

---

# Journal of Environmental Statistics

August 2014, Volume 6, Issue 3.

<http://www.jenvstat.org>

---

## Efficient Approximation of the Spatial Covariance Function for Large Datasets - Analysis of Atmospheric CO<sub>2</sub> Concentrations

**Patrick Vetter**

*Department of Statistics, European University Viadrina, Frankfurt (Oder), Germany*  
Wolfgang Schmid

*Department of Statistics, European University Viadrina, Frankfurt (Oder), Germany*  
Reimund Schwarze

*Department of Economics, European University Viadrina, Frankfurt (Oder), Germany*

---

### Abstract

Linear mixed effects models have been widely used in the spatial analysis of environmental processes. However, parameter estimation and spatial predictions involve the inversion and determinant of the  $n \times n$  dimensional spatial covariance matrix of the data process, with  $n$  being the number of observations. Nowadays environmental variables are typically obtained through remote sensing and contain observations of the order of tens or hundreds of thousands on a single day, which quickly leads to bottlenecks in terms of computation speed and requirements in working memory. Therefore techniques for reducing the dimension of the problem are required. The present work analyzes approaches to approximate the spatial covariance function in a real dataset of remotely sensed carbon dioxide concentrations, obtained from the Atmospheric Infrared Sounder of NASA's "Aqua" satellite on the 1st of May 2009. In a cross-validation case study it is shown how fixed rank kriging, stationary covariance tapering and the full-scale approximation are able to notably speed up calculations. However, the loss in predictive performance caused by the approximation strongly differs. The best results were obtained for the full-scale approximation, which was able to overcome the individual weaknesses of the fixed rank kriging and the covariance tapering.

**Keywords:** spatial covariance function, fixed rank kriging, covariance tapering, full-scale approximation, large spatial data sets, mid-tropospheric CO<sub>2</sub>, remote sensing, efficient approximation.

---

## 1. Introduction

The monitoring of environmental processes has been revolutionized in the recent past through the upcoming of remotely sensed satellite measurements. The resulting spatial resolution is far superior compared to the traditional monitoring through networks of measurement stations. This of course constitutes a major improvement for scientific research, but also introduces the need for statistical models that can handle such large data sets, which often involve observations on the order of tens or hundreds of thousand per day. One such environmental data set of particular interest for the ongoing political discussion and the related negotiations on climate change and global warming is the remotely sensed measurement of carbon dioxide concentration in the mid-troposphere as measured by the Atmospheric Infrared Sounder (AIRS) of NASA's "Aqua" satellite. It has been contemplated that space-based observations of  $CO_2$  could complement the weakly enforceable system of national reporting of sources of  $CO_2$  emissions in a meaningful way (Mintzer, Leonard, and Valencia (2010, p. 28)). In fact, recent studies have reported a 'gap' in  $CO_2$  reporting from China of 1.4 gigatonnes per year (Guan, Liu, Geng, Lindner, and Hubacek (2012)), which would amount to 5% of the global total. This gives rise for an objective assessment of  $CO_2$  emissions based on measurements and a corresponding validation of national reporting standards. In that way statistical modeling of atmospheric  $CO_2$  concentrations can serve as an important input to climate projection projects and for the estimation of  $CO_2$  surface fluxes. Linear mixed effects models have been widely used in the spatial analysis of such environmental data sets. However, parameter estimation and spatial predictions involve the inversion and determinant of the  $n \times n$  dimensional spatial covariance matrix of the data process, with  $n$  being the number of observations. As mentioned above, environmental variables as measured through remote sensing contain observations of the order of tens or hundreds of thousand on a single day, which quickly leads to bottlenecks in terms of computation speed and requirements in working memory.

**Linear Mixed-Effects Models** Consider a real-valued spatial process  $\{Z(\mathbf{s}) : \mathbf{s} \in D \subset \mathbb{R}^d\}$  defined on the domain of interest (e.g. the globe as in the  $CO_2$  example). The process is observed at  $n$  locations and is a noisy version of the smooth process  $\{Y(\mathbf{s}) : \mathbf{s} \in D \subset \mathbb{R}^d\}$ , which we are interested in making inference on. This defines the process  $Z(\cdot)$  at location  $\mathbf{s}$  as

$$Z(\mathbf{s}) = Y(\mathbf{s}) + \epsilon(\mathbf{s}) \quad (1)$$

where  $\{\epsilon(\mathbf{s}) : \mathbf{s} \in D\}$  is a spatial white-noise process with zero mean and  $var(\epsilon(\mathbf{s})) = \sigma_\epsilon^2 v(\mathbf{s})$ .  $\epsilon(\cdot)$  covers the *nugget effect*, or alternatively the measurement error of the instrument. The smooth process  $Y(\cdot)$  contains two parts,

$$Y(\mathbf{s}) = \mathbf{x}(\mathbf{s})' \boldsymbol{\alpha} + \nu(\mathbf{s}) \quad (2)$$

where the first one covers fixed-effects from a deterministic large-scale trend, modeled here as a linear function of  $p$  spatial covariates  $\mathbf{x}(\cdot)$ . The second term  $\nu(\cdot)$  models small-scale spatial random variations through a zero-mean process with positive and finite variance and (generally non-stationary) covariance function

$$cov(\nu(\mathbf{u}), \nu(\mathbf{v})) \equiv C(\mathbf{u}, \mathbf{v}) \quad \mathbf{u}, \mathbf{v} \in D \quad . \quad (3)$$

For the process  $Z(\cdot)$  at the  $n$  observed locations  $\mathbf{Z} = (Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n))$  this becomes

$$\mathbf{Z} = \mathbf{X}\boldsymbol{\alpha} + \boldsymbol{\nu} + \boldsymbol{\epsilon} \quad (4)$$

with  $\mathbf{X}$  being the  $n \times p$  matrix of covariate values at the observed data locations. Assuming  $\boldsymbol{\epsilon}$  and  $\boldsymbol{\nu}$  to be independent the resulting  $n \times n$  covariance matrix of  $\mathbf{Z}$  is

$$\boldsymbol{\Sigma} = \text{var}(\boldsymbol{\nu}) + \text{var}(\boldsymbol{\epsilon}) = \mathbf{C} + \sigma_\epsilon^2 \mathbf{V}_\epsilon \quad (5)$$

where  $\mathbf{C}$  is the covariance matrix of  $\boldsymbol{\nu}$  generated by the covariance function in (3) and  $\mathbf{V}_\epsilon = \text{diag}\{v_\epsilon(\mathbf{s}_1), \dots, v_\epsilon(\mathbf{s}_n)\}$ . The model described in (1)-(5) is also called a *linear mixed-effects model*.

To obtain an optimal linear spatial prediction of the smooth process  $Y(\cdot)$  at a specific location  $\mathbf{s}_0$ , *universal kriging* can be applied, as described for example in Cressie and Wikle (2011). Universal Kriging solves for the homogeneously linear combination of the data  $\boldsymbol{\lambda}'\mathbf{Z}$ , that minimizes the mean squared prediction error

$$MSPE(\boldsymbol{\lambda}) = E(Y(\mathbf{s}_0) - \boldsymbol{\lambda}'\mathbf{Z})^2.$$

In a purely gaussian setting this is also equivalent to deriving the posterior distribution  $[Y(\mathbf{s}_0)|\mathbf{Z}]$  and its first two moments  $E(Y(\mathbf{s}_0)|\mathbf{Z})$  and  $\text{var}(Y(\mathbf{s}_0)|\mathbf{Z})$ . The resulting universal kriging predictor and kriging variance are given in (6) and (7) (Cressie and Wikle (2011, p. 148))

$$\hat{Y}(\mathbf{s}_0) = \mathbf{x}(\mathbf{s}_0)' \hat{\boldsymbol{\alpha}}_{gls} + \mathbf{c}_Y(\mathbf{s}_0)' \boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\alpha}}_{gls}) \quad (6)$$

$$\begin{aligned} \hat{\sigma}^2(\mathbf{s}_0) &= \text{Var}(Y(\mathbf{s}_0)) \\ &= \mathbf{c}_Y(\mathbf{s}_0)' \boldsymbol{\Sigma}^{-1} \mathbf{c}_Y(\mathbf{s}_0) \\ &\quad + (\mathbf{x}(\mathbf{s}_0) - \mathbf{X}\boldsymbol{\Sigma}^{-1}\mathbf{c}_Y(\mathbf{s}_0))'(\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}(\mathbf{x}(\mathbf{s}_0) - \mathbf{X}\boldsymbol{\Sigma}^{-1}\mathbf{c}_Y(\mathbf{s}_0))), \end{aligned} \quad (7)$$

where  $\mathbf{c}_Y(\mathbf{s}_0) = \text{cov}(Y(\mathbf{s}_0), \mathbf{Z})$  describes the cross-covariance between  $Y(\mathbf{s}_0)$  and the observed data  $\mathbf{Z}$ , generated through the covariance function in (3), and  $\hat{\boldsymbol{\alpha}}_{gls} = (\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{Z}$ , is the generalized-least squares estimator of  $\boldsymbol{\alpha}$ . Computational problems of speed and storage may arise in the calculation of the inverse of the  $n \times n$  covariance matrix  $\boldsymbol{\Sigma}$ , which is needed for kriging predictions and variances in (6) and (7) and requires  $O(n^3)$  computations. This becomes even more difficult in iterative Maximum-Likelihood parameter estimations, where  $\boldsymbol{\Sigma}^{-1}$  has to be calculated in each iteration step. Another potential shortage can be identified for the case of a large number of prediction locations, for which the  $m \times n$  cross-covariance matrix might become very large and needs huge amounts of storage in the current workspace. The following approaches have been recently developed to tackle the large-matrix-problem by applying a low-rank approximation of the spatial process  $\nu(\cdot)$  (e.g. Cressie and Johannesson (2006), Shi and Cressie (2007), Cressie and Johannesson (2008) and Katzfuss and Cressie (2009)), by introducing sparseness to  $\boldsymbol{\Sigma}$  (Furrer, Genton, and Nychka (2006)) and a combination of both approaches (Sang and Huang (2012)).

This paper focuses on efficient inference on linear mixed-effects models from a frequentists perspective, however another class of methods in the recent literature is the Integrated Nested Laplace Approximation (INLA) approach proposed by Rue, Martino, and Chopin (2009) imbedded in a Bayesian framework. Here, an efficient approximation of the posterior marginals of the elements of the latent field in a latent gaussian model is introduced, that is clearly superior in terms of computation speed compared to traditional simulation based MCMC schemes. In Eidsvik, Martino, and Rue (2009) this procedure was also successfully applied in a spatial generalized linear mixed model. Lindgren, Rue, and Lindström

(2011) show that, by using stochastic partial differential equations (SPDE), an explicit link between Gaussian fields and Gaussian Markov random fields (GRMF) can be established, at least when using the Matérn class of covariance functions. With the Markov property the involved precision matrix becomes sparse and sparse matrix algorithms can be applied. In Lasinio, Mastrantonio, and Pollice (2013) the SPDE approach is extended and includes an INLA approximation to further enhance computational feasibility. The authors also carry out a simulation study, in which they compare the SPDE/INLA approach with the covariance tapering applied in this paper in terms of their respective predictive performance. It is shown that the tapering approach can catch up with the SPDE/INLA approach. The INLA approach can be even applied for predictive process models, as in Eidsvik, O. Finley, Banerjee, and Rue (2012), in the spirit of the Fixed Rank Kriging. First a reduced rank spatial process is established, which aims at reducing the dimensionality of the model, and the INLA approximation is used to conduct Bayesian inference. In the recent literature the research is now focused on translating the outlined approximation approaches into a spatio-temporal context. In Stroud, Stein, Lesht, Schwab, and Beletsky (2010) a dynamic state-space model is proposed using an Ensemble Kalman Filter. The authors apply the covariance tapering in order to speed up computations and to reduce storage requirements for the full ensemble covariance and Kalman Gain matrices, as well for the approximation of the gaussian error random field.

## 2. Approximating the Spatial Covariance Function

### 2.1. Fixed Rank Kriging

As a way of dealing with the inversion of the  $n \times n$  covariance matrix in a large data setting, Cressie and Johannesson (2006, 2008) proposed to approximate the spatial process  $\nu(\cdot)$  in (2) by a vector  $\boldsymbol{\eta}$  of  $r$  random effects with  $r \ll n$  and a corresponding set of spatial basis functions  $\mathbf{S}(\cdot)$ . The model for  $\nu(\cdot)$ , which the authors call *spatial random-effects model*, is

$$\nu(\mathbf{s}) = \mathbf{S}(\mathbf{s})'\boldsymbol{\eta} + \xi(\mathbf{s}), \quad , \mathbf{s} \in D \quad (8)$$

and the corresponding smooth process  $Y(\cdot)$  becomes

$$Y(\mathbf{s}) = \mathbf{x}(\mathbf{s})'\boldsymbol{\alpha} + \mathbf{S}(\mathbf{s})'\boldsymbol{\eta} + \xi(\mathbf{s}) \quad , \mathbf{s} \in D \quad (9)$$

which results in a *spatial mixed effects model*, where  $\mathbf{S}(\mathbf{s}) \equiv (S_1(\mathbf{s}), \dots, S_r(\mathbf{s}))'$  is the set of  $r$  basis functions evaluated at location  $\mathbf{s} \in D$ , and  $\boldsymbol{\eta}$  is a  $r$ -dimensional zero-mean vector of random effects with  $r \times r$  dimensional covariance matrix  $\text{var}(\boldsymbol{\eta}) = \mathbf{K}$ . The zero-mean micro-scale variation process  $\xi(\cdot)$  with variance  $\sigma_\xi^2 v_\xi(\cdot)$  accounts for the spatial variation not explained by the dimension reduced model. Assuming that the micro-scale variation  $\xi(\cdot)$  is white-noise in space and that  $\boldsymbol{\eta}$  and  $\xi(\mathbf{s})$  are independent the covariance function of  $\nu(\cdot)$  is consequently

$$C(\mathbf{u}, \mathbf{v}) = \mathbf{S}(\mathbf{u})'\mathbf{K}\mathbf{S}(\mathbf{v}), \quad \mathbf{u}, \mathbf{v} \in D. \quad (10)$$

It is important to note, that no assumptions about stationarity or isotropy are made in the spatial random-effects model. The resulting theoretical covariance matrix  $\boldsymbol{\Sigma}$  of the data process is

$$\boldsymbol{\Sigma} = \mathbf{SKS}' + \sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi, \quad (11)$$

with  $\mathbf{S}$  being the  $n \times r$  dimensional matrix of basis functions evaluated at each observed location  $\mathbf{s}_1, \dots, \mathbf{s}_n$  and  $\mathbf{V}_\xi \equiv \text{diag}\{v_\xi(\mathbf{s}_1), \dots, v_\xi(\mathbf{s}_n)\}$  covering the heterogeneity of the small-scale spatial variation. This representation of the covariance matrix allows for the application of the Sherman-Morrison-Woodsbury formula as in (12) (Henderson and Searle (1981, p. 53))

$$(\mathbf{A} + \mathbf{U}\mathbf{B}\mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{B}\mathbf{V}\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{B}\mathbf{V}\mathbf{A}^{-1} \quad (12)$$

and consequently the inverse of (11) can be written as

$$\begin{aligned} \boldsymbol{\Sigma}^{-1} &= (\sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi)^{-1} \\ &\quad - (\sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi)^{-1} \mathbf{S} \{ \mathbf{K}^{-1} + \mathbf{S}'(\sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi)^{-1} \mathbf{S} \}^{-1} \mathbf{S}'(\sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi)^{-1}. \end{aligned} \quad (13)$$

Obviously using this representation only the inverse of the fixed-rank  $r \times r$  matrix  $\mathbf{K}$  and the diagonal  $n \times n$  matrix  $(\sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi)$  describing the nugget effect are needed. In that way great savings in terms of storage and reductions of computing time can be achieved. As stated in Cressie and Johannesson (2008, p. 214) the computational cost reduces from  $O(n^3)$  to  $O(nr^2)$  and accordingly rises only linear with the size of the dataset. With the setting of (8)-(13) the corresponding kriging prediction and variance for the prediction location  $\mathbf{s}_0 \in D$  are

$$\hat{\mathbf{Y}}(\mathbf{s}_0) = \mathbf{x}(\mathbf{s}_0)' \hat{\boldsymbol{\alpha}}_{gls} + \mathbf{c}_Y(\mathbf{s}_0)' \boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X} \hat{\boldsymbol{\alpha}}_{gls}) \quad (14)$$

and

$$\begin{aligned} \hat{\sigma}^2(\mathbf{s}_0) &= \mathbf{S}(\mathbf{s}_0)' \mathbf{K} \mathbf{S}(\mathbf{s}_0) - \mathbf{c}_Y(\mathbf{s}_0)' \boldsymbol{\Sigma}^{-1} \mathbf{c}_Y(\mathbf{s}_0) \\ &\quad + (\mathbf{x}(\mathbf{s}_0) - \mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{c}_Y(\mathbf{s}_0))' (\mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} (\mathbf{x}(\mathbf{s}_0) - \mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{c}_Y(\mathbf{s}_0)) \end{aligned} \quad (15)$$

with

$$\mathbf{c}_Y(\mathbf{s}_0) = \mathbf{S} \mathbf{K} \mathbf{S}(\mathbf{s}_0) + \sigma_\xi^2 v_\xi(\mathbf{s}_0) I(\mathbf{s}_0 \in \{\mathbf{s}_1, \dots, \mathbf{s}_n\}), \quad (16)$$

and  $\mathbf{S}(\mathbf{s}_0)$  is the  $r$  dimensional vector of basis functions evaluated at the prediction location  $\mathbf{s}_0$ .

### Basis Function Selection

Important for the spatial random-effects model in (8) is the specification of the basis functions. The application of basis functions to approximate non-stationary covariance functions has already been discussed in literature, e.g. in Nychka, Wikle, and Royle (2002). Popular classes of functions are smoothing spline basis functions (e.g. in Wahba (1990)), wavelet basis functions (e.g. in Vidakovic (1999)) and radial basis functions (e.g. in Hastie, Tibshirani, and Friedman (2003)). An overview of available classes is also given in Wikle (2010). Since the predictions in Section 3 are required for a sphere, namely the globe, the class of multi-resolutional local bisquare functions is used (as suggested in Cressie and Johannesson (2008))

$$S_{i,l}(\mathbf{s}) = \begin{cases} [1 - (\|\mathbf{s} - \mathbf{m}_{i,l}\| / r_l)^2]^2 & \text{if } \|\mathbf{s} - \mathbf{m}_{i,l}\| \leq r_l \\ 0 & \text{otherwise} \end{cases}, \quad (17)$$

where  $\mathbf{m}_{i,l}$  is the center point of the  $i$ th basis function in resolution level  $l$  and  $r_l$  is defined through Cressie and Johannesson (2006) as

$$r_l = (1.5) \times (\text{shortest distance between the center points in resolution level } l).$$

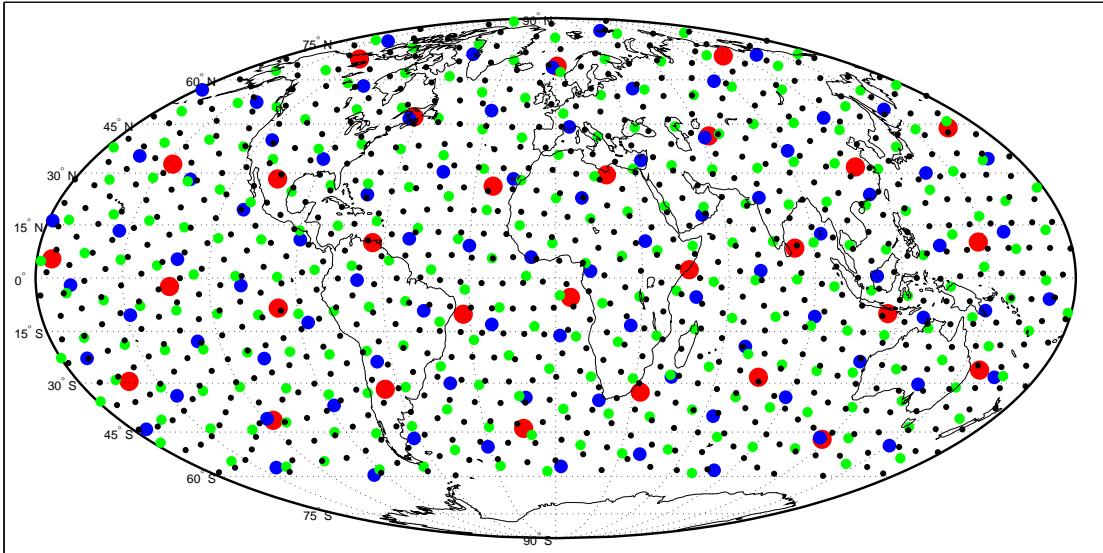


Figure 1: Basis function center points for 4 different resolutions of a Discrete Global Grid

By specifying multiple resolutions the covariance model is able to cover different scales of spatial variation. Along with the type of basis function the locations of the center points have to be specified. Ideally they should cover the entire domain and be equidistant. This can be achieved through the application of a multi-resolutional *Discrete Global Grid (DGG)* of hexagons (see [Sahr, White, and Kimerling \(2003\)](#)). In Figure 1 a DGG (ISEA3H<sup>1</sup>) was generated for the globe with 4 different resolutions, which is later used for the analysis of atmospheric carbon dioxide concentrations. Center points below 60 degrees South have been deleted, since no measurements of the satellite can be obtained from that area. The corresponding inter-cell distances, measured in great-arc distances, are 4430.85 km for resolution 1 (red dots) and 2558.15 km, 1476.95 km and 852.71 km for resolution 2 (blue dots), 3 (green dots) and 4 (black dots) respectively. Depending on how much resolutions are considered, there are 29, 116, 370 or 1127 basis functions to evaluate for each of the  $n$  observations, resulting in an increasing dimension of  $\mathbf{K}$ . Hence, there is a trade-off between the explained spatial variation and the computational advantages obtained through the basis function approximation.

### Parameter Estimation

The Fixed Rank Kriging approach requires estimates for the parameters  $\sigma_\epsilon^2, \sigma_\xi^2, \mathbf{K}$  and  $\boldsymbol{\alpha}$ . As already mentioned, a suitable estimator for  $\boldsymbol{\alpha}$  is the generalized-least squares estimator  $\hat{\boldsymbol{\alpha}}_{gls} = (\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{Z}$ . Since  $\sigma_\epsilon^2$  and  $\sigma_\xi^2$  are not individually identifiable,  $\sigma_\epsilon^2$  is assumed to be known and can be estimated through the semi-variogram at spatial lags close to zero using robust variogram estimates (see [Cressie and Hawkins \(1980\)](#)). The intercept of a fitted line using Weighted Least Squares (see [Cressie \(1985\)](#)) represents the estimate for  $\sigma_\epsilon^2$ . The parameters  $\mathbf{K}$  and  $\sigma_\xi^2$  of the spatial random-effects model can be estimated either through a

<sup>1</sup>Characteristics of different DGGs can be found on URL:<http://webpages.sou.edu/sahrk/dgg/isea/tables/tables.html>

Binned Method-of-Moments (BMoM) estimation procedure (as described in Cressie and Johannesson (2008)) or by Maximum Likelihood Estimation using an Expectation-Maximization (EM) algorithm (see Katzfuss and Cressie (2009)). In the BMoM approach the parameters  $\mathbf{K}$  and  $\sigma_\xi^2$  are estimated through minimizing a weighted Frobenius Norm between the theoretical covariance matrix  $\boldsymbol{\Sigma}$  and an empirical counterpart obtained by binning the data. However, as the authors in Katzfuss and Cressie (2009) state, BMoM estimation is inferior in providing accurate estimates of prediction uncertainty, is much more complicate to apply and requires many subjective decisions. Therefore our focus is on the EM algorithm and the BMoM procedure will not be described in detail. The interested reader is referred to Cressie and Johannesson (2008).

### *Maximum Likelihood Estimation via EM-Algorithm*

First, some distributional assumptions for the data have to be made in order to apply Maximum Likelihood estimation. For simplicity  $\mathbf{Z}$  represents the vector of detrended data in this case and it is assumed that  $\mathbf{Z}$  follows a multivariate normal distribution

$$\mathbf{Z} \sim N_n(\mathbf{0}, \mathbf{SKS}' + \sigma_\epsilon^2 \mathbf{V}_\epsilon + \sigma_\xi^2 \mathbf{V}_\xi).$$

Together with  $\sigma_\epsilon^2$ ,  $\mathbf{V}_\epsilon$  and  $\mathbf{V}_\xi$  assumed known, the Log-Likelihood becomes

$$\ell(\mathbf{K}, \sigma_\xi^2; \mathbf{Z}) \equiv \log f(\mathbf{Z}; \mathbf{K}, \sigma_\xi^2) = -\frac{1}{2} \log \det \boldsymbol{\Sigma} - \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{Z} \mathbf{Z}'). \quad (18)$$

However, as Katzfuss and Cressie (2009, p. 3381) state, finding estimates for  $\mathbf{K}$  and  $\sigma_\xi^2$  that maximize the likelihood equations is complicated. For that reason an EM algorithm (Dempster, Laird, and Rubin (1977)) is suggested. Instead of maximizing the likelihood of the observed data, it is assumed that knowing the distribution of some unobserved random variables, in this case  $\boldsymbol{\eta} \sim N_r(\mathbf{0}, \mathbf{K})$  and  $\boldsymbol{\xi} \sim N_n(\mathbf{0}, \sigma_\xi^2 \mathbf{V}_\xi)$  independently distributed, would result in the joint distribution function of both, the observed and the missing data, which in turn can be maximized much easier. The algorithm consists of two steps. The first one is the Expectation step, in which the conditional expectation of the complete-data likelihood at a certain value of the parameter vector  $\boldsymbol{\theta}^{[t]}$  at the  $t$ -th iteration, given the observed data, has to be calculated

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{[t]}) = E_{\boldsymbol{\theta}^{[t]}} \{ \log f(\boldsymbol{\eta}, \boldsymbol{\xi}; \boldsymbol{\theta}) | \mathbf{Z} \}. \quad (19)$$

In the Maximization step the parameters are updated, so that (19) is maximized, which results in the updating scheme

$$\mathbf{K}^{[t+1]} = \mathbf{K}^{[t]} + \mathbf{K}^{[t]} \left( \mathbf{S}' \boldsymbol{\Sigma}^{[t]-1} [\mathbf{Z} \mathbf{Z}' \boldsymbol{\Sigma}^{[t]-1} - \mathbf{I}_n] \mathbf{S} \right) \mathbf{K}^{[t]} \quad (20)$$

$$\sigma_\xi^{2[t+1]} = \sigma_\xi^{2[t]} + \sigma_\xi^{2[t]} \text{tr} \left( \frac{1}{n} \boldsymbol{\Sigma}^{[t]-1} [\mathbf{Z} \mathbf{Z}' \boldsymbol{\Sigma}^{[t]-1} - \mathbf{I}_n] \mathbf{V}_\xi \right) \sigma_\xi^{2[t]}. \quad (21)$$

For a thorough derivation of the Expectation and Maximization steps the reader is referred to Katzfuss and Cressie (2009). Both steps are repeated until a convergence criterion based either on the change in the maximized likelihood or on the change in the parameter values is fulfilled. The estimates found are solutions to the likelihood equations. However, the user has to make subjective decisions concerning the convergence criterion and the starting values for the iteration procedure, which might influence both the efficiency and the accuracy of

the algorithm. Furthermore the algorithm might lead to a local maximum depending on the choice of the initial values.

The outlined fixed rank kriging approach is very suitable for dealing with datasets of massive size. Due to the low dimensional spatial random effects vector  $\boldsymbol{\eta}$  inverting the data covariance matrix  $\boldsymbol{\Sigma}$  requires operations that rise only linear with the size of the dataset. In addition if the number of random effects  $r$  is sufficiently small no assumptions on the form of  $\text{var}(\boldsymbol{\eta}) = \mathbf{K}$  are necessary and the restricting assumptions of stationarity and/or isotropy can be avoided. Fixed rank kriging is able to cover spatial variations on larger scales with a small number of basis functions. However in order to capture many scales of spatial variation of the phenomenon finer resolutions of basis functions, and consequently a larger  $r$  are needed. This reduces the computational efficiency and introduces the need of parametric covariance functions for  $\mathbf{K}$  and simplifying assumptions. Obviously the choice of the basis functions is critical for the fixed rank kriging approach. Through choosing the number, the location and the type of the basis functions, the user is left with many subjective decisions, which possibly affect the outcome.

## 2.2. Covariance Tapering

Another way of efficiently dealing with  $\boldsymbol{\Sigma}^{-1}$  is to introduce sparseness. With  $\boldsymbol{\Sigma}$  and  $\mathbf{c}_Y(\mathbf{s}_0)$  being sparse, significant computational savings in calculating the kriging predictions and variances in (6) and (7) can be achieved. The operation  $\mathbf{u} = \boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\alpha}}_{\text{gls}})$  in (6) or  $\mathbf{u} = \boldsymbol{\Sigma}^{-1}\mathbf{c}_Y(\mathbf{s}_0)$  in (7) can be solved efficiently through sparse matrix techniques based on the Cholesky factorization  $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}'$ . Then solving the triangular systems  $\mathbf{Aw} = \mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\alpha}}_{\text{gls}}$  or  $\mathbf{Aw} = \mathbf{c}_Y(\mathbf{s}_0)$  respectively and  $\mathbf{A}'\mathbf{u} = \mathbf{w}$  yields the desired quantity. The computational complexity of the preceding operations is of order  $O(nk^2)$ , where  $k$  denotes the average number of non-zero elements in each row of  $\boldsymbol{\Sigma}$ , and consequently rises only linear with the size of the dataset.

Based on Furrer *et al.* (2006) sparseness can be introduced by setting covariances to zero for observations more than a specific distance apart. The intuition behind this, is that observations far from the prediction location are not expected to have a large influence on the prediction and can therefore be neglected. Another argument for restricting to a local neighborhood is that even if the process inhibits long-range spatial dependence, the conditional correlation is expected to be very small after observing a closely located neighbor, since most of the information was already covered in the correlation with the neighboring observation. Let  $C_{\boldsymbol{\theta}}(h)$  be a second order stationary and isotropic covariance function, with  $h = \|\mathbf{s}_i - \mathbf{s}_j\|$  and parameter vector  $\boldsymbol{\theta}$ . Using a taper function  $T(h, \gamma)$ , which is an isotropic and second order stationary covariance function with compact support, being equal to zero for  $h \geq \gamma$ , the tapered covariance function is the Schur product of  $C_{\boldsymbol{\theta}}(\cdot)$  and  $T(\cdot)$

$$C_{\text{tap}}(h, \gamma) = C_{\boldsymbol{\theta}}(h) \circ T(h, \gamma). \quad (22)$$

The tapered covariance function will also be a valid covariance function, since the Schur product of two positive definite matrices is again positive definite according to Horn and Johnson (1994, Theorem 5.2.1). An overview and some suggestions on choosing the type of taper function can be found in Furrer *et al.* (2006), including the spherical covariance function

and functions from the Wendland family

$$T_{spherical}(h, \gamma) = \left(1 - \frac{h}{\gamma}\right)_+^2 \left(1 + \frac{h}{2\gamma}\right), \quad h > 0, \quad (23)$$

$$T_{wendland,1}(h, \gamma) = \left(1 - \frac{h}{\gamma}\right)_+^4 \left(1 + 4\frac{h}{\gamma}\right), \quad h > 0, \quad (24)$$

$$T_{wendland,2}(h, \gamma) = \left(1 - \frac{h}{\gamma}\right)_+^6 \left(1 + 6\frac{h}{\gamma} + \frac{35h^2}{2\gamma^2}\right), \quad h > 0. \quad (25)$$

Furrer *et al.* (2006) also investigated the asymptotic behavior of the kriging estimators with the tapered covariance function and proved that under certain conditions the estimator is asymptotically equivalent to the one obtained by using the original covariance function.

### Parameter Estimation

The concept of tapering the covariance function not only improves the computational efficiency in kriging applications, in fact it can also be used in maximum likelihood estimation procedures, as it is demonstrated in Kaufman, Schervish, and Nychka (2008). In assuming multivariate normality a simple approximation of the log-likelihood function is obtained through replacing the model covariance matrix  $\Sigma(\boldsymbol{\theta})$  by its tapered counterpart  $\Sigma_{tap} = \Sigma(\boldsymbol{\theta}) \circ \mathbf{T}(\gamma)$ , where  $\mathbf{T}(\gamma)_{i,j} = T(||\mathbf{s}_i - \mathbf{s}_j||, \gamma)$ . However this may lead to biased estimates in practice for small values of  $\gamma$ , as Kaufman *et al.* (2008, p. 1546) points out. For that reason the authors suggest to apply a two-taper approximation, where both the model and the sample covariance matrix are tapered. In this case  $\mathbf{Z}$  represents the vector of detrended data and the one- and two-taper likelihood functions become

$$\ell_{1,taper}(\boldsymbol{\theta}) = -\frac{1}{2} \log \det(\Sigma_{tap}) - \frac{1}{2} \mathbf{Z}' \Sigma_{tap}^{-1} \mathbf{Z} \quad (26)$$

$$\ell_{2,taper}(\boldsymbol{\theta}) = -\frac{1}{2} \log \det(\Sigma_{tap}) - \frac{1}{2} \mathbf{Z}' (\Sigma_{tap}^{-1} \circ \mathbf{T}(\gamma)) \mathbf{Z}. \quad (27)$$

Maximizing  $\ell_{2,taper}(\boldsymbol{\theta})$  leads to unbiased estimators, but at the cost of an increased computational complexity, as the two-taper approximation involves calculating the full inverse  $\Sigma_{tap}^{-1}$ , whereas the simple approach (one-taper approximation) only requires solving the sparse system of equations  $\Sigma_{tap}^{-1} \mathbf{Z}$ . For the functional form of  $C_Y(h, \boldsymbol{\theta})$  used to generate  $\Sigma(\boldsymbol{\theta})$  there are many choices in the literature (see e.g. Cressie and Wikle (2011) for an overview) including the popular class of Matérn covariance functions (Matérn (1986)) defined as

$$C_Y(h, \sigma^2, \rho, \nu) = \frac{\sigma^2 (h/\rho)^\nu}{\Gamma(\nu) 2^{\nu-1}} K_\nu(h/\rho), \quad h \geq 0, \sigma^2, \rho, \nu > 0, \quad (28)$$

with  $K_\nu$  being the modified Bessel function of order  $\nu$  (see Abramowitz and Stegun (1964)),  $\sigma^2$  is the sill of the semi-variogram,  $\rho$  is a range parameter and  $\nu$  controls for the smoothness of the process.

Despite Maximum Likelihood Estimation there is also the possibility for Variogram-model fitting using the empirical semi-variogram. Assuming stationarity and isotropy of the semi-variogram and the covariance function of the process  $Y$ , their relationship can be established through

$$\gamma_Y(||h||, \boldsymbol{\theta}) = C_Y(0, \boldsymbol{\theta}) - C_Y(||h||, \boldsymbol{\theta}). \quad (29)$$

An empirical estimate of  $\gamma_Y(\cdot)$  (see Cressie and Wikle (2011, p. 131)) can be computed through

$$\begin{aligned}\hat{\gamma}_Y(h) &= \frac{1}{2} (\hat{\gamma}_Z(h) - \sigma_\epsilon^2) \\ &= \frac{1}{2} \left( \text{ave}\{(Z(\mathbf{s}_i) - Z(\mathbf{s}_j))^2 : \|\mathbf{s}_i - \mathbf{s}_j\| \in T(h); i, j = 1, \dots, n\} - \sigma_\epsilon^2 \right),\end{aligned}\quad (30)$$

where  $T(h)$  is a small tolerance region around  $h$ , resulting in a binning of the data. Using Least-Squares variogram fitting approaches a theoretical covariance function (e.g. (28)) can be fitted to the empirical semi-variogram by minimizing a weighted or non-weighted loss function (Cressie (1993)), given as

$$\text{loss}(\boldsymbol{\theta})_{OLS} = \sum_k (\hat{\gamma}_k - \gamma_k(\boldsymbol{\theta}))^2 \quad (31)$$

$$\text{loss}(\boldsymbol{\theta})_{WLS} = \sum_k n_k \left( \frac{\hat{\gamma}_k - \gamma_k(\boldsymbol{\theta})}{\gamma_k(\boldsymbol{\theta})} \right)^2, \quad (32)$$

where  $\hat{\gamma}_k$  denotes the value of the empirical semi-variogram for the  $k$ th bin,  $\gamma_k(\boldsymbol{\theta})$  is the value of the theoretical semi-variogram in the  $k$ th bin and  $n_k$  is the number of pairs in the  $k$ th bin. Through tapering the covariance function many zeroes are introduced to  $\boldsymbol{\Sigma}$  and the required operations rise only linear with the size of the dataset through the application of sparse matrix techniques. Disregarding covariances for locations with large distances only leads to a slight loss in predictive performance, since much of the information in the distant location is already covered in the correlation with closely located observations. In that way covariance tapering is able to efficiently capture spatial dependence at small spatial scales. However, depending on the choice of the taper range  $\gamma$  dependences at large scales are ignored and this might lead to a decline in predictive performance in regions with few data points. Another problem arises when using stationary covariance tapers on a non-stationary process. In this case small taper ranges are recommended in order to keep the bias small.

### 2.3. Full-Scale Approximation

Both approaches, the Fixed Rank Kriging and the Covariance Tapering, can be combined in a way, so that their advantages are fully exploited, as in Sang and Huang (2012). The former is able to efficiently capture the large-scale spatial dependence, since for that purpose only a small number of basis functions is needed (i.e. a small choice of  $r$ ), however in order to describe local behaviour the dimension of  $\mathbf{S}$  and  $\mathbf{K}$  has to be increased accordingly. In contrast the Covariance Tapering is efficient for the spatial dependence at small spatial scales (i.e. a small choice of  $\gamma$ ), whereas larger scales require larger taper ranges. In combining both approaches even small choices of  $r$  and  $\gamma$  are sufficient for providing a good approximation of the spatial dependence at the full scale. Consider the spatial random effects model in (8)

$$\nu(\mathbf{s}) = \mathbf{S}(\mathbf{s})' \boldsymbol{\eta} + \xi(\mathbf{s}), \quad , \mathbf{s} \in D,$$

where the residual process  $\xi(\cdot)$  from the Fixed Rank Kriging is no longer assumed to be independent in space, but instead has the following covariance function approximated through Covariance Tapering

$$C_\xi(\mathbf{u}, \mathbf{v}) = \{C_Y(\mathbf{u}, \mathbf{v}) - \mathbf{S}(\mathbf{u})' \mathbf{K} \mathbf{S}(\mathbf{v})\} \circ T(\mathbf{u}, \mathbf{v}, \gamma), \quad \mathbf{u}, \mathbf{v} \in D. \quad (33)$$

Consequently the covariance matrix  $\Sigma$  becomes

$$\Sigma = \mathbf{SKS}' + \mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon, \quad (34)$$

where  $\mathbf{C}_\xi$  denotes the sparse covariance matrix of the residual process  $\xi(\cdot)$  generated through (33). The approximate log-likelihood function, ignoring the constant term and assuming normality on  $\mathbf{Z}$  representing the vector of the detrended data, is

$$\ell(\boldsymbol{\theta}) = -\frac{1}{2} \log \det \{ \mathbf{SKS}' + \mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon \} - \frac{1}{2} \mathbf{Z}' \{ \mathbf{SKS}' + \mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon \}^{-1} \mathbf{Z}. \quad (35)$$

In evaluating the log-likelihood function, the inverse and the determinant of the  $n \times n$  covariance matrix  $\Sigma$  are needed. The form of (34) allows for efficiently inverting  $\Sigma$  using the Sherman-Morrison-Woodsbury formula.

$$\Sigma^{-1} = (\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon)^{-1} - (\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon)^{-1} \mathbf{S} \{ \mathbf{K}^{-1} + \mathbf{S}' (\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon)^{-1} \mathbf{S} \}^{-1} \mathbf{S}' (\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon)^{-1} \quad (36)$$

The determinant of (34) can be computed through

$$\det(\Sigma) = \det(\mathbf{K}^{-1} + \mathbf{S}' (\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon)^{-1} \mathbf{S}) \det(\mathbf{K}^{-1})^{-1} \det(\mathbf{C}_\xi + \sigma_\epsilon^2 \mathbf{V}_\epsilon) \quad (37)$$

In computing (36) and (37) only inverses and determinants of sparse  $n \times n$  and of  $r \times r$  matrices are needed, which have a computational complexity of  $O(nr^2 + nk^2)$ . Alternatively the model can be fitted in a two-step procedure, so that the fixed rank kriging is estimated first through the efficient EM-Algorithm outlined in Section 2.2 and covariance tapering is applied to the residual process afterwards by applying variogram-model fitting as in Section 2.3. This has the advantage of avoiding the computational demanding maximization of the full likelihood. With this setting the full-scale approximation is able to combine the capabilities of the fixed rank kriging and the covariance tapering and to overcome their individual weaknesses. Whereas the  $r$  dimensional spatial random effect  $\boldsymbol{\eta}$  captures large scale spatial dependence, the residual process  $\xi(\cdot)$  with tapered covariance function  $C_\xi$  efficiently describes local behavior. Importantly non-stationary and anisotropic behavior can be captured by the fixed rank part. However, an open problem remains the choice of the approximation parameters  $\gamma$  and  $r$  that lead to the most efficient outcome.

### 3. Efficiency Evaluation - Analysis of Atmospheric $CO_2$ Concentrations

#### 3.1. Data Description

The spatial dataset used for the comparative study consists of 12842 measurements of mid-tropospheric  $CO_2$  concentrations obtained from the Atmospheric InfraRed Sounder (AIRS) on board NASA's Aqua satellite. The unit of measurement is *ppm* corresponding to  $10^{-6}$  and denotes the number of  $CO_2$  molecules in one million parts of air. This Level-2 product (AIRX2STC)<sup>2</sup> contains observations at  $90 \times 90$  km nominal horizontal resolution at nadir, measured between  $-180^\circ$  and  $180^\circ$  longitude and  $-60^\circ$  and  $90^\circ$  latitude. However, in order to avoid change-of-support issues, it is assumed that measurements are at point support. The

<sup>2</sup>Level 2 and 3 products are freely downloadable at <http://disc.sci.gsfc.nasa.gov/AIRS/data-holdings>

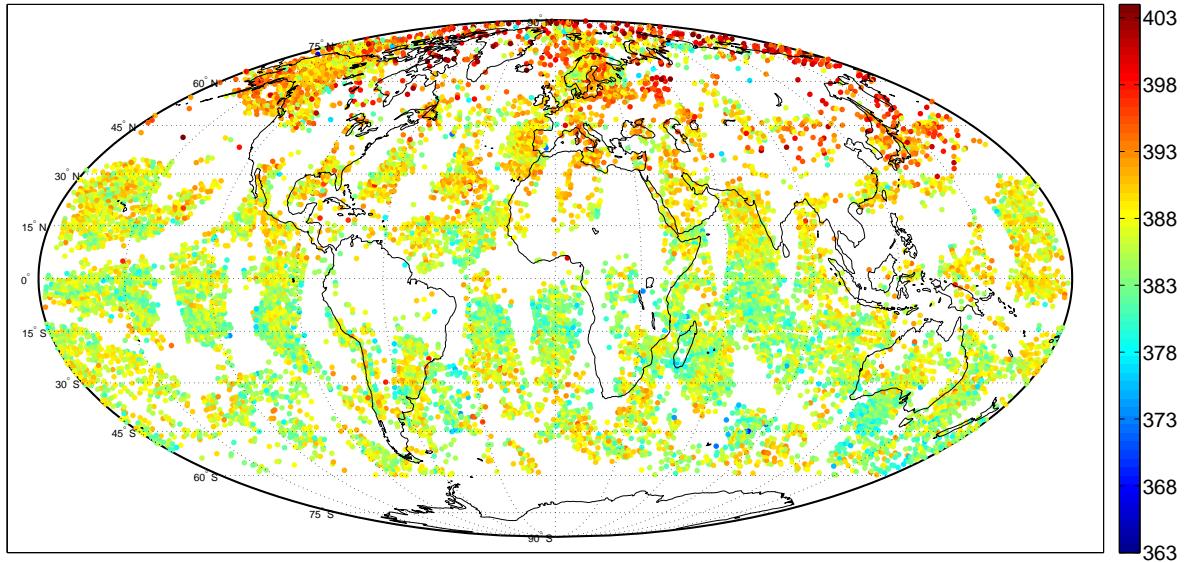


Figure 2: Mid-tropospheric CO<sub>2</sub> concentrations on the 1<sup>st</sup> of May 2009

Aqua satellite reaches global coverage twice a day. Since this study considers spatial-only processes, the data are treated as if they were taken at one specific time point, neglecting the time discrepancy between measurements. The dataset consists of observations taken on the 1<sup>st</sup> of May in 2009 and it already reveals some characteristic patterns of the natural carbon dioxide process. Particularly the data exhibits higher CO<sub>2</sub> concentrations and volatility in the northern hemisphere corresponding to a seasonal pattern caused by the growth stage of plants occurring in springtime.

### 3.2. Preliminary Steps

In order to evaluate the efficiency of the different approaches in approximating the spatial covariance function, it has to be ensured that the respective model assumptions are fulfilled. In Figures 3a, 3b and 3c the data is plotted against the degrees in longitude and/or latitude direction. Obviously the CO<sub>2</sub> process evolves differently in space depending on the orientation. Whereas a trend pattern can hardly be identified in the East-West direction, CO<sub>2</sub> concentrations tend to rise with decreasing distance to the north pole. Figure 3b also indicates that this latitude trend is of a non-linear kind. For that reason a non-linear regression was performed using polynomials up to order 3 for the latitude direction as covariates and a linear trend is assumed for the longitude direction. The subsequent analysis will be based on the detrended data. In addition the process variance appears to be higher in the northern hemisphere, as can be seen in Figure 3b. In Figures 3d and 3e empirical directional variograms of the detrended data with orientation 0° (North), 45° (North-East), 90° (East) and 135° (South-East) for the northern and southern hemisphere are shown. The empirical directional variograms were generated by using a tolerance angle of 22.5°. Importantly, since the spatial process evolves over the globe (for simplicity it is assumed that the earth is a perfect sphere with radius  $R = 6371\text{km}$ ) great-circle distances have to be used. As can be seen in Figure 3e the empirical directional variograms for the data in the southern hemisphere are quite similar, indicating that deviations from isotropy can be regarded as small. However due to the seasonal effect in the northern hemisphere caused by the growth of plants in springtime

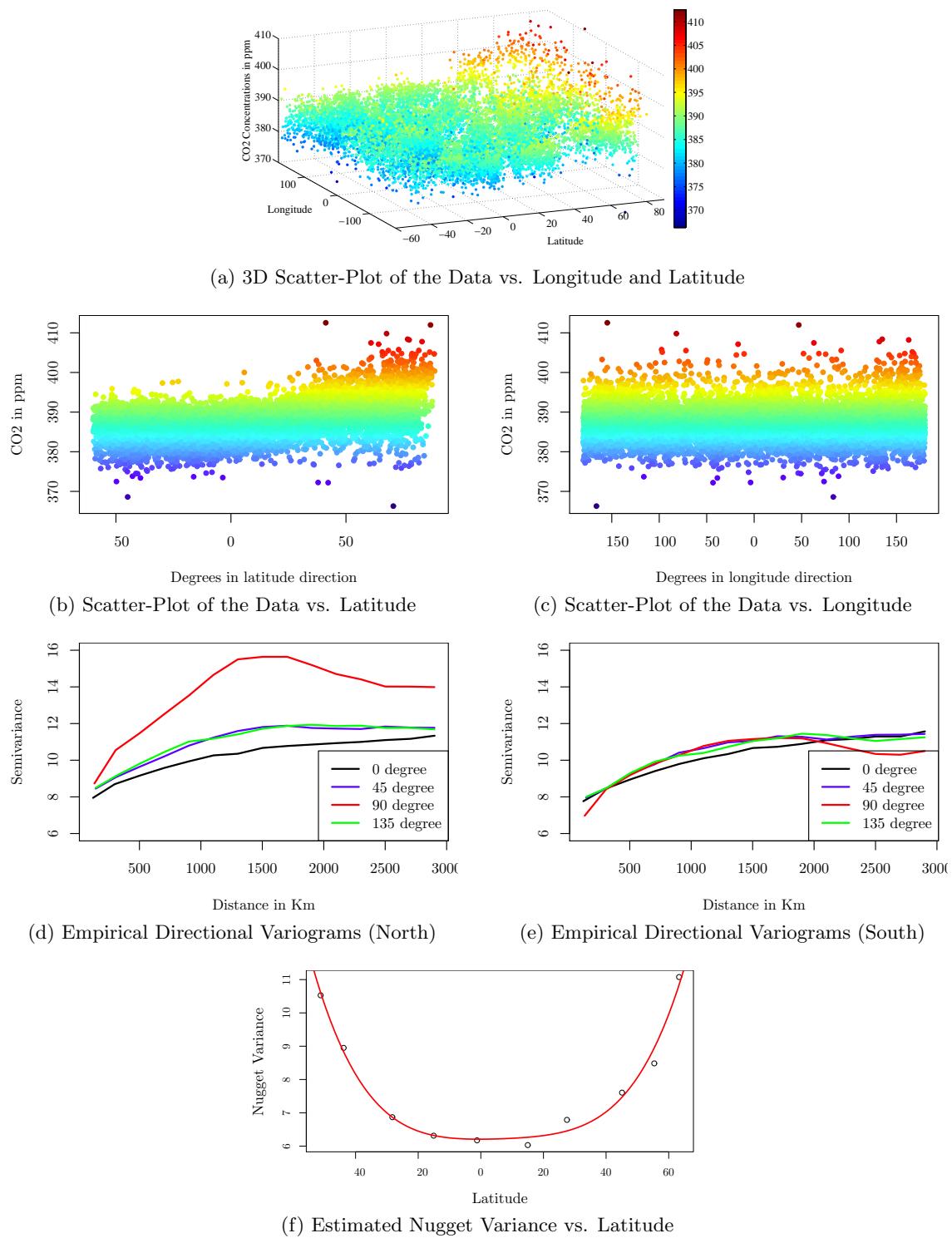


Figure 3: Exploratory Data Analysis

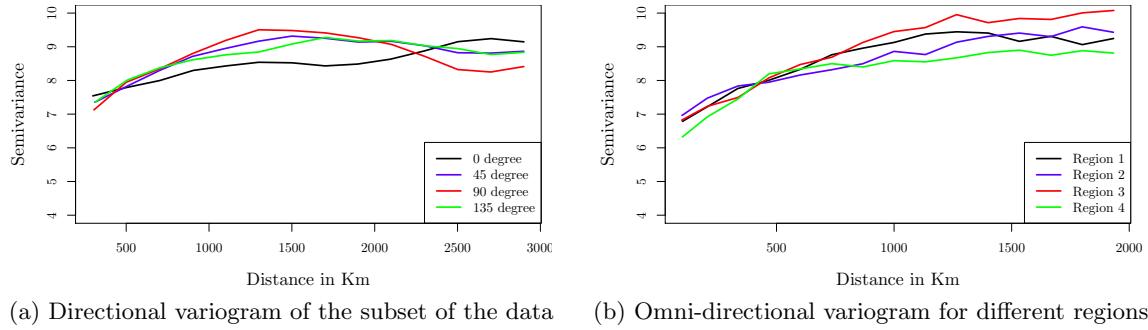


Figure 4: Spatial dependence structure in the subset of the data

the volatility is much higher. This can also be seen in the variograms in Figure 3d, where the sill is much higher in the East-West direction. In effect, the variogram becomes increasingly anisotropic with increasing degrees in latitude direction, leading to a non-stationary behavior of the process. Another indicator for non-stationarity is apparent in Figure 3f showing the estimated variance of the measurement error process depending on the degrees in latitude direction. Estimates for the nugget variance were computed for several subregions along the latitude direction using the approach mentioned in Section 2.2. Obviously measurements of the satellite become increasingly noisy the closer they are located to the poles. To account for this heterogeneity of the nugget variance a nonlinear regression was performed using polynomials up to order 4. This was used to provide values for  $\mathbf{V}_\epsilon$  in (5). As has been shown, the spatial process is characterized by a non-stationary dependence structure and a variogram that varies with the orientation, depending on the degrees in latitude direction. Consequently the stationary covariance tapering can be regarded as inappropriate and adjustments to non-stationary and anisotropic covariance functions are needed for spatial predictions on a global scale. In contrast, the fixed rank kriging approximation is able to work without these assumptions and is therefore suitable for this problem. Likewise the full-scale approximation can handle non-stationary and anisotropic processes through its fixed rank part. However, to ensure comparability of the outlined approaches, a subset of the data consisting of 5073 measurements between  $-20^\circ$  and  $20^\circ$  latitude is considered first. For this region around the equator deviations from stationarity and isotropy can be neglected, as can be seen in Figure 4. The empirical directional variograms of the subset of the data (Figure 4a) indicate a comparable spatial dependence structure for all directions and as Figure 4b suggests, this property holds irrespectively of the location. The empirical omni-directional variograms were calculated at four equally spaced reference regions within the subset of the data. In addition, the variance of the measurement error of the instrument can be considered as constant in the subset of the data as can be seen in Figure 3f. In that way the subset serves as a stationary scenario for evaluating the efficiency of the approximation approaches, which will be compared to the case, when these assumptions are not fulfilled.

### 3.3. Comparative Study

The focus of the study is to compare the efficiency of the approaches outlined in Section 2 in approximating the spatial covariance function by relating their predictive performance with the corresponding demand in computational resources. Both quantities are directly affected

by the choice of the number of basis functions  $r$  and/or the taper range  $\gamma$ . Increasing values result in a higher approximation quality but at the same time in higher computing times and storage requirements. Consequently it is of interest, which approach is able to solve this trade-off best. For the choice of the basis functions either 1, 2 or 3 resolutions from the DGG in Figure 1 are selected, resulting in 10, 42 or 132 basis functions for the subset and in 29, 166 or 370 basis functions for the complete dataset, respectively. The taper range  $\gamma$  is varied between 50km and 1500km. This results in 3 Fixed Rank Kriging models, 10 Covariance Tapering models and in another 30 Full-scale approximations covering each parameter combination. The predictive performance is evaluated by a series of cross-validation experiments. For each model a 10-fold cross-validation is performed, where the dataset has been divided randomly into 10 subsamples. In each round one subsample is retained as a validation set for testing purposes and the remaining subsamples are used to fit the model. This procedure is repeated 10 times, so that each observation was part of the validation set once. The predictions of the validation set can then be compared to the original data to construct out-of-sample performance measures, whereas the MSPE will be used in this study

$$MSPE \left( \hat{Y}(\mathbf{s}_0) \right) = \frac{1}{m} \sum_{i=1}^m \left( \hat{Y}(\mathbf{s}_i) - Y(\mathbf{s}_i) \right)^2 \quad , \mathbf{s}_0 = (\mathbf{s}_1, \dots, \mathbf{s}_m) . \quad (38)$$

However, the MSPE has to be adjusted, since predictions are based on the smooth process  $Y(\cdot)$  but only the noisy process  $Z(\cdot)$  is observed and consequently the squared residuals would be affected by the measurement error variance. Recall that  $\mathbf{Z} = \mathbf{Y} + \boldsymbol{\epsilon}$  and that  $var(\boldsymbol{\epsilon}) = \sigma_\epsilon^2 \mathbf{V}_\epsilon$  was assumed known and fitted through a polynomial function of order 4, so that the correct representation of the MSPE in the presence of measurement error can be obtained by subtracting the location specific nugget variance from the squared residual (see Cressie (1993, p. 128))

$$MSPE \left( \hat{Y}(\mathbf{s}_0) \right) = \frac{1}{m} \sum_{i=1}^m \left\{ \left( \hat{Y}(\mathbf{s}_i) - Z(\mathbf{s}_i) \right)^2 - \sigma_\epsilon^2 v_\epsilon(\mathbf{s}_i) \right\} \quad , \mathbf{s}_0 = (\mathbf{s}_1, \dots, \mathbf{s}_m) . \quad (39)$$

Besides the predictive performance, it is also of interest how much computational resources have been used by the models. In particular the computing time needed for calculating the important quantities in kriging predictions and in likelihood maximizations, which are the solution of the system of linear equations  $\boldsymbol{\Sigma}^{-1} \mathbf{Z}$  and the determinant of  $\boldsymbol{\Sigma}$ , is monitored. Furthermore, the maximum amount of working memory used in these calculations is recorded, disregarding all preliminary calculations. However, it has to be noted that depending on how much prediction locations  $\mathbf{s}_0$  are considered, the operation  $\mathbf{c}_Y(\mathbf{s}_0) \boldsymbol{\Sigma}^{-1}$  might also need significant amounts of working memory, especially for smooth prediction surfaces.

### 3.4. Subset Results - Stationary scenario

For the subset of the data, which serves as a stationary scenario, the trade-off between predictive performance and demand in computational resources is visualized in Figure 5. In Figure 5a the time in seconds needed to calculate the important quantities  $\boldsymbol{\Sigma}^{-1} \mathbf{Z}$  and  $\det \boldsymbol{\Sigma}$  is plotted against the MSPE and in Figure 5b the maximum amount of working memory in GB is shown. In comparing the fixed rank kriging (black dots) with the covariance tapering (blue line) it can be seen that the latter approach is more efficient in approximating the spatial covariance function, since for every level of the MSPE less or equal time and memory is

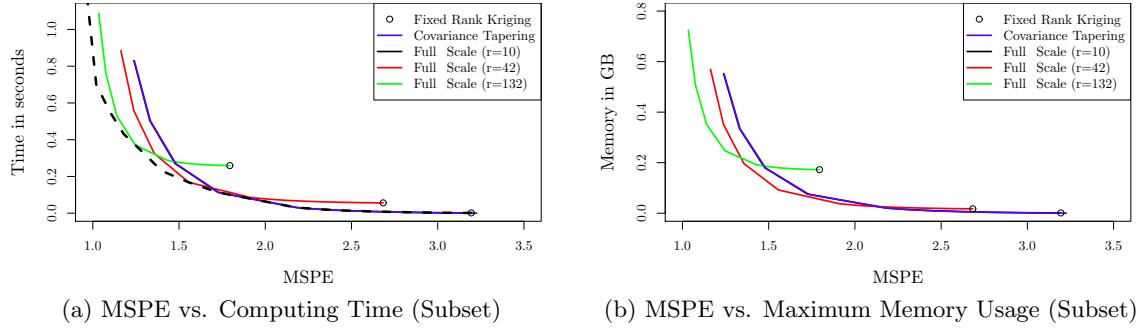


Figure 5: Efficiency Evaluation of the Covariance Approximation Approaches (Subset)

	Full Model	Full-scale Approximation r=132, $\gamma=1500$	Full-scale Approximation r=42, $\gamma=625$	Fixed Rank Kriging r=132	Fixed Rank Kriging r=42	Covariance Tapering $\gamma=750$	Covariance Tapering $\gamma=625$
MSPE	0.975	1.034	1.734	1.795	2.684	1.734	1.963
Time $\Sigma^{-1}\tilde{\mathbf{Z}}$	11.086	0.76104	0.06663	0.18632	0.01967	0.06561	0.04159
Time det $\Sigma$	12.726	0.32702	0.03035	0.07271	0.00656	0.04635	0.02916
Memory $\Sigma^{-1}\tilde{\mathbf{Z}}$	2.464	0.72400	0.06453	0.17236	0.01745	0.07450	0.04708
Memory det $\Sigma$	2.365	0.23800	0.02121	0.05666	0.00574	0.02449	0.01548

Table 1: Efficiency Evaluation - Subset

needed. However it has to be noted, that this result strongly depends on the range of spatial dependence relative to the total extent of the spatial domain and the spatial distribution of the data locations. As denoted earlier, covariance tapering has advantages in describing local and fixed rank kriging in large-scale dependencies. Accordingly having a process with a small range of spatial dependence in relative terms will result in efficiency advantages for the covariance tapering. In contrast a high proportion of clustered data decreases the sparsity of  $\Sigma$  and increases the demand in computational resources for the covariance tapering. To overcome the individual weaknesses and to exploit the advantages of the fixed rank kriging and the covariance tapering their combination in a full-scale approximation leads to further efficiency gains, as can be seen in Figures 5a and 5b. For lower approximation qualities the full-scale approximation (black line) is slightly more efficient as the covariance tapering (blue line), however in order to achieve lower values of the MSPE it is worth including higher resolutions of basis functions in the full-scale approximation (red and green line) to further reduce the computational complexity. The complete summary of results for the 43 different models is shown in Table 3 in the Appendix, whereas a characteristic snapshot is shown in Table 1. To evaluate the overall quality of the approximations, the results for the full model without any approximation, i.e. the model with an untapered Matérn covariance function, as a baseline are shown. The full model achieved a MSPE of 0.975 and the approximation that came closest to that level is the full-scale approximation with 132 basis functions and a taper range of 1500km. As can be seen, almost the same predictive performance can be accomplished, but about 20 times faster and with only around 30% of the maximum working memory required. Table 1 also compares the efficiency of the approximations for a fixed level of the MSPE of about 1.75. Clearly the full-scale approximation outperforms the other approaches in terms of speed and storage, whereas the advantage over the covariance tapering is rather small for

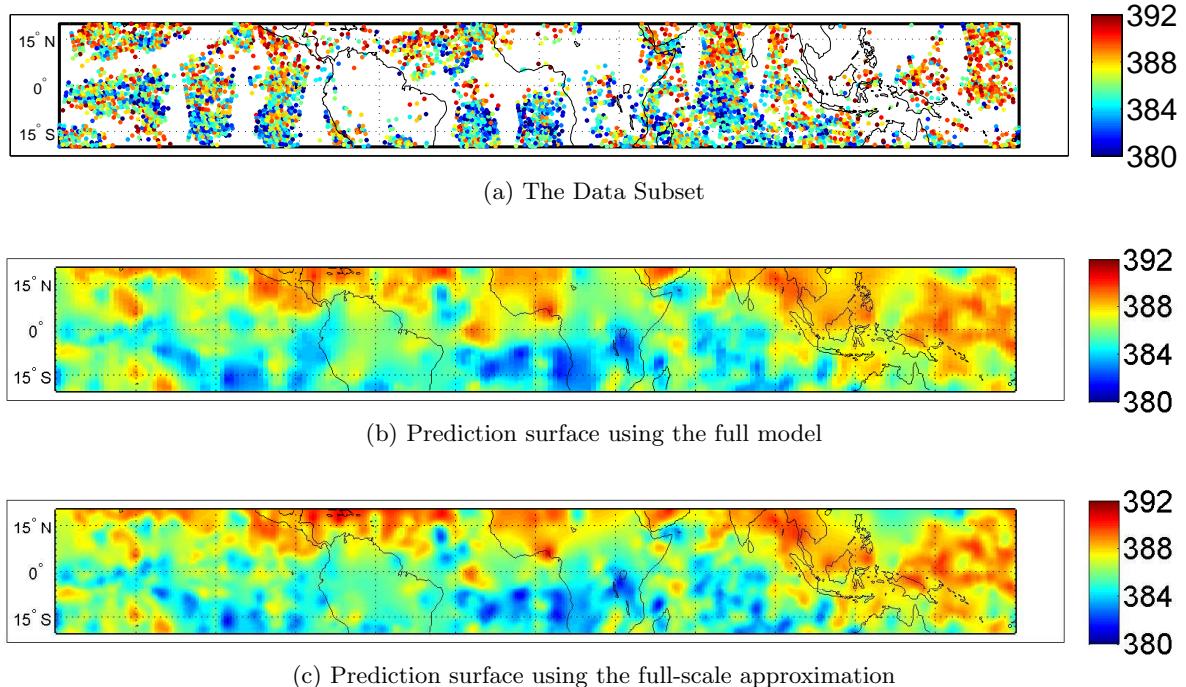


Figure 6: Prediction Surfaces of the Subset

this level of the MSPE, but becomes much larger for better approximation qualities. Finally Figure 6 shows kriging surfaces of the full model (Figure 6b) and the full-scale approximated model with  $r = 132$  and  $\gamma = 1500\text{km}$  (Figure 6c), whereas 250000 prediction locations were used to produce the latter and, due to memory restrictions, only 40000 pixels can be produced for the full model. As the similarity of both plots indicate, the quality of the approximation is very good and comes together with remarkably high computational savings.

### 3.5. Global Dataset Results - Non-stationary Scenario

As outlined in Section 3.2 the global dataset is characterized by a non-stationary and anisotropic dependence structure, which affects the efficiency of the stationary covariance tapering. Compared to the stationary case, the efficiency curves (blue lines) are shifted to the right in Figures 7a and 7b. The stationary covariance tapering is not able to provide high quality approximations of the non-stationary covariance function, because the increasing process variance in the northern hemisphere is not captured. Consequently estimated prediction errors will not yield reliable estimates of the prediction uncertainty. In contrast the Fixed Rank Kriging (black dots) is able to account for the spatially varying dependence structure and yields lower values of the MSPE. Nevertheless the stationary covariance tapering is still more efficient for low approximation qualities. The directional semi-variograms in Figure 3d already revealed the anisotropic and non-stationary character of the dependence structure by the increased process variance in the northern hemisphere. However, at small spatial scales the spatial dependence patterns are very similar, despite the differing nugget variances. Consequently a stationary covariance tapering is still capable of providing good approximations of the covariance function in a local neighborhood although a non-stationary behavior is apparent at larger spatial

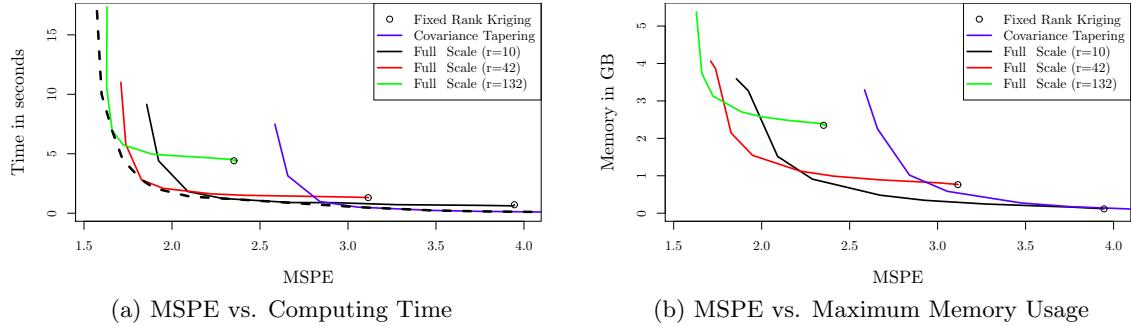


Figure 7: Efficiency Evaluation of the Covariance Approximation Approaches

	Full Model	Full-scale Approximation $r=29, \gamma=625$	Approximation $r=29, \gamma=750$	Fixed Rank Kriging $r=370$	Kriging $r=29$	Covariance Tapering $\gamma=1500$	Covariance Tapering $\gamma=750$
MSPE	2.556	2.480	2.288	2.352	3.946	2.658	3.055
Time $\Sigma^{-1}\tilde{\mathbf{Z}}$	155.6	0.620	0.736	2.911	0.396	1.881	0.317
Time $\det \Sigma$	186.0	0.399	0.453	1.499	0.264	1.259	0.194
Memory $\Sigma^{-1}\tilde{\mathbf{Z}}$	14.18	0.696	0.911	2.347	0.115	2.254	0.584
Memory $\det \Sigma$	13.79	0.210	0.272	1.058	0.046	0.815	0.230

Table 2: Efficiency Evaluation - Global Dataset

scales. In that way, a full-scale approximation is again able to further increase efficiency and to supply high quality approximations of the spatial covariance function at all spatial scales, as indicated by the efficiency curves (black, red and green lines) in Figure 7. A complete summary of the results of all models can be found in Table 4 in the Appendix, whereas a characteristic snapshot is shown in Table 2. As in the subset, the results of the full model with an untapered stationary Matérn covariance function are provided for comparative purpose. However, analogously to the stationary covariance tapering, it does not show a high predictive performance and results in a MSPE of 2.556 accompanied by high computation times and a huge amount of 14 GB of working memory. Using the Full-scale approximation with  $r = 29$  and  $\gamma = 625$  a comparable value can be obtained about 335 times faster and with only about 5% of the memory required at maximum. For comparing the efficiency of the approximation approaches a fixed level of about 2.3 of the MSPE is considered. Again the full-scale approximation was superior in terms of efficiency compared to both single approaches and the lead is even more pronounced in the non-stationary scenario than for the subset. Finally the full-scale approximation can be used to compute a high quality prediction surface for the process of atmospheric  $CO_2$  concentrations over the globe, as it is shown in Figure 8, where a kriging surface containing 250000 prediction locations for the full-scale approximation with  $r = 370$  and  $\gamma = 1500\text{km}$  was produced.

### 3.6. Choice of the covariance function

The predictive performance of both the covariance tapering and the full-scale approximation approach is directly affected by the choice of the underlying covariance model. Hence, the robustness of the obtained results has to be checked. For that purpose the cross-validation

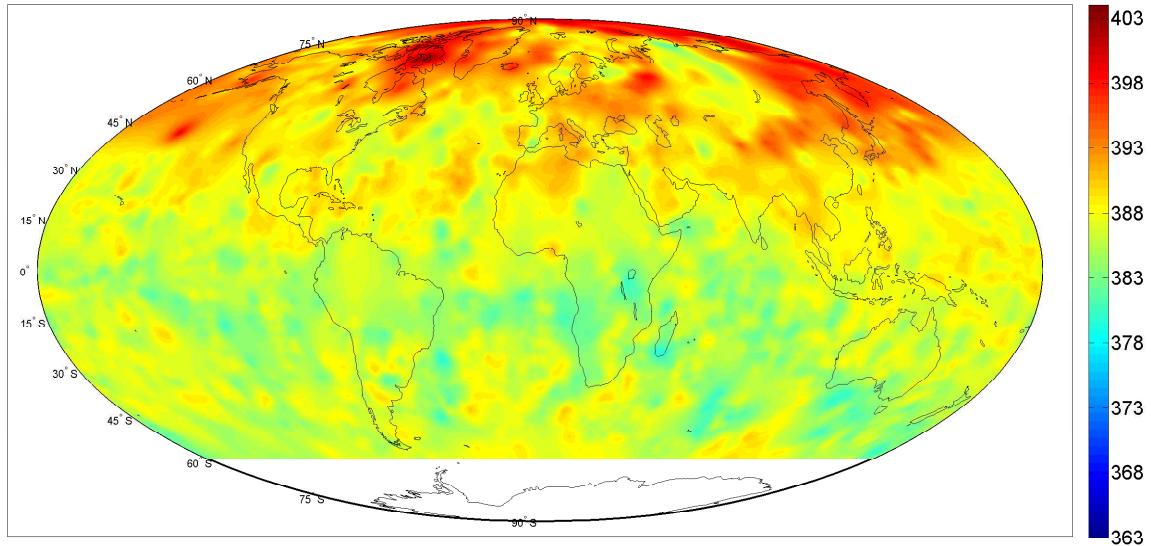


Figure 8: Mid-tropospheric CO<sub>2</sub> concentrations on the 1<sup>st</sup> of May 2009

study for the global dataset was also performed for two other popular choices of covariance functions, namely the spherical and the exponential model, which can be found for example in Cressie (1993, p. 61). In Figure 9 the trade-off between computation time and predictive performance for all approximation approaches is shown, whereas the colors black (Matérn), red (Exponential) and green (Spherical) represent the different underlying covariance models. As can be seen, the efficiency of the approximation approaches is hardly changed by altering the covariance model. The main results still hold true. In effect, the covariance tapering is still more efficient at lower and fixed rank kriging at higher approximation qualities and a combination of both approaches in a full-scale approximation is always superior in terms of efficiency. However, some differences can be identified, with the spherical model yielding better results than the exponential model and the Matérn showing the worst performance. Obviously parameter parsimony in the covariance model is more important than flexibility in the variogram fit for this dataset.

### 3.7. Choice of the taper function

Another factor influencing the efficiency of the approximation approaches is the choice of the taper function. In Figure 10 the corresponding results of the comparative study are shown for 3 types of taper functions, which were already introduced in Section 2.2, namely the Spherical and the Wendland taper functions of order 1 and 2. For this analysis an exponential model was used for calculating the covariances. As can be seen, changing the taper function only leads to small changes in the overall efficiency of the approximation approaches. Again the main results are still valid and there is a tendency for parsimony to be more important than flexibility of fitting the covariance function. Using the spherical taper function was always superior in terms of efficiency than the Wendland type functions, whereas the higher order Wendland function performed worst.

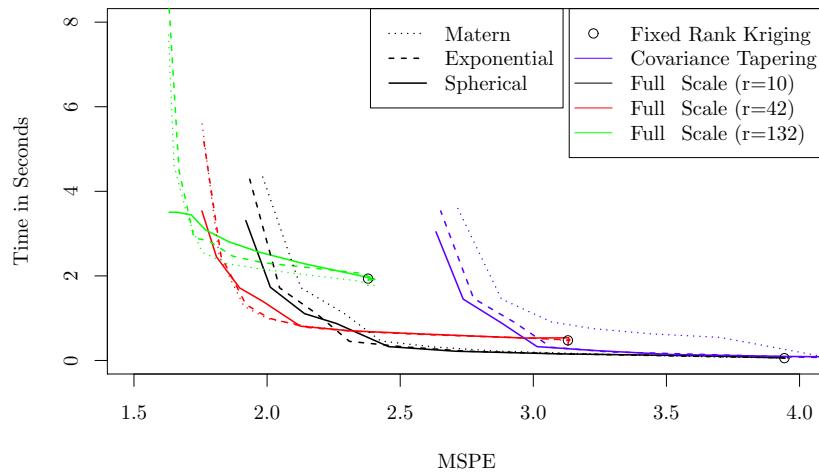


Figure 9: Efficiency evaluation for different choices of the covariance function

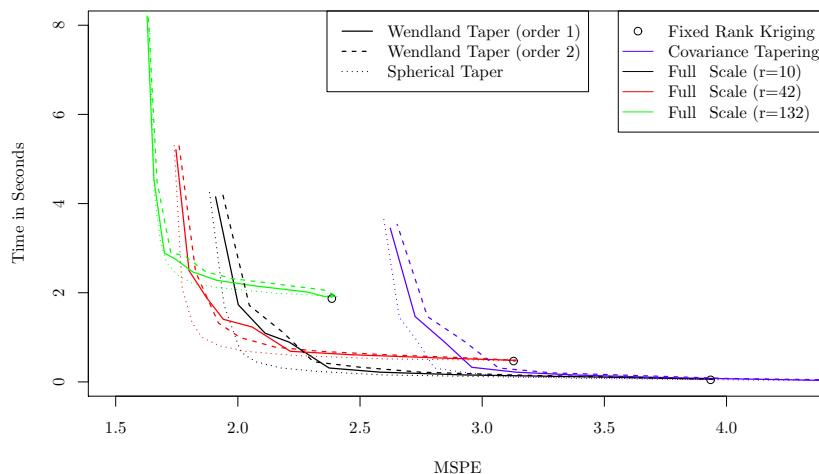


Figure 10: Efficiency evaluation for different choices of the taper function

### 3.8. Choice of the approximation parameters

The trade-off between predictive performance and computational complexity can be directly controlled through the choice of the approximation parameters  $r$  and  $\gamma$  in the full-scale approximation. In Figures 5a and 7a this trade-off was illustrated for a fixed number of basis functions  $r$  through the black, red and green lines. However these Figures can also give an idea on how the overall efficiency curve of the full-scale approximation would look like, as the enveloping black dotted lines sketch. These curves depict the optimal combinations of  $r$  and  $\gamma$  yielding the lowest achievable computational complexity at all levels of the MSPE. Interestingly the tangency points of the black, red and green lines on the hypothetical black dotted line correspond to a certain taper range  $\gamma^*$  between 750km and 1000km. Consequently Full-scale approximations with a taper range higher/lower than  $\gamma^*$  are always dominated through models with more/less basis functions. Moreover  $\gamma^*$  appears to coincide with the estimated effective range of the fitted Matérn covariance function. Intuitively this also makes sense, since this is exactly the scale of spatial dependence where the covariance tapering has advantages over the fixed rank kriging.

## 4. Conclusions

This paper investigated approaches to approximate the spatial covariance function and analyzed the trade-off between the loss in information due to the approximation and the reductions in computational complexity. Based on a remotely sensed data set of carbon dioxide concentrations in the mid-troposphere an efficiency evaluation was conducted, monitoring the predictive performance through the MSPE and the computational complexity through computation speed and storage requirements. All outlined approaches, namely fixed rank kriging (Cressie and Johannesson (2006, 2008)), covariance tapering (Furrer *et al.* (2006)) and the full-scale approximation (Sang and Huang (2012)) were able to notably speed up the calculations of the important quantities in maximum likelihood estimation and in kriging predictions, which are the determinant of the covariance matrix of the observed data  $\Sigma$  and the solution of the system of linear equations  $\Sigma^{-1}\mathbf{Z}$ . The required computations rise only linear with the size of the data set, instead of cubic. However, depending on the degree of the approximation, controlled by parameters  $r$  as the number of random effects in the fixed rank kriging approach and  $\gamma$  as the taper range in the covariance tapering approach, the loss in predictive performance differs substantially. In a subset of the data, where the process can be regarded as stationary, it was shown that covariance tapering outperformed the fixed rank kriging. However through combining both approaches in a full-scale approximation even more efficient approximations can be generated. The individual weaknesses, namely the inefficiency of the fixed rank kriging to describe local spatial dependence and of the covariance tapering to cover large-scale spatial dependence, can be overcome. In the full data set, involving a strong non-stationary behavior in the latitude direction, the advantage of the fixed rank kriging is apparent, since no assumptions on stationarity and/or isotropy have to be made and therefore the increased process variance in the northern hemisphere can be easily captured. This feature also translates into the full-scale approximation. Interestingly the analysis gives an idea on how to choose the approximation parameters  $r$  and  $\gamma$  optimally. For each level of the MSPE the most efficient combination of parameters involves a certain taper range  $\gamma^*$ , which coincides with the effective range of the fitted Matérn covariance function for the  $CO_2$  example. However a thorough investigation of the optimal choice of the approximation parameters, i.e.

model selection, is left open for future research. The spatial random-effects model (see Section 2.1) has also been extended into a spatio-temporal random effects model. In Cressie, Shi, and Kang (2010),Katzfuss and Cressie (2011) and González-Manteiga, Crujeiras, Katzfuss, and Cressie (2012) the low-dimensional latent random process is represented by a state-space model following a vector autoregressive process of order 1, whereas the first paper takes a filtering perspective, the second paper introduces a smoothing algorithm and the last paper gives a Bayesian solution to the smoothing problem. Finally an extension for the full-scale approximation is currently under investigation in Zhang, Sang, and Huang (2013), where the space-time covariance function is approximated through a large-scale spatio-temporal process of low rank and a small-scale spatio-temporal remainder component, which is subject to covariance tapering.

## A. Tables

Table 3: Efficiency Evaluation - Subset

Model Type	Parameters		MSPE	Time in sec $\Sigma^{-1}\tilde{\mathbf{Z}}$	Time in sec $\det \Sigma$	Memory in GB $\Sigma^{-1}\tilde{\mathbf{Z}}$	Memory in GB $\det \Sigma$
Full Model			0.975	11.0866	12.7267	2.4645	2.3657
Fixed Rank Kriging	10	0	3.194	0.00111	0.00037	0.00099	0.00033
	42	0	2.684	0.01967	0.00656	0.01745	0.00574
	132	0	1.795	0.18632	0.07271	0.17236	0.05666
Covariance Tapering	0	50	3.216	0.00001	0.00001	0.00001	0.00001
	0	100	3.186	0.00005	0.00003	0.00005	0.00002
	0	200	3.055	0.00067	0.00041	0.00072	0.00024
	0	300	2.783	0.00304	0.00185	0.00325	0.00107
	0	400	2.466	0.00826	0.00526	0.00900	0.00296
	0	500	2.192	0.01758	0.01197	0.01966	0.00646
	0	625	1.963	0.04159	0.02916	0.04708	0.01548
	0	750	1.734	0.06561	0.04635	0.07450	0.02449
	0	1000	1.482	0.15570	0.11324	0.17895	0.05883
	0	1500	1.241	0.47393	0.35510	0.55164	0.18134
Full-scale Approximation (r=10)	10	50	3.184	0.00113	0.00038	0.00100	0.00033
	10	100	3.155	0.00115	0.00041	0.00104	0.00034
	10	200	3.026	0.00182	0.00074	0.00171	0.00056
	10	300	2.759	0.00438	0.00200	0.00424	0.00140
	10	400	2.447	0.01292	0.00209	0.00999	0.00328
	10	500	2.177	0.02048	0.01055	0.02065	0.00679
	10	625	1.951	0.04693	0.02531	0.04807	0.01580
	10	750	1.726	0.07337	0.04008	0.07549	0.02481
	10	1000	1.477	0.16990	0.10052	0.17994	0.05915
	10	1500	1.239	0.51520	0.31532	0.55263	0.18167
Full-scale Approximation (r=42)	42	50	2.676	0.01968	0.00656	0.01746	0.00574
	42	100	2.654	0.01958	0.00673	0.01750	0.00575
	42	200	2.558	0.02024	0.00706	0.01817	0.00597
	42	300	2.358	0.02266	0.00845	0.02070	0.00681
	42	400	2.121	0.02801	0.01173	0.02645	0.00869
	42	500	1.912	0.03929	0.01648	0.03711	0.01220
	42	625	1.734	0.06663	0.03035	0.06453	0.02121
	42	750	1.556	0.09396	0.04422	0.09195	0.03023
	42	1000	1.356	0.19425	0.10091	0.19640	0.06456
	42	1500	1.163	0.55367	0.30158	0.56909	0.18708
Full-scale Approximation (r=132)	132	50	1.791	0.18665	0.07240	0.17237	0.05666
	132	100	1.780	0.18651	0.07260	0.17241	0.05668
	132	200	1.736	0.18797	0.07214	0.17308	0.05690
	132	300	1.644	0.19186	0.07206	0.17561	0.05773
	132	400	1.531	0.19894	0.07361	0.18136	0.05962
	132	500	1.429	0.21101	0.07757	0.19202	0.06312
	132	625	1.338	0.23900	0.09079	0.21944	0.07214
	132	750	1.247	0.26698	0.10401	0.24686	0.08115
	132	1000	1.140	0.37446	0.15351	0.35131	0.11549
	132	1500	1.034	0.76104	0.32702	0.72400	0.23800

Table 4: Efficiency Evaluation - Global Dataset

Model Type	Parameters		MSPE	Time	Time	Memory in GB	Memory in GB
	r	$\gamma$		in sec $\Sigma^{-1}\tilde{\mathbf{Z}}$	in sec $\det \Sigma$		
Full Model			2.556	155.6	186.0	14.18	13.79
Covariance Tapering	29	0	3.946	0.396	0.264	0.115	0.046
	166	0	3.115	0.620	0.385	0.762	0.341
	370	0	2.352	2.911	1.499	2.347	1.058
	0	50	5.069	0.008	0.005	0.009	0.003
	0	100	5.034	0.011	0.007	0.013	0.005
	0	200	4.726	0.030	0.019	0.045	0.017
	0	300	4.149	0.059	0.035	0.102	0.039
	0	400	3.752	0.098	0.055	0.175	0.068
	0	500	3.474	0.153	0.085	0.274	0.106
	0	625	3.264	0.235	0.140	0.429	0.168
Full-scale Approximation (r=29)	0	750	3.055	0.317	0.194	0.584	0.230
	0	1000	2.840	0.600	0.380	1.017	0.392
	0	1500	2.658	1.881	1.259	2.254	0.815
	29	50	3.946	0.396	0.264	0.115	0.046
	29	100	3.940	0.303	0.266	0.120	0.048
	29	200	3.752	0.330	0.278	0.165	0.060
	29	300	3.277	0.373	0.294	0.244	0.082
	29	400	2.925	0.519	0.315	0.345	0.111
	29	500	2.672	0.503	0.345	0.482	0.148
	29	625	2.480	0.620	0.399	0.696	0.210
Full-scale Approximation (r=166)	29	750	2.288	0.736	0.453	0.911	0.272
	29	1000	2.091	1.130	0.640	1.517	0.435
	29	1500	1.924	2.839	1.518	3.273	0.858
	166	50	3.115	0.620	0.385	0.762	0.341
	166	100	3.116	0.632	0.388	0.767	0.343
	166	200	3.006	0.659	0.399	0.811	0.356
	166	300	2.677	0.722	0.419	0.889	0.377
	166	400	2.411	0.769	0.436	0.988	0.406
	166	500	2.220	0.864	0.465	1.123	0.444
	166	625	2.085	1.043	0.520	1.334	0.506
Full-scale Approximation (r=370)	166	750	1.949	1.221	0.574	1.545	0.568
	166	1000	1.827	1.754	0.759	2.141	0.730
	166	1500	1.738	3.838	1.639	3.867	1.152
	370	50	2.352	2.911	1.499	2.347	1.058
	370	100	2.368	2.928	1.496	2.352	1.060
	370	200	2.337	3.008	1.515	2.396	1.073
	370	300	2.160	3.157	1.552	2.474	1.094
	370	400	2.002	3.276	1.566	2.574	1.123
	370	500	1.885	3.401	1.577	2.708	1.161
	370	625	1.804	3.726	1.632	2.919	1.222
	370	750	1.724	4.050	1.686	3.129	1.284
	370	1000	1.660	5.047	1.884	3.734	1.447
	370	1500	1.629	7.899	2.707	5.372	1.869

## B. Source Code

The MATLAB source code for the fixed rank kriging, basis function calculation and EM-Estimation was partially taken from the web tutorial on [http://www.stat.osu.edu/~sses/collab\\_co2.html](http://www.stat.osu.edu/~sses/collab_co2.html).

```

1 function h = distance_spherical(x,y)
2 %DISTANCE_ Determine distances between locations using the Spherical law of cosines
3 % formula
4 %
5 % This function produces a matrix that describes the
6 % distances between two sets of locations.
7 %
8 % INPUT PARAMETERS
9 % x - location coordinates (degrees) for data set #1 [n1 x D] -(long,lat)
10 % y - location coordinates (degrees) for data set #2 [n2 x D] -(long,lat)
11 % OUTPUT PARAMETERS
12 % h - distance (km) between points in x from points in y [n1 x n2]
13
14 [n1,D] = size(x);
15 [n2,D2] = size(y);
16
17 if D~=D2
18     error('ERROR in DISTANCE_: locations must have same number of dimensions (columns)')
19 end
20
21 h = zeros(n1,n2);
22 if D==1
23     for id = 1:D
24         h = h + (x(:,id)*ones(1, n2)-ones(n1,1)*y(:,id)').^2;
25     end
26     h = sqrt(h);
27 else
28     r=6371.0087714; %WGS84 mean radius
29     x=x*pi/180;
30     y=y*pi/180;
31     h = r*acos(sin(x(:,2)*ones(1, n2)).*sin(ones(n1,1)*y(:,2)')+cos(x(:,2)*ones(1, n2)
32             ).*cos(ones(n1,1)*y(:,2)')...
33             .*cos(x(:,1)*ones(1, n2)-ones(n1,1)*y(:,1)' ));
34     h=real(h);
35 end
36 return;

```

Listing 1: Compute geographic distances

```

1 function S=Basis(loc,BF_loc)
2
3 for i=1:3
4     [B,IX] = sort(BF_loc{i},1);
5     BF_loc{i}=BF_loc{i}(IX(:,2),:);
6 end
7
8 S=zeros(size(loc,1),size(BF_loc{1},1)+size(BF_loc{2},1)+size(BF_loc{3},1));
9 count=0;
10 for i=1:3
11     hrl=distance_spherical(BF_loc{i},BF_loc{i});
12     rl(i,1)=1.5*min(hrl(hrl>1e-3));
13     for j=1:length(BF_loc{i})
14         count=count+1;
15         h=distance_spherical(BF_loc{i}(j,:),loc);
16         s=(1-(h./rl(i,1)).^2).^2;
17         s(h>rl(i,1))=0;
18         S(:,count)=s;

```

```

19    end
20 end

```

Listing 2: Evaluate basis functions at specific locations

```

1 function [K sig_xi]=EM(S,z,V,sig_eps,V2)
2
3 n=size(S,1);
4
5 if nargin<5, V2=sparse(1:n,1:n,1); end
6
7 diagV=diag(V);
8 diagV2=diag(V2);
9
10 % initial values
11 varest=var(z,1);
12 K_old=.9*varest*eye(size(S,2));
13 sig2=.1*varest;
14 t=1;
15 done=0;
16
17 while done==0,
18
19     % update help terms
20     diagDinv=(sig2(t)*diagV2+sig_eps*diagV).^( -1);
21     DInv=sparse(1:n,1:n,diagDinv);
22     tempt=inv(inv(K_old)+S'*DInv*S);
23
24     % update K
25     SigInv2=(tempt*S')*DInv;
26     KSDInv=K_old*S'*DInv;
27     KSSigInv=KSDInv-KSDInv*S*SigInv2;
28     muEta=KSSigInv*z;
29     SigEta=K_old-KSSigInv*S*K_old;
30     K_new=SigEta+muEta*muEta';
31
32     % update sigma_xi (sig2)
33     muEps=sig2(t)*(DInv*z-DInv*S*(SigInv2*z));
34     trSigInv=trace(DInv)-trace(SigInv2*DInv*S);
35     sig2(t+1)=1/n*(n*sig2(t)-(sig2(t))^2*trSigInv+muEps'*muEps);
36
37     % check for convergence
38     diff=sum(sum((K_new-K_old).^2,1),2)+(sig2(t+1)-sig2(t))^2;
39     if diff<avgtol*(size(S,2))^2, done=1; end
40     if t>maxit,
41         done=1;
42         disp(strcat('Algorithm did not converge after ',num2str(maxit),' iterations'))
43         ;
44     end
45
46     t=t+1
47     K_old=K_new;
48 end
49
50 K=K_new;
51 sig_xi=sig2(t);

```

Listing 3: EM algorithm for the parameters in the Fixed Rank Kriging approach

```

1 function [y_hat]=FRK(data,s_pred,K,sig_xi,sig_eps,BF_loc,V_eps,V_xi)
2
3 % data is a n x 3 matrix : first column is degrees in Latitude, %
% second is degrees in Longitude and %
% third is the data values

```

```

4 % s_pred is a n x 2 matrix of the prediction locations
5 % K      is the k x k dimensional covariance matrix of the
6 %       S_predatial random effect
7 % sig_xi is a scalar, representing the variance parameter of the %      micro-scale
8 %       variation process
9 % sig_eps is a scalar, representing the variance parameter of the %      measurement error process
10 % BF_loc cell array with the locations of the basis function
11 %       centers for each spatial resolution in a different cell
12 % V_eps n x n matrix describing the S_predatial heterogeneity of %      the
13 %       measurement error process
14 % V_xi n x n matrix describing the S_predatial heterogeneity of %
15 %       the micro-scale variation process
16
17 if nargin<8, V2=sparse(1:length(data),1:length(data),1); end
18
19 % Observations
20 lon=data(:,2);
21 lat=data(:,1);
22 z=data(:,3);
23
24 % Prediction locations
25 lon_pred=s_pred(:,2);
26 lat_pred=s_pred(:,1);
27
28 n=size(z,1);
29 m=length(lon_pred);
30
31 % Evaluate spatial basis functions at observation locations
32 S_obs=Basis(data(:,[2 1]),BF_loc);
33
34 % Evaluate spatial basis functions at prediction locations
35 S_pred=Basis([lon_pred lat_pred],BF_loc);
36
37 % Fixed Rank Kriging
38 temp=inv(sig_xi*V_xi + sig_eps*V_eps);
39
40 temp2=inv(inv(K)+S_obs'*temp*S_obs);
41
42 E=sparse(m,n);
43 for i=1:n,
44     E((lon_pred==lon(i) & lat_pred==lat(i)),i)=1;
45 end;
46
47 temp3=temp*z-temp*S_obs*(temp2*S_obs'*temp*z);
48
49 y_hat=S_pred*(K*S_obs'*temp3)+sig_xi*E*temp3;

```

Listing 4: Fixed Rank Kriging predictions

```

1 function [predtap]=tapped(data,s_pred,sig_eps,V,Cs,cross)
2
3 % Cs is n x n covariance matrix of the spatial random effect at
4 %       the observed location
5 % cross is the n x m matrix of cross covariances between Y at the %      prediction
6 %       locations and the observed data Z
7
8 lon=data(:,2);
9 lat=data(:,1);
10 z=data(:,3);
11 % coordinates of prediction locations
12 lon_pred=s_pred(:,2);
13 lat_pred=s_pred(:,1);

```

```

14
15 n=size(z,1); % number of observations
16 m=length(lon_pred); % number of prediction locations
17
18 SigInvz=(Cs+sig_eps*V)\z;
19
20 %predictions
21 predtap=cross'*SigInvz;

```

Listing 5: Covariance Tapering predictions

```

1 function [predFSA]=FSA(data,s_pred,K,sig_xi,sig_eps,BF_loc,V,Cs,cross)
2
3 % observed data
4 lon=data(:,2);
5 lat=data(:,1);
6 z=data(:,3);
7
8 % coordinates of prediction locations
9 lon_pred=s_pred(:,2);
10 lat_pred=s_pred(:,1);
11
12 n=size(z,1); % number of observations
13 m=length(lon_pred); % number of prediction locations
14
15 % Evaluate spatial basis functions at observation locations
16 S=Create_S(data(:,[2 1]),BF_loc);
17
18 % Evaluate spatial basis functions at prediction locations
19 Sp=Create_S([lon_pred lat_pred],BF_loc);
20
21 Dinvz=(Cs+sig_eps*V)\z;
22 Dinvz=(Cs+sig_eps*V)\S;
23 SigInvz=Dinvz-Dinvz*((inv(K)+S'*Dinvz)\S'*Dinvz);
24
25 %predictions
26 predFSA=cross'*SigInvz+Sp*(K*S'*SigInvz);

```

Listing 6: FSA predictions

```

1 library(matlab)
2 library('R.matlab')
3 library(fields)
4 library(geoR)
5 library(tcltk)
6 library(gstat)
7 library(rgdal)
8 library(Matrix)
9
10 source('D:/C02/FRK_auflösungen/frkres.R')
11 x<-cbind(lonbtR,latbtR) # observation locations
12 x2<-cbind(lonrausR,latrausR) # prediction locations
13
14
15 # Compute empirical variogram for the detrended data
16
17 data=data.frame(ztilde,lonbtR,latbtR)
18 coordinates(data)=~lonbtR+latbtR
19 proj4string(data)="+proj=longlat"
20 breaks = seq(0, 1500, l = 0.1*1500)
21 variogstat<-variogram(ztilde~1,data,boundaries=breaks)
22 variogstat$gamma=variogstat$gamma-sig2_eps
23
24 # Fit Exponential Model
25 vmodel=vgm(psill=5.5,model="Exp",range=700,nuget=0)

```

```

26 fitgstat=fit.variogram(variogstat,model=vmodel,fit.sills=T,fit.ranges=T,fit.method=1)
27 sill=fitgstat$psill
28 range=fitgstat$range
29
30 # Create tapered covariance matrices based on the detrended data
31 dist<-rdist.earth(x,miles=F)
32 gc()
33
34 Cst<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.cov(x
35 ,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
36 gc()
37 Cst=as.spam(Cst)
38 gc()
39 rm('dist')
40
41 dist<-rdist.earth(x,x2,miles=F)
42
43 gc()
44
45 cross<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.cov(
46 (x,x2,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
47 gc()
48 cross=as.spam(cross)
49 gc()
50 rm('dist')
51
52 # Obtain computing times
53
54 V_eps=diag(nrow=length(x[,1]))
55 a<-as.spam(Cst+sig2_eps*V_eps)
56 start <- Sys.time()
57 blubb=chol(a)
58 b=backsolve(blubb,ztilde)
59 f=backsolve(t(blubb),b,upper.tri=FALSE)
60 timetap<-as.numeric((Sys.time() - start),units='secs')
61
62 a<-as.spam(Cst+sig2_eps*V_eps)
63 start <- Sys.time()
64 deta=det(a)
65 timedettap=as.numeric((Sys.time() - start),units='secs')
66
67 start <- Sys.time()
68 pred=t(cross) %*% f
69 timepredtap=as.numeric((Sys.time() - start),units='secs')
70
71 dist<-rdist.earth(x2,miles=F)
72
73 gc()
74
75 Csp<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.cov(
76 (x2,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
77 gc()
78 Csp=as.spam(diag(Csp))
79 gc()
80 rm('dist')
81
82 writeMat('D:/CO2/R/covmatt.mat',Cst=as.matrix(Cst),Csp=as.matrix(Csp),timecholtap=
83 timetap,timedettap=timedettap,timepredtap=timepredtap)
83 rm('Cst')
84 rm('Csp')
85 cross=triplet(cross,tri=T)
86 gc()
87 writeMat('D:/CO2/R/crosst.mat',crosstri=cross)

```

```

88 rm('cross')
89 gc()
90
91
92 # Compute empirical variogram for the residuals of the fixed rank kriging
93
94 data=data.frame(resFRK,lonbtR,latbtR)
95 coordinates(data)=~lonbtR+latbtR
96 proj4string(data)="+proj=longlat"
97 breaks = seq(0, 1500, 1 = 0.1*1500)
98 variogstat<-variogram(resFRK~1,data,boundaries=breaks)
99 variogstat$gamma=variogstat$gamma-sig2_eps
100
101 # Fit exponential Model
102 vmodel=vgm(psill=5.5,model="Exp",range=700,nuget=0)
103 fitgstat=fit.variogram(variogstat,model=vmodel,fit.sills=T,fit.ranges=T,fit.method=1)
104 sill=fitgstat$psill
105 range=fitgstat$range
106
107 # Compute covariance matrices for the tapering part of the full-scale approximation
108
109 dist<-rdist.earth(x,miles=F)
110
111 gc()
112
113 CsFSA<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.cov
114   (x,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
115 gc()
116 CsFSA=as.spam(CsFSA)
117 gc()
118 rm('dist')
119
120 dist<-rdist.earth(x,x2,miles=F)
121
122 gc()
123
124 crossFSA<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.
125   cov(x,x2,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
126 gc()
127 crossFSA=as.spam(crossFSA)
128 gc()
129 rm('dist')
130
131 # Computation time
132
133 V_eps=diag(nrow=length(x[,1]))
134 a<-as.spam(CsFSA+sig2_eps*V_eps)
135 start <- Sys.time ()
136 blubb=chol(a)
137 b=backsolve(blubb,ztilde)
138 dinvz=backsolve(t(blubb),b,upper.tri=FALSE)
139 b=backsolve(blubb,S)
140 dinvs=backsolve(t(blubb),b,upper.tri=FALSE)
141 blubb=chol(solve(K)+t(S)%*%dinvz)
142 b=backsolve(blubb,t(S))
143 temp=backsolve(t(blubb),b,upper.tri=FALSE)
144 Siginvz=dinvz-dinvs%*%(temp%*%dinvz)
145 timecholFSA=as.numeric((Sys.time () - start),units='secs')
146
147 start <- Sys.time ()
148 deta=det(solve(K)+t(S)%*%dinvs)%*%(det(solve(K))^{(-1)})%*%det(a)
149 timedetFSA=as.numeric((Sys.time () - start),units='secs')
150
151 start <- Sys.time ()

```

```

152 pred=t(crossFSA)%%Siginvz+Sp1%*%(K%*%t(S)%*%Siginvz)
153 timepredFSA=as.numeric((Sys.time () - start),units='secs')
154
155 dist<-rdist.earth(x2,miles=F)
156
157 gc()
158
159 CspFSA<-cov.spatial(dist, cov.model= "exponential",cov.pars=c(sill,range))*wendland.
160 cov(x2,Dist.args=list(method="greatcircle",miles=F),theta=tapper,k=2)
160 gc()
161 CspFSA=as.spam(diag(Csp1))
162 gc()
163
164 rm('dist')
165 writeMat('covmatFSA.mat',CsFSA=as.matrix(CsFSA),CspFSA=as.matrix(CspFSA),timecholFSA=
166 timecholFSA,timedetFSA=timedetFSA,timepredFSA=timepredFSA)
166 crossS=triplet(crossS,tri=T)
167 gc()
168 writeMat('crossFSA.mat',crossStri=crossS)
169 rm('CsFSA')
170 rm('CspFSA')
171 rm('crossFSA')
172 gc()
173
174
175 # Computation time for Fixed Rank Kriging
176 V_eps=diag(nrow=length(x[,1]))
177 a<-as.spam(sig2_eps*V_eps)
178 start <- Sys.time()
179 blubb=chol(a)
180 b=backsolve(blubb,ztilde)
181 dinvz=backsolve(t(blubb),b,upper.tri=FALSE)
182 b=backsolve(blubb,S)
183 dinvs=backsolve(t(blubb),b,upper.tri=FALSE)
184 blubb=chol(solve(K)+t(S)%*%dinvz)
185 b=backsolve(blubb,t(S))
186 temp=backsolve(t(blubb),b,upper.tri=FALSE)
187 Siginvz=dinvz-dinvs%*%(temp%*%dinvz)
188 timecholfrk=as.numeric((Sys.time () - start),units='secs')
189
190 start <- Sys.time()
191 data=det(solve(K)+t(S)%*%dinvz)%*%(det(solve(K))^{(-1)})%*%det(a)
192 timedetfrk=as.numeric((Sys.time () - start),units='secs')
193
194 start <- Sys.time()
195 pred=Sp1%*%(K%*%t(S)%*%Siginvz)
196 timepredfrk=as.numeric((Sys.time () - start),units='secs')
197
198
199 writeMat('timefrk.mat',timepredfrk=timepredfrk,timedetfrk=timedetfrk,timecholfrk=
    timecholfrk)

```

Listing 7: Variogram-Fitting and performance evaluation

```

1 b=10; % b-fold crossvalidation
2 Indices = crossvalind('Kfold', length(ztilde), b);
3
4 load('data.mat') % data, locations and basis function centers
5
6 for k=1:b,
7 lonbt=lon(Indices~=k);
8 latbt=lat(Indices~=k);
9 zbt=ztilde(Indices~=k);
10 zraus{k}=ztilde(Indices==k);
11 latraus=lat(Indices==k);
12 lonraus=lon(Indices==k);

```

```

13 S_obs=Create_S([lonbt latbt],BF_loc);
14
15 n=length(zbt);
16 V_eps=sparse(1:n,1:n,1);
17 data=[latbt lonbt zbt];
18 s_pred=[latraus lonraus];
19
20 % EM-Estimation of K and sigma^2_xi
21 [K sig_xi]=EMestimation(S_obs,zbt,V_eps,sig_eps);
22
23 % FRK-Predictions
24 [pred_FRK{k}]=FRK(data,s_pred,K,sig_xi,sig_eps,BF_loc,V_eps);
25
26 % Calculate FRK-Residuals
27 s_obs=[latbt lonbt];
28
29 [pred_frk_obs]=FRK(data,s_obs,K,sig_xi,sig_eps,BF_loc,V_eps);
30 resFRK=zbt-pred_frk_obs;
31
32 lonbtR=lonbt;
33 latbtR=latbt;
34 lonrausR=lonraus;
35 latrausR=latraus;
36
37 i=0;
38 for tapper=[50 100:100:500 625 750 1000 1500],
39 i=i+1
40 delete 'frkres.R'
41 delete 'crosst.mat'
42 delete 'covmatt.mat'
43 delete 'crossFSA.mat'
44 delete 'covmatFSA.mat'
45
46 % Estimation of Covariance matrices based on the Covariance
47 % Tapering / Full-scale approximation and the recording of the % computing times is
        done in R-Software
48
49 saveR('frkres.R','resFRK','lonbtR','latbtR','latrausR','lonrausR','zbt','tapper',
       'sig2_eps');
50 eval(['!C:/PROGRA~1/R/R-2.15.2/bin/x64/Rscript ' 'Variogram_fitting
      .R']);
51
52 load('crosst.mat')
53 load('covmatt.mat')
54 Cst=sparse(Cst);
55 cross=sparse(double(crosstri.i),double(crosstri.j),crosstri.values,n,length(lonraus));
56 clear crosstri
57
58 [predtap{i}]=tapped(data,s_pred,sig_eps,V_eps,Cst,cross);
59
60 clear Cst
61 clear cross
62
63 load('crossFSA.mat')
64 load('covmatFSA.mat')
65 CsFSA=sparse(CsFSA);
66 crossFSA=sparse(double(crossFSAttri.i),double(crossFSAttri.j),crossFSAttri.values,n,
                  length(lonraus));
67 clear crossFSAttri
68
69 [predFSA{i}]=FSA(data,pred_locs,K,sig_xi,sig_eps,V_eps,CsFSA,crossFSA);
70
71 clear CsFSA
72 clear crossFSA
73
74
75

```

```

76 tapperm(i)=tapper;
77
78 %Calculate MSE
79
80 MSE_FSA(i)=(1/length(zraus{k}))*(((zraus{k}-predFSA{i}).^2) '*ones(length(zraus{k}),1))
81 ;
81 MSEtap(i)=(1/length(zraus{k}))*(((zraus{k}-predtap{i}).^2) '*ones(length(zraus{k}),1));
82
83 end
84
85 MSE_FSAs{k}=MSE_FSA-sig2_eps;
86 MSEtaps{k}=MSEtap-sig2_eps;
87
88 end
89
90 for k=1:b,
91 MSE_FRK(k)=(1/length(zraus{k}))*(((zraus{k}-pred_FRK{k}).^2) '*ones(length(zraus{k}),1)
91 )-sig2_eps;
92 end
93
94 MSE_FRK=mean(MSE_FRK);
95
96 MSE_FSA=MSE_FSAs{1};
97 MSEtap=MSEtaps{1};
98 for k=2:b,
99 MSE_FSA=[MSE_FSA; MSE_FSAs{k}];
100 MSEtap=[MSEtap; MSEtaps{k}];
101 end
102
103 MSE_FSA=mean(MSE_FSA,1);
104 MSEtap=mean(MSEtap,1);

```

Listing 8: Cross-Validation study

## Acknowledgements

This work was supported by the German Ministry of Education and Research (BMBF) under its funding program 'Economics of Climate Change' [grant number 01LA1139A]

## References

- Abramowitz M, Stegun IA (1964). *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Dover, New York.
- Cressie N (1985). “Fitting variogram models by weighted least squares.” *Journal of the International Association for Mathematical Geology*, **17**(5), 563–586.
- Cressie N (1993). *Statistics for Spatial Data*. Revised edition edition. Wiley-Interscience.
- Cressie N, Hawkins D (1980). “Robust estimation of the variogram: I.” *Mathematical Geology*, **12**(2), 115–125.
- Cressie N, Johannesson G (2006). “Spatial prediction of massive datasets.” In *Proceedings of the Australian Academy of Science Elizabeth and Frederick White Conference*, volume 1247, pp. 1–11.
- Cressie N, Johannesson G (2008). “Fixed rank kriging for very large spatial data sets.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **70**(1), 209–226.
- Cressie N, Shi T, Kang E (2010). “Fixed Rank filtering for spatio-temporal data.” *Journal of Computational and Graphical Statistics*, **19**, 724–745.
- Cressie N, Wikle C (2011). *Statistics for Spatio-Temporal Data*. Wiley Series in Probability and Statistics. Wiley.
- Dempster A, Laird N, Rubin D (1977). “Maximum likelihood from incomplete data via the EM algorithm.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **39**(1), 1–38.
- Eidsvik J, Martino S, Rue H (2009). “Approximate Bayesian Inference in Spatial Generalized Linear Mixed Models.” *Scandinavian Journal of Statistics*, **36**(1), 1–22.
- Eidsvik J, O’Finley A, Banerjee S, Rue H (2012). “Approximate Bayesian inference for large spatial datasets using predictive process models.” *Computational Statistics & Data Analysis*, **56**(6), 1362 – 1380.
- Furrer R, Genton MG, Nychka D (2006). “Covariance tapering for interpolation of large spatial datasets.” *Journal of Computational and Graphical Statistics*, **15**(3), 502–523.
- González-Manteiga W, Crujeiras RM, Katzfuss M, Cressie N (2012). “Bayesian hierarchical spatio-temporal smoothing for very large datasets.” *Environmetrics*, **23**(1), 94–107.
- Guan D, Liu Z, Geng Y, Lindner S, Hubacek K (2012). “The gigatonne gap in China’s carbon dioxide inventories.” *Nature Climate Change*, **2**, 672–675.

- Hastie T, Tibshirani R, Friedman J (2003). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- Henderson H, Searle S (1981). “On deriving the inverse of a sum of matrices.” *Siam Review*, **23**(1), 53–60.
- Horn RA, Johnson CR (1994). *Topics in Matrix Analysis*. Cambridge University Press.
- Katzfuss M, Cressie N (2009). “Maximum likelihood estimation of covariance parameters in the spatial-random-effects model.” In *Proceedings of the Joint Statistical Meetings*, pp. 3378–3390. American Statistical Association.
- Katzfuss M, Cressie N (2011). “Spatio-temporal smoothing and EM estimation for massive remote-sensing data sets.” *Journal of Time Series Analysis*, **32**, 430–446.
- Kaufman CG, Schervish MJ, Nychka DW (2008). “Covariance tapering for likelihood-based estimation in large spatial data sets.” *Journal of the American Statistical Association*, **103**(484), 1545–1555.
- Lasinio GJ, Mastrantonio G, Pollice A (2013). “Discussing the ”big n problem”.” *Statistical Methods and Applications*, **22**(1), 97–112.
- Lindgren F, Rue H, Lindström J (2011). “An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **73**(4), 423–498.
- Matérn B (1986). *Spatial Variation*. Lecture Notes in Statistics. Springer-Verlag.
- Mintzer I, Leonard J, Valencia I (2010). “Counting the Gigatonnes: Building trust in greenhouse gas inventories from the United States and China.” *World Wildlife Federation*.
- Nychka D, Wikle C, Royle JA (2002). “Multiresolution models for nonstationary spatial covariance functions.” *Statistical Modelling*, **2**(4), 315–331.
- Rue H, Martino S, Chopin N (2009). “Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **71**(2), 319–392.
- Sahr K, White D, Kimerling AJ (2003). “Geodesic discrete global grid systems.” *Cartography and Geographic Information Science*, **30**(2), 121–134.
- Sang H, Huang J (2012). “A full scale approximation of covariance functions for large spatial data sets.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **74**(1), 111–132.
- Shi T, Cressie N (2007). “Global statistical analysis of MISR aerosol data: a massive data product from NASA’s Terra satellite.” *Environmetrics*, **18**(7), 665–680.
- Stroud JR, Stein ML, Lesht BM, Schwab DJ, Beletsky D (2010). “An Ensemble Kalman Filter and Smoother for Satellite Data Assimilation.” *Journal of the American Statistical Association*, **105**(491), 978–990.

- Vidakovic B (1999). *Statistical Modeling by Wavelets*. Wiley Series in Probability and Statistics. John Wiley & Sons, Inc., Hoboken, NJ, USA.
- Wahba G (1990). *Spline Models for Observational Data*. Society for Industrial and Applied Mathematics.
- Wikle C (2010). “Low-rank representations for spatial processes.” In *Handbook of Spatial Statistics*, Chapman & Hall/CRC Handbooks of Modern Statistical Methods, pp. 107–118. CRC Press.
- Zhang B, Sang H, Huang JZ (2013). “Full-Scale Approximations of Spatio-Temporal Covariance Models for Large Datasets.” *Statistica Sinica*, **in press**.

**Affiliation:**

Patrick Vetter  
 Department of Statistics  
 Europa-Universität Viadrina Frankfurt (Oder), Germany  
 Postal Code 15230  
 E-mail: [vetter@europa-uni.de](mailto:vetter@europa-uni.de)