

Big models for structured environmental big data

Finn Lindgren



London 2014-09-17

Data

Spatial models

Spatial

Matérn/SPDE

Basis connections

Markov

Non-stationarity

Examples

Precipitation

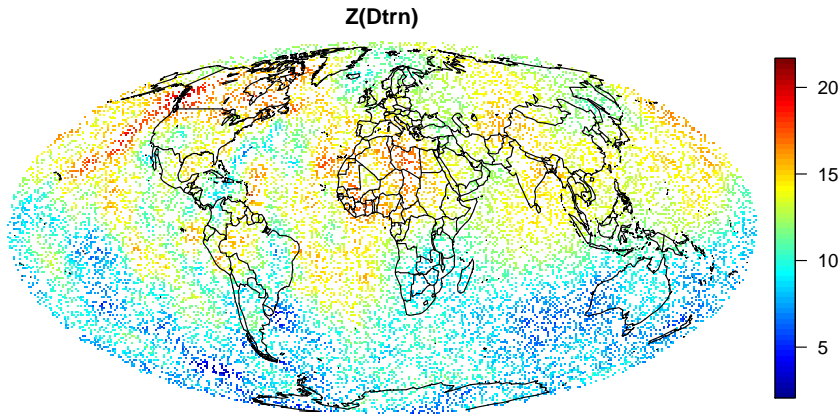
LGCP

CO₂

Temperature

End

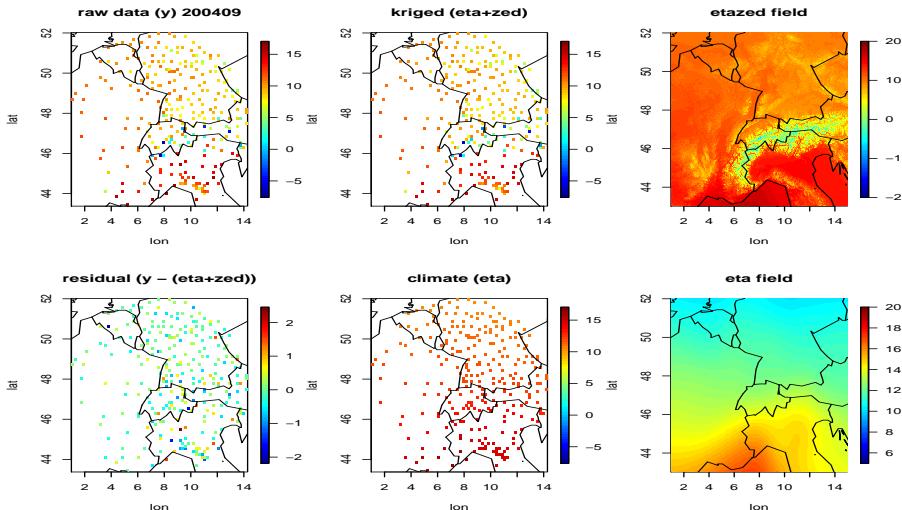
“Big” data



Synthetic data mimicking the spatial variability of satellite based CO₂ measurements.

Lots of data, decent spatial coverage, but more “missing” than “observed”.

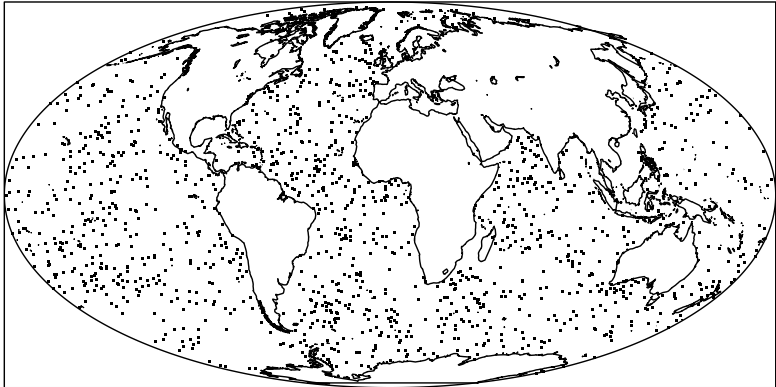
Sparse spatial coverage of temperature measurements



Total observations: $\approx 20,000,000$ from daily timeseries over 160 years

Spatial information: 350 observations

Point process data are very weakly informative



Big data, information, and stochastic models

Summary

- ▶ Space-time data is nearly always sparse and spatially uneven
- ▶ Nearly *all* potential observables are “missing”
- ▶ Fortunately, “missingness” is a problem only for specific estimation algorithms, not for statistical analysis as such
- ▶ Hierarchical/generative/Bayesian/regularization models provide a general and useful framework for analysis of multi data source data sets

Challenges:

- ▶ Even larger output than input
- ▶ Quantifying uncertainty is harder than reasonable point estimates
- ▶ Can we use physical world knowledge?
- ▶ Representing multi-scale stochastic dependence
- ▶ ...in a numerics friendly manner
- ▶ Hierarchical models are general, but MCMC is very slow
- ▶ We *can* use a lot more direct computation

Brief introduction to spatial modelling

Spatial statistics framework

- ▶ Spatial domain Ω , or space-time domain $\Omega \times \mathbb{T}$, $\mathbb{T} \subset \mathbb{R}$.
- ▶ Random field $u(\mathbf{s})$, $\mathbf{s} \in \Omega$, or $u(\mathbf{s}, t)$, $(\mathbf{s}, t) \in \Omega \times \mathbb{T}$.
- ▶ Observations y_i . Usually $y_i = u(\mathbf{s}_i) + \epsilon_i$ (simple georeferenced data in geostatistics) or more generally $y_i \sim \text{GLMM}$, with $u(\cdot)$ as a structured random effect.
- ▶ We'll focus on Bayesian hierarchical models in the form of spatial latent Gaussian models.

Two basic model and method components

- ▶ We need stochastic models for $u(\cdot)$.
- ▶ We need computationally efficient (Bayesian) inference methods for the posterior distribution of $u(\cdot)$ given data \mathbf{y} .

Covariance functions and stochastic PDEs

The Matérn covariance family on \mathbb{R}^d

$$\text{Cov}(u(\mathbf{0}), u(\mathbf{s})) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} (\kappa \|\mathbf{s}\|)^\nu K_\nu(\kappa \|\mathbf{s}\|)$$

Scale $\kappa > 0$, smoothness $\nu > 0$, variance $\sigma^2 > 0$



Whittle (1954, 1963): Matérn as SPDE solution

Matérn fields are the stationary solutions to the SPDE

$$(\kappa^2 - \nabla \cdot \nabla)^{\alpha/2} u(\mathbf{s}) = \mathcal{W}(\mathbf{s}), \quad \alpha = \nu + d/2$$

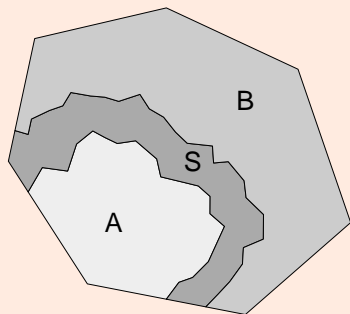
$\mathcal{W}(\cdot)$ white noise, $\nabla \cdot \nabla = \sum_{i=1}^d \frac{\partial^2}{\partial s_i^2}$, $\sigma^2 = \frac{\Gamma(\nu)}{\Gamma(\alpha) \kappa^{2\nu} (4\pi)^{d/2}}$



Continuous and discrete Markov properties

Markov properties

S is a separating set for A and B : $u(A) \perp u(B) \mid u(S)$



Solutions to

$$(\kappa^2 - \nabla \cdot \nabla)^{\alpha/2} u(s) = \mathcal{W}(s)$$

are Markov when α is an integer.

(Rozanov, 1977)

Discrete representations ($Q = \Sigma^{-1}$):

$$Q_{AB} = 0$$

$$Q_{A|S,B} = Q_{AA}$$

$$\mu_{A|S,B} = \mu_A - Q_{AA}^{-1} Q_{AS} (u_S - \mu_S)$$

Continuous domain Markov approximations

Continuous Markovian spatial models (Lindgren et al, 2011)

Local basis: $u(\mathbf{s}) = \sum_k \psi_k(\mathbf{s}) u_k$, (compact, piecewise linear)

Basis weights: $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1})$, sparse \mathbf{Q} based on an SPDE

Special case: $(\kappa^2 - \nabla \cdot \nabla)u(\mathbf{s}) = \mathcal{W}(\mathbf{s})$, $\mathbf{s} \in \Omega$

Precision: $\mathbf{Q} = \kappa^4 \mathbf{C} + 2\kappa^2 \mathbf{G} + \mathbf{G}_2$ ($\kappa^4 + 2\kappa^2|\boldsymbol{\omega}|^2 + |\boldsymbol{\omega}|^4$)

Conditional distribution in a Gaussian model

$\mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}_u, \mathbf{Q}_u^{-1})$, $\mathbf{y}|\mathbf{u} \sim \mathcal{N}(\mathbf{A}\mathbf{u}, \mathbf{Q}_{y|\mathbf{u}}^{-1})$ ($A_{ij} = \psi_j(\mathbf{s}_i)$)

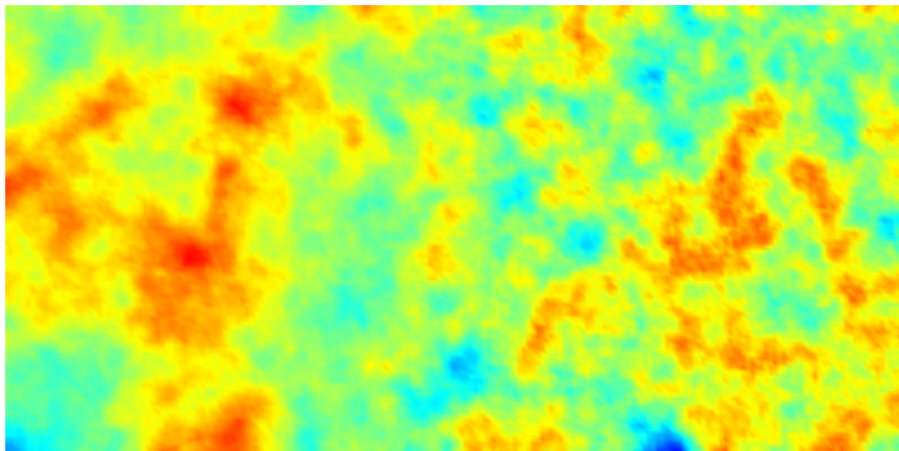
$\mathbf{u}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{u|\mathbf{y}}, \mathbf{Q}_{u|\mathbf{y}}^{-1})$

$\mathbf{Q}_{u|\mathbf{y}} = \mathbf{Q}_u + \mathbf{A}^T \mathbf{Q}_{y|\mathbf{u}} \mathbf{A}$ (~"Sparse iff ψ_k have compact support")

$\boldsymbol{\mu}_{u|\mathbf{y}} = \boldsymbol{\mu}_u + \mathbf{Q}_{u|\mathbf{y}}^{-1} \mathbf{A}^T \mathbf{Q}_{y|\mathbf{u}} (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_u)$

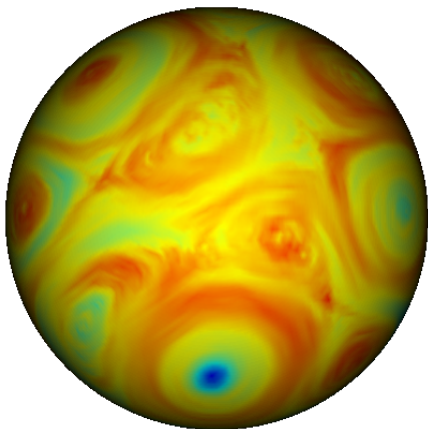
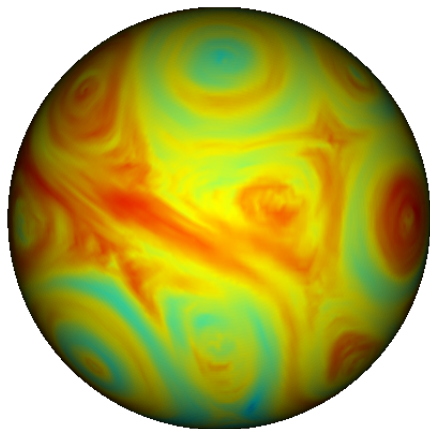
We've translated the spatial inference problem into sparse numerical linear algebra similar to finite element PDE solvers

Non-stationary field



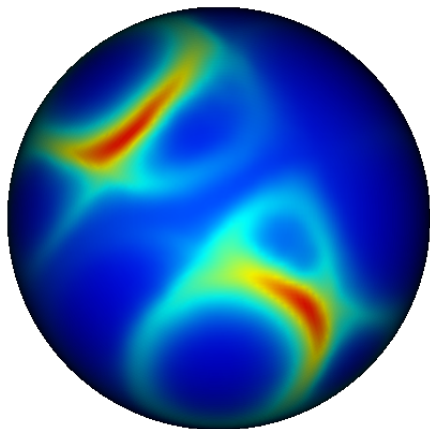
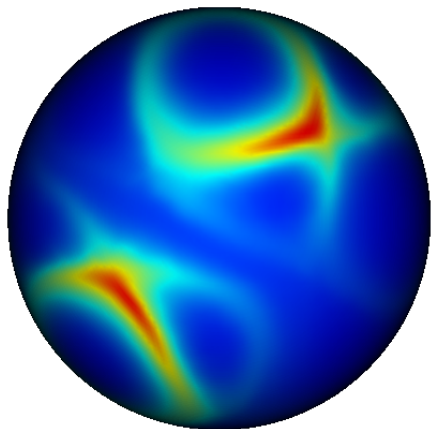
$$(\kappa(\mathbf{s}))^2 - \nabla \cdot \nabla)u(\mathbf{s}) = \kappa(\mathbf{s})\mathcal{W}(\mathbf{s}), \quad \mathbf{s} \in \Omega$$

Anisotropic field on a globe via vector parameter field

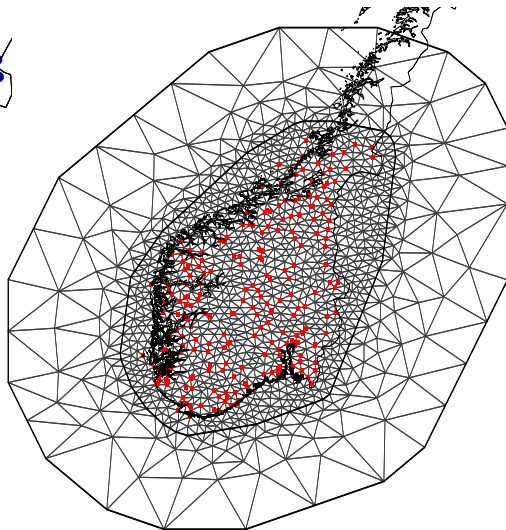
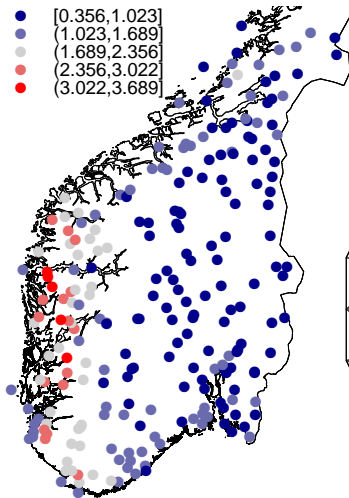


$$(\kappa(\mathbf{s})^2 - \nabla \cdot \mathbf{H}(\mathbf{s})\nabla)u(\mathbf{s}) = \kappa(\mathbf{s})\mathcal{W}(\mathbf{s}), \quad \mathbf{s} \in \Omega$$

Four covariance functions



Example: Precipitation (Ingebrigtsen et al., 2013)



Non-stationary precision construction

Finite element construction of basis weight precision

Non-stationary SPDE:

$$(\kappa(\mathbf{s}))^2 - \nabla \cdot \nabla (\tau(\mathbf{s})u(\mathbf{s})) = \mathcal{W}(\mathbf{s})$$

The SPDE parameters are constructed via spatial covariates:

$$\log \tau(\mathbf{s}) = b_0^\tau(\mathbf{s}) + \sum_{j=1}^p b_j^\tau(\mathbf{s})\theta_j, \quad \log \kappa(\mathbf{s}) = b_0^\kappa(\mathbf{s}) + \sum_{j=1}^p b_j^\kappa(\mathbf{s})\theta_j$$

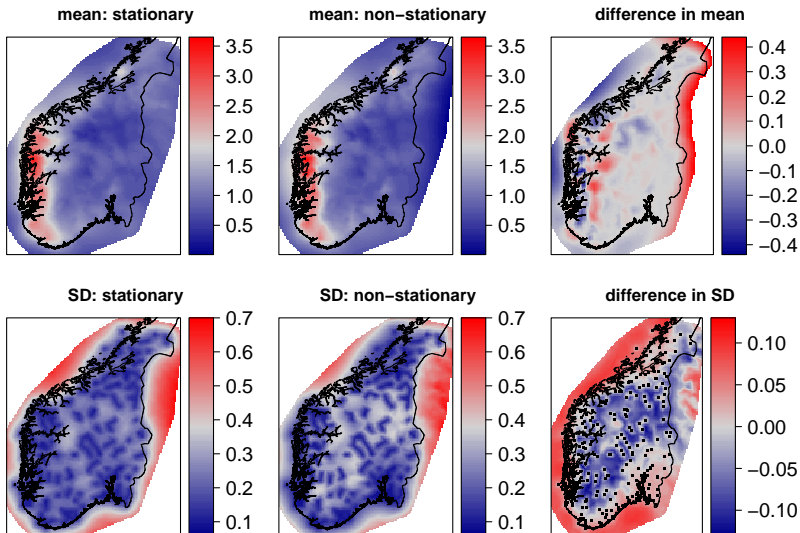
Finite element calculations give

$$\mathbf{T} = \text{diag}(\tau(\mathbf{s}_i)), \quad \mathbf{K} = \text{diag}(\kappa(\mathbf{s}_i))$$

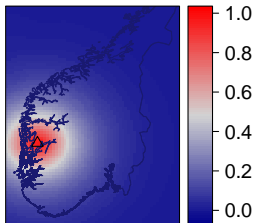
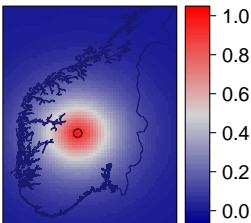
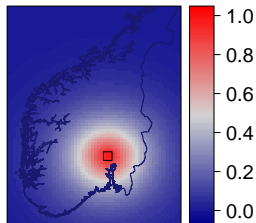
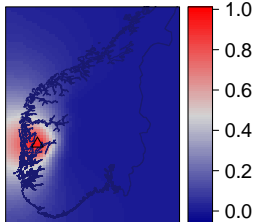
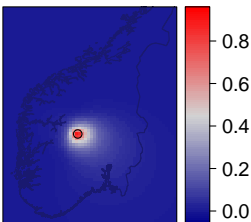
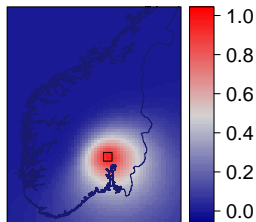
$$C_{ii} = \int \psi_i(\mathbf{s}) d\mathbf{s}, \quad G_{ij} = \int \nabla \psi_i(\mathbf{s}) \cdot \nabla \psi_j(\mathbf{s}) d\mathbf{s}$$

$$\mathbf{Q} = \mathbf{T} (\mathbf{K}^2 \mathbf{C} \mathbf{K}^2 + \mathbf{K}^2 \mathbf{G} + \mathbf{G} \mathbf{K}^2 + \mathbf{G} \mathbf{C}^{-1} \mathbf{G}) \mathbf{T}$$

Results for stationary and non-stationary models



Correlations for stationary and non-stationary models

Kvamskogen: stationary**Hemsedal: stationary****Hønefoss: stationary****Kvamskogen: non-stationary****Hemsedal: non-stationary****Hønefoss: non-stationary**

Example: Point process data

Log-Gaussian Cox processes

Point intensity:

$$\lambda(\mathbf{s}) = \exp \left(\sum_i b_i(\mathbf{s}) \beta_i + u(\mathbf{s}) \right)$$

Inhomogeneous Poisson process log-likelihood:

$$\ln p(\{\mathbf{y}_k\} | \boldsymbol{\lambda}) = |\Omega| - \int_{\Omega} \lambda(\mathbf{s}) d\mathbf{s} + \sum_{k=1}^n \ln \lambda(\mathbf{y}_k)$$

The likelihood can be approximated numerically, e.g.

$$\int_{\Omega} \lambda(\mathbf{s}) d\mathbf{s} \approx \sum_{j=1}^N \lambda(\mathbf{s}_j) w_j$$

Laplace approximations

Quadratic posterior log-likelihood approximation

$$p(\mathbf{u} \mid \boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\mu}_u, \mathbf{Q}_u^{-1}), \quad \mathbf{y} \mid \mathbf{u}, \boldsymbol{\theta} \sim p(\mathbf{y} \mid \mathbf{u})$$

$$p_G(\mathbf{u} \mid \mathbf{y}, \boldsymbol{\theta}) \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{Q}}^{-1})$$

$$\mathbf{0} = \nabla_{\mathbf{u}} \{ \ln p(\mathbf{u} \mid \boldsymbol{\theta}) + \ln p(\mathbf{y} \mid \mathbf{u}) \} \Big|_{\mathbf{u}=\tilde{\boldsymbol{\mu}}}$$

$$\tilde{\mathbf{Q}} = \mathbf{Q}_u - \nabla_{\mathbf{u}}^2 \ln p(\mathbf{y} \mid \mathbf{u}) \Big|_{\mathbf{u}=\tilde{\boldsymbol{\mu}}}$$

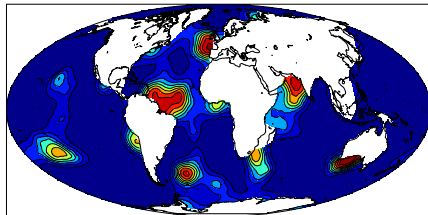
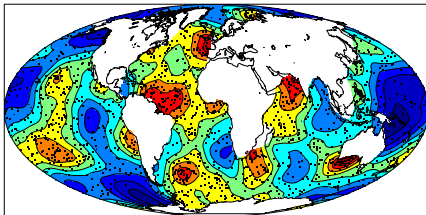
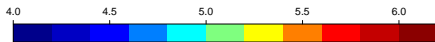
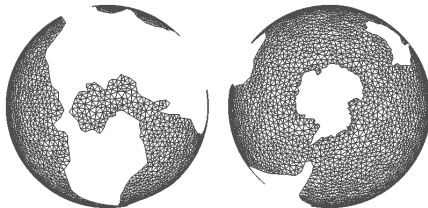
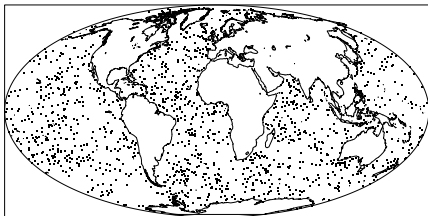
Direct Bayesian inference with INLA (r-inla.org)

$$\tilde{p}(\boldsymbol{\theta} \mid \mathbf{y}) \propto \frac{p(\boldsymbol{\theta})p(\mathbf{u} \mid \boldsymbol{\theta})p(\mathbf{y} \mid \mathbf{u}, \boldsymbol{\theta})}{p_G(\mathbf{u} \mid \mathbf{y}, \boldsymbol{\theta})} \Big|_{\mathbf{u}=\tilde{\boldsymbol{\mu}}}$$

$$\tilde{p}(\mathbf{u}_i \mid \mathbf{y}) \propto \int p_{GG}(\mathbf{u}_i \mid \mathbf{y}, \boldsymbol{\theta}) \tilde{p}(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}$$

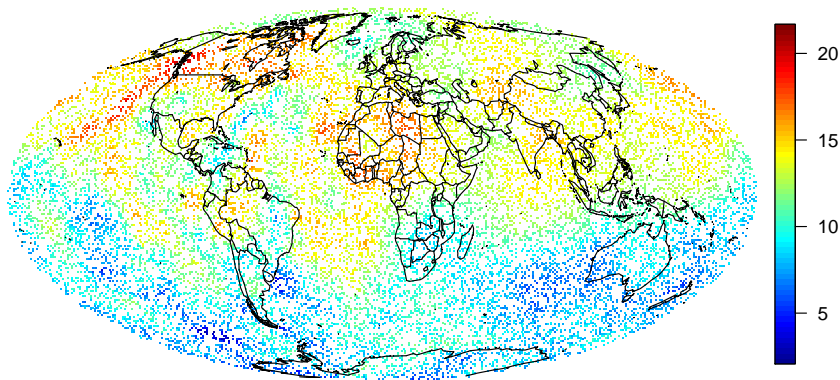
For a large model class this approach is much faster, as well as more accurate, than MCMC

SPDE based inference for point process data



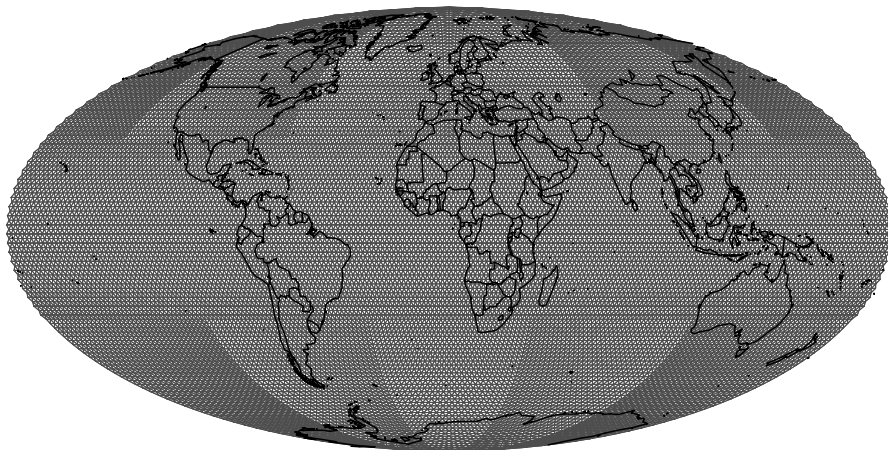
Synthetic CO2 data

Z(Dtrn)



Finite element mesh

Triangulation mesh

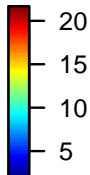
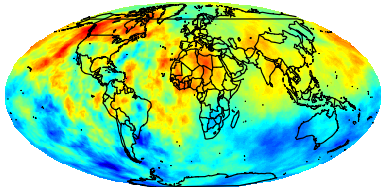


Finite element mesh

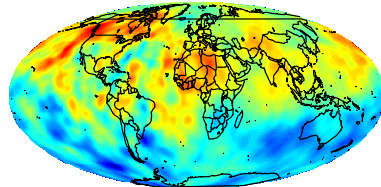


Estimates of CO2 from synthetic data

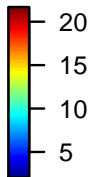
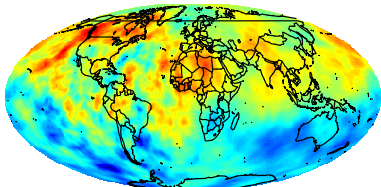
Ytrue



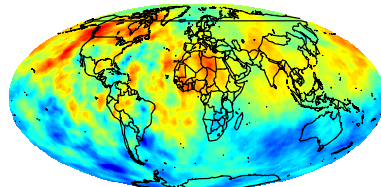
YLKR



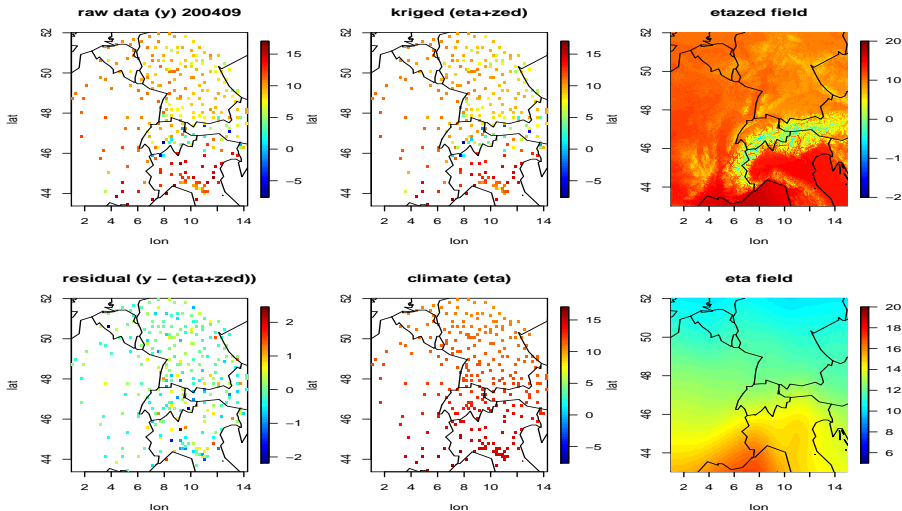
Y0



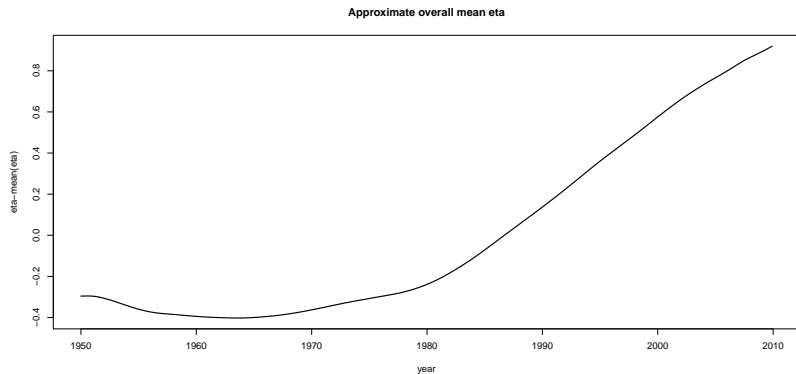
Y1



Temperature estimation with geographical covariates and stochastic fields



Spatial average climate estimate



Model:

weather = elevation effect + slow climate process + seasonal effects
+ fast weather process + observation noise

The computational work-horse

Cholesky decomposition (Cholesky, 1924)

$$Q = LL^T, \quad L \text{ lower triangular}$$

$$Q^{-1}x = L^{-T}L^{-1}x, \quad \text{via backward substitution}$$

$$\log \det Q = 2 \log \det L = 2 \sum_i \log L_{ii}$$

André-Louis Cholesky (1875–1918)

"He invented, for the solution of the condition equations in the method of least squares, a very ingenious computational procedure which immediately proved extremely useful, and which most assuredly would have great benefits for all geodesists, if it were published some day." (Euology by Commandant Benoit, 1922)



The end!