

# Posterior Inference for Sparse Hierarchical Non-stationary Models

Lassi Roininen

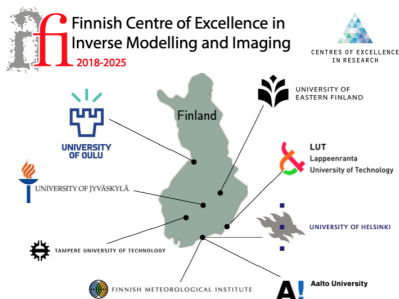
Lappeenranta University of Technology, Finland

with Karla Monterrubio-Gómez, Sara Wade, Neil Chada, Janne Huttunen, Kenneth Muhumuza,  
Sari Lasanen, Jarkko Suuronen, Alberto Mendoza, Theo Damoulas, Mark Girolami



Potsdam 2020

# Centre of Excellence in Inverse Modelling and Imaging



# Motivation

- Detect material interfaces, inhomogeneous structures, anisotropies, Gaussian and non-Gaussian features
- Gaussian and non-Gaussian hierarchical random field priors, parametric models for structures, noise models etc
- Metropolis-within Gibbs, elliptical slice sampling, Hamiltonian Monte Carlo, and optimisation methods
- Applications: Subsurface imaging (electrical impedance tomography, Darcy flow models) and near-space remote sensing (High-power radar experiments, satellite tomography and remote sensing)

# Non-Gaussian Priors

# TV and Besov space priors

- Lassas and Siltanen 2004 showed that TV are not discretisation-invariant
- Lassas, Saksman and Siltanen 2009 constructed Besov space priors
  - Often defined via wavelet expansions.
  - For edge-preserving inversion the Haar wavelet basis is often used
  - However due to the structure of the Haar basis, discontinuities are preferred on an underlying dyadic grid given by the discontinuities of the basis functions. For example, on the domain  $(0, 1)$ , discontinuity is vastly preferred at  $x = 1/4$  over  $x = 1/3$ .
  - Thus Besov priors make, in most practical cases, both a strong and unrealistic assumption.

# Non-Gaussian models – $\alpha$ -stable priors

- Markku Markkanen, Lassi Roininen, Janne M J Huttunen and Sari Lasanen, Cauchy difference priors for edge-preserving Bayesian inversion, *Journal of Ill-posed and Inverse Problems* (2019).
- Alberto Mendoza, Lassi Roininen, Mark Girolami, Jere Heikkinen and Heikki Haario, *Statistical Methods To Enable Practical On-Site Tomographic Imaging of Whole-Core Samples*, *Geophysics* (2019).
- Neil Chada, Sari Lasanen and Lassi Roininen, *Posterior Convergence Analysis of  $\alpha$ -Stable Sheets*, *arXiv* (2019).
- Kenneth Muhumuza, Lassi Roininen, Janne M. J. Huttunen, Timo Lähivaara, *A Bayesian-based approach to improving acoustic Born waveform inversion of seismic data for viscoelastic media*, *arXiv* (2019).

# Stable random walks

- Let  $U(t), t \in \mathbb{I} \subset \mathbb{R}^+$  be a stochastic process. We call it a Lévy  $\alpha$ -stable process starting from zero, or simply as stable process, if  $U(0) = 0$ ,  $U$  has independent increments and

$$U(t) - U(s) \sim S_\alpha \left( (t-s)^{1/\alpha}, \beta, 0 \right) \quad (1)$$

for any  $0 \leq s < t < \infty$  and for some  $0 < \alpha \leq 2, -1 \leq \beta \leq 1$ .

- For the continuous limit of the Cauchy walk, we apply independently scattered measures. We obtain random walk approximation

$$U_{t_i} - U_{t_{i-1}} \sim S_\alpha(h^{\frac{1}{\alpha}}, \beta, 0)$$

where  $t_i - t_{i-1} =: h$ . It is easy to see that such random walk approximations converge to the  $\alpha$ -stable Lévy motion as  $h \rightarrow 0$  in distribution on the Skorokhod space of functions that are right-continuous and have left limits.

# Chada, Lasanen, Roininen: $\alpha$ -stable sheet paper

- The paper lays the groundwork for Bayesian inverse problems with stable fields, specifically stable stochastic integrals  $U(x) = \int_E f(x, x') M(dx')$
- The paper has expository flavour: We study the very simple stable sheets as an illustrative and easy to follow example. For stable sheets,

$$f(x, x') = \begin{cases} 1 & \text{when } x'_i \leq x_i \text{ for all } i = 1, \dots, d \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

- Stable integral is defined like the usual Itô integral, but with stable random measure  $M$  in place of Brownian motion  $B$ /Brownian sheet.
- Note: no second moments means that instead of  $L^2$ , the integrals are limits of integrals of simple functions in probability.

# Stable Integrals

- Stable integrals  $U(x)$  can be presented in many equivalent ways (equivalent = equivalent in distribution) When it comes to Bayesian inverse problems, the best way seems to be through Lévy-LePage series representation,
- Lévy-LePage series representation is

$$U(x) = (C_\alpha |E|)^{\frac{1}{\alpha}} \sum_{k=1}^{\infty} \rho_k \Gamma_k^{1/\alpha} f(x, V_k), \quad (3)$$

where  $0 < \alpha < 2$ ,

$$C_\alpha = \left( \int_0^\infty x^{-\alpha} \sin(x) dx \right)^{-1}, \quad (4)$$

$\rho_k$  is a Rademacher sequence (i.i.d. with values  $\pm 1$  with equal probabilities),  $\Gamma_k$  are arrival times of a Poisson process with arrival rate 1, and  $V_k$  are i.i.d. uniformly distributed on  $E$ . The three sequences  $\rho_k$ ,  $\Gamma_k$  and  $V_k$  are mutually independent.

# Discretisation

- Lévy-LePage series representation gives
  - 1) Sample path regularity in  $L^p$ ,  $1 \leq p < \infty$  (also in the more general Sobolev space  $H_p^s$ ,  $s < 1/p$ ,  $p \geq 2$ ),
  - 2) Convergence in distributions on sample space.
- From 1) and 2), we proceed to posterior convergence in distribution for finite-dimensional data. The discretization of  $U$  on  $[0, 1]^d$  is taken to be

$$U^N(x) = U(h \lceil x/h \rceil), \quad (5)$$

where the ceiling function  $\lceil t \rceil = \min\{m \in \mathbb{Z}^d : t_j \leq m_j, j = 1, \dots, d\}$ .

- Computationally, the discretisation is determined from set of difference equations (here in 2D case)

$$\begin{aligned} U(hm_1, hm_2) - U(hm_1, h(m_2 - 1)) - U(h(m_1 - 1), hm_2) \\ + U(h(m_1 - 1), h(m_2 - 1)) \sim S_\alpha(|h|^{d/\alpha}, 0, 0) \end{aligned} \quad (6)$$

with i.i.d. right hand sides and zero boundary values on the coordinate axes.

# Convergence $h \rightarrow 0$

## Theorem

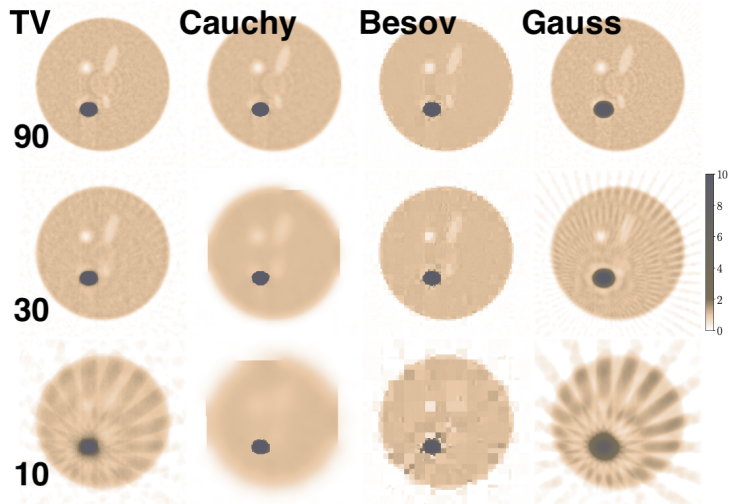
*Let  $1 \leq p < \infty$ . The approximations  $U^N(x) = U(h \lceil x/h \rceil)$  converge to  $U$  on  $L^p((0,1)^d)$  in distribution.*

Open questions:

- Can we do the same for infinite-dimensional data (e.g. with Gaussian noise)?
- How to obtain stronger posterior convergence for  $\alpha = 1$ ?

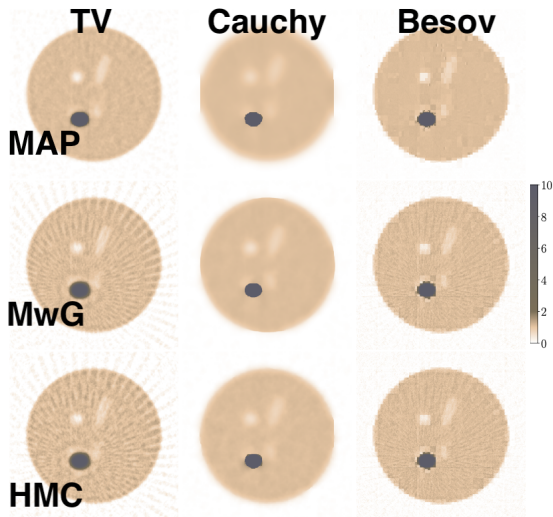
# MAP estimates

- Log tomography with different number of projections and prior models



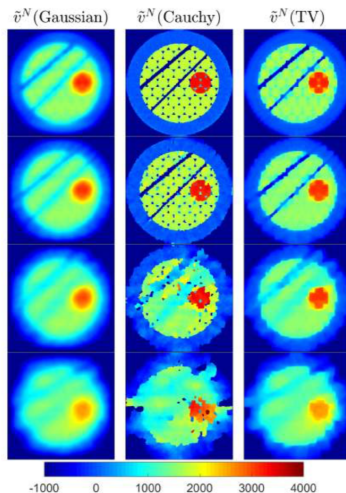
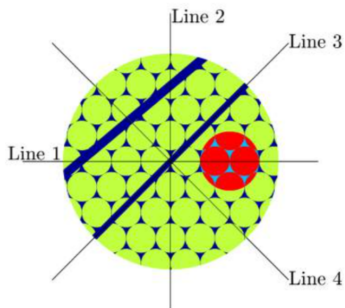
# MAP and CM estimates – 30 projections

- Log tomography



# Tomographic imaging of whole-core samples

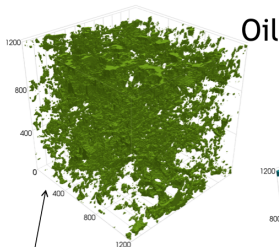
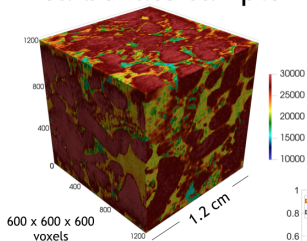
- 46, 23, 12, 6 projections with 10% noise



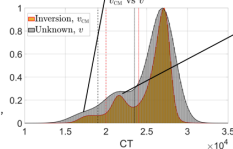
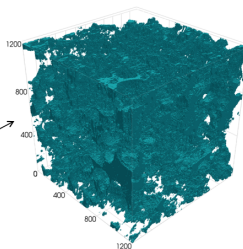
# Micro-CT

**Example 4:** mixed-wet carbonate reservoir rocks from the Middle-East\*.

Waterflood oil-bearing  
carbonate sample



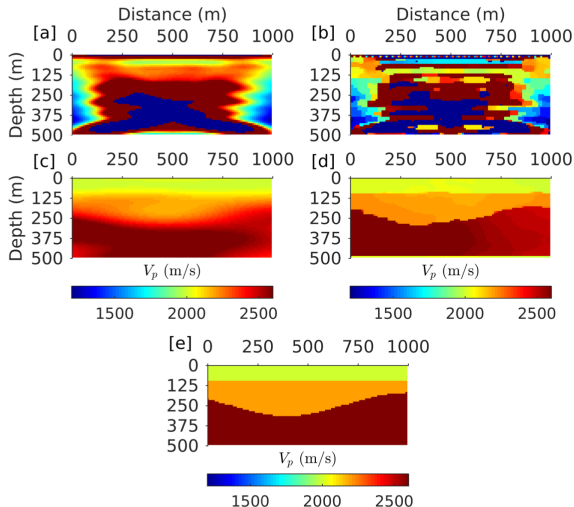
Water



\*Alhammadi, A., Alratrout, A., Bijeljic, B., & Blunt, M. (2018, May 17). X-ray micro-tomography datasets of mixed-wet carbonate reservoir rocks for in situ effective contact angle measurements. Retrieved August 16, 2018, from [www.digitalrockportal.org](http://www.digitalrockportal.org)

# Born waveform inversion of seismic data

- A Bayesian approach to improving acoustic Born waveform inversion of seismic data for viscoelastic media



# Sparse Hierarchical Non-stationary Models

Karla Monterrubio-Gómez, Lassi Roininen, Sara Wade, Theodoros Damoulas, and Mark Girolami, Posterior Inference for Sparse Hierarchical Non-stationary Models, arXiv 2019.

# Hierarchical GP model

- Based on the non-stationary Matérn kernel via varying length-scaling  $\ell(x_i)$ .
- Hierarchical model for 1- $d$  problems:

$$\begin{aligned}
 y_i &\sim \mathcal{N}(z(x_i), \sigma_\varepsilon^2), \quad i = 1, \dots, m, \\
 z(\cdot) &\sim \mathcal{GP}\left(0, C_\phi^{\text{NS}}(\cdot, \cdot)\right), \\
 \log \ell(\cdot) &\sim \mathcal{GP}\left(\mu_\ell, C_\varphi^{\text{S}}(\cdot, \cdot)\right), \\
 (\tau^2, \varphi, \sigma_\varepsilon^2, \mu_\ell) &\sim \pi(\tau^2)\pi(\varphi)\pi(\sigma_\varepsilon^2)\pi(\mu_\ell),
 \end{aligned} \tag{7}$$

where  $C_\phi^{\text{NS}}(\cdot, \cdot)$  denotes a non-stationary kernel,  $C_\varphi^{\text{S}}(\cdot, \cdot)$  is a stationary covariance function with parameters  $\varphi$ , and  $\mu_\ell$  the constant mean of  $\log \ell(\cdot)$ .

- Extremely flexible, 2-level improves predictive performance
- Fully Bayesian inference challenging:
  - Computationally expensive (2 nested GPs), latent processes and hyperparameters tend to be strongly coupled
  - Model is sensitive to the choice of hyperparameters.

# Sparse Hierarchical Non-stationary Models

- **Idea:** Use Gaussian Markov random fields - precision matrix equivalent to  $\mathbf{z} \sim \mathcal{N}(0, (Q_\phi^{\text{NS}})^{-1})$   $\mathbf{z} \sim \mathcal{N}(0, C_\phi^{\text{NS}})$
- How to create  $Q$ ?
  - Roininen et al. 2019 derive a SPDE formulation for non-stationary Matérn fields.
  - For  $d = 1$  and  $\nu = 2 - 1/2$ ,

$$(1 - \ell(\cdot)^2 \Delta) \mathbf{z} = \tau \sqrt{\ell(\cdot)} \mathbf{w}, \quad (8)$$

where  $\Delta$  is the Laplace operator,  $\mathbf{w}$  is white noise on  $\mathbb{R}$ ,  $\text{Var}(\mathbf{w}) = \Gamma(\nu + 1/2)(4\pi)^{1/2}/\Gamma(\nu)$ , and  $\ell(\cdot)$  is a spatially varying length-scale.

- A finite-dimensional approximation can be written as

$$L(\ell)\mathbf{z} = \mathbf{w},$$

where  $\mathbf{z} \in \mathbb{R}^n$  with  $n$  the discretisation size.  $L(\ell)$  is a sparse matrix depending on  $\ell_j := \ell(jh)$ , with  $h$  the discretisation step in a chosen finite difference approximation.

# The model

- GP regression model:  $\mathbf{y} = \mathbf{A}\mathbf{z} + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2 I_m)$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{z} \in \mathbb{R}^n$ .
- Hierarchical formulation

$$\begin{aligned}
 \mathbf{y} \mid \mathbf{z}, \sigma_\varepsilon^2 &\sim \mathcal{N}(\mathbf{A}\mathbf{z}, \sigma_\varepsilon^2 I_m), \\
 \mathbf{z} \mid \phi &\sim \mathcal{N}(0, Q_\phi^{-1}) \\
 \log \ell := \mathbf{u} \mid \varphi &\sim \mathcal{N}(\boldsymbol{\mu}_\ell, C_\varphi) \\
 (\tau^2, \sigma_\varepsilon^2, \varphi, \mu_\ell) &\sim \pi(\tau^2)\pi(\sigma_\varepsilon^2)\pi(\varphi)\pi(\mu_\ell)
 \end{aligned} \tag{9}$$

where  $\boldsymbol{\mu}_\ell$  is the  $n$ -dimensional constant mean vector.

- Key component:  $(C_\phi^{\text{NS}})^{-1} := Q_\phi = L(\phi)^\top L(\phi)$ , which depends on  $\mathbf{u}$  and  $\tau^2$ .
- $\varphi$  parameters of the covariance that describe properties of the length-scales.

# Hyperpriors

- **Stationary** assumption for spatially varying length-scale
- Explore two priors for  $\mathbf{u}$ :

## Squared Exponential:

- ▶ Strong prior smoothness assumptions on how the correlation of the non-stationary process changes with distance.
- ▶ Precision matrix is **dense** and depends on length-scale  $\lambda$  and magnitude  $\tau_\ell$ .

## AR(1):

- ▶ Ornstein-Uhlenbeck covariance
- ▶ Allows quick changes but is smoother than white noise.
- ▶ Precision is **sparse**  $Q_\varphi = L(\varphi)^\top L(\varphi)$ , where  $L(\varphi)$  is a banded matrix that depends on  $\lambda$  and  $\tau_\ell$ .
- To improve model identifiability, we fix  $\tau$ ,  $\mu_\ell$  and  $\tau_\ell$ .

# Inference for one-dimensional problems

Posterior of interest:

$$\pi(\mathbf{z}, \mathbf{u}, \lambda, \sigma_\varepsilon^2 \mid \mathbf{y}) \propto \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{z}, \sigma_\varepsilon^2 \mathbf{I}_m) \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_z, \mathbf{Q}_u^{-1}) \mathcal{N}(\mathbf{u} \mid \boldsymbol{\mu}_\ell, \mathbf{C}_\varphi) \pi(\lambda) \pi(\sigma_\varepsilon^2).$$

## • Metropolis-within-Gibbs (MWG)

- Length scale  $\mathbf{u}$  are updated individually.
- When proposing  $u_k^*$ , for  $k = 1, \dots, n$ , log-ratio of acceptance probability simplifies- ( $O(n)$  for SE and  $O(1)$  for AR).
- When proposing hyperparameter  $\varphi^*$ , we require:  $\log \left( \frac{\mathcal{N}(\mathbf{u} \mid \boldsymbol{\mu}_\ell, \mathbf{C}_\varphi^*)}{\mathcal{N}(\mathbf{u} \mid \boldsymbol{\mu}_\ell, \mathbf{C}_\varphi)} \right) - (O(n^3))$  for SE and  $O(n)$  for AR).
- Does not perform well for SE.

## • Whitened Elliptical Slice Sampling (w-ELL-SS)

$\mathbf{z} = \mathbf{L}(\mathbf{u})^{-1} \boldsymbol{\xi}$  with  $\boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{I}_n)$  and  $\mathbf{u} = \mathbf{R}_\varphi \boldsymbol{\zeta} + \boldsymbol{\mu}_\ell$  with  $\boldsymbol{\zeta} \sim \mathcal{N}(0, \mathbf{I}_n)$ .

$$\pi(\boldsymbol{\zeta}, \boldsymbol{\xi}, \lambda, \sigma_\varepsilon^2 \mid \mathbf{y}) \propto \mathcal{N}(\mathbf{y} \mid \mathbf{A} \mathbf{L}(\mathbf{R}_\varphi \boldsymbol{\zeta} + \boldsymbol{\mu}_\ell)^{-1} \boldsymbol{\xi}, \sigma_\varepsilon^2 \mathbf{I}_m) \mathcal{N}(\boldsymbol{\xi} \mid 0, \mathbf{I}_n) \mathcal{N}(\boldsymbol{\zeta} \mid 0, \mathbf{I}_n) \pi(\lambda) \pi(\sigma_\varepsilon^2).$$

- $\mathbf{u}$  updated jointly through  $\boldsymbol{\zeta}$ .
- Likelihood can be evaluated as a product of univariate Gaussians
- $\mathbf{z} = \mathbf{L}(\mathbf{u})^{-1} \boldsymbol{\xi}$  can be solved in  $O(n)$
- $\mathbf{u} = \mathbf{R}_\varphi \boldsymbol{\zeta} + \boldsymbol{\mu}_\ell$  - ( $O(n^2)$  for SE and  $O(n)$  for AR)
- Each iteration may require several likelihood evaluations.

## • Marginal Elliptical Slice Sampling (m-ELL-SS)

$$\pi(\zeta, \lambda, \sigma_\varepsilon^2 \mid \mathbf{y}) \propto \mathcal{N}(\mathbf{y} \mid 0, A Q_{\zeta, \varphi}^{-1} A^\top + \sigma_\varepsilon^2 I_n) \mathcal{N}(\zeta \mid 0, I_m) \pi(\lambda) \pi(\sigma_\varepsilon^2).$$

- $\mathbf{u}$  updated jointly through  $\zeta$ .
- $\mathbf{u} = R_\varphi \zeta + \boldsymbol{\mu}_\ell$  - ( $O(n^2)$  for SE and  $O(n)$  for AR)
- Marginal likelihood computation:

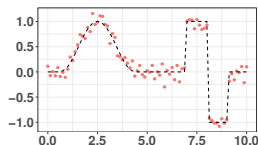
$$\log \pi(\mathbf{y} \mid \mathbf{u}, \lambda, \sigma_\varepsilon^2) = -\frac{m}{2} \log(2\pi) - \frac{1}{2} \log \det(\Psi) - \frac{1}{2} \mathbf{y}^\top \Psi^{-1} \mathbf{y}$$

where  $\Psi = A Q_{\mathbf{u}}^{-1} A^\top + \sigma_\varepsilon^2 I_m$ .

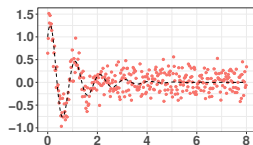
- Employ Woodbury identity for  $\Psi^{-1}$
- Quadratic term:  $\sigma_\varepsilon^{-2} \left( \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top A \left( L(\mathbf{u})^\top L(\mathbf{u}) + \sigma_\varepsilon^{-2} A^\top A \right)^{-1} A^\top \mathbf{y} \right)$
- Determinant computation is the dominant term ( $O(m^3)$  or  $O(nm)$ )
- Improved mixing

# Synthetic 1-d experiments

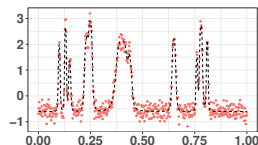
## 3 synthetic 1-dimensional examples



(a) Experiment 1

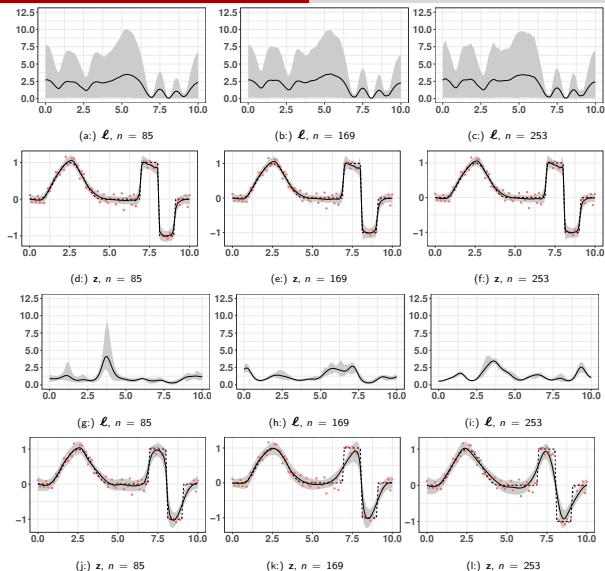


(b) Experiment 2

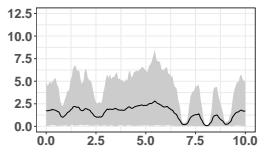
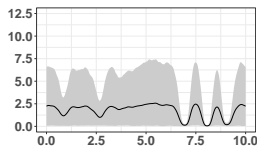
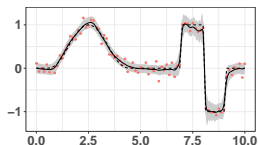
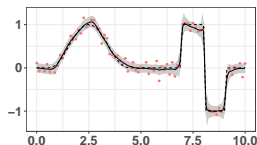


(c) Experiment 3

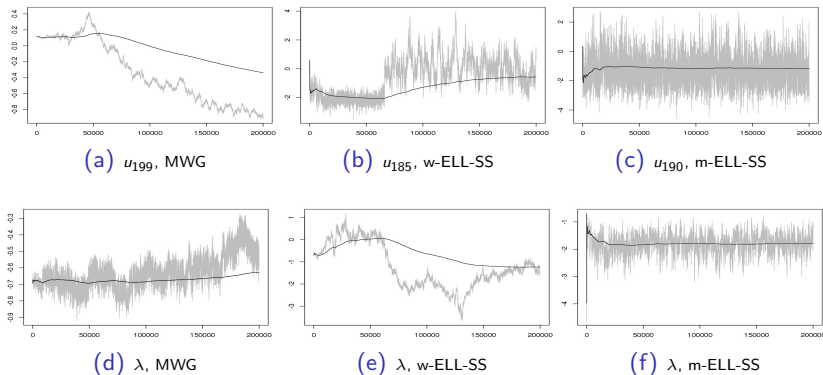
**Figure:** (a): 81 observations with domain  $[0, 10]$  and  $\sigma_{\varepsilon}^2 = 0.01$ . (b): 350 observations with domain  $[0, 8]$  and  $\sigma_{\varepsilon}^2 = 0.04$ . (c): 512 observations with domain  $[0, 1]$  and  $\sigma_{\varepsilon}^2 = 0.04$



**Figure:** MWG. (a)-(c): Estimated  $\ell$  process with 95% credible intervals for AR(1) on different grids. (d)-(f): Estimated  $z$  process with 95% credible intervals for AR(1) on different grids with observed data in red. (g)-(i): Estimated  $\ell$  process with 95% credible intervals for SE on different grids. (j)-(l): Estimated  $z$  process with 95% credible intervals for SE hyperprior on different grids with observed data in red.

(a) w-ELL-SS  $\ell$ ,  $n = 253$ (b) m-ELL-SS  $\ell$ ,  $n = 253$ (c) w-ELL-SS  $\ell$ ,  $n = 253$ (d) m-ELL-SS  $z$ ,  $n = 253$ 

**Figure:** Results for Experiment 1 at the highest resolution ( $n=253$ ) for SE hyperprior with (left column) w-ELL-SS algorithm and (right column) m-ELL-SS algorithm.



**Figure:** Example 1: Traceplots with cumulative averages of the chains for SE hyperprior with  $n = 253$ . (Top row:) element of  $\mathbf{u}$  with the lowest ESS. (Bottom row:) the hyperparameter.

OES = ESS/CPUtime

		MWG			w-ELL-SS			m-ELL-SS		
		$n = 85$	$n = 169$	$n = 253$	$n = 85$	$n = 169$	$n = 253$	$n = 85$	$n = 169$	$n = 253$
AR(1)	$\sigma_\varepsilon^2$	622.76	173.12	65.99	380.89	102.38	38.91	<b>661.20</b>	<b>257.81</b>	<b>116.35</b>
	$\ell_{min}$	<b>635.36</b>	114.02	41.05	30.90	8.99	2.94	287.16	<b>114.36</b>	<b>59.71</b>
	$z_{min}$	<b>203.80</b>	42.10	13.91	9.12	2.34	0.86	129.75	<b>52.16</b>	<b>22.30</b>
	$\lambda$	89.84	15.66	6.00	22.77	5.26	2.36	<b>111.80</b>	<b>45.54</b>	<b>21.53</b>
	MAE	0.041	0.051	0.054	0.041	0.051	0.054	0.041	0.051	0.053
	EC	0.988	0.975	0.971	0.988	0.975	0.975	0.988	0.975	0.975
SE	$\sigma_\varepsilon^2$	11.19	4.88	7.49	246.24	77.72	8.89	<b>856.15</b>	<b>253.91</b>	<b>125.97</b>
	$\ell_{min}$	1.22	0.73	0.64	21.69	10.22	2.79	<b>244.91</b>	<b>122.57</b>	<b>55.82</b>
	$z$	0.06	0.01	0.01	4.71	1.37	0.24	<b>76.80</b>	<b>24.11</b>	<b>9.87</b>
	$\lambda$	0.59	0.75	0.31	2.31	0.29	0.01	<b>16.59</b>	<b>4.15</b>	<b>2.21</b>
	MAE	0.078	0.100	0.133	0.040	0.050	0.054	0.039	0.049	0.052
	EC	0.889	0.826	0.763	0.988	0.975	0.971	0.988	0.975	0.979

**Table:** Experiment 1: OES with both hyperpriors under various discretisation schemes ( $n = 86, 169, 253$ ) and three different algorithms.  $\ell_{min}$  and  $z_{min}$  report OES for the minimum ESS across all dimensions. Highest values in boldface.

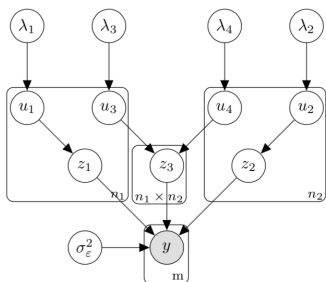
- AR(1) hypermodel adds further computational gains.
- MWG performs poorly for highly correlated hyperprior.
- MWG deteriorates efficiency as the number of observations or discretisation size increase.
- w-ELL-SS for weak likelihoods performs well regardless the hyperprior employed at the price of highly correlated chains.
- Marginal sampler converges to the stationary distribution faster.
- m-ELL-SS good compromise between computational complexity and efficiency of the chains.

# Extensions for two-dimensional problems

Employs additive Gaussian process models (AGP)

$$\mathbf{y} = A_1 \mathbf{z}_1 + A_2 \mathbf{z}_2 + A_3 \mathbf{z}_3 + \varepsilon,$$

- $A_1 \in \mathbb{R}^{m \times n_1}$ ,  $A_2 \in \mathbb{R}^{m \times n_2}$  and  $A_3 \in \mathbb{R}^{m \times (n_1 n_2)}$  known matrices.
- $z_1(\cdot)$  and  $z_2(\cdot)$  independent univariate non-stationary processes.
- $z_3(\cdot)$  is a bivariate, non-stationary, separable process -interaction term



Hierarchical model:

$$\mathbf{y} \mid \{\mathbf{z}_r\}_{r=1}^3, \sigma_\varepsilon^2 \sim \mathcal{N}(\mathbf{A}_1 \mathbf{z}_1 + \mathbf{A}_2 \mathbf{z}_2 + \mathbf{A}_3 \mathbf{z}_3, \sigma_\varepsilon^2 \mathbf{I}_m)$$

$$\mathbf{z}_r \mid \boldsymbol{\phi}_r \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\boldsymbol{\phi}_r}^{\text{NS}}), \quad r = 1, 2, 3$$

$$\mathbf{u}_s \mid \boldsymbol{\varphi}_s \sim \mathcal{N}(\boldsymbol{\mu}_{\ell_s}, \mathbf{C}_{\boldsymbol{\varphi}_s}^{\text{S}}), \quad s = 1, 2, 3, 4$$

$$(\sigma_\varepsilon^2, \boldsymbol{\varphi}) \sim \pi(\sigma_\varepsilon^2) \pi(\boldsymbol{\varphi}_1) \pi(\boldsymbol{\varphi}_2) \pi(\boldsymbol{\varphi}_3) \pi(\boldsymbol{\varphi}_4),$$

with  $\boldsymbol{\varphi} = (\boldsymbol{\varphi}_1, \dots, \boldsymbol{\varphi}_4)$ .

- AGP works based on one-dimensional kernels
- Posterior:

$$\begin{aligned} \pi(\{\mathbf{z}_r\}_{r=1}^3, \{\mathbf{u}_s, \lambda_s\}_{s=1}^4, \sigma_\varepsilon^2 \mid \mathbf{y}) &\propto \mathcal{N}(\mathbf{y} \mid \mathbf{A}_1 \mathbf{z}_1 + \mathbf{A}_2 \mathbf{z}_2 + \mathbf{A}_3 \mathbf{z}_3, \sigma_\varepsilon^2 \mathbf{I}_m) \mathcal{N}(\mathbf{z}_1 \mid \mathbf{0}, \mathbf{Q}_{\mathbf{u}_1}^{-1}) \\ &\quad \mathcal{N}(\mathbf{z}_2 \mid \mathbf{0}, \mathbf{Q}_{\mathbf{u}_2}^{-1}) \mathcal{N}(\mathbf{z}_3 \mid \mathbf{0}, \mathbf{Q}_{\mathbf{u}_{3,4}}^{-1}) \mathcal{N}(\mathbf{u}_1 \mid \boldsymbol{\mu}_{\ell_1}, \mathbf{C}_{\boldsymbol{\varphi}_1}) \cdots \mathcal{N}(\mathbf{u}_4 \mid \boldsymbol{\mu}_{\ell_4}, \mathbf{C}_{\boldsymbol{\varphi}_4}) \\ &\quad \pi(\lambda_1) \cdots \pi(\lambda_4) \pi(\sigma_\varepsilon^2), \end{aligned}$$

with  $\mathbf{Q}_{\mathbf{u}_{3,4}}^{-1} := \mathbf{Q}_{\mathbf{u}_3}^{-1} \otimes \mathbf{Q}_{\mathbf{u}_4}^{-1}$

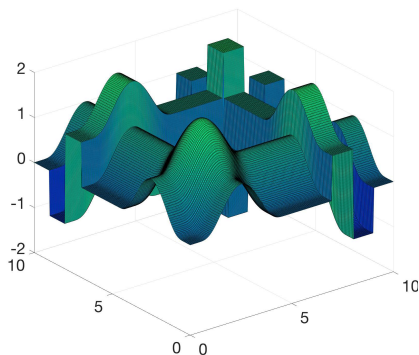
# Inference for two-dimensional problems

- Blocked Gibbs sampler, that updates the three blocks of parameters  $(\mathbf{z}_1, \mathbf{u}_1, \lambda_1)$ ;  $(\mathbf{z}_2, \mathbf{u}_2, \lambda_2)$ ; and  $(\mathbf{z}_3, \mathbf{u}_3, \mathbf{u}_4, \lambda_3, \lambda_4)$  from their full conditional distributions.
- Block marginal elliptical slice sampler (Block-m-ELL-SS)**
  - To sample  $(\mathbf{z}_1, \mathbf{u}_1, \lambda_1)$ , the full conditional is factorised:

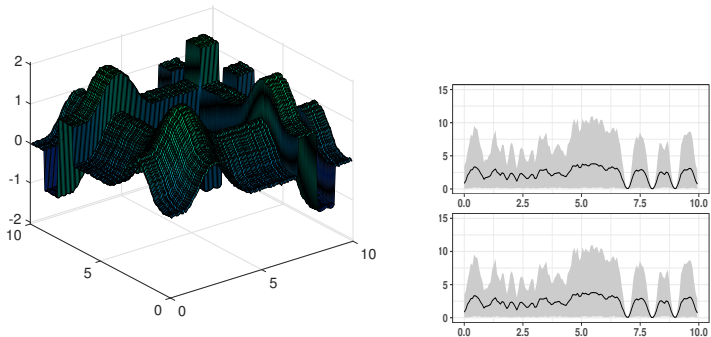
$$\pi(\mathbf{z}_1, \zeta_1, \lambda_1 \mid \mathbf{y}, \sigma_\varepsilon^2, \mathbf{z}_2, \mathbf{z}_3) = \pi(\zeta_1, \lambda_1 \mid \mathbf{y}, \sigma_\varepsilon^2, \mathbf{z}_2, \mathbf{z}_3) \pi(\mathbf{z}_1 \mid \zeta_1, \lambda_1, \mathbf{y}, \sigma_\varepsilon^2, \mathbf{z}_2, \mathbf{z}_3),$$

- Interaction term: use eigendecompositions and matrix-vector multiplications for Kronecker matrices!

# Synthetic 2-d experiment



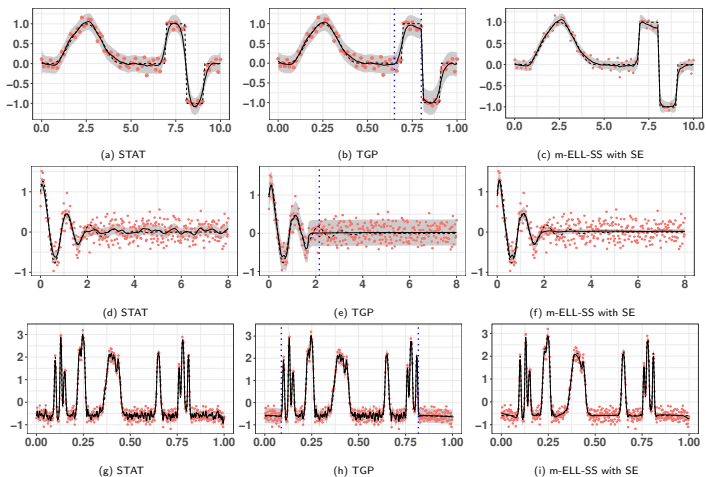
**Figure:** 2-dimensional synthetic data.  $m = 20,449$  noisy observations in an expanded grid of  $n_1 = n_2 = 143$  equally spaced points in  $[0, 10]$ , employing  $z(x_1, x_2) = z(x_1) + z(x_2)$ .



**Figure:** Posterior mean surface and one-dimensional length-scale processes with 95% credible intervals.

- Capture smooth areas and edges.
- 2-level AGP correctly learns the varying correlation along the surface.
- 99.26 minutes for 50,000 iterations

# Comparative evaluation



**Figure:** Each row shows one of the simulated experiments. Red dots depict observed data, dotted lines show the true signal, solid lines show the posterior mean, and grey areas depict 95% credible intervals. (a)(d)(g)(j): Stationary GP (b)(e)(h)(k): TGP, with blue dotted lines depicting MAP cut-off points. (c)(f)(i)(l): 2-level GP with m-ELL-SS algorithm and the hyperprior with lowest MAE.

# Thank you