



ALMA MATER STUDIORUM  
UNIVERSITÀ DI BOLOGNA

ARCHIVIO ISTITUZIONALE  
DELLA RICERCA

Alma Mater Studiorum Università di Bologna  
Archivio istituzionale della ricerca

P-spline smoothing for spatial data collected worldwide

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

*Published Version:*

Greco F, V.M. (2018). P-spline smoothing for spatial data collected worldwide. SPATIAL STATISTICS, 27, 1-17 [10.1016/j.spasta.2018.08.008].

*Availability:*

This version is available at: <https://hdl.handle.net/11585/647337> since: 2018-10-21

*Published:*

DOI: <http://doi.org/10.1016/j.spasta.2018.08.008>

*Terms of use:*

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).  
When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

**Greco, F., Ventrucci, M., Castelli, E., 2018. P-spline smoothing for spatial data collected worldwide. *Spat. Stat.* 27, 1–17.**

<https://doi.org/10.1016/j.spasta.2018.08.008>

The final published version is available online at:

<https://doi.org/10.1016/j.spasta.2018.08.008>

Rights / License:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

*This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>)*

***When citing, please refer to the published version.***

# P-spline smoothing for spatial data collected worldwide

Fedele Greco<sup>a</sup>, Massimo Ventrucci<sup>a,\*</sup>, Elisa Castelli<sup>b</sup>

<sup>a</sup>*Department of Statistical Sciences, University of Bologna, via delle Belle Arti n. 41, 40126 Bologna, Italy*

<sup>b</sup>*Institute of Atmospheric Sciences and Climate (CNR-ISAC), Via Piero Gobetti 101, 40129 Bologna, Italy*

---

## Abstract

Spatial data collected worldwide from a huge number of locations is frequently used in environmental and climate studies. Spatial modelling for this type of data presents both methodological and computational challenges. In this work we illustrate a computationally efficient non-parametric framework in order to model and estimate the spatial field while accounting for geodesic distances between locations. The spatial field is modelled via penalized splines (P-splines) using intrinsic Gaussian Markov Random Field (GMRF) priors for the spline coefficients. The key idea is to use the sphere as a surrogate for the Globe, then build the basis of B-spline functions on a geodesic grid system. The basis matrix is sparse as is the precision matrix of the GMRF prior, thus computational efficiency is gained by construction. We illustrate the approach with a real climate study, where the goal is to identify the Intertropical Convergence Zone using high-resolution remote sensing data.

*Keywords:* smoothing, intrinsic Gaussian Markov Random field, P-spline, geodesic, ITCZ

---

## 1. Introduction

High-resolution spatial data collected worldwide, usually by means of remote sensing techniques, is wide-spread in environmental and climate studies: most of the statistical methods developed in modelling this kind of data use the sphere as a surrogate for the Globe. Modelling data collected at a global scale presents both methodological and computational challenges. The traditional toolkit for a spatial data modeller when dealing with geostatistical datasets and aiming to make predictions at unmonitored locations would suggest to apply kriging techniques (see, e.g., Banerjee et al. (2014)). These rely on the assumption of a smooth Gaussian Random Field (GRF), continuous

---

\*Corresponding author: massimo.ventrucci@unibo.it

10 in space but only observed at a discrete set of points, any finite realization of it be-  
11 ing generated by a multivariate Gaussian distribution. The covariance structure of this  
12 distribution is specified via a spatial covariance function. The practice is largely domi-  
13 nated by spatial covariances defined on Euclidean distances, such as the Matérn family,  
14 thus a preliminary step in the analysis is the projection of the 3d Cartesian coordinates  
15 (from the Earth’s surface) over a 2d coordinate space. The standard choice is to work  
16 with geographic coordinates (latitude-longitude), but other mapping methods can be  
17 used. Banerjee (2005) provides a review of such mapping techniques and discusses the  
18 impact of the chosen metric on spatial prediction via kriging. The traditional toolkit  
19 outlined above presents two main difficulties when modelling high-resolution data ob-  
20 served over a spherical domain.

21 The first issue is that the process of spatial prediction needs to be coherent with  
22 the geometry of the sphere. Using a planar metric over a 2d projection is inappro-  
23 priate because it generates spurious anisotropy and non-stationarity of the covariance  
24 function (Banerjee (2005)). The geodesic (aka great circle) distance, i.e. the length  
25 of the shortest path between two points over the surface of a sphere, is a natural can-  
26 didate for measuring distances over a spherical domain. However, using great circle  
27 distances in a Matérn family does not necessarily guarantee a positive definite covari-  
28 ance (Gneiting, 2013). Banerjee (2005) used a simulation to study the behaviour of  
29 different metrics regarding estimation of the Matérn covariance parameters on a region  
30 as large as Colorado, finding a substantial impact of the chosen metric on the range of  
31 the correlation function. This means that with data collected on larger regions on Earth  
32 (e.g. the whole Globe), a biased estimation of the underlying field is to be expected to  
33 some extent, when covariance functions built on Euclidean distances are used. A large  
34 number of papers have tackled this issue by essentially proposing new models for data  
35 on a spherical domain, both in a parametric and non-parametric framework.

36 In the parametric setting, several papers focused on building valid stochastic pro-  
37 cesses for the sphere, see, e.g., Jun and Stein (2007); Jeong and Jun (2015); Heaton  
38 et al. (2014); Porcu et al. (2016) and references therein. The stochastic partial differen-  
39 tial equation (SPDE) approach by Lindgren et al. (2011) has gained a lot of attention  
40 recently. This approach builds a GRF with Matérn covariance as the finite element  
41 solution of a particular SPDE, an idea that can be generalized for different types of  
42 manifolds including the sphere.

43 In the non-parametric setting, Wahba (1981) was first to introduce smoothing splines  
44 onto the sphere, while analysing weather data collected from a large number of stations  
45 around the world. Outside the spline realm, Di Marzio et al. (2014) presented local lin-  
46 ear regression for spherical data, including the case of smoothing of a scalar response

47 on a spherical predictor. Wood (2017) discusses in detail the connection between spline  
48 smoothing and thin plate splines for a sphere, pointing out that low rank smoothers are  
49 also applicable to spherical data. Although low rank smoothers allows for a reduction  
50 in the number of parameters to estimate, the main role in alleviating the computational  
51 burden is played by the sparsity of the smoothing matrix, obtained by using local ba-  
52 sis functions, i.e. non-null over a limited domain. B-splines are local functions built  
53 upon joint polynomials connected at knots, which are applied in different contexts,  
54 such as in penalized spline (P-spline) regression (Eilers and Marx, 1996). With spa-  
55 tial data, bivariate B-splines over triangulations (Lai and Schumaker, 2007) provide  
56 a basis for piecewise polynomial surfaces and are used in spatial models (Lai et al.,  
57 2009; Baramidze et al., 2006). Finite Elements provide an alternative basis for piece-  
58 wise polynomial surfaces over triangulations (Sangalli et al., 2013). Also, more recent  
59 proposals deal with data distributed on two-dimensional general domains using finite  
60 elements (Duchamp and Stuetzle, 2003; Ettinger et al., 2016) or non-rational B-splines  
61 basis (Wilhelm et al., 2016).

62 The second difficulty concerning the application of kriging techniques to high-  
63 resolution global datasets is purely computational. Continuous covariance functions  
64 used in geostatistics involve a dense covariance structure for the underlying GRF. When  
65 the number of data locations  $n$  is large, this modelling framework becomes impractical  
66 because of the need to invert large dense matrices, with a computational cost increasing  
67 by cubic growth with  $n$ . Statistics literature on the *big n problem* has boomed in the last  
68 decade, mostly due to the increasing availability of high resolution remote sensing data  
69 for environmental studies. Some of the models for large data that can be implemented  
70 in a fully Bayesian hierarchical setting (for a review see Banerjee (2017)) are based  
71 on a low-rank representation of the field (Wikle and Cressie, 1999; Banerjee et al.,  
72 2008). Other proposals look to find a sparse representation of the covariance, like in  
73 tapering (Furrer et al., 2006), or of the precision matrix (Rue and Held, 2005). In this  
74 framework, the paper by Lindgren et al. (2011) derives an approximated solution to  
75 the SPDE in terms of a Gaussian Markov Random Field (GMRF), instead of a GRF,  
76 in order to gain computational efficiency. A recent approach that allows to deal with  
77 GRFs in a computationally efficient manner is Datta et al. (2016), where sparsity is  
78 introduced without the need for dimension reduction. The fully Bayesian framework  
79 presented in this paper follows both directions, in the sense that it is built on a low-  
80 rank representation using local B-splines and exploits the sparsity induced by a GMRF  
81 prior.

82 We propose a computationally efficient non-parametric approach to estimate the  
83 spatial field underlying data on the sphere that properly accounts for geodesic dis-

84 tances between locations. Our method is based on a low-rank P-spline smoother to  
85 gain flexibility w.r.t. parametric models. The main contribution of this work is the  
86 extension of the P-spline model for smoothing data collected over a spherical domain.  
87 The model is built on a set of bivariate B-splines computed on a Geodesic Discrete  
88 Global Grid (GDGG) system (Sahr et al., 2003), yielding a quasi-regular triangular  
89 mesh over the Globe. Geodesic grids have been used in spatial statistics to create flex-  
90 ible multi-resolution models implemented in a likelihood-based inferential framework  
91 (Cressie and Johannesson, 2008; Nychka et al., 2015). In contrast to the latter works,  
92 in this paper we follow a fully Bayesian approach and fit the model using an efficient  
93 Gibbs sampler, exploiting sparsity of the basis matrix and of the precision of the GMRF  
94 prior. We illustrate the method on a real climate study, where the goal is to identify the  
95 *Intertropical Convergence Zone* (ITCZ) from high-resolution remote sensing data col-  
96 lected worldwide over sea, with missing data occurring over land.

97 The rest of the paper is organized as follows. Section 2 describes the dataset and  
98 application goals. Section 3 presents our proposal for smoothing data over the sphere  
99 that we dub *Geodesic P-splines*. Section 4 illustrates the method used on a climate-  
100 related case study, focusing on the detection of the ITCZ. A discussion is provided in  
101 Section 5.

## 102 2. Motivating example

103 Our interest in geodesic P-splines is motivated by a climate-related case study  
104 aimed at investigating the location of the ITCZ using satellite data. The ITCZ is a  
105 region of the atmosphere broadly located within the tropical belt where the north-east  
106 and south-east trade winds converge, an area characterised by high cloudiness and se-  
107 vere convective precipitation (Holton and Hakim, 2013). An important aspect regards  
108 seasonal variability in the ITCZ position: the ITCZ is located roughly North of the  
109 equator in the boreal Spring and Summer, while it migrates to southern regions in Au-  
110 tumn and Winter. The location of the ITCZ affects duration and intensity of the wet  
111 and dry seasons at the tropics and plays a key role in the general circulation of the  
112 atmosphere: assessing its variability is crucial for improving global climate models.  
113 Moreover, understanding the long-term trend characterizing this phenomenon is cru-  
114 cial for monitoring changes in climate patterns on a global scale.

115 The phenomenon regulating the ITCZ behaviour cannot be measured directly, hence  
116 several studies have investigated it using some suitable proxy variables, like maximum  
117 precipitation (Zhang, 2001), wind field (Žagar et al., 2011), vorticity and reflectivity of  
118 the clouds (Waliser and Gautier, 1993). As a general feature, all these studies benefit

119 from the increasing availability of satellite measurements. In this paper we focus on  
120 data from the infrared channels of the Along Track Scanning Radiometer (ATSR) series  
121 of instruments, which were in orbit from 1991 to 2012 for the accurate retrieval of sea  
122 surface temperature. Recently, in the frame of the European Space Agency ATSR Long  
123 Term Stability project (<https://earth.esa.int/web/sppa/activities/multi-sensors-timeseries/alts/about>),  
124 Casadio et al. (2016) developed the Advanced Infra-Red Water Vapour Estimator al-  
125 gorithm (AIRWAVE) for the retrieval of the Total Column of Water Vapour (TCWV)  
126 from the ATSR measurements. In this work we use TCWV as a proxy variable for  
127 locating ITCZ.

128 Data on TCWV regarding year 2008 was provided by the National Research Coun-  
129 cil - Institute of Atmospheric Sciences and Climate (CNR-ISAC), Italy. Data comes  
130 as monthly averages of TCWV in a raster of dimension 360 columns (longitude val-  
131 ues) by 180 rows (latitude values), thus each cell covers one degree over latitude and  
132 longitude. In Figure 1, the data for January and July is displayed. The AIRWAVE  
133 algorithm provides reliable data over the sea and in clear sky conditions, thus TCWV  
134 observations over land are missing (roughly a third of the total number of cells), except  
135 in areas covered by lakes. The percentage of raster cells with missing observations is  
136 about 40%.

137 The application goal is to estimate the ITCZ position and its uncertainty. We con-  
138 sider the TCWV data on the fine raster grid as point-level data observed at the centroid  
139 of each cell. We are actually managing raster data as point data. This is a standard pro-  
140 cedure when modelling high resolution raster data, particularly when adopting splines  
141 that need to be evaluated at fixed points. The same rationale has been adopted in Eilers  
142 et al. (2006); Lee and Durbán (2009); Ugarte et al. (2012). We focus on modelling the  
143 latent field of TCWV separately for each month, deferring spatio-temporal modelling  
144 to future work. The statistical challenges we tackle in this paper are related to efficient  
145 smoothing of large data to remove measurement error and to allow for rapid predic-  
146 tions at unmonitored locations. We believe that the extension of Bayesian P-Splines  
147 to a spherical domain can be a valuable strategy because of its efficiency and compu-  
148 tational stability. Bayesian inference provides immediate tools for ITCZ location, by  
149 analysing the joint posterior distribution of the latent field. One issue concerning ITCZ  
150 detection is that there is no definition in terms of a fixed threshold. This situation calls  
151 for methods to search for peaks in the latent field, bearing in mind that the ITCZ is  
152 expected to be located at the Equator.

153 In Section 4, the ITCZ detection problem is addressed by searching for the latitudes  
154 where the TCWV latent field shows the highest values. We provide a graphical output,  
155 by plotting the posterior probability that a point on the Earth belongs to the ITCZ.

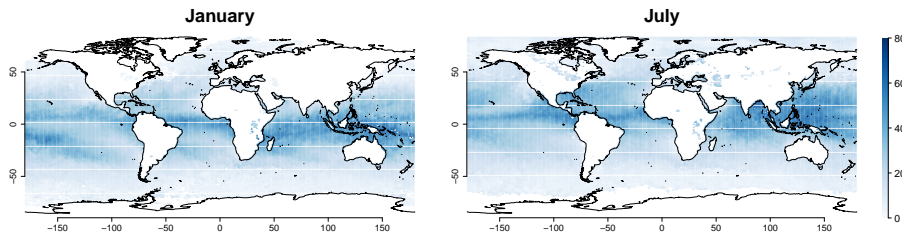


Figure 1: TCWV data for January and July (unit measure,  $Kg/m^2$ ). In general, TCWV measurements are available only over sea, as the data cannot be accurately retrieved over land; however, note that observations are still present in correspondence of wide lakes, e.g. the Great Lakes of North America and the Victoria lake.

### 156 3. Geodesic P-splines

#### 157 3.1. Background on P-splines for spatial data

158 In the one-dimensional setting, P-splines (Eilers and Marx, 1996) are usually adopted  
 159 to model the smooth effect of a covariate on the response as a linear combination of  
 160 B-splines scaled by spline coefficients. Key features of this method are (a) equally-  
 161 spaced univariate B-splines of a certain degree  $d$ , these being non zero over a limited  
 162 interval of the covariate domain, and (b) a penalty on the  $r^{th}$  order differences between  
 163 adjacent spline coefficients to control smoothness. The popularity of P-splines is due  
 164 to numerical stability and flexibility in the modelling choices; e.g., the penalty order  
 165 and the degree of the B-splines can be decided according to the application at hand.  
 166 Higher-dimensional smoothers, suitable for modelling spatial data, can be constructed  
 167 as tensor product P-splines (Eilers et al., 2006). In a frequentist framework, estimation  
 168 is obtained via penalized maximum likelihood or iterative re-weighted least squares,  
 169 with the smoothing parameter selected via cross validation or optimized over some  
 170 information criterion. This method has become increasingly popular and is currently  
 171 implemented in R packages such as `mgcv` (Wood, 2017).

In order to build the ground for our proposal we next revise spatial P-splines for data observed on a two-dimensional latitude-longitude plane following Eilers et al. (2006). Let us assume  $y_i$  is a Gaussian observation at location  $(lat_i, lon_i)$ ,  $i = 1, \dots, n$ , the model is

$$y_i = \mu(lat_i, lon_i) + \epsilon_i \quad ; \quad \epsilon_i \sim \mathcal{N}(0, \tau_\epsilon^{-1}),$$

where  $\mu(lat_i, lon_i)$  is a two-dimensional function, with no parametric assumptions on it and  $\tau_\epsilon$  is the noise precision. We can think of  $\mu(lat_i, lon_i)$  as a smooth surface representing the latent field which is modelled as a linear combination of bivariate B-spline



basis functions:

$$\mu(\text{lat}_i, \text{lon}_i) = \sum_{q=1}^Q \sum_{l=1}^L b_l(\text{lat}_i) b_q(\text{lon}_i) \beta_{l,q},$$

172 where  $b_l(\text{lat}_i) b_q(\text{lon}_i)$  is the tensor product of marginal B-splines, evaluated at  $(\text{lat}_i, \text{lon}_i)$ ,  
 173 and  $\beta_{l,q}$  is the associated spline coefficient. The marginal B-splines  $b_l, l = 1, \dots, L$   
 174 ( $b_q, q = 1, \dots, Q$ ), are defined on a set of knots that are chosen to be equally-spaced  
 175 over the latitude (longitude) domain. Taking the tensor product of the marginal bases  
 176 comes to  $K = QL$  bivariate B-splines built on a regular grid over the plane; see Figure  
 177 2, left panel. In this sense, P-splines give a low-rank representation of the latent field,  
 178 as  $K$  is typically chosen to be much lower than  $n$ . In matrix notation,  $\boldsymbol{\mu} = \mathbf{B}\boldsymbol{\beta}$ , where  $\mathbf{B}$   
 179 is a basis matrix of dimension  $n \times K$  and  $\boldsymbol{\beta}$  the vector of spline coefficients. When data  
 180 is organized in a regular grid with no missing values, the basis matrix can be computed  
 181 by the Kronecker product  $\mathbf{B} = \mathbf{B}_{\text{lat}} \otimes \mathbf{B}_{\text{lon}}$ . When data are irregularly scattered over  
 182 the plane, efficient row-wise Kronecker operations can still be used to compute  $\mathbf{B}$ , as  
 183 this is equivalent to having data organized on a fine regular grid with missing values.  
 184 We suggest the reader see Eilers et al. (2006) for details on P-splines for spatial data  
 185 and to Lee (2010) for insights into the mixed model formulation of P-splines within a  
 186 spatio-temporal setting.

187 P-splines have been framed in a fully hierarchical Bayesian context by Lang and  
 188 Brezger (2004). The hierarchical model can be cast starting from the following likeli-  
 189 hood:

$$y|\alpha, \boldsymbol{\beta}, \tau_\epsilon \sim \mathcal{N}(\boldsymbol{\mu}, \tau_\epsilon^{-1} \mathbf{I}) \quad ; \quad \boldsymbol{\mu} = \alpha \mathbf{1} + \mathbf{B}\boldsymbol{\beta}$$

The penalty is reproduced by an  $r^{\text{th}}$  order random walk (RW) prior on the spline coef-  
 ficients, that in general can be expressed as

$$\pi(\boldsymbol{\beta}|\tau_\beta) = (2\pi)^{-r \text{rank}(\mathbf{R})/2} (|\tau_\beta \mathbf{R}|^*)^{1/2} \exp\left\{-\frac{\tau_\beta}{2} \boldsymbol{\beta}^\top \mathbf{R} \boldsymbol{\beta}\right\}, \quad (1)$$

190 where  $\tau_\beta$  is a scalar precision hyper-parameter and  $\mathbf{R}$  is the structure matrix of dimen-  
 191 sion  $K \times K$ . The non-zero entries in  $\mathbf{R}$  impose conditional dependencies among the  
 192 spline coefficients, thus encoding the type of penalty. Formally, the RW is a particular  
 193 type of Intrinsic Gaussian Markov Random Field (IGMRF). The smoothing properties  
 194 of an IGMRF are determined by the pattern of non-zero entries of  $\mathbf{R}$  and by its rank  
 195 deficiency. Any vector in the null space of  $\mathbf{R}$  can be added to  $\boldsymbol{\beta}$  and density (1) remains  
 196 unchanged. For this reason, IGMRF priors are appropriate to model local deviations  
 197 around an overall mean or, in general, a polynomial trend, with  $\tau_\beta$  controlling the size

198 of such deviations. For spatial smoothing, we will focus on a prior that leaves the  
 199 overall mean unspecified, therefore  $\text{rank}(\mathbf{R}) = K - 1$ .

The precision matrix for P-spline smoothing over a plane proposed in Eilers et al. (2006) is constructed as the Kronecker sum

$$\mathbf{R} = (\mathbf{I}_L \otimes \mathbf{R}_{lon}) + (\mathbf{R}_{lat} \otimes \mathbf{I}_Q) \quad (2)$$

where  $\mathbf{R}_{lat}$  and  $\mathbf{R}_{lon}$  are the (marginal) structure matrices of a RW on latitudinal and longitudinal knots, respectively. If we take  $\mathbf{R}_{lat}$  and  $\mathbf{R}_{lon}$  as the structure of a 1<sup>st</sup> order RW, this is equivalent to assuming an intrinsic Conditional Autoregressive (ICAR) model (Besag, 1974), with structure

$$R_{ij} = \begin{cases} k_i & i = j \\ -1 & i \sim j \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

200 where  $k_i$  is the number of knots adjacent to the  $i^{\text{th}}$  knot; e.g.  $k_i = \{2, 3, 4\}$  according to  
 201 whether  $i$  is a knot on the vertex, the border, or the interior of the regular grid. Usually  
 202 an ICAR prior is assumed on a set of  $n$  random effects, one for each data location,  
 203 but here the ICAR is on the spline coefficients. In this sense, the basis  $\mathbf{B}$  allows the  
 204 stochastic field on the  $K$  spline coefficients to be expanded at a much larger number  
 205 of locations like  $n$ . This strategy allows for a substantial reduction in the number of  
 206 parameters to estimate. Choosing higher order random walks in each dimension is  
 207 possible: this will yield an higher-order IGMRF prior, having a structure matrix with  
 208 larger rank-deficiency; e.g. taking a 2<sup>nd</sup> order RW on latitude and longitude comes to  
 209 an IGMRF that models deviations from a plane. For a discussion of the properties of  
 210 IGMRFs and their applications see Rue and Held (2005).

### 211 3.2. P-splines on Geodesic Discrete Global Grid Systems

212 The assumption of equally-spaced knots is convenient for building Bayesian penal-  
 213 ized spline models, because it allows to create a suitable smoothing prior by simply  
 214 using an IGMRF model for regularly spaced locations on the spline coefficients. Fol-  
 215 lowing this idea, knot placement must take into account the geometry of the data's  
 216 support. Thus, building an equally spaced basis on the latitude-longitude plane is not  
 217 a sensible choice when the data covers the whole Globe or a large region thereof. Fig-  
 218 ure 2 highlights that equally spaced B-splines in terms of Euclidean distances over  
 219 the latitude-longitude plane (left panel) are not equally-spaced over the sphere (right

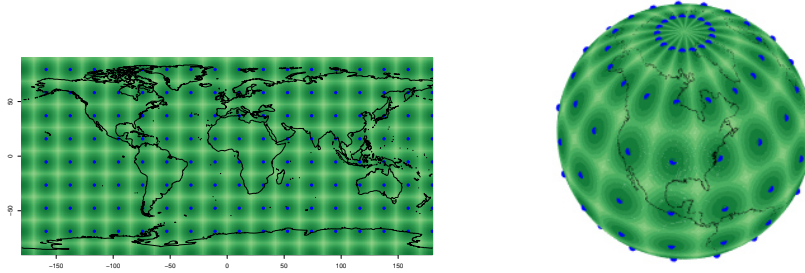


Figure 2: Cubic B-splines equally-spaced in terms of Euclidean distances over the latitude and longitude plane (left panel; computed as the tensor product of marginal B-spline basis, see Section 3.1). The right panel displays how these bases appear on the sphere.

220 panel). The spacing between the knots and the shape of the basis varies substan-  
 221 tially latitude-wise: in such a knot-grid, imposing an IGMRF with structure (3) and  
 222 a single precision parameter  $\tau_\beta$  on the spline coefficients would generate the spurious  
 223 anisotropy discussed in Banerjee (2005). Of course, this would be a naive approach to  
 224 spatial smoothing over the sphere, since it does not introduce conditional dependence  
 225 between knots located at extreme longitudes, which are actually close on the sphere  
 226 surface. A circular penalty imposing conditional correlations among these knots seems  
 227 a more sensible choice, however the irregular knot distribution over the sphere would  
 228 still generate spurious non-stationarity, as this paper will discuss at a later moment. In  
 229 what follows we propose an approach for (a) building geodesic knot-grids which are  
 230 almost equally spaced in terms of geodesic distances, (b) building basis functions and  
 231 penalty matrices on such grids.

### 232 3.2.1. Building the geodesic grid

233 Although building *exactly* equally spaced grids over the sphere surface is an impos-  
 234 sible task, GDGGs offer a close approximation to equal spacing and their architecture  
 235 provides immediate solutions to build basis functions and penalty matrices. Details on  
 236 the spatial configuration of GDGGs can be found in Randall et al. (2002). Sahr et al.  
 237 (2003) outline five design choices that need to be undertaken for GDGGs construction:  
 238 our choices are listed below.

- 239 1. Choice of a *base regular polyhedron*: we choose the icosahedron, which is a  
 240 polyhedron made of 20 equilateral triangles and 12 nodes and consider this as a  
 241 rough representation of a unit sphere. An icosahedron is displayed in Figure 3,  
 242 left panel.

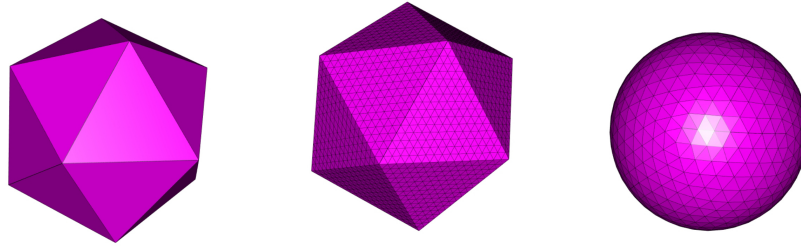


Figure 3: On the left panel, the icosahedron. On the central panel, the icomesh, i.e. the regular triangular mesh after the *split* operation is repeated four times ( $\nu = 4$ ). On the right panel, the icosphere, i.e. the mesh obtained from normalizing the icomesh nodes of the central panel.

- 243 2. Choice of a fixed *orientation* for the base regular polyhedron relative to the Earth:  
 244 we set one node of the icosahedron at coordinates  $(0, 0, 1)$ , assuming this to be  
 245 the North Pole.
- 246 3. Choice of a *hierarchical spatial partitioning* method defined symmetrically on  
 247 each face of the base regular polyhedron. At this step, we split each triangle of  
 248 the icosahedron in four equal triangles. By repeating this operation an arbitrary  
 249 number of times we obtain a refined mesh, which we denote as *icomesh*. In  
 250 Figure 3, central panel, the reader can see the icomesh resulting from four split  
 251 iterations.
- 252 4. Transforming the base polyhedron partition into the corresponding spherical sur-  
 253 face. This is achieved by simply normalizing the icomesh nodes, so that they lay  
 254 on the sphere; we denote this mesh as *icosphere*, see Figure 3, right panel. The  
 255 icosphere is a refined icosahedron, hence a much better representation of the  
 256 sphere.
- 257 5. Choice of a method to *assign points to grid cells*. The ability to assign points  
 258 to the grid cells composing the tessellation can be useful for several purposes.  
 259 In our case, it is fundamental for determining which triangle each data location  
 260 falls into when it comes to the computation of the basis functions, as discussed  
 261 in the following section.

262 Following the above five steps, we obtain a geodesic grid of knots which are almost  
 263 equidistant in terms of great-circle distances. To summarize, the GDGG is constructed  
 264 by splitting each icosahedron face into four triangles, in a recursive way. Note that,  
 265 while the icosphere is a sphere tessellated into spherical triangles, the icomesh is a  
 266 regular mesh made by equilateral triangles.

267 *3.2.2. Building the basis and the penalty matrix*

268 The number of split iterations determines the dimension of the basis, i.e. the num-  
 269 ber of columns of the basis matrix  $\mathbf{B}$ . Let  $n$  be the number of observations and  $\nu$  be  
 270 the number of split iterations, the basis has dimension  $n \times K$ , with  $K = 5 \cdot 2^{(2\nu+1)} + 2$   
 271 (Randall et al., 2002). For  $\nu = 0$  we have  $K = 12$ , which is the number of vertices  
 272 of the icosahedron, by increasing  $\nu$  we obtain an icomesh with higher resolution. We  
 273 adopt B-spline basis functions centred at the knots, each basis spanning six triangles  
 274 (thus assuming the six closest nodes as neighbours) except for those centred at the 12  
 275 icosahedron vertices (that have five neighbours).

276 The next step consists of evaluating the  $K$  B-splines, of a certain degree  $d$ , at an  
 277 arbitrary data point on the sphere. Once that the triangle containing such a point is  
 278 determined, B-splines can be evaluated using Bernstein polynomials (Lai and Schu-  
 279 maker, 2007). To this aim, we find it convenient to work on the icomesh instead of  
 280 the icosphere, as it is simpler to deal with planar than with spherical triangles. There-  
 281 fore, we first project the 3d data location from the icosphere onto the icomesh do-  
 282 main, obtaining a point,  $\mathbf{v}$ , which falls inside a planar triangle (that lies on one of the  
 283 icosahedron faces) and, second, we evaluate the  $K$  B-splines at this 2d point. Fol-  
 284 lowing Lai and Schumaker (2007), any point  $\mathbf{v} = (x, y)$  inside a triangle of vertices  
 285  $\mathbf{v}_1 = (x_1, y_1)$ ,  $\mathbf{v}_2 = (x_2, y_2)$ ,  $\mathbf{v}_3 = (x_3, y_3)$  has a unique representation as

$$\mathbf{v} = \mathbf{v}_1 b_1 + \mathbf{v}_2 b_2 + \mathbf{v}_3 b_3,$$

where  $(b_1, b_2, b_3)$  are called barycentric coordinates and are such that  $b_1 + b_2 + b_3 = 1$ .  
 The Bernstein polynomial of degree  $d$  is

$$H_{tjk}^d = \frac{d!}{t!j!k!} b_1^t b_2^j b_3^k \quad (4)$$

with  $t, j, k$  integer numbers summing to  $d$ . The following property

$$\sum_{t+j+k=d} H_{tjk}^d = 1$$

286 guarantees that for each location on the sphere the basis functions add up to 1. This is  
 287 a desirable property for any smoothing model, giving a flat spatial field when there is  
 288 no variation around the overall level, i.e. all spline coefficients are equal.

289 Let  $\mathbf{z}_i = (z_{i1}, z_{i2})$  denote the location for observation  $i$  projected on the icomesh,  
 290  $\mathbf{B}[i, ]$  the row entry of  $\mathbf{B}$  with the B-splines evaluated at  $\mathbf{z}_i$  and  $\{k_1, k_2, k_3\}$  the indices  
 291 for the three knots closest to  $\mathbf{z}_i$  (note that these are the vertices of the triangle containing

$d = 1$	$d = 2$	$d = 3$
$\begin{pmatrix} H_{100}^1 \\ H_{010}^1 \\ H_{001}^1 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$	$\begin{pmatrix} H_{200}^2 + H_{100}^2 \\ H_{020}^2 + H_{010}^2 \\ H_{002}^2 + H_{001}^2 \end{pmatrix} = \begin{pmatrix} b_1^2 + b_1 \\ b_2^2 + b_2 \\ b_3^2 + b_3 \end{pmatrix}$	$\begin{pmatrix} H_{300}^3 + H_{200}^3 + H_{100}^3 \\ H_{030}^3 + H_{020}^3 + H_{010}^3 \\ H_{003}^3 + H_{002}^3 + H_{001}^3 \end{pmatrix} = \begin{pmatrix} b_1^3 + b_1^2 + b_1 \\ b_2^3 + b_2^2 + b_2 \\ b_3^3 + b_3^2 + b_3 \end{pmatrix}$

Table 1: Non-zero elements of  $\mathbf{B}[i, \cdot]$ , for B-splines of degree  $d = \{1, 2, 3\}$ .

292 observation  $i$ ). It is important to note that only the B-splines centred at  $\{k_1, k_2, k_3\}$  are  
293 non-zero at  $\mathbf{z}_i$ , whereas the B-splines centred at the remaining knots in the icomesh  
294 are zero at  $\mathbf{z}_i$ . The three non-zero elements of  $\mathbf{B}[i, \{k_1, k_2, k_3\}]$  can be expressed as  
295 Bernstein polynomials (4), i.e. polynomials in the barycentric coordinates. Table 1  
296 reports the non zero elements of  $\mathbf{B}[i, \cdot]$  for linear ( $d = 1$ ), quadratic ( $d = 2$ ) and cubic  
297 ( $d = 3$ ) B-splines.

298 The resulting basis matrix  $\mathbf{B}$  is sparse because the B-splines are non-zero over a  
299 domain spanning over only six triangles on the icomesh. Figure 4, left panel, shows  
300 how the new basis functions appear when projected over latitude and longitude. This  
301 plot suggests that a fairly similar degree of smoothness is applied everywhere using this  
302 new basis, avoiding the kind of spurious anisotropy introduced by the basis in Figure  
303 2. The Geodesic P-splines setting is completed by specifying the matrix  $\mathbf{R}$ , that we  
304 choose as the ICAR structure (3) with rank-deficiency 1. The number of neighbouring  
305 knots is  $k_i = 5$ , if  $i$  is one of the 12 nodes of the icosahedron, and  $k_i = 6$ , if  $i$  is one of  
306 the remaining  $K - 12$  nodes.

### 307 3.2.3. Model properties

308 When using IGMRF priors with precision matrix  $\tau_\beta \mathbf{R}$  on the spline coefficients  $\boldsymbol{\beta}$ ,  
309 the structure of conditional dependence imposed by  $\mathbf{R}$  determines the structure of the  
310 marginal variances of each coefficient,  $\text{Var}(\beta_i) = \tau_\beta^{-1} R_{ii}^-$ ,  $i = 1, \dots, K$ ,  $\mathbf{R}^-$  being the  
311 generalised inverse of  $\mathbf{R}$ . Different structures can lead to extremely different marginal  
312 variances. To overcome this problem, Sørbye and Rue (2014) suggest scaling the pre-  
313 cision matrix so that the hyperprior for  $\tau_\beta$  can be selected to give the same degree of  
314 smoothness, a priori, starting from different structure matrices. The scaled precision  
315 matrix can be obtained as  $\mathbf{R}^* = \kappa \mathbf{R}$ , where  $\kappa$  is the geometric mean of the diagonal en-  
316 tries of  $\mathbf{R}^-$ . IGMRFs with scaled precision matrices, although being characterised by a  
317 different correlation structure, have a common feature: the average marginal variance  
318 is equal to one.

319 Figure 5 compares the marginal variances for three models corresponding to a naive  
320 penalty (top-left), longitude-wise circular penalty (bottom) and a geodesic penalty (top-  
321 right). For the sake of comparison, the precision matrices associated with the three

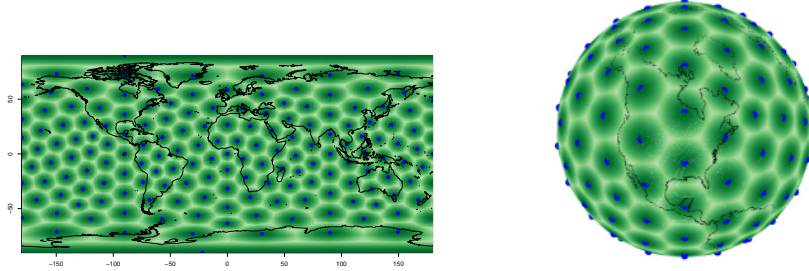


Figure 4: Cubic B-splines equally-spaced in terms of geodesic distances over the sphere (right panel; computed using Bernstein polynomials on a GDGG, see Section 3.2). The left panel displays how these bases appear on the latitude longitude plane.

322 models were scaled. For naive penalty, we mean an IGMRF prior for the spline co-  
 323 efficients laying on a planar grid, using the ICAR structure (3). The longitude-wise  
 324 circular penalty is an IGMRF on a planar grid with structure (2), but assuming  $\mathbf{R}_{lon}$  as  
 325 the structure of a circular 1<sup>st</sup> order RW. For geodesic penalty, we mean an IGMRF on  
 326 a GDGG using the ICAR structure as described in Section 3.2.2. In the left panel, the  
 327 non-stationarity in the marginal variances implied by using the ICAR structure on a  
 328 regular planar grid (naive penalty) is evident. In the bottom panel, marginal variances  
 329 obtained by building a circular penalty longitude-wise show a variation latitude-wise  
 330 as expected. The IGMRF prior on the geodesic grid with the ICAR structure implies  
 331 stability in the marginal variances that is not achieved with the other specifications. As  
 332 a matter of fact, the geodesic grid is almost a torus since all knots, except the icosaha-  
 333 dron nodes, have six neighbours; we believe this is a desirable feature of our model as  
 334 it mimics the idea of second-order stationarity typical of Matérn correlation functions.

### 325 3.2.4. Hyperpriors

336 To complete the fully Bayesian model we need to set priors for the hyper-parameters  
 337  $\tau_\beta$  and  $\tau_\epsilon$ . The precision  $\tau_\beta$  regulates the amount of smoothing. When  $\tau_\beta$  goes to in-  
 338 finity,  $\mu$  is a constant (because the rank deficiency of  $\mathbf{R}$  is 1), while  $\tau_\beta \in (0, +\infty)$  gives  
 339 a more flexible surface. A standard approach is to use a Gamma,  $\text{Ga}(a, b)$ , with shape  
 340  $a$  and rate  $b$ , for both random walk and noise precisions. Usual parametrizations are  $a$   
 341 equal to 1 and  $b$  small (e.g.  $\text{Ga}(1, 5e-5)$ ), or  $a$  and  $b$  small (e.g.  $\text{Gamma}(1e-3, 1e-3)$ ),  
 342 as an attempt of non informativeness on the variance scale. Several papers in the lit-  
 343 erature have discussed issues related to the Gamma conjugate priors in hierarchical  
 344 additive models and proposed alternatives (Gelman, 2006; Simpson et al., 2017). Typ-  
 345 ically, the main impact regards the prior for the random walk precision, whereas the

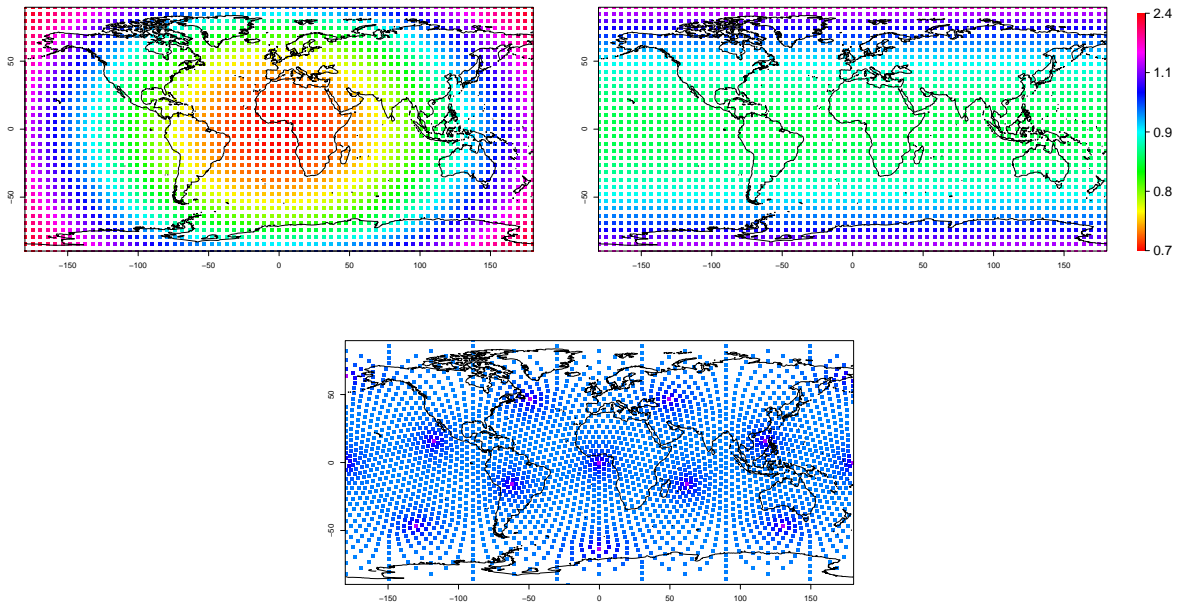


Figure 5: Marginal variances with naive penalty (top left panel), longitude-wise circular penalty (top right panel) and geodesic penalty (bottom panel). (The colour bar on the right is valid for all the three panels).



346 prior for the noise precision is negligible. In general, choice about the prior  $\pi(\tau_\beta)$  will  
 347 be relevant in situations where we have a poor sample size compared to the number  
 348 of parameters which require to be estimated. In the case study under examination the  
 349 large sample size available for estimating each spline coefficient makes the impact of  
 350  $\pi(\tau_\beta)$  negligible.

### 351 3.2.5. Computations

Model estimation does not raise particular issues with respect to planar P-spline  
 models, once matrices  $\mathbf{B}$  and  $\mathbf{R}$  have been built. Indeed, the model belongs to the class  
 of Latent Gaussian Markov Models and approximate Bayesian inference can be per-  
 formed efficiently using the R-INLA package (Rue et al., 2009). In our case study, we  
 find it more appropriate to use a Gibbs sampling algorithm as the tools developed for  
 detecting the ITCZ (see section 4.3) require a sample from the joint posterior distribu-  
 tion of the model. The most expensive step is to sample from the full conditional for  
 the spline coefficients

$$\beta|\tau_\beta, \tau_\epsilon, \mathbf{y} \sim N(\mathbf{Q}^{-1}\mathbf{B}^\top \mathbf{y}, \mathbf{Q}^{-1}) \quad \mathbf{Q} = \left( \mathbf{B}^\top \mathbf{B} + \frac{\tau_\beta}{\tau_\epsilon} \mathbf{R} \right) \quad (5)$$

352 under the linear constraint  $\mathbf{1}_K^\top \mathbf{B} \beta = 0$  needed for intercept identifiability. We use an ef-  
 353 ficient Gibbs sampler coded in R with the use of sparse matrix algebra as implemented  
 354 in the spam package (Furrer and Sain, 2010) to exploit sparsity of  $\mathbf{Q}$  in (5). The spam  
 355 package contains routines to perform an efficient Cholesky decomposition of  $\mathbf{Q}$ , which  
 356 is important for fast sampling from a GMRF under linear constraints like the full condi-  
 357 tional  $\pi(\beta|\tau_\beta, \tau_\epsilon, \mathbf{y})$  in (5). The full conditionals for all the parameters in the model and  
 358 the code for implementing the Gibbs sampler in R can be found in the supplementary  
 359 material.

## 360 4. Application

### 361 4.1. Modelling TCWV data

362 The goal of our application is to detect the ITCZ location by using the TCWV  
 363 dataset described in Section 2. The operative definition of ITCZ that we use, as sug-  
 364 gested by researchers from ISAC-CNR, Italy, is “the strip surrounding the Earth surface  
 365 where TCWV shows highest values”.

366 To this aim, we first apply Geodesic P-splines for smoothing observed TCWV data,  
 367 which is affected by noise and does not provide measurements over the land, in order

368 to predict the latent field all over the world. Then, we exploit the model output for lo-  
 369 cating ITCZ by sampling from the joint posterior distribution of the latent field. A set of  
 370  $m = 1000$  data randomly scattered over the Earth’s surface is held out from model es-  
 371 timation for validation purposes (see section 4.2). We denote with  $\mathbf{y} = (y_1, \dots, y_n)^\top$  the  
 372 vector of TCWV observations used for model estimation and with  $\mathbf{y}^* = (y_1^*, \dots, y_m^*)^\top$   
 373 the vector of validation data. We fitted the Geodesic P-spline model to the data dis-  
 374 played in Figure 1, referred to January and July, 2008.

375 The hierarchical model is

Likelihood:

$$\begin{aligned} \mathbf{y}|\alpha, \boldsymbol{\beta}, \tau_\epsilon &\sim \mathcal{N}(\boldsymbol{\mu}, \tau_\epsilon^{-1} \mathbf{I}) \\ \boldsymbol{\mu} &= \alpha \mathbf{1} + \mathbf{B}\boldsymbol{\beta} \end{aligned} \quad (6)$$

Prior:

$$\begin{aligned} \alpha &\sim \mathcal{N}(0, \tau_\alpha^{-1}) \\ \boldsymbol{\beta}|\tau_\beta &\sim \mathcal{N}(\mathbf{0}, \tau_\beta^{-1} \mathbf{R}^*) \quad \boldsymbol{\beta} \text{ is subject to } \mathbf{1}_K^\top \mathbf{B}\boldsymbol{\beta} = 0 \end{aligned} \quad (7)$$

Hyper-prior:

$$\begin{aligned} \tau_\beta &\sim \text{Ga}(1, 5e - 5) \\ \tau_\epsilon &\sim \text{Ga}(1, 5e - 5) \end{aligned}$$

376 At the likelihood level, the matrix  $\mathbf{B}$  in (6) is the B-spline basis on a GDGG as  
 377 described in Section 3.2. The latent field  $\boldsymbol{\mu}$  is a surface varying smoothly over the  
 378 sphere, with  $\alpha$  the global spatial mean and  $\boldsymbol{\beta}$  the spline coefficients. At the prior level,  
 379 we have a diffuse Gaussian prior, with  $\tau_\alpha$  fixed at a small value for the intercept and an  
 380 ICAR prior, with precision  $\tau_\beta \mathbf{R}^*$ , for the spline coefficients. Using the scaled matrix  
 381  $\mathbf{R}^*$  is a fundamental step: this allows us to select the same prior for  $\tau_\beta$  and  $\tau_\epsilon$ , as both  $\mathbf{I}$   
 382 in (6) and  $\mathbf{R}^*$  in (7) have average marginal variance equal to 1. The results presented in  
 383 this section are obtained using a  $\text{Ga}(1, 5e - 5)$  for both hyperparameters, after checking  
 384 that the results were non sensitive to other choices for  $a$  and  $b$ .

385 To compute  $\mathbf{B}$ , the latitude and longitude coordinates are converted into spherical  
 386 coordinates, then projected on the icomesh and finally cubic B-splines are evaluated on  
 387 a GDGG with  $K$  knots, where  $K$  depends on  $\nu$ , the number of split iterations perfor-  
 388 med on the icosahedron: the choice of  $\nu$  is a critical aspect of the method and will be  
 389 discussed in section 4.3, where we compare results obtained with  $\nu = 1, \dots, 6$ .

390 Model estimation is performed by Gibbs sampling: we draw a total of 5000 samples  
 391 after convergence (achieved after a quick burnin due to the large sample size available

392 for each model parameter). Regarding computational time, it takes about ten seconds  
 393 (in a laptop with Intel core i7, 2.5GHz, 16 GB ram) to run a hundred iterations when  
 394  $\nu = 5$ .

395 In a Bayesian framework, spatial prediction is naturally based on the joint posterior  
 396 predictive distribution of the latent field: sampling from this distribution is particularly  
 397 efficient when using Bayesian P-splines. Once the posterior distribution of the latent  
 398 field  $\pi(\boldsymbol{\mu}|\mathbf{y})$  has been obtained, prediction  $\tilde{\mu}$  at an arbitrary location  $\tilde{\mathbf{x}}$  can be performed,  
 399 after evaluating the basis functions at  $\tilde{\mathbf{x}}$ , using the posterior predictive distribution:

$$\pi(\tilde{\mu}|\mathbf{y}) = \int \pi(\tilde{\mu}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}|\mathbf{y})d\boldsymbol{\theta} \quad (8)$$

400 where  $\boldsymbol{\theta} = (\alpha, \boldsymbol{\beta}, \tau_\alpha, \tau_\beta, \tau_\epsilon)$ . This is achieved by composite sampling once  $G$  sam-  
 401 ples from the posterior distribution are available. Samples from distribution (8) are  
 402 obtained by sampling from  $\pi(\tilde{\mu}|\boldsymbol{\theta}^g)$ , where  $\boldsymbol{\theta}^g$  is an MCMC sample from the posterior  
 403 distribution of  $\boldsymbol{\theta}$ .

#### 404 4.2. Model checking

405 In this section the goal is twofold: firstly, we investigate the predictive perfor-  
 406 mance of the proposed Geodesic P-spline (G-Pspline) model for different choices of  
 407  $\nu$ . Secondly, we compare the G-Pspline and the stochastic partial differential equa-  
 408 tion (SPDE) approach by Lindgren et al. (2011), both in terms of computational and  
 409 predictive performance, using the TCWV data.

The predictive performance is evaluated by first estimating the model on training  
 dataset and then computing error measures on a validation dataset. Let  $\hat{y}_j^* = E(y_j^*|\mathbf{y})$   
 denote the mean of the posterior predictive distribution at the validation location  $j$ ,  $j =$   
 $1, \dots, m$ , as measures of predictive performances we consider both the relative mean  
 absolute prediction error (RMAPE),

$$RMAPE = \frac{1}{m} \sum_{j=1}^m \left| \frac{\hat{y}_j^* - y_j^*}{y_j^*} \right|, \quad (9)$$

410 where the average is taken over the validation locations, and the relative mean square  
 411 prediction error (RMSPE), which is the same as (9) except for averaging squares, in-  
 412 stead of absolute, relative errors.

##### 413 4.2.1. G-Pspline model performance for varying $\nu$

414 The first three lines of Table 2 report the RMAPE and the RMSPE for the G-Pspline  
 415 model for different number of knots,  $K$  (note that the different  $K$ 's are associated to

416  $\nu = 1, \dots, 6$ ). The RMSPE shows a minimum at  $\nu = 5$  (0.109) while the RMAPE is  
 417 almost unchanged when increasing  $\nu$  from 5 to 6 as it decays from 12.8% ( $\nu = 5$ ) to  
 418 12.2% ( $\nu = 6$ ). Based on these results we find appropriate to select  $\nu = 5$  for ITCZ  
 419 location, as this allows us to gain computational speed in the procedure described in  
 420 section 4.3.

<b>G-Pspline</b>	<i>K</i>	42	162	642	2562	10242	40962
	RMAPE	0.263	0.213	0.169	0.146	0.128	0.122
	RMSPE	0.218	0.206	0.166	0.136	0.109	0.122
<b>SPDE</b>	<i>K</i>	43	164	644	2580	10243	40841
	RMAPE	0.518	0.204	0.166	0.148	0.128	0.118
	RMSPE	1.619	0.186	0.163	0.151	0.106	0.110

Table 2: Relative mean absolute prediction error (RMAPE) and relative mean square prediction error (RMSPE) obtained with Geodesic P-Splines (G-Pspline) and SPDE.  $K$  denotes the number of knots of the geodesic grid for G-Pspline (for  $\nu = 1, \dots, 6$ ), or the number of knots of the triangular mesh for SPDE (obtained by tuning the `max.edge` argument in the R-INLA function `inla.mesh.2d`).

421 The choice of  $\nu$  is the starting step of the modelling process, similarly to the choice  
 422 of a suitable triangulation in the SPDE approach. Literature on P-splines recomm-  
 423 mends using a number of knots  $K$  large enough to describe the spatial variation of the  
 424 data and let the penalty determine the right amount of smoothing. Under different  $K$   
 425 levels, provided that  $K$  is large enough, the same degree of smoothing is obtained by  
 426 rescaling the smoothing parameter accordingly. This is confirmed in Table 2, where  
 427 measures of predictive performance of the G-Pspline model remain almost unchanged  
 428 for  $K = 10242$  ( $\nu = 5$ ) and  $K = 40962$  ( $\nu = 6$ ). In a Bayesian P-spline setting, this  
 429 rescaling is reflected in the posterior distribution of  $\tau_\beta$ ; if  $K$  changes, the location of  
 430  $\pi(\tau_\beta|\mathbf{y})$  shifts accordingly.

#### 431 4.2.2. Comparing G-Pspline and SPDE

432 SPDE is implemented in the R-INLA package (Martins et al., 2013) and, as a start-  
 433 ing point, needs the definition of a triangular mesh covering the study region, analo-  
 434 gously to the definition of a geodesic grid in our framework. For the sake of comparison  
 435 of the prediction performance the SPDE mesh and the G-Pspline geodesic grid should  
 436 have similar size. For each column in Table 2 (i.e., for each  $\nu$  from 1 to 6) the triang-  
 437 ular mesh is built using the R-INLA function `inla.mesh.2d`, by tuning the `max.edge`  
 438 argument (the largest allowed triangle edge length) in order to have roughly the same  
 439 number of knots  $K$ .

440 Looking at Table 2 column-wise, G-Pspline and SPDE perform similarly in terms  
 441 of RMAPE and RMSPE. The boxplots in Figure 6 show the variability of the (log)

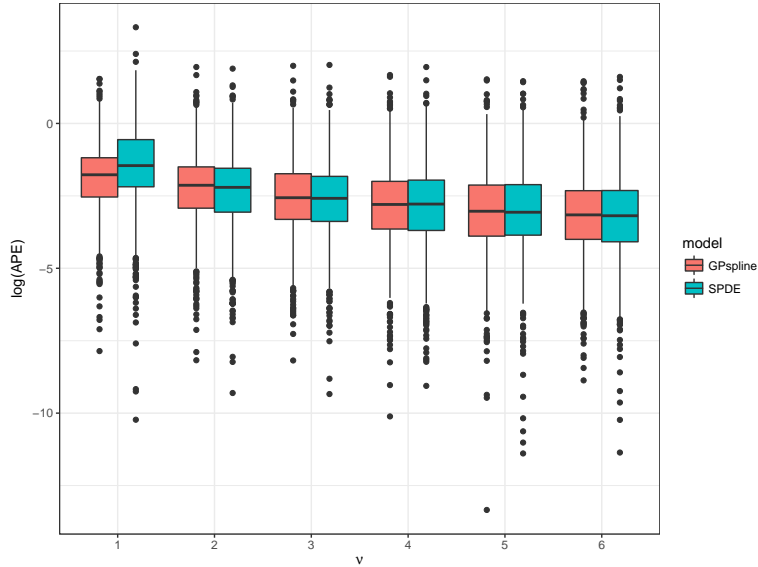


Figure 6: Boxplots of the (log) relative absolute prediction errors (APE),  $\log(|\hat{y}_j^* - y_j^*|/y_j^*)$ , measured on a validation set of  $m = 1000$  locations. Prediction performance for the G-Pspline and the SPDE models is very similar for each  $\nu$ .

442 relative absolute prediction errors over the  $m$  validation locations, for both models and  
 443 varying  $\nu$ . For each  $\nu$ , the variability of the relative absolute prediction errors is prac-  
 444 tically the same for G-Pspline and SPDE. Boxplots for the squared prediction errors  
 445 present a similar pattern and are not shown here. Based on these results we can con-  
 446 clude that, in our case study, prediction performance measured on a validation set of  
 447  $m = 1000$  locations is overall very similar for G-Pspline and SPDE.

448 As a final note on computation time, the SPDE model fitted within R-INLA is faster  
 449 than the G-Pspline fitted via Gibbs sampling: SPDE takes around four minutes, while  
 450 our Gibbs sampler takes around ten minutes to run 5000 iterations. Nonetheless, the  
 451 procedure described in Section 4.3 for ITCZ location requires MCMC samples from  
 452 the model posterior. Sampling from the posterior of the latent field within R-INLA  
 453 (using `inla.posterior.sample()`) is computationally intensive for our model, as it  
 454 takes around 30 seconds to run 10 samples. Therefore, to the purpose of ITCZ location,  
 455 the proposed Geodesic P-spline approach using Gibbs sampling is overall faster than  
 456 SPDE within R-INLA, while maintaining the same predictive performance. Results on  
 457 the ITCZ location obtained with SPDE for January and July (not shown here) were very  
 458 similar to those presented in Figure 8 which are obtained using the procedure discussed  
 459 next.

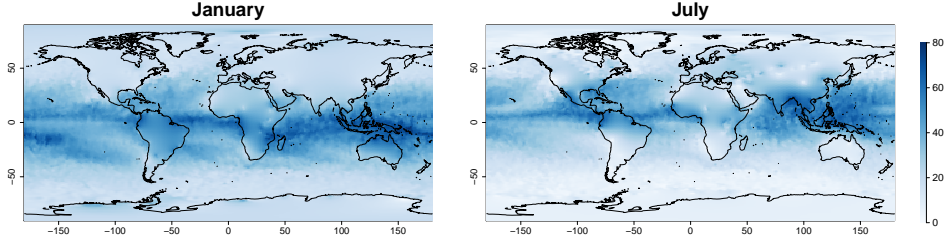


Figure 7: Model prediction of the latent field for January and July.

#### 460 4.3. Locating the ITCZ

461 In Figure 7 we illustrate the maps of the TCWV posterior means obtained with  
 462  $\nu = 5$ : this accomplishes our first task, i.e. to remove random noise from data and to  
 463 reconstruct the latent field on the whole of Earth’s surface.

464 The problem of ITCZ location is addressed by summarising the posterior predic-  
 465 tive distribution of the TCWV latent field. The procedure outlined below requires the  
 466 specification of a reasonable guess concerning the width of the ITCZ region denoted  
 467 as  $W$ ; we based our choice on expert knowledge by ISAC-CNR researchers and set  
 468  $W = 1000 \text{ km}$ . The ITCZ width relative to the length of a Meridian (which is about  
 469  $20000 \text{ km}$ ) is around  $w = W/20,000 = 0,05$ .

470 Our algorithm to locate the ITCZ consists of a discrete search performed longitude-  
 471 wise (i.e. at each meridian). Let  $m = 1, \dots, M$  indicate a set of  $M$  meridians: for a  
 472 given  $m$ , we sample from the posterior predictive distribution of the latent field at a fine  
 473 grid over latitude. Then, we compute the posterior probability that a point at a given  
 474 latitude belongs to the region where the TCWV shows the highest values (i.e. the point  
 475 falls into the ITCZ region), integrating out uncertainty about model parameters.

476 Let  $\tilde{\boldsymbol{\mu}}_m = (\tilde{\mu}_{1m}, \dots, \tilde{\mu}_{lm}, \dots, \tilde{\mu}_{Lm})$  be the vector of the latent field predicted at loca-  
 477 tions  $l = 1, \dots, L$ , where  $(lat_{1m}, \dots, lat_{lm}, \dots, lat_{Lm})$  is a regular sequence from  $90^\circ$  to  
 478  $-90^\circ$ . The algorithm proceeds as follows. For  $m = 1, \dots, M$ :

- 479 • evaluate the bases at locations  $l = 1, \dots, L$ , this gives a meridian-specific  $L \times K$   
 480 dimensional basis matrix  $\tilde{\mathbf{B}}_m$ ;
- 481 • sample  $G$  realizations from the posterior predictive distribution (8) by computing  
 482  $\tilde{\boldsymbol{\mu}}_m^g = \alpha^g \mathbf{1} + \tilde{\mathbf{B}}_m \hat{\boldsymbol{\beta}}^g$ ,  $g = 1, \dots, G$ ;
- 483 • for  $g = 1, \dots, G$ , rank the vector  $\tilde{\boldsymbol{\mu}}_m^g$ . This gives a posterior sample of the  
 484 ranks, indicated by vector  $\boldsymbol{\phi}_m^g = (\phi_{1m}^g, \dots, \phi_{lm}^g, \dots, \phi_{Lm}^g)$ , e.g.  $\phi_{lm}^g = L$  if  $l =$   
 485  $\text{argmax}_l(\tilde{\mu}_{lm}^g)$ , while  $\phi_{lm}^g = 1$  if  $l = \text{argmin}_l(\tilde{\mu}_{lm}^g)$ .

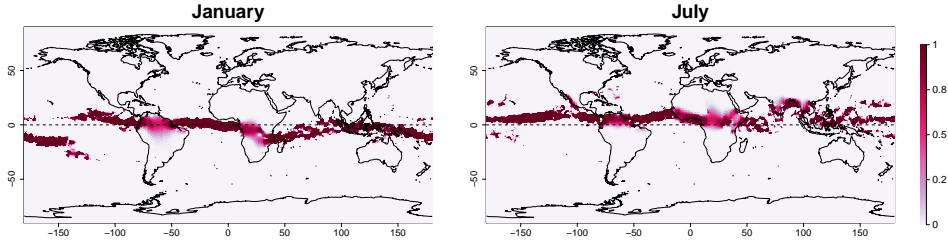


Figure 8: ITCZ location for January and July.

The probability that a point  $l$  belonging to meridian  $m$  falls into the ITCZ is computed as

$$Pr(lat_{lm} \in ITCZ|y) = \frac{1}{G} \sum_{g=1}^G \mathbb{1} \left( 1 - \frac{\phi_{lm}^g}{L} < w \right) \quad (10)$$

486 where  $\mathbb{1}$  is the indicator function and  $\phi_{lm}^g/L$  is the normalised rank. To sum up the  
 487 above, (10) is the probability that the point with geographical coordinates  $(lat_l, long_m)$   
 488 falls inside the ITCZ, where the length of the ITCZ is fixed according to  $W$ . Results are  
 489 displayed in Figure 8 for the two months under examination: this Figure is obtained  
 490 running the algorithm with  $L = 1000$  and  $M = 360$ . The ITCZ is mostly located in  
 491 the south (north) of the Equator in January (July), as expected on the basis of prior  
 492 knowledge concerning its seasonal behaviour. The map for January shows the double  
 493 ITCZ, which is typical of the Central Pacific region in some period during the year  
 494 (Waliser and Gautier, 1993). The proposed method allows to locate the ITCZ even  
 495 over land (in particular in Africa and South America) where data is not available, this  
 496 being reflected by higher posterior uncertainty. Of course, the width of ITCZ reported  
 497 in Figure 8 is strictly dependent on the choice of  $W$ : although this is very relevant  
 498 when studying a single month, we believe that it is not such a crucial choice if the  
 499 method is used for studying the spatio-temporal trend of the phenomenon. Indeed, in  
 500 this case it would be important to keep  $W$  fixed along the study period in order to ensure  
 501 comparability among results.

## 502 5. Discussion

503 We presented a Bayesian hierarchical framework for smoothing data collected world-  
 504 wide at a large number of locations. With respect to traditional methods, the proposed  
 505 model accounts for geodesic distances between the data, thus overcoming the limita-  
 506 tions of covariance functions for Euclidean spaces when applied to global datasets.

507 The non-parametric model formulation proposed extends the Bayesian P-spline ap-  
508 proach for smoothing worldwide collected data. Using a sphere as a representation of  
509 the Globe, the idea is to build a new basis of B-splines on a suitable geodesic grid while  
510 keeping the hierarchical model formulation of Bayesian P-splines, with the associated  
511 advantages in terms of flexibility and computation. Two key features of P-splines are  
512 maintained in the Geodesic P-spline model: (a) the use of local bell-shaped functions,  
513 e.g. the B-splines on the icomesh, that yield a sparse basis matrix; (b) the use of B-  
514 splines centred at equidistant knots, i.e. the nodes of the icomesh. Point (b) suggests  
515 that an IGMRF for regular locations is a sensible prior distribution for the spline co-  
516 efficients, giving stable marginal variances as opposed to the standard P-spline model  
517 construction. Computational efficiency is due to (a) reduction of the latent field dimen-  
518 sion, as the smoothing prior operates on the spline coefficients (low-rank smoother)  
519 and (b) fast MCMC based on sparse Cholesky factorization of the structure matrix of  
520 the full conditional for the latent field. These advantages allow for a fast fitting of the  
521 model to data collected worldwide in high-resolution.

522 We applied the Geodesic P-spline model to the TCWV data retrieved with the AIR-  
523 WAVE algorithm at a huge number of locations on Earth. The smoothing approach in  
524 this example is desirable as it allows field estimation at unmonitored locations. We  
525 provided inferential tools to locate the ITCZ based on ranking samples from the poste-  
526 rior distribution of the latent field, estimated at a fine grid over the Globe. Results are  
527 coherent with prior knowledge concerning ITCZ, indicating a shift towards southern  
528 regions in autumn and winter.

529 A critical aspect is the choice of hyperprior for the random walk precision,  $\pi(\tau_\beta)$ .  
530 We expect a large impact of  $\pi(\tau_\beta)$  in situations where sample size is small compared  
531 to the number of parameters required to be estimated. In the case study on TCWV,  
532 the sample size available for estimating each spline coefficients is large enough, which  
533 makes the impact of  $\pi(\tau_\beta)$  very small. In the results presented in Section 4 we used a  
534 Gamma with shape  $a = 1$  and rate  $b = 5e - 5$  for both  $\tau_\beta$  and  $\tau_\epsilon$ , after checking that the  
535 posterior  $\pi(\tau_\beta|\mathbf{y})$  remained unchanged under different choices of  $a$  and  $b$ . We believe  
536 that controlling that the posterior learns from the data in the same way for different  
537 choices of the prior is a reasonable approach to test the robustness of the Bayesian  
538 specification. On the topic of prior selection for variance parameters the literature has  
539 shown rapid growth over the past decade; see, e.g., Gelman (2006); Simpson et al.  
540 (2017) and references therein.

541 The model can be extended in several directions, both in a methodological and ap-  
542 plied sense. In this paper we focused on an IGMRF structure for the spline coefficients  
543 equivalent to the ICAR used for lattice data, using the six surrounding knots as neigh-



544 bours. Investigation of geodesic grids suitable for higher order IGMRF priors would  
545 be interesting. Another attractive research line would be to look into a model based on  
546 nested B-splines, defined on a set of geodesic grids of different resolution, following  
547 Nychka et al. (2015). In a fully Bayesian framework this requires careful hyperprior  
548 specification, as it is not clear how to prevent confounding between nested components.

549 As for application, a future research line worthy of investigation is modelling the  
550 ITCZ based on different proxy variables, focusing the analysis on a wide temporal  
551 range, following the ideas in Castelli et al. (2017). The application of Geodesic P-  
552 spline models to the 20 years of ATSR data will allow the investigation of ITCZ merid-  
553 ional migration trends. Moreover, joint modelling of TCWV and other ITCZ related  
554 phenomena, possibly available at misaligned locations, will result in more reliable es-  
555 timates of the ITCZ latent field, especially at locations where TCWV retrieval is not  
556 possible with current ATSR technology.

#### 557 **Acknowledgements**

558 The work by Fedele Greco and Massimo Ventrucci is funded by the PRIN2015-  
559 supported project *Environmental processes and human activities: capturing their inter-*  
560 *actions via statistical methods* (EPHASTAT) by MIUR (Italian Ministry of Education,  
561 University and Scientific Research). We thank Bianca Maria Dinelli (ISAC-CNR),  
562 Enzo Papandrea and Stefano Casadio (Serco s.p.a.) for guidance in the interpretation  
563 of the results. ATSR TCWV dataset was developed in the frame of the ESA ALTS  
564 project ESA Contract No. 4000108531/13/I-NB.

#### 565 **References**

- 566 Banerjee, S., 2005. On Geodetic Distance Computations in Spatial Modeling. *Biomet-*  
567 *rics* 61 (2), 617–625.  
568 URL <http://dx.doi.org/10.1111/j.1541-0420.2005.00320.x>
- 569 Banerjee, S., 2017. High-dimensional bayesian geostatistics. *Bayesian Analysis* 12 (2),  
570 583–614.  
571 URL <https://doi.org/10.1214/17-BA1056R>
- 572 Banerjee, S., Gelfand, A., Carlin, B., 2014. *Hierarchical Modeling and Analysis for*  
573 *Spatial Data*, Second Edition. Chapman & Hall/CRC Monographs on Statistics &  
574 *Applied Probability*. Taylor & Francis.

- 575 Banerjee, S., Gelfand, A. E., Finley, A. O., Sang, H., 2008. Gaussian predictive process  
576 models for large spatial data sets. *Journal of the Royal Statistical Society: Series B*  
577 (Statistical Methodology) 70 (4), 825–848.  
578 URL <http://dx.doi.org/10.1111/j.1467-9868.2008.00663.x>
- 579 Baramidze, V., Lai, M. J., Shum, C. K., 2006. Spherical splines for data interpolation  
580 and fitting. *SIAM Journal on Scientific Computing* 28 (1), 241–259.  
581 URL <https://doi.org/10.1137/040620722>
- 582 Besag, J., 1974. Spatial interaction and the statistical analysis of lattice systems (with  
583 discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodol-*  
584 *ogy)* 36(2), 192–225.
- 585 Casadio, S., Castelli, E., Papandrea, E., Dinelli, B., Pisacane, G., Burini, A., Bojkov,  
586 B., 2016. Total column water vapour from along track scanning radiometer series  
587 using thermal infrared dual view ocean cloud free measurements: The advanced  
588 infra-red water vapour estimator (airwave) algorithm. *Remote Sensing of Environ-*  
589 *ment* 172, 1–14.
- 590 Castelli, E., Ventrucci, M., Greco, F., Valeri, M., Dinelli, B., Papandrea, E., Casa-  
591 dio, S., 2017. On the contribution of 20 years of ATSR data and geodesic P-spline:  
592 Efficient spatial smoothing method to ITCZ trend analysis. In: *Proceedings of the*  
593 *2017 conference on Big Data from Space (BiDS17)*, 28-30 October 2017, Toulouse,  
594 France.
- 595 Cressie, N., Johannesson, G., 2008. Fixed rank kriging for very large spatial data sets.  
596 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 70 (1),  
597 209–226.  
598 URL <http://dx.doi.org/10.1111/j.1467-9868.2007.00633.x>
- 599 Datta, A., Banerjee, S., Finley, A. O., Gelfand, A. E., 2016. Hierarchical nearest-  
600 neighbor gaussian process models for large geostatistical datasets. *Journal of the*  
601 *American Statistical Association* 111 (514), 800–812.  
602 URL <https://doi.org/10.1080/01621459.2015.1044091>
- 603 Di Marzio, M., Panzera, A., Taylor, C. C., 2014. Nonparametric regression for spheri-  
604 cal data. *Journal of the American Statistical Association* 109 (506), 748–763.  
605 URL <http://dx.doi.org/10.1080/01621459.2013.866567>

- 606 Duchamp, T., Stuetzle, W., 2003. Spline smoothing on surfaces. *Journal of Computa-*  
607 *tional and Graphical Statistics* 12 (2), 354–381.  
608 URL <https://doi.org/10.1198/1061860031743>
- 609 Eilers, P., Currie, I., Durbán, M., 2006. Fast and compact smoothing on large multid-
- 610 *imensional grids. Computational Statistics & Data Analysis* 5, 61–76.
- 611 Eilers, P., Marx, B., 1996. Flexible Smoothing with B-splines and Penalties. *Statistical*  
612 *Science* 11, 89–121.
- 613 Ettinger, B., Perotto, S., Sangalli, L. M., 2016. Spatial regression models over two-
- 614 *dimensional manifolds. Biometrika* 103 (1), 71–88.  
615 URL <http://dx.doi.org/10.1093/biomet/asv069>
- 616 Furrer, R., Genton, M. G., Nychka, D., 2006. Covariance tapering for interpolation of
- 617 *large spatial datasets. Journal of Computational and Graphical Statistics* 15, 502–
- 618 523.
- 619 Furrer, R., Sain, S., 2010. spam: A Sparse Matrix R Package with Emphasis on MCMC
- 620 *Methods for Gaussian Markov Random Fields. Journal of Statistical Software, Arti-*  
621 *cles* 36 (10), 1–25.  
622 URL <https://www.jstatsoft.org/v036/i10>
- 623 Gelman, A., 2006. Prior distributions for variance parameters in hierarchical models
- 624 *(comment on article by Browne and Draper). Bayesian Analysis* 1 (3), 515–534.  
625 URL <http://dx.doi.org/10.1214/06-BA117A>
- 626 Gneiting, T., 2013. Strictly and non-strictly positive definite functions on spheres.
- 627 *Bernoulli* 19 (4), 1327–1349.  
628 URL <http://dx.doi.org/10.3150/12-BEJSP06>
- 629 Holton, J. R., Hakim, G. J., 2013. *An Introduction to Dynamic Meteorology (Fifth*  
630 *Edition)*. Academic Press, Boston.  
631 URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/B9780123848666000301)  
632 [B9780123848666000301](https://www.sciencedirect.com/science/article/pii/B9780123848666000301)
- 633 Jeong, J., Jun, M., 2015. A class of matérn-like covariance functions for smooth
- 634 *processes on a sphere. Spatial Statistics* 11, 1 – 18.  
635 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S2211675314000554)  
636 [S2211675314000554](http://www.sciencedirect.com/science/article/pii/S2211675314000554)

- 637 Jun, M., Stein, M. L., 2007. An approach to producing space: Time covariance func-  
638 tions on spheres. *Technometrics* 49 (4), 468–479.  
639 URL <http://www.jstor.org/stable/25471392>
- 640 Lai, M., Schumaker, L., 2007. Spline Functions on Triangulations. No. v. 13 in Ency-  
641 clopedia of Mathematics an. Cambridge University Press.  
642 URL <https://books.google.it/books?id=6hvvqGgbBmEoC>
- 643 Lai, M. J., Shum, C. K., Baramidze, V., Wenston, P., Aug 2009. Triangulated spherical  
644 splines for geopotential reconstruction. *Journal of Geodesy* 83 (8), 695–708.  
645 URL <https://doi.org/10.1007/s00190-008-0283-0>
- 646 Lang, S., Brezger, A., 2004. Bayesian p-splines. *Journal of Computational and Graph-*  
647 *ical Statistics* 13, 183–212.
- 648 Lee, D., 2010. Smothing mixed model for spatial and spatio-temporal data. Tech. rep.,  
649 Ph.D. Thesis, Department of Statistics, Universidad Carlos III de Madrid, Spain.
- 650 Lee, D., Durbán, M., 2009. Smooth-car mixed models for spatial count data. *Compu-*  
651 *tational Statistics and data Analysis* 53, 2968–2977.
- 652 Lindgren, F., Rue, H., Lindström, J., 2011. An explicit link between Gaussian fields  
653 and Gaussian Markov random fields: the stochastic partial differential equation ap-  
654 proach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*  
655 73 (4), 423–498.  
656 URL <http://dx.doi.org/10.1111/j.1467-9868.2011.00777.x>
- 657 Martins, T. G., Simpson, D., Lindgren, F., Rue, H., 2013. Bayesian computing with  
658 inla: New features. *Computational Statistics & Data Analysis* 67 (0), 68 – 83.  
659 URL <http://www.sciencedirect.com/science/article/pii/S0167947313001552>
- 660
- 661 Nychka, D., Bandyopadhyay, S., Hammerling, D., Lindgren, F., Sain, S., 2015. A  
662 multiresolution gaussian process model for the analysis of large spatial datasets.  
663 *Journal of Computational and Graphical Statistics* 24 (2), 579–599.  
664 URL <http://dx.doi.org/10.1080/10618600.2014.914946>
- 665 Porcu, E., Bevilacqua, M., Genton, M. G., 2016. Spatio-temporal covariance and cross-  
666 covariance functions of the great circle distance on a sphere. *Journal of the American*  
667 *Statistical Association* 111 (514), 888–898.  
668 URL <https://doi.org/10.1080/01621459.2015.1072541>

- 669 Randall, D. A., Ringler, T. D., Heikes, R., Jones, P., Baumgardner, J., 2002. Climate  
670 modeling with spherical geodesic grids. *Computing in Science and Engineering* 4,  
671 32–41.
- 672 Rue, H., Held, L., 2005. *Gaussian Markov Random Fields*. Chapman and Hall/CRC.
- 673 Rue, H., Martino, S., Chopin, N., 2009. Approximate Bayesian inference for latent  
674 Gaussian models by using integrated nested Laplace approximations. *Journal of the*  
675 *Royal Statistical Society: Series B (Statistical Methodology)* 71 (2), 319–392.  
676 URL <http://dx.doi.org/10.1111/j.1467-9868.2008.00700.x>
- 677 Sahr, K., White, D., Kimerling, A. J., 2003. Geodesic Discrete Global Grid Systems.  
678 *Cartography and Geographic Information Science* 30 (2), 121–134.  
679 URL <http://dx.doi.org/10.1559/152304003100011090>
- 680 Sangalli, L. M., Ramsay, J. O., Ramsay, T. O., 2013. Spatial spline regression models.  
681 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 75 (4),  
682 681–703.  
683 URL <http://dx.doi.org/10.1111/rssb.12009>
- 684 Simpson, D., Rue, H., Riebler, A., Martins, T. G., Sørbye, S. H., 02 2017. Penalising  
685 model component complexity: A principled, practical approach to constructing pri-  
686 ors. *Statist. Sci.* 32 (1), 1–28.  
687 URL <https://doi.org/10.1214/16-STS576>
- 688 Sørbye, S., Rue, H., 2014. Scaling intrinsic gaussian markov random field priors in  
689 spatial modelling. *Spatial Statistics* 8, 39 – 51.  
690 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S2211675313000407)  
691 [S2211675313000407](http://www.sciencedirect.com/science/article/pii/S2211675313000407)
- 692 Ugarte, M., Goicoa, T., Etxeberria, J., Militino, A., 2012. A p-spline anova type model  
693 in space-time disease mapping. *Stochastic Environmental Research and Risk As-*  
694 *essment* 26 (6), 835–845.
- 695 Wahba, G., 1981. Spline interpolation and smoothing on the sphere. *SIAM Journal on*  
696 *Scientific and Statistical Computing* 2 (1), 5–16.  
697 URL <https://doi.org/10.1137/0902002>
- 698 Waliser, D. E., Gautier, C., 1993. A satellite-derived climatology of the ITCZ. *Journal*  
699 *of Climate* 6 (11), 2162–2174.

- 700 Wikle, C., Cressie, N., 1999. A dimension-reduced approach to space-time Kalman  
701 filtering. *Biometrika* 86 (4), 815–829.  
702 URL +<http://dx.doi.org/10.1093/biomet/86.4.815>
- 703 Wilhelm, M., Ded, L., Sangalli, L. M., Wilhelm, P., 2016. Igs: An isogeometric  
704 approach for smoothing on surfaces. *Computer Methods in Applied Mechanics and*  
705 *Engineering* 302, 70 – 89.  
706 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0045782516000025)  
707 [S0045782516000025](http://www.sciencedirect.com/science/article/pii/S0045782516000025)
- 708 Wood, S., 2017. *Generalized Additive Models: An Introduction with R*. Chapman and  
709 Hall/CRC.
- 710 Žagar, N., Skok, G., Tribbia, J., 2011. Climatology of the ITCZ derived from ERA  
711 Interim reanalyses. *Journal of Geophysical Research: Atmospheres* 116 (D15).
- 712 Zhang, C., 2001. Double ITCZs. *Journal of Geophysical Research: Atmospheres*  
713 106 (D11), 11785–11792.
- 714 Heaton, M., Katzfuss, M., Berrett, C., Nychka, D., 2014. Constructing valid spatial  
715 processes on the sphere using kernel convolutions. *Environmetrics* 25 (1), 2–15.  
716 URL <http://dx.doi.org/10.1002/env.2251>