

Spatio-Temporal Smoothing for Very Large Datasets

Satellite Based Vegetation Measurements

Johan Lindström¹

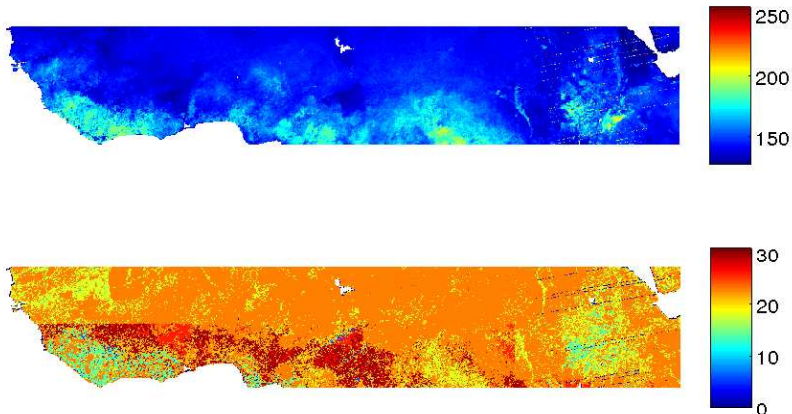
¹Centre for Mathematical Sciences
Lund University

Pan-American Advanced Study Institute
Búzios, June 25, 2014



LUND
UNIVERSITY

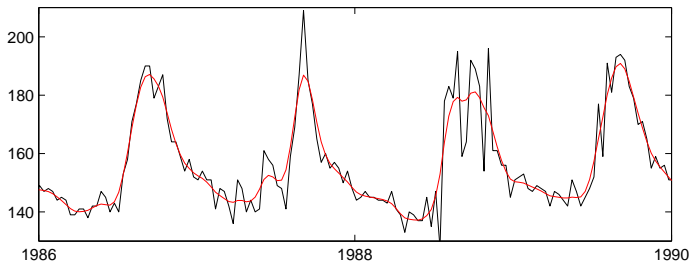
Satellite measured vegetation — The African Sahel



NDVI [0, 255] and cloud [0, 31] data from February 1985.

Smoothed version of the NDVI Data

Smooth the data to fill in missing values and remove noise due to cloud cover, etc.



Important ecological questions:

- ▶ Plant phenology (start and end of season)
- ▶ Plant productivity (integral)

What can the methods handle?

Assuming a standard desktop computer.

	N
Standard Kriging	$\sim 10^3$
Big N-methods	$\sim 10^5$
Satellite data	$\sim 10^8$

Big N-methods: process convolution (Higdon, 2001), tapering (Furrer et al., 2006), predictive processes (Banerjee et al., 2008), block composite (Eidsvik et al., 2014), SPDE/GMRF (Lindgren et al., 2011).

Space-time is often **worse**. For the SPDE-methods:

Dimension	Time for $R = \text{chol}(Q)$	$\text{nnz}(Q)$	$\text{nnz}(R)$
\mathbb{R}^1	$\mathcal{O}(N^1)$	$c^1 N$	$\sim N^1$
\mathbb{R}^2	$\mathcal{O}(N^{1.5})$	$c^2 N$	$\sim N^{1+1/2}$
\mathbb{R}^3	$\mathcal{O}(N^2)$	$c^3 N$	$\sim N^{1+2/3}$
\mathbb{R}^d	???	$c^d N$	$\sim N^{1+(d-1)/d}$

Can't we just wait for better computers?

Current Satellite data (1982–)

10 day temporal resolution

~ 8 km spatial resolution

New EU earth observation platform (launch 2015–)

3 – 5 day temporal resolution

≤ 1 km spatial resolution

Available data grows faster than processing power!

The standard likelihood

Log-likelihood

$$l(\theta|\mathbf{y}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}(\theta)| - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu}(\theta))^\top \boldsymbol{\Sigma}(\theta)^{-1} (\mathbf{y} - \boldsymbol{\mu}(\theta)),$$

with derivatives

$$\begin{aligned} \frac{\partial l(\theta|\mathbf{y})}{\partial \theta} &= -\frac{1}{2} \operatorname{tr} \left(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \theta} \right) + (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\mu}}{\partial \theta} \\ &\quad + \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \theta} \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \end{aligned}$$

Spline smoothing (Wahba, 1981)

$$\operatorname{GCV}(\theta) = \frac{\|\mathbf{y} - \mathbf{S}(\theta)\mathbf{y}\|_2^2}{(n - \operatorname{tr}(\mathbf{S}(\theta)))^2}, \quad \mathbf{S}(\theta) \propto (\mathbf{A}^\top \mathbf{A} + \mathbf{Q}(\theta))^{-1}$$

Parameter estimates

Parameters are estimated by one of:

1. $\operatorname{argmax}_{\theta} l(\theta|\mathbf{y})$.
2. Solving: $\nabla l(\theta|\mathbf{y}) = 0$.
3. $\operatorname{argmin}_{\theta} \operatorname{GCV}(\theta)$.

and we need to compute:

$\log \Sigma $	For case 1.
$\Sigma^{-1}\mathbf{y}$	For all cases.
$\operatorname{tr}(\Sigma^{-1})$	For cases 2, 3.

These can be approximated to high accuracy (Aune et al., 2014; Anitescu et al., 2012; Stein et al., 2013)

The Hutchinson (1990) trace estimator

The trace of a matrix Σ^{-1} is given by

$$\text{tr}(\Sigma^{-1}) \approx \frac{1}{K} \sum_{k=1}^K \mathbf{e}_k^\top \Sigma^{-1} \mathbf{e}_k$$

where the elements of \mathbf{e}_k are iid $\{-1, 1\}$ with probability $1/2$.

$$\mathbb{E}(\mathbf{e}^\top \Sigma^{-1} \mathbf{e}) = \text{tr}(\Sigma^{-1} \mathbb{E}(\mathbf{e}\mathbf{e}^\top)) = \text{tr}(\Sigma^{-1} \mathbf{I})$$

$$\text{V}(\mathbf{e}^\top \Sigma^{-1} \mathbf{e}) = 2 \left(\|\Sigma^{-1}\|_F^2 - \sum_i (\Sigma_{ii}^{-1})^2 \right) = 2 \sum_{i \neq j} (\Sigma_{ij}^{-1})^2$$

Usefull matrix algebra Harville (1997); Petersen and Pedersen (2012)

Probing vectors

A time series example ($K = 5$):

$$\mathcal{V} = \begin{bmatrix} \pm 1 & 0 & 0 & 0 & 0 \\ 0 & \pm 1 & 0 & 0 & 0 \\ 0 & 0 & \pm 1 & 0 & 0 \\ 0 & 0 & 0 & \pm 1 & 0 \\ 0 & 0 & 0 & 0 & \pm 1 \\ \pm 1 & 0 & 0 & 0 & 0 \\ 0 & \pm 1 & 0 & 0 & 0 \\ \vdots & & & & \end{bmatrix} \quad \mathcal{V}(\mathcal{V}\mathcal{V}^\top) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ \vdots & & & & & & \ddots \end{bmatrix} \dots$$

\mathbf{v}_k columns of \mathcal{V} .

Similar argument in \mathbb{R}^2 .

In general a **NP-complete graph-colouring** problem, but **near optimal greedy algorithms** exist.

Probing vectors II

Variance for the original estimator:

$$V \left(\frac{1}{K} \sum_{k=1}^K \mathbf{e}_k^\top \Sigma^{-1} \mathbf{e}_k \right) = \frac{2}{K} \sum_{i \neq j} (\Sigma_{ij}^{-1})^2$$

Variance for the new estimator:

$$\begin{aligned} V \left(\sum_{k=1}^K \mathbf{v}_k^\top \Sigma^{-1} \mathbf{v}_k \right) &= 2 \left(\left\| \Sigma^{-1} \circ V(\mathbf{v}\mathbf{v}^\top) \right\|_F^2 - \sum_{i=1}^n (\Sigma_{ii}^{-1})^2 \right) \\ &= 2 \sum_{\substack{i=j+mK \\ m \in \mathbb{Z}^+}} (\Sigma_{ij}^{-1})^2 \end{aligned}$$

This is better if elements in Σ^{-1} decrease away from the diagonal

Parameter estimates (again)

Parameters are estimated by one of:

1. $\operatorname{argmax}_{\theta} l(\theta|\mathbf{y})$.
2. Solving: $\nabla l(\theta|\mathbf{y}) = 0$.
3. $\operatorname{argmin}_{\theta} \operatorname{GCV}(\theta)$.

and we need to compute:

$\log \Sigma $	For case 1.
$\Sigma^{-1}\mathbf{y}$	For all cases.
$\operatorname{tr}(\Sigma^{-1})$	For cases 2, 3.

We can reduce $\operatorname{tr}(\Sigma^{-1})$ to computations of $\Sigma^{-1}\mathbf{b}$.

Aune et al. (2014) uses $\log |\Sigma| = \operatorname{tr} \log \Sigma$, and a series expansion of $\log \Sigma$.

Using the trace estimator the problem reduces to computing

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \quad \text{for different } \mathbf{b}$$

Preconditioned Conjugate Gradient

Iterative method for solving

$$\mathbf{Ax} = \mathbf{b}$$

using a **preconditioner**, \mathbf{M} , and computing only \mathbf{Ax} and $\mathbf{M}^{-1}\mathbf{x}$.

A suitable preconditioner \mathbf{M} has

- ▶ $\mathbf{M}^{-1}\mathbf{x}$ easy to compute
- ▶ $\mathbf{M}^{-1}\mathbf{A} \approx \mathbf{I}$.

Classical solutions

- ▶ Use $\text{diag}(\mathbf{A})^{-1}$.
- ▶ For sparse matrices use the (modified) incomplete Cholesky.

Comments

The Good

- ▶ No need to store the full matrices.
- ▶ $\mathcal{O}(N^1)$ for sparse matrices.
- ▶ Anitescu et al. can be seen as GEE \rightarrow parameter uncertainties

The Bad

- ▶ Still $\mathcal{O}(N^2)$ for covariances matrices.
- ▶ Only Aune et al. evaluates the likelihood.

The Ugly

- ▶ Numerical linear solvers in [MATLAB](#) and [R](#) are reasonable, but not outstanding.
- ▶ Sparse matrix support is “lacking” in [MATLAB](#) and bad in [R](#).
- ▶ Optimal preconditioner depends on the problem.

Smoothing splines (Wahba, 1981)

Find $\hat{\mu}(s)$ such that

$$\hat{\mu} = \underset{\mu}{\operatorname{argmin}} \frac{1}{n} \sum_i \|y_i - \mu_i\|_2^2 + \lambda \int \left(\frac{\partial^\alpha \mu(s)}{\partial s^\alpha} \right)^2 ds$$

$\hat{\mu}$ is given as the solutions to an SPDE (Wahba, 1981):

$$\Delta^{\alpha/2} \mathbf{x}(s) = \mathcal{W}$$

This corresponds to a Matérn (1960) covariance with **range** $\rightarrow \infty$, and the spline can be seen as a hierarchical model

$$\mathbf{y} | \mathbf{x} \in \mathcal{N}(\mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I}) \quad \mathbf{x} \in \mathcal{N}(\mathbf{0}, (\tau^2 \mathbf{Q})^{-1}).$$

Smoothing splines \approx Gaussian processes

The spline smoothing of \mathbf{y} can be seen as a conditional expectation

$$\begin{aligned}\hat{\boldsymbol{\mu}} &= \mathbf{E}(\mathbf{x}|\mathbf{y}; \tau^2, \sigma^2) = \left(\tau^2\mathbf{Q} + \mathbf{A}^\top(\sigma^2\mathbf{I})^{-1}\mathbf{A}\right)^{-1}\left(\mathbf{A}^\top(\sigma^2\mathbf{I})^{-1}\mathbf{y}\right) \\ &= \left(\lambda\mathbf{Q} + \mathbf{A}^\top\mathbf{A}\right)^{-1}\mathbf{A}^\top\mathbf{y},\end{aligned}$$

where $\lambda = \tau^2\sigma^2$.

λ found by minimizing the GCV-error

$$\text{GCV}(\lambda) = \frac{\sum_i \|\mathbf{y}_i - \hat{\boldsymbol{\mu}}_i\|_2^2}{(n - \text{tr}(\mathbf{S}(\lambda)))^2}$$

with $\mathbf{S}(\lambda) = \mathbf{A} \left(\lambda\mathbf{Q} + \mathbf{A}^\top\mathbf{A}\right)^{-1}\mathbf{A}^\top$

Smoothing splines on a regular grid

For our splines ($\alpha = 2$) we have

$$\hat{\mu} = (\lambda \mathbf{Q} + \mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{y}$$

$$\mathbf{Q} = \mathbf{G}^\top \mathbf{G}$$

where $\mathbf{G} \approx \Delta$.

For a regular grid in 1D:

$$\mathbf{G} = \begin{bmatrix} 1 & -1 & 0 & \dots & \dots \\ -1 & 2 & -1 & 0 & \dots \\ 0 & -1 & 2 & -1 & \dots \end{bmatrix}$$

Neumann boundary condition

Smoothing splines on a regular grid

For our splines ($\alpha = 2$) we have

$$\hat{\mu} = (\lambda \mathbf{Q} + \mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{y}$$

$$\mathbf{Q} = \mathbf{G}^\top \mathbf{G}$$

where $\mathbf{G} \approx \Delta$.

For a regular grid in 1D:

$$\mathbf{G} = \begin{bmatrix} \mathbf{2} & -1 & 0 & \dots & \mathbf{1} \\ -1 & \mathbf{2} & -1 & 0 & \dots \\ 0 & -1 & \mathbf{2} & -1 & \dots \end{bmatrix}$$

Cyclic boundary condition (Thon et al., 2012)

Preconditioner on the grid

Assume complete observations, $\mathbf{A} = \mathbf{I}$, and take $\mathbf{M} = \lambda \mathbf{Q} + \mathbf{I}$.

$$\mathbf{M} = \lambda \underbrace{\mathbf{G}^\top \mathbf{G}}_{=\mathbf{U}\mathbf{\Gamma}\mathbf{U}^\top} + \underbrace{\mathbf{I}}_{=\mathbf{U}\mathbf{U}^\top} = \mathbf{U} \underbrace{(\lambda \mathbf{\Gamma}^2 + \mathbf{I})}_{=\mathbf{\Xi}} \mathbf{U}^\top,$$

where

$$\mathbf{U} = \begin{cases} \text{DCT,} & \text{Neumann} \\ \text{FFT,} & \text{Cyclic} \end{cases}$$

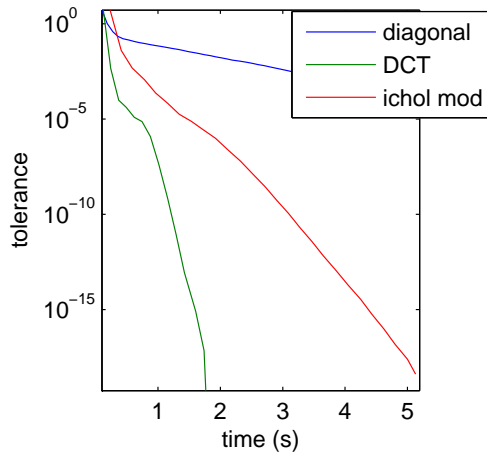
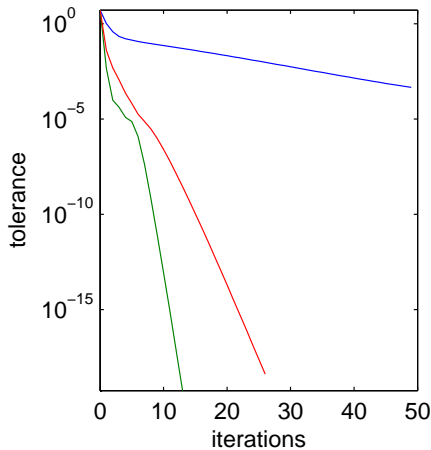
$$\mathbf{\Gamma}_{ii} = \begin{cases} 2 - 2 \cos \left((i-1) \frac{\pi}{n} \right), & \text{Neumann} \\ 2 - 2 \cos \left((i-1) \frac{2\pi}{n} \right), & \text{Cyclic} \end{cases}$$

For Neumann boundary conditions

$$\mathbf{M}^{-1} \mathbf{b} = \text{iDCT} \left(\mathbf{\Xi}^{-1} \text{DCT}(\mathbf{b}) \right), \quad \mathbf{\Xi}_{ii} = \lambda \mathbf{\Gamma}_{ii}^2 + 1$$

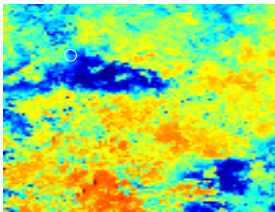
Convergence speed for PCG

Test run on a $2D$ 256×256 grid

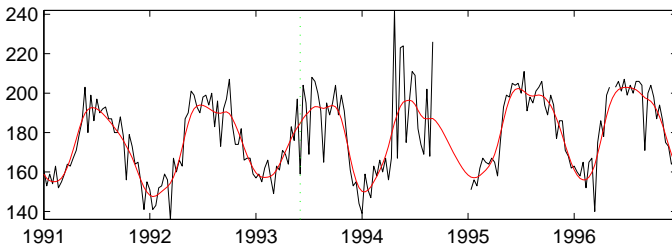
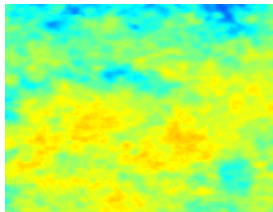


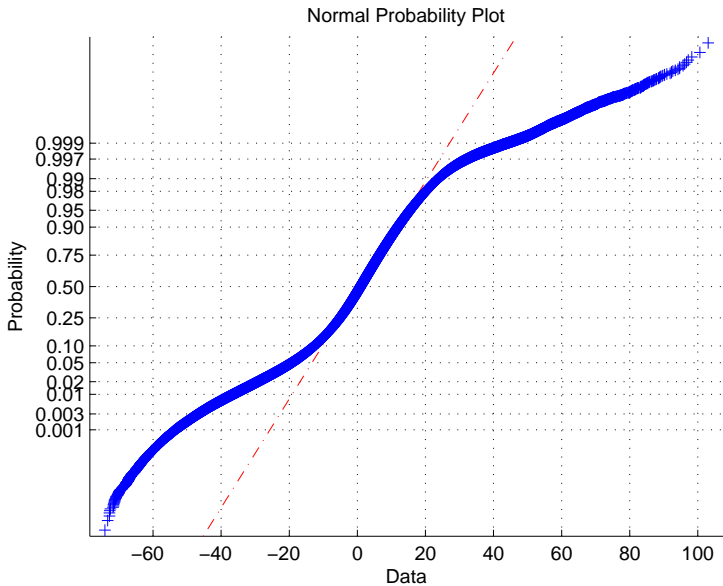
Preliminary results

NDVI May 1993

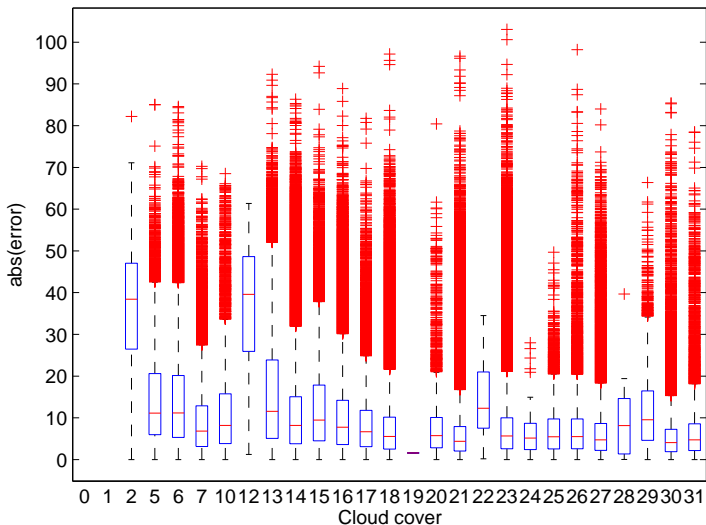


Smooth



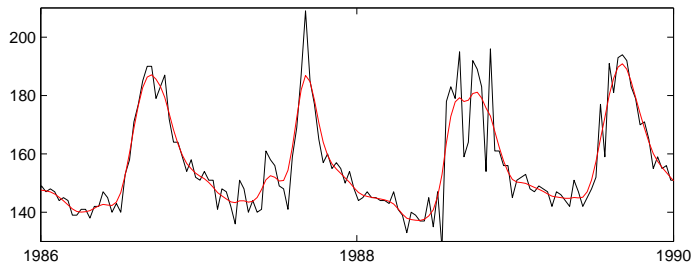


Can we use the cloud cover?



Further problems

- ▶ Non-Gaussian errors
- ▶ Axis anisotropy
- ▶ Non-stationarity in time
- ▶ Spatially correlated errors (due to clouds)
- ▶ Non-stationarity in space?



*“When optimising function that involve high-dimensional determinant approximations, it is important to **use whatever structure is available** in order to speed up computations.”*

— Aune et al. (2014, pp. 256)

Bibliography I

- Anitescu, M., Chen, J., and Wang, L. (2012), "A Matrix-free Approach for Solving the Parametric Gaussian Process Maximum Likelihood Problem," *SIAM J. Sci. Comput.*, 34, A240–A262.
- Aune, E., Simpson, D. P., and Eidsvik, J. (2014), "Parameter estimation in high dimensional Gaussian distributions," *Stat. Comput.*, 24, 247–263.
- Banerjee, S., Gelfand, A. E., Finley, A. O., and Sang, H. (2008), "Gaussian predictive process models for large spatial data sets," *J. R. Stat. Soc. B*, 70, 825–848.
- Eidsvik, J., Shaby, B. A., Reich, B. J., Wheeler, M., and Niemi, J. (2014), "Estimation and Prediction in Spatial Models With Block Composite Likelihoods," *J. Comput. Graphical Stat.*, 23, 295–315.
- Furrer, R., Genton, M. G., and Nychka, D. (2006), "Covariance Tapering for Interpolation of Large Spatial Datasets," *J. Comput. Graphical Stat.*, 15, 502–523.
- Harville, D. A. (1997), *Matrix Algebra From a Statistician's Perspective*, Springer, 1st ed.
- Higdon, D. (2001), "Space and space time modeling using process convolutions," Tech. rep., Institute of Statistics and Decision Sciences, Duke University, Durham, NC, USA.

Bibliography II

- Hutchinson, M. (1990), "A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines," *Commun. Stat. Simul. Comput.*, 19, 433–450.
- Lindgren, F., Rue, H., and Lindström, J. (2011), "An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach," *J. R. Stat. Soc. B*, 73, 423–498.
- Matérn, B. (1960), "Spatial variation. Stochastic models and their application to some problems in forest surveys and other sampling investigations." Ph.D. thesis, Stockholm University, Stockholm, Sweden.
- Petersen, K. B. and Pedersen, M. S. (2012), "The Matrix Cookbook," <http://matrixcookbook.com>.
- Stein, M. L., Chen, J., and Anitescu, M. (2013), "Stochastic approximation of score functions for Gaussian processes," *Ann. Stat.*, 7, 1162–1191.
- Thon, K., Rue, H., Skrøvseth, S. O., and Godtliebsen, F. (2012), "Bayesian multiscale analysis of images modeled as Gaussian Markov random fields," *Comput. Stat. Data Anal.*, 56, 49–61.
- Wahba, G. (1981), "Spline interpolation and smoothing on the sphere," *SIAM J. Sci. Comput.*, 2, 5–16.