigma_d$  represents the *decay-time* for the periodic component, and the period has been fixed to one year. Therefore, our final covariance is expressed as

$$k([\mathbf{x}_i, t_i], [\mathbf{x}_j, t_j]) = k_1(\mathbf{x}_i, \mathbf{x}_j) + k_2(t_i, t_j), \quad (9)$$

which is parameterized by only three more hyperparameters collected in  $\boldsymbol{\theta} = \{\nu, \sigma_1, \dots, \sigma_F, \sigma_n, \sigma_t, \sigma_d, \gamma\}$ . Note that this kernel function allows us to incorporate time easily, but the relations between time  $t_i$  and signal  $\mathbf{x}_i$  samples is missing. Some approximations to deal with this issue exist in the literature, such as cross-kernel composition [4, 6] or latent force models [1].

We show the advantage of encoding such prior knowledge and structure in the relevant problem of solar irradiation prediction using GPR. Noting the non-stationary temporal behaviour of the signal, we develop a particular time-based composite covariance to account for the relevant seasonal signal variations. Data from the AEMET radiometric observatory of Murcia (Southern Spain, 38.0° N, 1.2° W) were used. Table 2 reports the obtained results with GPR models and several statistical regression

**Table 2.** Results for the estimation of the daily solar irradiation of linear and nonlinear regression models. Subscript  $\text{METHOD}_t$  indicates that the  $\text{METHOD}$  includes time as input variable. Best results are highlighted in bold, the second best in italics.

METHOD	ME	RMSE	MAE	R
RLR	0.27	4.42	3.51	0.76
RLR <sub>t</sub>	0.25	4.33	3.42	0.78
SVR [26]	0.54	4.40	3.35	0.77
SVR <sub>t</sub>	0.42	4.23	3.12	0.79
RVM [28]	0.19	4.06	3.25	0.80
RVM <sub>t</sub>	0.14	3.71	3.11	0.81
GPR [21]	0.14	3.22	2.47	<i>0.88</i>
GPR <sub>t</sub>	<i>0.13</i>	<i>3.15</i>	<i>2.27</i>	<i>0.88</i>
<b>TGPR</b>	<b>0.11</b>	<b>3.14</b>	<b>2.19</b>	<b>0.90</b>

methods: regularized linear regression (RLR), support vector regression (SVR), relevance vector machine (RVM) and GPR. All methods were run with and without using two additional dummy time features containing the year and day-of-year (DOY). We will indicate the former case with a subscript, like e.g. SVR<sub>t</sub>. First, including time information improves all baseline models. Second, the best overall results are obtained by the GPR models, when including time information or not. Third, in particular, the proposed temporal GPR (TGPR) outperforms the rest in accuracy (root-mean-square error, RMSE, and mean absolute error, MAE) and goodness-of-fit ( $R$ ), and closely follows the elastic net in bias (ME). TGPR performs better than GPR and GPR<sub>t</sub> in all quality measures.

### 3.3 Heteroscedastic GPR: Learning the noise model

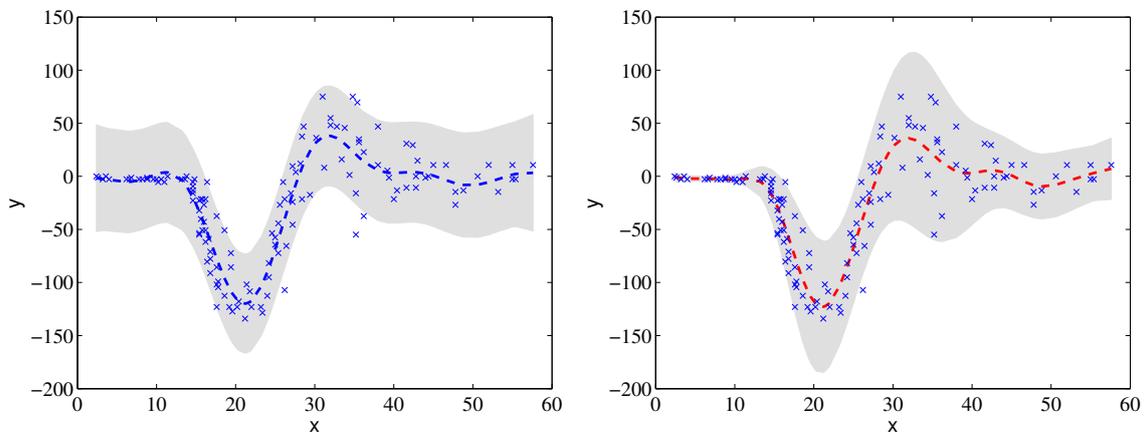
The standard GPR is essentially homoscedastic, i.e., assumes constant noise power  $\sigma^2$  for all observations. This assumption can be too restrictive for some problems. Heteroscedastic GPs, on the other hand, let noise power vary smoothly throughout input space, by changing the prior over  $\varepsilon_n$  to  $\varepsilon_n \sim \mathcal{N}(0, e^{g(\mathbf{x}_n)})$ , and placing a GP prior over  $g(\mathbf{x}) \sim \mathcal{GP}(\mu_0 \mathbf{1}, k_{\theta_g}(\mathbf{x}, \mathbf{x}'))$ . Note that the exponential is needed<sup>2</sup> in order to describe the non-negative variance. The hyperparameters of the covariance functions of both GPs are collected in  $\boldsymbol{\theta}_f$  and  $\boldsymbol{\theta}_g$ , accounting for the signal and the noise relations, respectively.

Relaxing the homoscedasticity assumption into heteroscedasticity yields a richer, more flexible model that contains the standard GP as a particular case corresponding to a constant  $g(\mathbf{x})$ . Unfortunately, this also hampers analytical tractability, so approx-

<sup>2</sup> Of course, other transformations are possible, just not as convenient.

imate methods must be used to obtain posterior distributions for  $f(\mathbf{x})$  and  $g(\mathbf{x})$ , which are in turn required to compute the predictive distribution over  $y_*$ .

As an alternative to the costly classic heteroscedastic GP approaches, variational techniques allow to approximate intractable integrals arising in Bayesian inference and machine learning. A sophisticated variational approximation called *Marginalized Variational (MV)* approximation was introduced in [16]. The MV approximation renders (approximate) Bayesian inference in the heteroscedastic GP model both fast and accurate. We will refer to this variational approximation for heteroscedastic GP regression as VHGP. A simple comparison between the homoscedastic canonical GP and the VHGP model is shown in Fig. 3.



**Fig. 3.** Predictive mean and variance of the standard GP (left) and the heteroscedastic GP (right). It is noticeable that in the low noise regime the VHGP produces tighter confidence intervals as expected, while high noise variance associated to high signal variance (middle of the observed signal) the predictive variance is more reasonable too.

### 3.4 Warped GPR: Learning the output transformation

Very often, in practical applications, one transforms the observed variable to better pose the problem. Actually, it is a standard practice to ‘linearize’ or ‘uniformize’ the distribution of the observations (which is commonly skewed due to the sampling strategies in *in situ* data collection) by applying non-linear link functions like the logarithmic, the exponential or the logistic functions.

*Warped* GPR [27] essentially warps observations  $\mathbf{y}$  through a nonlinear parametric function  $g$  to a latent space  $z_i = g(y_i) = g(f(\mathbf{x}_i) + \varepsilon_i)$ , where  $f$  is a possibly noisy latent function with  $d$  inputs, and  $g$  is a function with scalar inputs parametrized by  $\psi$ . The function  $g$  must be *monotonic*, otherwise the probability measure will not be conserved in the transformation, and the distribution over the targets may not be

**Table 3.** Results using both raw and empirically-transformed observation variables.

	ME	RMSE	MAE	R
<b>Raw</b>				
GPR	0.02	1.74	0.33	0.82
VHGPR	0.29	2.51	0.46	0.65
WGPR	0.08	1.71	0.30	0.83
<b>Empirically-based</b>				
GPR	0.15	1.69	0.29	0.86
VHGPR	0.15	1.70	0.29	0.85
WGPR	0.17	1.75	0.30	0.86

valid [27]. It can be shown that replacing  $y_i$  by  $z_i$  into the standard GP model leads to an extended problem that can be solved by taking derivatives of the negative log likelihood function in (5), but now with respect to both  $\boldsymbol{\theta}$  and  $\boldsymbol{\psi}$  parameter vectors.

For both the GPR and WGPR models we need to define the covariance (kernel, or Gram) function  $k(\cdot, \cdot)$ , which should capture the similarity between samples. We used the standard Automatic Relevance Determination (ARD) covariance [21]. Model hyperparameters are collectively grouped in  $\boldsymbol{\theta} = \{\nu, \sigma_n, \sigma_1, \dots, \sigma_d\}$ . In addition, for the WGPR we need to define a parametric smooth and monotonic form for  $g$ , which can be defined as:

$$g(y_i; \boldsymbol{\psi}) = \sum_{\ell=1}^L a_{\ell} \tanh(b_{\ell} y_i + c_{\ell}), \quad a_{\ell}, b_{\ell} \geq 0,$$

where  $\boldsymbol{\psi} = \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ . Recently, flexible non-parametric functions have replaced such parametric forms [18], thus placing another prior for  $g(\mathbf{x}) \sim \mathcal{GP}(f, c(f, f'))$ , whose model is learned via variational inference.

For illustration purposes, we focus on the estimation of chlorophyll-a concentrations from remote sensing upwelling radiance just above the ocean surface. We used the SeaBAM dataset [19, 20], which gathers 919 *in situ* pigment measurements around the United States and Europe. The dataset contains coincident *in situ* chlorophyll concentration and remote sensing reflectance measurements ( $R_{rs}(\lambda)$ , [ $\text{sr}^{-1}$ ]) at some wavelengths (412, 443, 490, 510 and 555 nm) that are present in the SeaWiFS ocean colour satellite sensor. The chlorophyll concentration values range from 0.019 to 32.79  $\text{mg}/\text{m}^3$  (revealing a clear exponential distribution).

Table 3 shows different scores –bias (ME), accuracy (RMSE, MAE) and goodness-of-fit ( $R$ )– between the observed and predicted variable when using the raw data (no *ad hoc* transform at all) and the empirically adjusted transform. Results are shown for three flavours of GPs: the standard GPR [21], the variational heteroscedastic GP (VHGPR) [17], and the proposed warped GP regression (WGPR) [18, 27] for different

rates of training samples. Empirically-based warping slightly improves the results over working with raw data for the same number of training samples, but this requires prior knowledge about the problem, time and efforts to fit an appropriate function. On the other hand, WGPR outperforms the rest of GPs in all comparisons over standard GPR and VHGPR ( $\sim +1 - 10\%$ ). Finally, WGPR nicely compensates the lack of prior knowledge about the (possibly skewed) distribution of the observation variable.

### 3.5 Source code and toolboxes

The most widely known sites to obtain free source code on GP modeling are GPML<sup>3</sup> and GPstuff<sup>4</sup>. The former website centralizes the main activities in GP modeling and provides up-to-date resources concerned with probabilistic modeling, inference and learning based on GPs, while the latter is a versatile collection of GP models and computational tools required for inference, sparse approximations and model assessment methods. We also recommend to the interested reader in regression in general, our MATLAB SimpleR<sup>5</sup> toolbox that contains many regression tools organized in families: tree-based, bagging and boosting, neural nets, kernel regression methods, and several Bayesian nonparametric models like GPs.

## 4 Analysis of Gaussian Process Models

An interesting possibility in GP models is to extract knowledge from the trained model. We will show in what follows two different approaches: 1) feature ranking exploiting the automatic relevance determination (ARD) covariance and 2) uncertainty estimation looking at the predictive variance estimates.

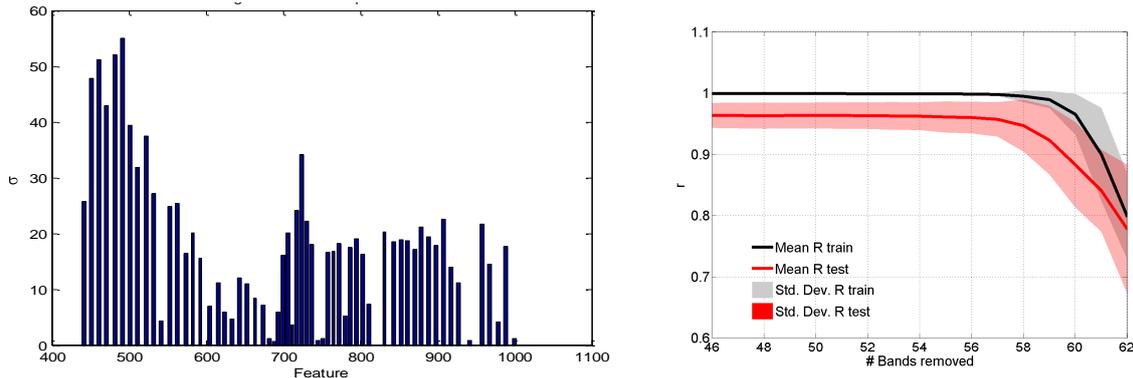
### 4.1 Ranking features through the ARD covariance

One of the advantages of GPs is that during the development of the GP model the predictive power of each single band is evaluated for the parameter of interest through calculation of the ARD. Specifically, band ranking through  $\sigma_b$  may reveal the bands that contribute the most to the development of a GP model. An example of the  $\sigma_b$ 's for one GP model trained with field leaf chlorophyll content (*Chl*) data and with 62 CHRIS bands is shown in Fig. 4 (left). The band with highest  $\sigma_b$  is the least contributing to the model. It can be noted that a relatively few bands (about 8) were evaluated as crucial for *Chl* estimation, while the majority of bands were evaluated as less contributing.

<sup>3</sup> <http://www.gaussianprocess.org/>

<sup>4</sup> <http://becs.aalto.fi/en/research/bayes/gpstuff/>

<sup>5</sup> <http://isp.uv.es/soft.htm>



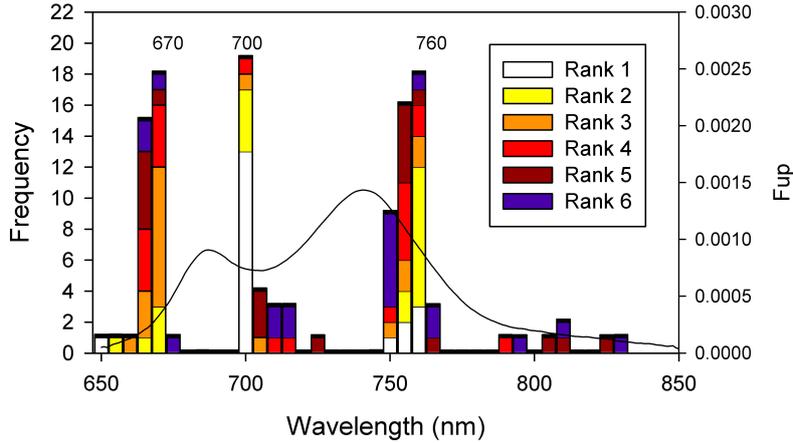
**Fig. 4.** Estimated  $\sigma_b$  values for one GP model using 62 CHRIS bands (left). The lower the  $\sigma_b$  the more important the band is for regression. *Chl r* and standard deviation (SD) of training and validation for GP fittings using backward elimination (right).

This does not necessarily mean that other bands are obstructing optimized accuracies. Only when less than 4 bands were left accuracies started to degrade rapidly Fig. 4 (right). The figure suggests that the most relevant spectral region is to be found between 550 and 1000 nm. Most contributing bands were positioned around the red edge, at 680 and 730 nm respectively, but not all bands within the red edge were evaluated as relevant. This is due to when having a large number of bands available then neighbouring bands do not provide much additional information and can thus be considered as redundant.

Consequently, the  $\sigma_b$  proved to be a valuable tool to detect most sensitive bands of a sensor towards a biophysical parameter. A more systematic analysis was applied by sorting the bands on their relevance and counting the band rankings over 50 repetitions. In [32] the four most relevant bands were tracked for *Chl*, LAI and fCOVER and for different Sentinel-2 settings. It demonstrated the potential of Sentinel-2, with its new band in the red-edge, for vegetation properties estimation. Also in [34]  $\sigma_b$  were used to analyze band sensitivity of Sentinel-2 towards LAI. A similar approach was pursued on analyzing leaf *Chl* based on tracking the most sensitive spectral regions of sun-induced fluorescence data [29], as displayed in Fig. 5.

## 4.2 Uncertainty intervals

In this section, we use GP models for retrieval and portability in space and time. For this, we will exploit the associated predictive variance (i.e. uncertainty interval) provided by GP models. Consequently, retrievals with high uncertainties refer to pixel spectral information that deviates from what has been represented during the training phase. In turn, low uncertainties refer to pixels that were well represented in the training phase. The quantification of variable-associated uncertainties is a strong re-

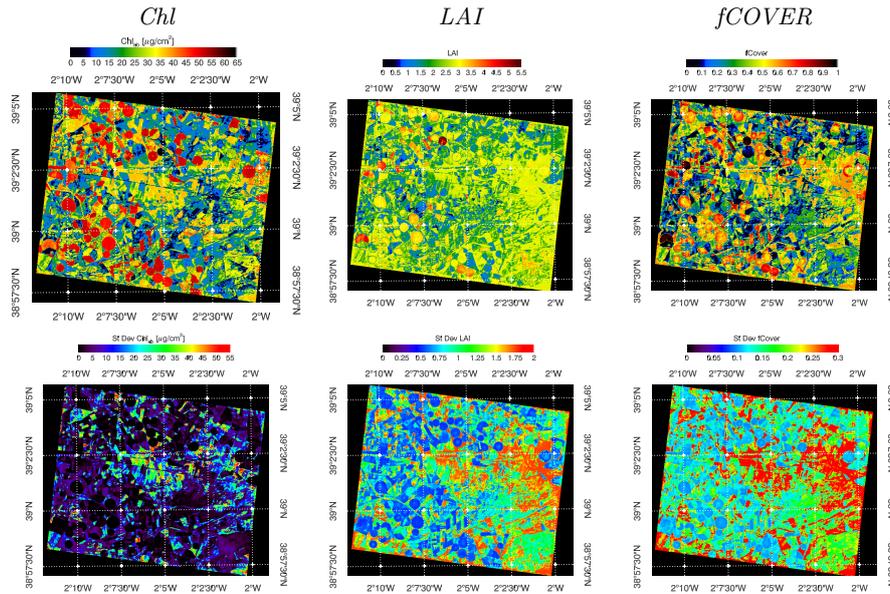


**Fig. 5.** Frequency plots of the top eight ranked bands with lowest  $\sigma_b$  values in 20 runs of GPR prediction of *Chl* based on upward fluorescence ( $F_{up}$ ) emission. An emission curve is given as illustration.

quirement when remote sensing products are ingested in higher level processing, e.g. to estimate ecosystem respiration, photosynthetic activity, or carbon sequestration [14].

The application of GPs for the estimation of biophysical parameters was initially demonstrated in [30]. A locally collected field dataset called SPARC-2003 at Barrax (Spain) was used for training and validation of GPs for the vegetation parameters of LAI, *Chl* and fCOVER. Sufficiently high validation accuracies were obtained ( $R^2 > 0.86$ ) for processing a CHRIS image into these parameters, as shown in Fig. 6. Within the uncertainty maps, areas with reliable retrievals are clearly distinguished from areas with unreliable retrievals. Low uncertainties were found on irrigated areas and harvested fields. High uncertainties were found on areas with remarkably different spectra, such as bright, whitish calcareous soils, or harvested fields. This indicates that the input spectrum deviates from what has been presented during the training stage, thereby imposing uncertainties to the retrieval.

GP models were subsequently applied to the SPARC dataset that was re-sampled to different Sentinel-2 band settings and then uncertainties were inspected [32]. On the whole, adding spectral information led to reduction of uncertainties and thus more meaningful biophysical parameter maps. The locally-trained GP models were applied to simulated Sentinel-2 images in a follow-up study [33]. Time series over the local Barrax site as well images across the world were processed. Also the role of an extended training dataset by adding spectra of non-vegetated surfaces were evaluated. Subsequently the uncertainty values were analyzed. By using the extended training dataset not only further improved performances but also allowed a decrease in theoretical uncertainties. The GP models were successfully applied to simulated Sentinel-2 images covering various sites; associated relative uncertainties were on the same order as those generated by the reference image.

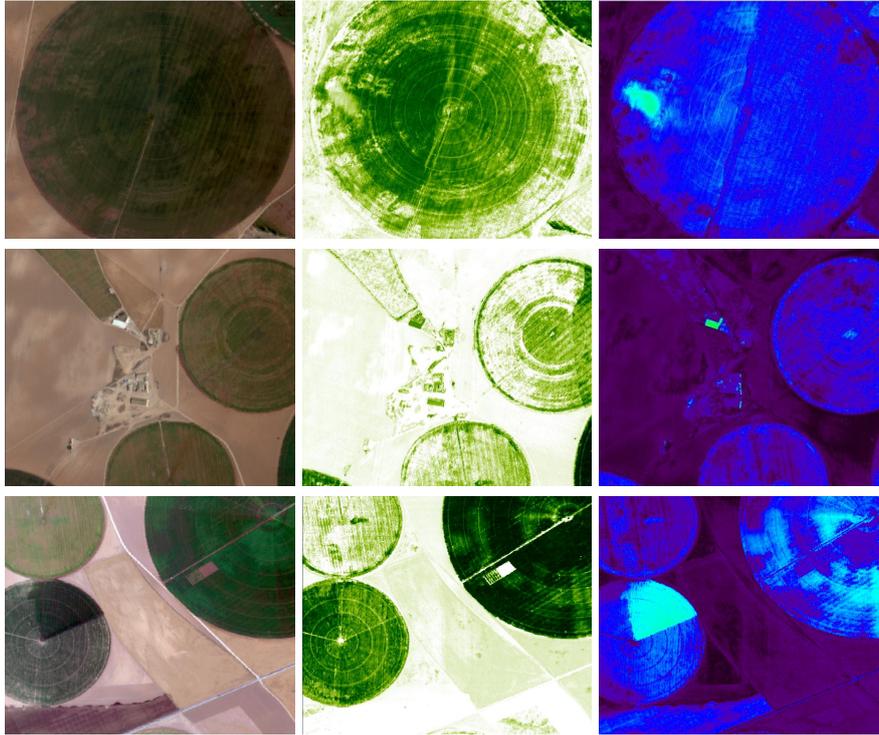


**Fig. 6.** Prediction maps (top) and associated uncertainty intervals (bottom), generated with GP and four bands of the CHRIS 12-07-2003 nadir image.

As a final example, uncertainty estimates were exploited to assess the robustness of the retrievals at multiple spatial scales. In [31], retrievals from hyperspectral airborne and spaceborne data over the Barrax area were compared. Based on the spaceborne SPARC-2003 dataset, GP developed a model that was excellently validated ( $r^2$ : 0.96). The SPARC-trained GP model was subsequently applied to airborne CASI flightlines (Barrax, 2009) to generate *Chl* maps. The accompanying uncertainty maps provided insight in the robustness of the retrievals. In general similar uncertainties were achieved by both sensors, which is encouraging for upscaling estimates from field to landscape scale. The high spatial resolution of CASI in combination with the uncertainties allows us to observe the spatial patterns of retrievals in more detail. Some examples of mean estimates and associated uncertainties are shown in Fig. 7.

## 5 Conclusions and further work

This paper provided a comprehensive survey to the field of Gaussian Processes (GPs) in the context of remote sensing data analysis for Earth observation applications, and in particular for biophysical parameter estimation. We summarized the main properties of GPs and the advantages over other methods for estimation: essentially GPs can provide competitive predictive power, give error-bars for the estimations, allows to design and optimize sensible kernel functions, and also to analyze the encoded knowledge in the model via automatic relevance determination kernel functions.



**Fig. 7.** Three examples [top, middle, bottom] of CASI RGB snapshots [left], *Chl* estimates [middle], and related uncertainty intervals [right].

GP models offer as well a solid Bayesian framework to formulate new algorithms well-suited to the signal characteristics. We have seen for example that by incorporating proper priors, we can encompass signal-dependent noise, and infer parametric forms of warping the observations as an alternative to either *ad hoc* filtering. On the downside, we need to mention the scalability issue: essentially, the optimization of GP models require computing determinants and invert matrices of size  $n \times n$ , which runs cubically in computational time and quadratically in memory storage. In the last years, however, great advances have appeared in machine learning and now it is possible to train GPs with several thousands of points.

All the developments were illustrated at a local scales through a full set of illustrative examples in the field of geosciences and remote sensing. In particular, we treated important problems of ocean and land sciences: from accurate estimation of oceanic chlorophyll content and pigments, to vegetation properties from multi- and hyperspectral sensors.

## Acknowledgments

The authors wish to deeply acknowledge the collaboration, comments and fruitful discussions with many researchers during the last decade on GP models for remote sens-

ing and geoscience applications: Miguel Lázaro-Gredilla (Vicarious), Robert Jenssen (Univ. Tromsø, Norway), Martin Jung (MPI, Jena, Germany), and Sancho Salcedo-Saez (Univ. Alcalá, Madrid, Spain).

This paper has been partially supported by the Spanish Ministry of Economy and Competitiveness under projects TIN2012-38102-C03-01 and ESP2013-48458-C4-1-P, and by the the European Research Council (ERC) consolidator grant entitled SEDAL with grant agreement 647423. AG is thankful to Marie Curie International Incoming Fellowship for supporting this work.

## References

1. Álvarez, M.A., Luengo, D., Lawrence, N.D.: Linear latent force models using gaussian processes. *IEEE Trans. Pattern Anal. Mach. Intell.* 35(11), 2693–2705 (2013), <http://dx.doi.org/10.1109/TPAMI.2013.86>
2. Baret, F., Weiss, M., Lacaze, R., Camacho, F., Makhmara, H., Pacholczyk, P., Smets, B.: Geov1: LAI and FAPAR essential climate variables and FCOVER global time series capitalizing over existing products. part1: Principles of development and production. *Rem. Sens. Env.* 137(0), 299 – 309 (2013)
3. Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rödenbeck, C., Arain, M.A., Baldocchi, D., Bonan, G.B., Bondeau, A., Cescatti, A., Lasslop, G., Lindroth, A., Lomas, M., Luyssaert, S., Margolis, H., Oleson, K.W., Rouspard, O., Veenendaal, E., Viovy, N., Williams, C., Woodward, F.I., Papale, D.: Terrestrial gross carbon dioxide uptake: Global distribution and covariation with climate. *Science* 329(834) (2010)
4. Camps-Valls, G., Gómez-Chova, L., Muñoz-Marí, J., Vila-Francés, J., Calpe-Maravilla, J.: Composite kernels for hyperspectral image classification. *IEEE Geosci. Remote Sens. Lett.* 3(1), 93–97 (2006)
5. Camps-Valls, G., Gómez-Chova, L., Vila-Francés, J., Amorós-López, J., Muñoz-Marí, J., Calpe-Maravilla, J.: Retrieval of oceanic chlorophyll concentration with relevance vector machines. *Rem. Sens. Env.* 105(1), 23–33 (2006)
6. Camps-Valls, G., Martínez-Ramón, M., Rojo-Álvarez, J.L., Muñoz-Marí, J.: Non-linear system identification with composite relevance vector machines. *IEEE Signal Processing Letters* 14(4), 279–282 (April 2007)
7. Camps-Valls, G., Tuia, D., Gómez-Chova, L., Malo, J. (eds.): *Remote Sensing Image Processing*. Morgan & Claypool (Sept 2011)
8. Donlon, C., Berruti, B., Buongiorno, A., Ferreira, M.H., Féménias, P., Frerick, J., Goryl, P., Klein, U., Laur, H., Mavrocordatos, C., Nieke, J., Rebhan, H., Seitz, B., Stroede, J., Sciarra, R.: The Global Monitoring for Environment and Security (GMES) Sentinel-3 mission. *Remote Sensing of Environment* 120, 37–57 (2012)
9. Dorigo, W.A., Zurita-Milla, R., de Wit, A.J.W., Brazile, J., Singh, R., Schaepman, M.E.: A review on reflective remote sensing and data assimilation techniques for enhanced agroecosystem modeling. *International Journal of Applied Earth Observation and Geoinformation* 9(2), 165–193 (2007)
10. Drusch, M., Del Bello, U., Carlier, S., Colin, O., Fernandez, V., Gascon, F., Hoersch, B., Isola, C., Laberinti, P., Martimort, P., Meygret, A., Spoto, F., Sy, O., Marchese, F., Bargellini, P.: Sentinel-2: ESA’s Optical High-Resolution Mission for GMES Operational Services. *Rem. Sens. Env.* 120, 25–36 (2012)
11. Durbha, S., King, R., Younan, N.: Support vector machines regression for retrieval of leaf area index from multiangle imaging spectroradiometer. *Rem. Sens. Env.* 107(1-2), 348–361 (2007)
12. Guanter, L., Zhang, Y., Jung, M., Joiner, J., Voigt, M., Berry, J.A., Frankenberg, C., Huete, A., Zarco-Tejada, P., Lee, J.E., Moran, M.S., Ponce-Campos, G., Beer, C., Camps-Valls, G., Buchmann, N., Gitanelle, D., Klumpp, K., Cescatti, A., Baker, J.M., Griffis, T.J.: Global and time-resolved monitoring

- of crop photosynthesis with chlorophyll fluorescence. *Proceedings of the National Academy of Sciences, PNAS* (2014)
13. Hastie, T., Tibshirani, R., Friedman, J.H.: *The elements of statistical learning: data mining, inference, and prediction*. Springer-Verlag, New York, USA, 2nd edn. (2009)
  14. Jagermeyr, J., Gerten, D., Lucht, W., Hostert, P., Migliavacca, M., Nemani, R.: A high-resolution approach to estimating ecosystem respiration at continental scales using operational satellite data. *Global Change Biology* 20(4), 1191–1210 (2014), cited By 2
  15. Jung, M., Reichstein, M., Margolis, H.A., Cescatti, A., Richardson, A.D., Arain, M.A., Arneeth, A., Bernhofer, C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law, B.E., Lindroth, A., Merbold, L., Montagnani, L., Moors, E.J., Papale, D., Sottocornola, M., Vaccari, F., Williams, C.: Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible heat derived from eddy covariance, satellite, and meteorological observations. *Journal of Geophysical Research: Biogeosciences* 116(G3), 1–16 (2011)
  16. Lázaro-Gredilla, M., Titsias, M.K.: Variational heteroscedastic gaussian process regression. In: *28th International Conference on Machine Learning, ICML 2011*. pp. 841–848. ACM, Bellevue, WA, USA (2011)
  17. Lázaro-Gredilla, M., Titsias, M.K., Verrelst, J., Camps-Valls, G.: Retrieval of biophysical parameters with heteroscedastic gaussian processes. *IEEE Geosc. Rem. Sens. Lett.* 11(4), 838–842 (2014)
  18. Lázaro-Gredilla, M.: Bayesian warped gaussian processes. In: *NIPS*. pp. 1628–1636 (2012)
  19. Maritorena, S., O’Reilly, J.: *OC2v2: Update on the initial operational SeaWiFS chlorophyll algorithm*, vol. 11, pp. 3–8. John Wiley & Sons, NASA Goddard Space Flight Center, Greenbelt, Maryland, USA (2000)
  20. O’Reilly, J.E., Maritorena, S., Mitchell, B.G., Siegel, D.A., Carder, K., Garver, S.A., Kahru, M., McClain, C.: Ocean color chlorophyll algorithms for SeaWiFS. *Journal of Geophysical Research* 103(C11), 24937–24953 (Oct 1998)
  21. Rasmussen, C.E., Williams, C.K.I.: *Gaussian Processes for Machine Learning*. The MIT Press, New York (2006)
  22. Roelofsen, H., Kooistra, L., Van Bodegom, P., Verrelst, J., Krol, J., Witte, J.c.: Mapping a priori defined plant associations using remotely sensed vegetation characteristics. *Rem. Sens. Env.* 140, 639–651 (2014)
  23. Sampson, P., Guttorp, P.: Nonparametric estimation of nonstationary spatial covariance structure. *Journal of the American Statistical Association Publication* 87(417), 108–119 (Mar 1992)
  24. Sarker, L.R., Nichol, J.E.: Improved forest biomass estimates using ALOS AVNIR-2 texture indices. *Rem. Sens. Env.* 115(4), 968–977 (2011)
  25. Schaepman, M., Ustin, S., Plaza, A., Painter, T., Verrelst, J., Liang, S.: Earth system science related imaging spectroscopy—An assessment. *Rem. Sens. Env.* 113(1), S123–S137 (2009)
  26. Smola, A.J., Schölkopf, B.: A tutorial on support vector regression. *Statistics and Computing* 14, 199–222 (2004)
  27. Snelson, E., Rasmussen, C., Ghahramani, Z.: Warped gaussian processes. In: *Advances in Neural Information Processing Systems, NIPS*. MIT Press (2004)
  28. Tipping, M.E.: *The Relevance Vector Machine*. In: Solla, S.A., Leen, T.K., Müller, K.R. (eds.) *Advances in Neural Information Processing Systems 12*. Cambridge, Mass: MIT Press (2000)
  29. Van Wittenberghe, S., Verrelst, J., Rivera, J., Alonso, L., Moreno, J., Samson, R.: Gaussian processes retrieval of leaf parameters from a multi-species reflectance, absorbance and fluorescence dataset. *Journal of Photochemistry and Photobiology B: Biology* 134, 37–48 (2014)
  30. Verrelst, J., Alonso, L., Camps-Valls, G., Delegido, J., Moreno, J.: Retrieval of vegetation biophysical parameters using Gaussian process techniques. *IEEE Trans. Geosc. Rem. Sens.* 50(5 PART 2), 1832–1843 (2012)
  31. Verrelst, J., Alonso, L., Rivera Caicedo, J., Moreno, J., Camps-Valls, G.: Gaussian process retrieval of chlorophyll content from imaging spectroscopy data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 6(2), 867–874 (2013)

32. Verrelst, J., Muñoz, J., Alonso, L., Delegido, J., Rivera, J., Moreno, J., Camps-Valls, G.: Machine learning regression algorithms for biophysical parameter retrieval: Opportunities for Sentinel-2 and -3. *Rem. Sens. Env.* 118(0), 127–139 (2012)
33. Verrelst, J., Rivera, J., Moreno, J., Camps-Valls, G.: Gaussian processes uncertainty estimates in experimental Sentinel-2 LAI and leaf chlorophyll content retrieval. *ISPRS Journal of Photogrammetry and Remote Sensing* 86, 157–167 (2013)
34. Verrelst, J., Rivera, J., Veroustraete, F., Muñoz Mari, J., Clevers, J., Camps-Valls, G., Moreno, J.: Experimental Sentinel-2 LAI estimation using parametric, non-parametric and physical retrieval methods - A comparison. *ISPRS Journal of Photogrammetry and Remote Sensing* (2015)
35. Yang, F., White, M., Michaelis, A., Ichii, K., Hashimoto, H., Votava, P., Zhu, A.X., Nemani, R.: Prediction of Continental-Scale Evapotranspiration by Combining MODIS and AmeriFlux Data Through Support Vector Machine. *IEEE Trans. Geosc. Rem. Sens.* 44(11), 3452–3461 (nov 2006)