

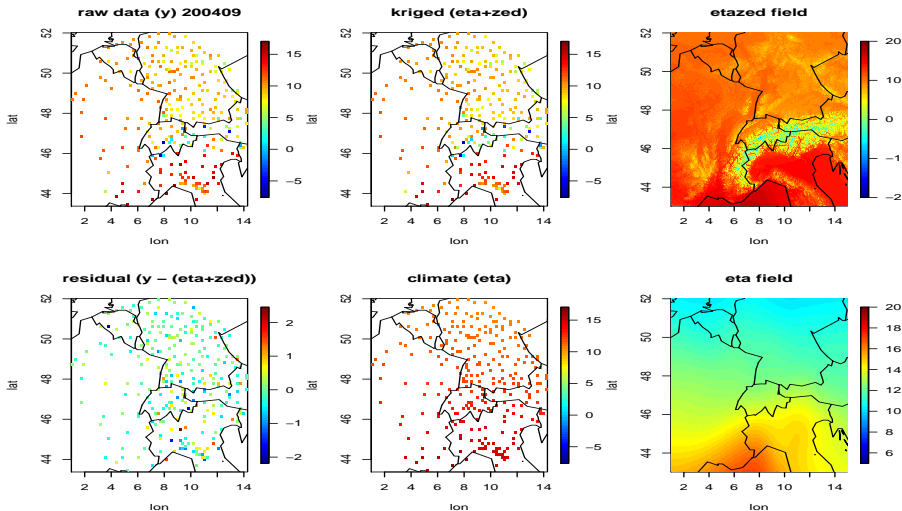
Towards realistic stochastic modelling of global *daily* temperatures

Finn Lindgren



Big Data in Environmental Science
University of British Columbia, PIMS
Vancouver, 12 May 2015

Sparse spatial coverage of temperature measurements



300 regional stations: $\approx 20,000,000$ from daily timeseries over 160 years

Full data: 70000 stations, several satellites, many ships, some lakes

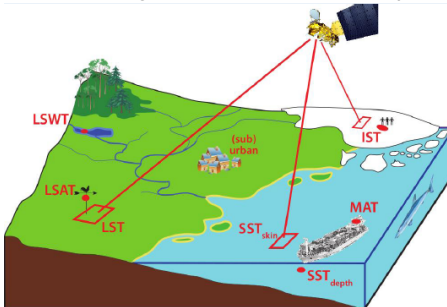
EUSTACE (*EU Surface Temperatures for All Corners of Earth*)



EUSTACE has received funding from the European Union's Horizon 2020 Programme for Research and Innovation, under Grant Agreement no 640171



EUSTACE will give publicly available daily estimates of surface air temperature since 1850 across the globe for the first time by combining surface and satellite data using novel statistical techniques.



Nick A. Rayner, Renate Auchmann, Janette Bessembinder, Stefan Brönnimann, Yuri Brugnara, Laura Carrea, Darren Ghent, Elizabeth Good, Katie Herring, Jacob Høyer, John Kennedy, Albert Klein Tank, Finn Lindgren, Colin Morice, Chris Merchant, John Remedios, Ag Stephens and Rasmus Tonboe

contact: nick.rayner@metoffice.gov.uk

EUSTACE data and modelling workflow

- ▶ *WP1: Translate sensor data into local estimates of air temperatures*
Relationships vary over land, ocean, ice, and lakes, and with season
The estimation errors are spatially correlated
Daily station data shows timeseries breaks in mean and distribution
- ▶ *WP2: Spatio-temporal blending of all data sources*
Two approaches: “Advanced traditional” and “Ambitious”
Point estimates of air temperature and uncertainty, and a random sample/ensemble from the posterior distribution
- ▶ *WP3: Validation of calibration models and data products*

Modelling challenges:

- ▶ Each data sources generates daily values of one (or more) of T_{\max} , T_{\min} , T_{avg} , and T_{range} at irregular locations
- ▶ T_{\max} and T_{\min} are strongly non-Gaussian and dependent, but T_{avg} and T_{range} are less dependent.
- ▶ After compensating for the seasonal cycle, T_{avg} is close to Gaussian, but T_{range} is non-Gaussian.

Preliminary model framework

- ▶ Build a priori independent multi-component models for T_{avg} and T_{range} , in space (\mathbf{s}) and time (year t and time within year τ):

$$T_{\text{avg}}(\mathbf{s}, t, \tau) = \text{seasonal}_{\text{avg}}(\mathbf{s}, \tau) + \text{longterm}_{\text{avg}}(\mathbf{s}, t, \tau) + \text{shortterm}_{\text{avg}}(\mathbf{s}, t, \tau)$$

A transformation or copula model is needed for the range, non-linear $h(\cdot)$:

$$\log(\text{scaling}_{\text{range}}(\mathbf{s}, t, \tau)) = \text{seasonal}_{\text{range}}(\mathbf{s}, \tau) + \text{longterm}_{\text{range}}(\mathbf{s}, t, \tau)$$

$$T_{\text{range}}(\mathbf{s}, t, \tau) = \text{scaling}_{\text{range}}(\mathbf{s}, t, \tau) \cdot h(\text{shortterm}_{\text{range}}(\mathbf{s}, t, \tau))$$

Each component is a Gaussian process with different spatial and temporal covariance properties.

- ▶ For observations of daily maxima and minima, define

$$T_{\text{min}} = T_{\text{avg}} - T_{\text{range}}/2$$

$$T_{\text{max}} = T_{\text{avg}} + T_{\text{range}}/2$$

Multiple data sources, all with their own issues

- ▶ Ground station air temperatures, with temporally persistent systematic deviations, e.g.

$$y_i^{\text{air,max}}(t) = T_{\text{max}}(s_i, t) + \sum_{k=1}^{K_i} H_k(t)\beta_{ik} + \epsilon_i(t),$$
 where $H_k(\cdot)$ typically are step functions around time series breakpoints.
- ▶ Ship measurements and buoys. Similar to ground stations, but moving around in space! And they measure water temperature as well. Sometimes with buckets.
- ▶ Geostationary and polar orbiting satellite measurements of land surface, ocean surface, and ice sheet surface temperatures. Different measurement footprints, spatially correlated errors, complicated links to air temperatures.
- ▶ Lake temperatures. Water temperatures linked with spatio-temporal averages of past air temperatures.

WP1 of EUSTACE will build calibration models for each data source, with uncertainties tracked into pseudo-observation models for air temperature, which is then fed into a unified Bayesian spatio-temporal estimate in WP2, via

$$p(T_{\text{avg}}, T_{\text{range}}, \beta \mid Y_1, \dots, Y_N) \propto p(T_{\text{avg}}, T_{\text{range}}, \beta) \prod_{k=1}^N p(Y_k \mid T_{\text{avg}}, T_{\text{range}}, \beta)$$

Laplace approximations for non-linear observations

By building the latent random field models as discretisations of stochastic PDEs collecting all the basis coefficients into a single a priori Gaussian vector, the spatio-temporal inference problem is turned into a sparse numerical algebra problem.

The non-linear transformations lead to a non-linear problem that can be approximated with a Gaussian posterior:

Quadratic posterior log-likelihood approximation

$$\begin{aligned} p(\mathbf{u} \mid \boldsymbol{\theta}) &\sim \mathcal{N}(\boldsymbol{\mu}_u, \mathbf{Q}_u^{-1}), \quad \mathbf{y} \mid \mathbf{u}, \boldsymbol{\theta} \sim p(\mathbf{y} \mid \mathbf{u}) \\ p_G(\mathbf{u} \mid \mathbf{y}, \boldsymbol{\theta}) &\sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{Q}}^{-1}) \\ \mathbf{0} &= \nabla_{\mathbf{u}} \{ \ln p(\mathbf{u} \mid \boldsymbol{\theta}) + \ln p(\mathbf{y} \mid \mathbf{u}) \} \Big|_{\mathbf{u}=\tilde{\boldsymbol{\mu}}} \\ \tilde{\mathbf{Q}} &= \mathbf{Q}_u - \nabla_{\mathbf{u}}^2 \ln p(\mathbf{y} \mid \mathbf{u}) \Big|_{\mathbf{u}=\tilde{\boldsymbol{\mu}}} \end{aligned}$$

The product structure for T_{range} doesn't necessarily generate positive definite precisions away from the mode, but the optimisation can use Fisher scoring, where $\mathbf{E}_{\mathbf{y} \mid \mathbf{u}} \left(-\nabla_{\mathbf{u}}^2 \ln p(\mathbf{y} \mid \mathbf{u}) \right)$ is positive definite.

Products of transformed processes

Assume that \mathbf{u} is a large scale process and \mathbf{v} is a small scale process, so that they are statistically identifiable from observations of the form

$$y_i = h_u(u_i) \cdot h_v(v_i) + \epsilon_i, \quad h_u \text{ and } h_v \text{ non-linear transformations.}$$

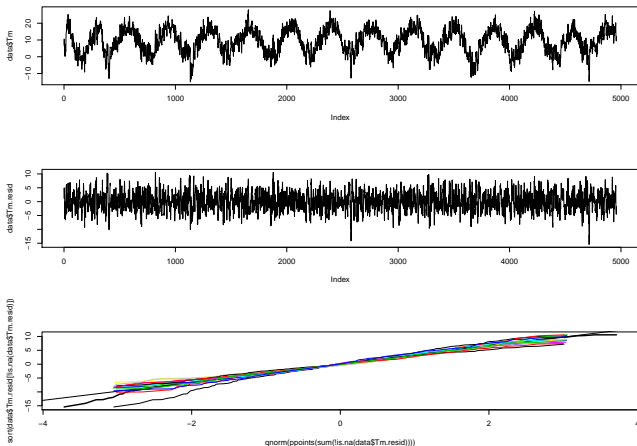
Write \mathbf{h}_u , \mathbf{h}'_u , \mathbf{h}''_u for the vectors of transformed values and derivatives of h_u at the u_i values, and similarly for \mathbf{v} . Then

$$\begin{aligned} C - \log p(\mathbf{y} \mid \mathbf{u}, \mathbf{v}) &= \frac{1}{2} (\mathbf{y} - \mathbf{h}_u \odot \mathbf{h}_v)^\top \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{h}_u \odot \mathbf{h}_v) \\ - \frac{\partial}{\partial \mathbf{v}} \log p(\mathbf{y} \mid \mathbf{u}, \mathbf{v}) &= - \text{diag}(\mathbf{h}_u \odot \mathbf{h}'_v) \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{h}_u \odot \mathbf{h}_v) \\ - \frac{\partial^2}{\partial \mathbf{v}^2} \log p(\mathbf{y} \mid \mathbf{u}, \mathbf{v}) &= \text{diag}(\mathbf{h}_u \odot \mathbf{h}'_v) \mathbf{Q}_\epsilon \text{diag}(\mathbf{h}_u \odot \mathbf{h}'_v) \\ &\quad - \text{diag}(\text{diag}(\mathbf{h}_u \odot \mathbf{h}''_v) \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{h}_u \odot \mathbf{h}_v)) \end{aligned}$$

and similarly for $\frac{\partial}{\partial \mathbf{u}}$, $\frac{\partial^2}{\partial \mathbf{u} \partial \mathbf{v}}$, and $\frac{\partial^2}{\partial \mathbf{u}^2}$. The problematic term in the Hessian involving \mathbf{y} disappears in Fisher scoring.

Seasonal effects and shortterm distribution for T_{avg}

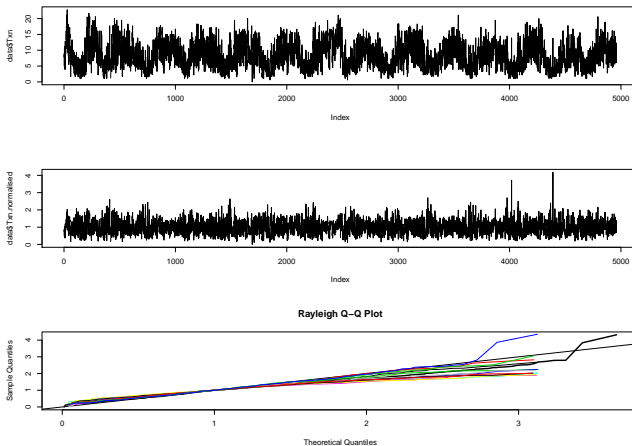
Raw data, estimated short term process, and Normal Q-Q plot for T_{avg}



The residual process is close to Normal, possibly with some seasonal variations.

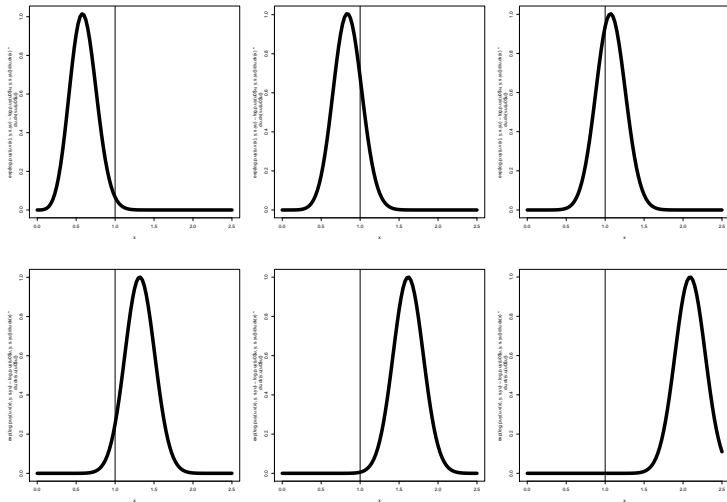
Seasonal effects and shortterm distribution for T_{range}

Raw data, estimated short term process, and Rayleigh Q-Q plot for T_{range}

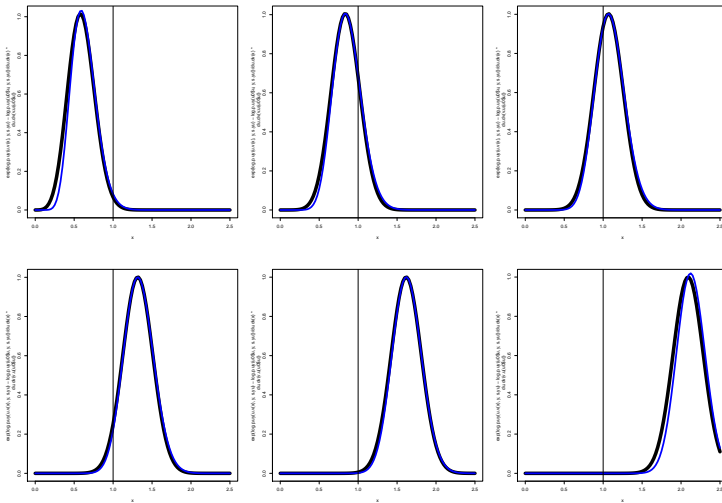


Under some theory, the residual process should have Rayleigh marginal distribution, but the true distribution has lighter tails, with a seasonal pattern.

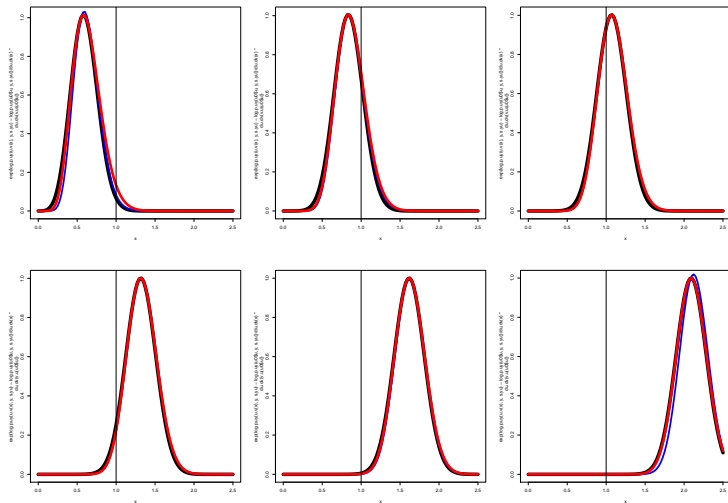
Posterior density approximations with Laplace



Posterior density approximations with Laplace



Posterior density approximations with Laplace



Illustrative modelling and computation principle, revisited

Ignoring the non-linearities, and the seasonal effects:

- ▶ A temporally slow, simplified stochastic heat equation (non-separable) for the longterm processes for T_{mean}

$$\frac{\partial}{\partial t} z(\mathbf{s}, t) - \gamma_z \nabla \cdot \nabla z(\mathbf{s}, t) = \mathcal{E}(\mathbf{s}, t)$$

$$(1 - \gamma_\mathcal{E} \nabla \cdot \nabla) \mathcal{E}(\mathbf{s}, t) = \mathcal{W}_\mathcal{E}(\mathbf{s}, t)$$

- ▶ A temporally quick, spatially non-stationary SPDE/GMRF (separable)

$$\left(\frac{\partial}{\partial t} + \gamma_t\right) (\kappa(\mathbf{s})^2 - \nabla \cdot \nabla) (\tau(\mathbf{s}) a(\mathbf{s}, t)) = \mathcal{W}_a(\mathbf{s}, t)$$
- ▶ Direct and linear observations:

$y_i = a(\mathbf{s}_i, t_i) + z(\mathbf{s}_i, t_i) + \epsilon_i$, discretised into

$$\mathbf{y} = \mathbf{A}(\mathbf{a} + (\mathbf{B} \otimes \mathbf{I})\mathbf{z}) + \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_\epsilon^{-1})$$

where \mathbf{B} maps from long-term basis functions to short-term, and \mathbf{A} maps from short-term basis functions to the observations.

The posterior precision can be formulated for $(\mathbf{a} + \mathbf{z}, \mathbf{z}) | \mathbf{y}$:

$$\mathbf{Q}_{(\mathbf{a}+\mathbf{z}, \mathbf{z}) | \mathbf{y}} = \begin{bmatrix} \mathbf{Q}_t \otimes \mathbf{Q}_a + \mathbf{A}^\top \mathbf{Q}_\epsilon \mathbf{A} & -\mathbf{Q}_t \mathbf{B} \otimes \mathbf{Q}_a \\ -\mathbf{B}^\top \mathbf{Q}_t \otimes \mathbf{Q}_a & \mathbf{Q}_z + \mathbf{B}^\top \mathbf{Q}_t \mathbf{B} \otimes \mathbf{Q}_a \end{bmatrix}$$

Locally isotropic 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

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

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

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

Combining this with an AR(1) discretisation of the temporal operator, we get

$$\mathbf{Q}_t \otimes \mathbf{Q}_a.$$

GMRF precision for the simplified stochastic heat equation

The precision is a quint-diagonal block matrix, with further structure

$$\begin{aligned} \mathbf{Q}_z &= \mathbf{M}_2^{(t)} \otimes \mathbf{M}_0^{(s)} + \mathbf{M}_1^{(t)} \otimes \mathbf{M}_1^{(s)} + \mathbf{M}_0^{(t)} \otimes \mathbf{M}_2^{(s)} \\ \mathbf{M}_0^{(s)} &= \mathbf{C} + \gamma_\varepsilon \mathbf{G} \\ \mathbf{M}_1^{(s)} &= \mathbf{G} + \gamma_\varepsilon \mathbf{G} \mathbf{C}^{-1} \mathbf{G} \\ \mathbf{M}_2^{(s)} &= \mathbf{G} \mathbf{C}^{-1} \mathbf{G} + \gamma_\varepsilon \mathbf{G} \mathbf{C}^{-1} \mathbf{G} \mathbf{C}^{-1} \mathbf{G} \end{aligned}$$

Ignoring the degenerate aspect of the model, the precision can be pseudo Cholesky factorised as $\mathbf{Q}_z = \tilde{\mathbf{L}}_z \tilde{\mathbf{L}}_z^\top$, where

$$\begin{aligned} \tilde{\mathbf{L}}_z &= \left[\left[\mathbf{L}_2^{(t)} \otimes \mathbf{L}_C, \quad \mathbf{L}_1^{(t)} \otimes \mathbf{L}_G, \quad \mathbf{L}_0^{(t)} \otimes \mathbf{G} \mathbf{L}_C^{-\top} \right], \right. \\ &\quad \left. \gamma_\varepsilon^{1/2} \left[\mathbf{L}_2^{(t)} \otimes \mathbf{L}_G, \quad \mathbf{L}_1^{(t)} \otimes \mathbf{G} \mathbf{L}_C^{-\top}, \quad \mathbf{L}_0^{(t)} \otimes \mathbf{G} \mathbf{C}^{-1} \mathbf{L}_G \right] \right] \end{aligned}$$

Since the kronecker products do not need to be explicitly stored, the pseudo Cholesky factors require very little time and memory, and can be computed even for a very large number of spatial basis functions (up to a million).

Posterior calculations, recalled from yesterday

Write $x = (a + z, z)$ for the full latent field.

$$Q_{x|y} = \begin{bmatrix} Q_t \otimes Q_a + A^\top Q_\epsilon A & -Q_t B \otimes Q_a \\ -B^\top Q_t \otimes Q_a & Q_z + B^\top Q_t B \otimes Q_a \end{bmatrix}$$

can be pseudo-Cholesky-factorised:

$$Q_{x|y} = \tilde{L}_{x|y} \tilde{L}_{x|y}^\top, \quad \tilde{L}_{x|y} = \begin{bmatrix} L_t \otimes L_a & \mathbf{0} & A^\top L_\epsilon \\ -B^\top L_t \otimes L_a & \tilde{L}_z & \mathbf{0} \end{bmatrix}$$

Posterior expectation, samples, and marginal variances (with $\tilde{A} = [A \quad \mathbf{0}]$):

$$Q_{x|y}(\mu_{x|y} - \mu_x) = \tilde{A}^\top Q_\epsilon (y - \tilde{A}\mu_x),$$

$$Q_{x|y}(x - \mu_{x|y}) = \tilde{L}_{x|y} w, \quad w \sim \mathcal{N}(\mathbf{0}, I), \quad \text{or}$$

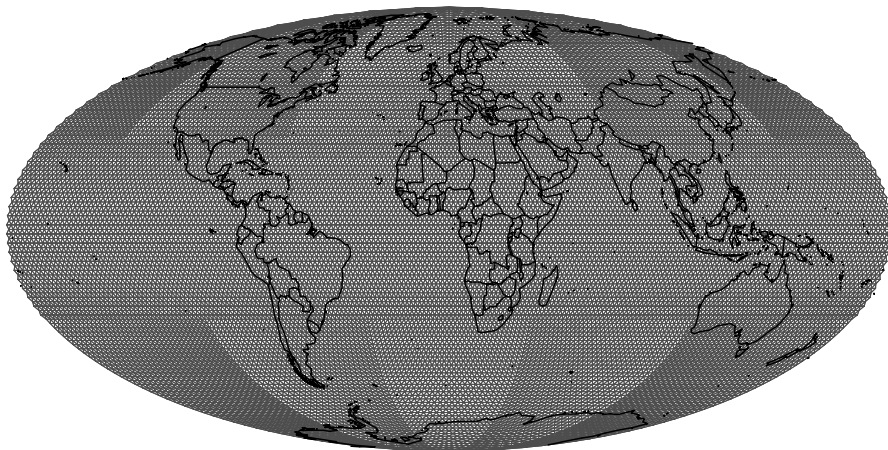
$$Q_{x|y}(x - \mu_x) = \tilde{A}^\top Q_\epsilon (y - \tilde{A}\mu_x) + \tilde{L}_{x|y} w, \quad w \sim \mathcal{N}(\mathbf{0}, I),$$

$$\text{Var}(x_i|y) = \text{diag}(\text{inla.qinv}(Q_{x|y})) \quad (\text{requires Cholesky})$$

The preconditioners for iterative solvers from yesterday are not fully satisfactory.

Finite element mesh

Triangulation mesh



Smaller subregions are partially self-contained sub-problems.

Domain decomposition

- ▶ Divide the domain into a collection of overlapping subdomain blocks
- ▶ Solve a local problem, e.g. the conditional solution, maintaining coherence by enforcing constraints on overlapping nodes.

Monte Carlo variance reduction for posterior variances

$$E(\mathbf{x}_i | \mathbf{y}) = E(E(\mathbf{x}_i | \mathbf{y}, \mathbf{x}_{\notin \text{subblock}}))$$

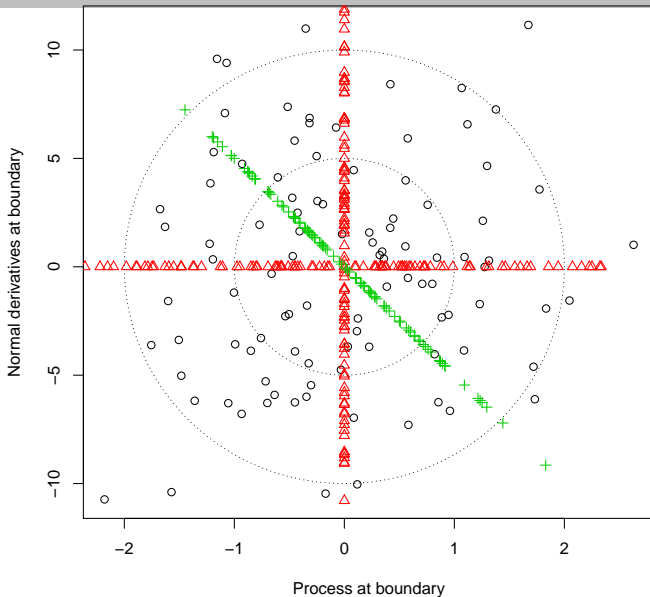
$$\text{Var}(\mathbf{x}_i | \mathbf{y}) = \text{Var}(\mathbf{x}_i | \mathbf{y}, \mathbf{x}_{\notin \text{subblock}}) + \text{Var}(E(\mathbf{x}_i | \mathbf{y}, \mathbf{x}_{\notin \text{subblock}}))$$

Also works for linear combinations, with some complications

Subdomain boundary adjustment (new idea)

- ▶ Apply *stochastic boundary correction* for each subdomain
- ▶ Solve the full local problem, reusing the appropriate randomness for overlapping subdomains
- ▶ Blend the results for overlapping domains.
- ▶ Apply this as a preconditioner in an iterative solver

All deterministic boundary conditions are 'inappropriate'



Stationary stochastic boundary adjustment

Recall the Matérn generating SPDE

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

RKHS inner product on a domain Ω for integer α :

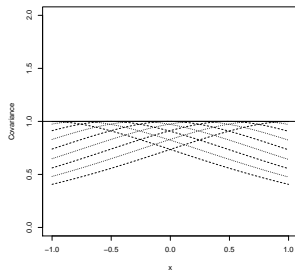
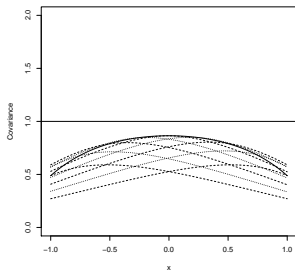
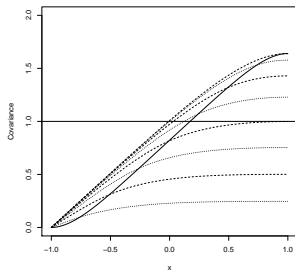
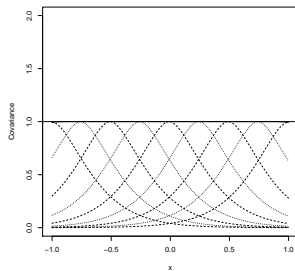
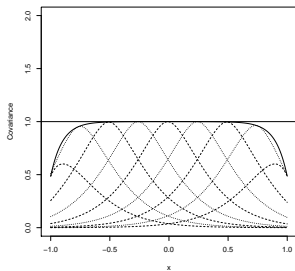
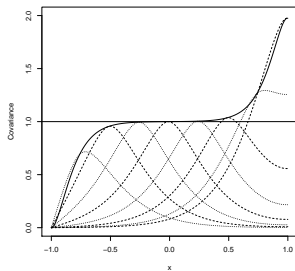
$$\langle f, g \rangle_{H(\Omega)} = \sum_{k=0}^{\alpha} \binom{\alpha}{k} \kappa^{2\alpha-2k} \langle \nabla^k f, \nabla^k g \rangle_{\Omega}$$

Boundary adjusted precision operator on a compact subdomain, where \mathcal{P} projects onto the operator null-space:

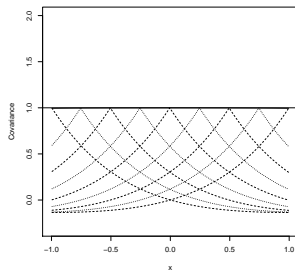
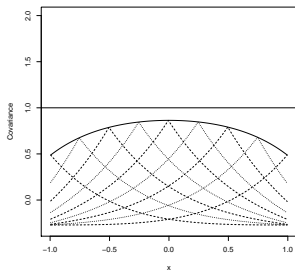
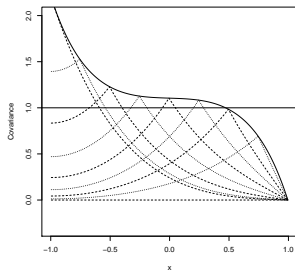
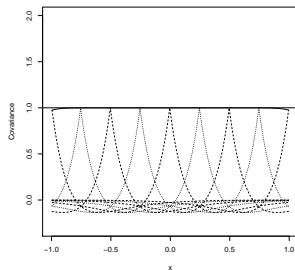
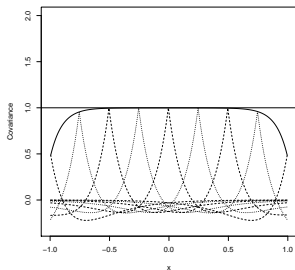
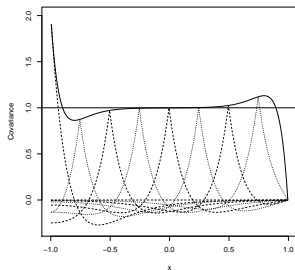
$$\begin{aligned} \mathcal{Q}_{\Omega}(f, g) &= \langle f, g \rangle_{H(\Omega)} - \langle \mathcal{P}f, \mathcal{P}g \rangle_{H(\Omega)} + \mathcal{Q}_{\mathcal{P};\partial\Omega}(\mathcal{P}f, \mathcal{P}g) \\ &= \langle f - \mathcal{P}f, g - \mathcal{P}g \rangle_{H(\Omega)} + \mathcal{Q}_{\mathcal{P};\partial\Omega}(\mathcal{P}f, \mathcal{P}g) \end{aligned}$$

Note that $\mathcal{Q}_{\mathcal{P};\partial\Omega}(\mathcal{P}f, \mathcal{P}g)$ may involve normal derivatives at the boundary.

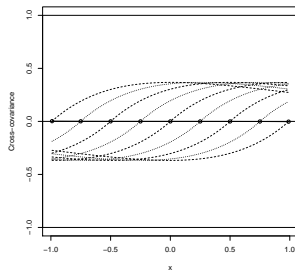
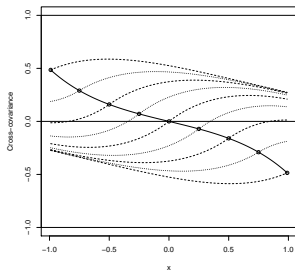
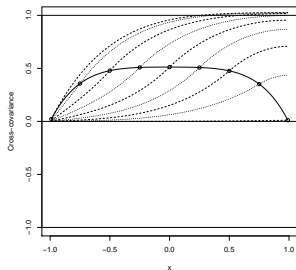
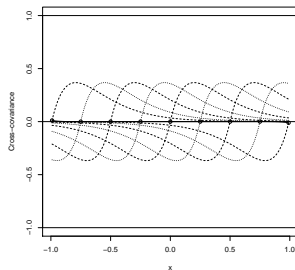
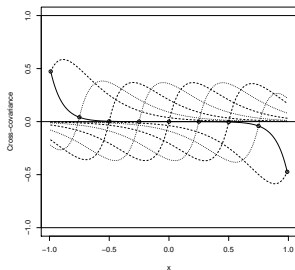
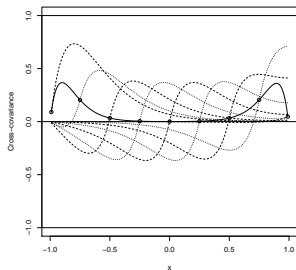
Covariances (D&N, Robin, Stoch) for $\kappa = 5$ and 1



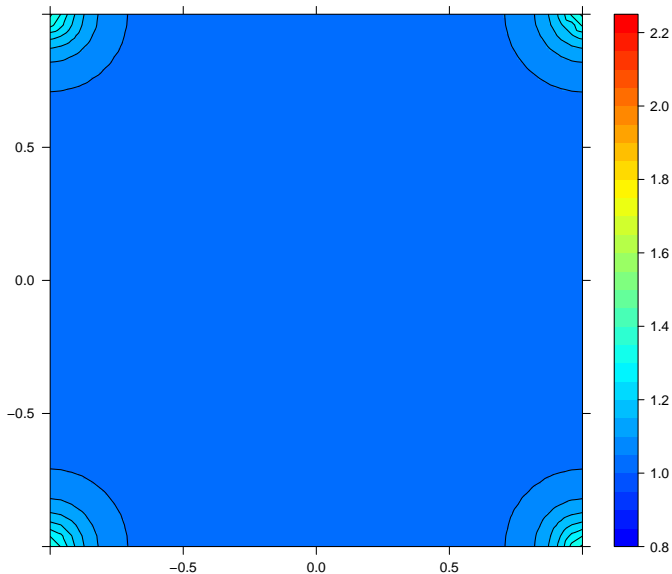
Derivative covariances (D&N, Robin, Stoch) for $\kappa = 5$ and 1



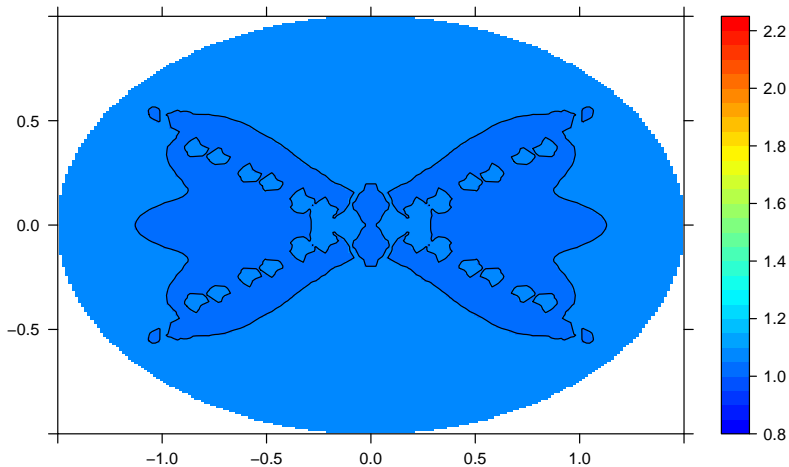
Process-derivative cross-covariances (D&N, Robin, Stoch)



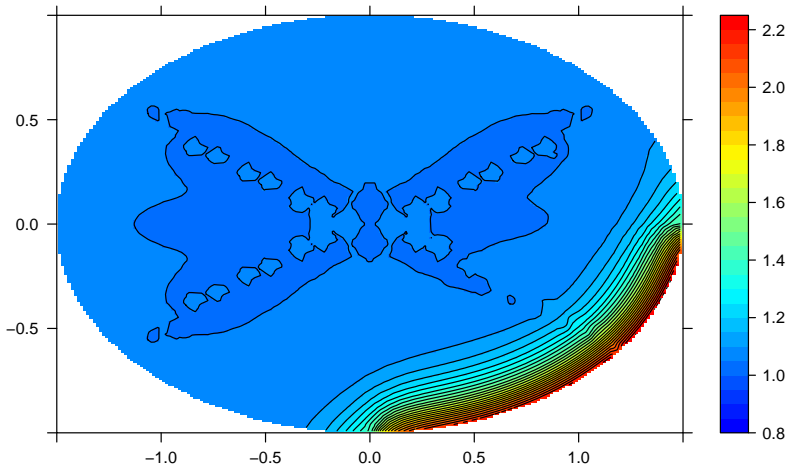
Square domain, stochastic boundary (variances)



Elliptical domain, stochastic boundary (variances)



Elliptical domain, mixed boundary (variances)



References

- ▶ Rue, H. and Held, L.: Gaussian Markov Random Fields; Theory and Applications; *Chapman & Hall/CRC*, 2005
- ▶ Lindgren, F.: Computation fundamentals of discrete GMRF representations of continuous domain spatial models; preliminary book chapter manuscript, 2014 (covers direct numerical methods)
<http://people.bath.ac.uk/fl353/tmp/gmrf.pdf>
- ▶ Lindgren, F., Rue, H., and Lindström, J.: An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach (with discussion); *JRSS Series B*, 2011
Non-CRAN package: R-INLA at <http://r-inla.org/>
- ▶ Bolin, D. and Lindgren, F.: Excursion and contour uncertainty regions for latent Gaussian models; *JRSS Series B*, 2014, in press. Accepted version at arXiv:1211.3946 and on journal web page.
CRAN package: `excursions`
Development: <http://bitbucket.org/davidbolin/excursions>