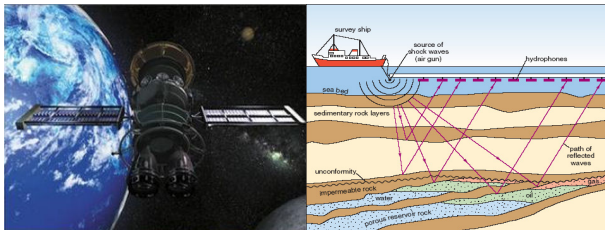


# Estimation and prediction in spatial models with block composite likelihoods

Jo Eidsvik,  
Department of Mathematical Sciences, NTNU, Norway

With Ben Shaby (Colorado State Univ), Brian Reich (North Carolina State Univ), Matt Wheeler and Jarad Niemi (Iowa State Univ)

# Large spatial (spatio-temporal) datasets



# Spatial Gaussian model

Model:  $Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$ .

- ▶  $\mathbf{s}$  is a spatial location on a domain of interest.
- ▶  $Y(\mathbf{s})$  is the response variable.
- ▶  $\mathbf{X}(\mathbf{s})$  are covariates.
- ▶  $\boldsymbol{\beta}$  are regression parameters.
- ▶  $w(\mathbf{s})$  is a spatially smooth process, covariance parameters;  $\phi, \sigma^2$ .
- ▶  $\epsilon(\mathbf{s})$  is independent noise  $\epsilon(\mathbf{s}) \sim N(0, \tau^2)$
- ▶  $\boldsymbol{\theta} = (\sigma^2, \phi, \tau^2)$ .

## Gaussian likelihood

Model:  $Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})\boldsymbol{\beta} + w(\mathbf{s}) + \epsilon(\mathbf{s})$ .

Sampling locations  $\mathbf{s}_1, \dots, \mathbf{s}_n$ .

Data:  $\mathbf{Y} = (Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n))'$ ,  $\mathbf{X} = (\mathbf{X}(\mathbf{s}_1), \dots, \mathbf{X}(\mathbf{s}_n))'$ .

Log likelihood:

$$l(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$$

Challenges for **parameter estimation** and **prediction**:

1. Build and store  $n \times n$  matrix  $\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \boldsymbol{\Sigma} = \mathbf{C} + \tau^2 \mathbf{I}_n$
2. Compute  $\log |\boldsymbol{\Sigma}|$  and  $\boldsymbol{\Sigma}^{-1}$  or  $(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})$

*Some approaches: tapering, fixed rank Kriging, Markov random fields, predictive process, approximate likelihood methods.*

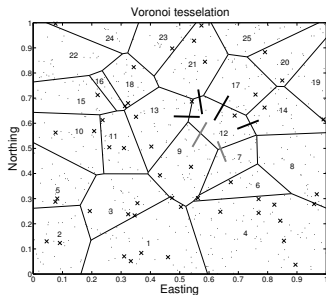
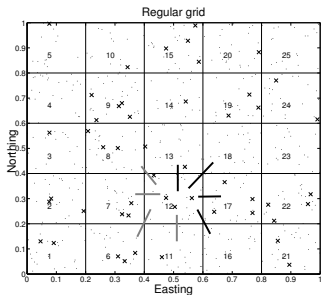
## Use model for composites

- ▶ Lindsay (1988), Curriere and Lele (1999), Varin (2008).
- ▶ Scales with dimension. Useful for parallel computing.
- ▶ Use joints for subsets of data, not full joint.  $M$  blocks.  
 $n = \sum_{k=1}^M n_k$ .
- ▶ Block-pair variables:  $\mathbf{Y}_{kl} = (\mathbf{Y}'_k, \mathbf{Y}'_l)'$ ,  $\mathbf{X}_{kl} = (\mathbf{X}'_k, \mathbf{X}'_l)'$ , size  $(n_k + n_l) \times (n_k + n_l)$  covariance  $\boldsymbol{\Sigma}_{kl} = \boldsymbol{\Sigma}_{kl}(\boldsymbol{\theta})$ .
- ▶ Log composite likelihood:

$$l_{CL}(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = \sum_{k=1}^{M-1} \sum_{l>k} \left\{ -\frac{1}{2} \log |\boldsymbol{\Sigma}_{kl}| - \frac{1}{2} (\mathbf{Y}_{kl} - \mathbf{X}_{kl}\boldsymbol{\beta})' \boldsymbol{\Sigma}_{kl}^{-1} (\mathbf{Y}_{kl} - \mathbf{X}_{kl}\boldsymbol{\beta}) \right\}$$

## Split and conquer - examples of blocking

$$l_{CL}(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}) = \sum_{k=1}^{M-1} \sum_{I \in \mathcal{N}_k^{\rightarrow}} \left\{ -\frac{1}{2} \log |\boldsymbol{\Sigma}_{kl}| - \frac{1}{2} (\mathbf{Y}_{kl} - \mathbf{X}_{kl} \boldsymbol{\beta})' \boldsymbol{\Sigma}_{kl}^{-1} (\mathbf{Y}_{kl} - \mathbf{X}_{kl} \boldsymbol{\beta}) \right\}$$



## Estimation : Maximum composite likelihood

$$(\hat{\theta}, \hat{\beta}) = \operatorname{argmax}_{\theta, \beta} \{l_{CL}(\mathbf{Y}; \beta, \theta)\}.$$

- ▶ Iterative algorithm  $p = 0, 1, \dots$
- ▶ Estimating equations:  $\frac{dl_{CL}(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\beta} = 0, \frac{dl_{CL}(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\theta} = 0.$
- ▶  $\hat{\beta}_p = \mathbf{A}^{-1} \mathbf{b},$   
 $\mathbf{A} = \mathbf{A}(\hat{\theta}_p; X_{kl} \text{pairs}), \mathbf{b} = \mathbf{b}(\hat{\theta}_p; X_{kl} \text{pairs}, Y_{kl} \text{pairs}).$
- ▶ Fisher-scoring:  $\hat{\theta}_{p+1} = \hat{\theta}_p - E \left( \frac{d^2 l_{CL}(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\theta^2} \right)^{-1} \frac{dl_{CL}(\mathbf{Y}; \hat{\beta}_p, \hat{\theta}_p)}{d\theta}$

## Analytical derivatives

Computed for each block-pair.

$$\begin{aligned} \frac{d \log |\boldsymbol{\Sigma}_{kl}|}{d\theta_r} &= \text{trace}(\boldsymbol{\Sigma}_{kl}^{-1} \frac{d\boldsymbol{\Sigma}_{kl}}{d\theta_r}) \\ \frac{d \mathbf{Z}'_{kl} \boldsymbol{\Sigma}_{kl}^{-1} \mathbf{Z}_{kl}}{d\theta_r} &= -\mathbf{Z}'_{kl} \boldsymbol{\Sigma}_{kl}^{-1} \frac{d\boldsymbol{\Sigma}_{kl}}{d\theta_r} \boldsymbol{\Sigma}_{kl}^{-1} \mathbf{Z}_{kl} \\ \mathbf{Z}_{kl} &= (\mathbf{Y}_{kl} - \mathbf{X}_{kl} \boldsymbol{\beta}) \end{aligned}$$



## Asymptotic properties: Godambe sandwich

$\hat{\theta} \rightarrow N(\theta, G^{-1})$ . (Varin, 2008).

$$G = G(\hat{\theta}) = H(\hat{\theta})J^{-1}(\hat{\theta})H(\hat{\theta}),$$

$$H(\hat{\theta}) = -E\left(\frac{d^2 l_{CL}}{d\theta^2}(\hat{\theta})\right), \quad J(\hat{\theta}) = \text{Var}\left(\frac{dl_{CL}}{d\theta}(\hat{\theta})\right)$$



# Synthetic data

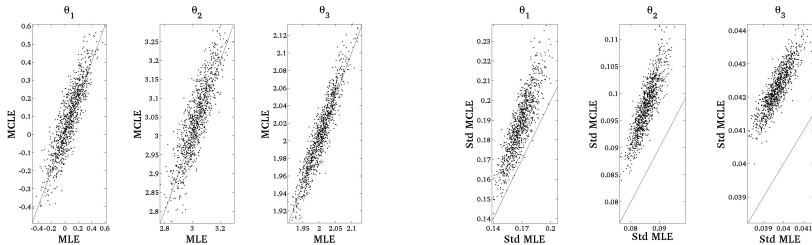
Unit square domain.  $n = 2000$ . 2000 replicates of data.

Covariance  $\Sigma(h) = \tau^2 I(h = 0) + \sigma^2(1 + \phi h) \exp(-\phi h)$ ,  $h = |\mathbf{s}_i - \mathbf{s}_j|$ .

Effective range 1/4th of domain.

# Likelihood vs composite with $M = 100$ blocks.

Parameter estimates: log precision, log range, log nugget precision.



## Prediction from composite likelihood

$\mathbf{Y}_k^a = (\mathbf{Y}'_{k0}, \mathbf{Y}'_k)'$ .  $\mathbf{Y}_{k0}$  unobserved variables of size  $n_{k0}$  in block  $k$ .  
 $\mathbf{X}_k^a = (\mathbf{X}'_{k0}, \mathbf{X}'_k)'$  associated covariates.

$$l_{CL}^k(\mathbf{Y}_{k0}) = - \sum_{l \in N_k} \frac{1}{2} [(\mathbf{Y}_k^a, \mathbf{Y}_l^a)' - (\mathbf{X}_k^a, \mathbf{X}_l^a)' \boldsymbol{\beta}]' \boldsymbol{\Sigma}_{0kl}^{-1} [(\mathbf{Y}_k^a, \mathbf{Y}_l^a)' - (\mathbf{X}_k^a, \mathbf{X}_l^a)' \boldsymbol{\beta}]$$

$\boldsymbol{\Sigma}_{0kl}$  is  $(n_{k0} + n_k + n_l) \times (n_{k0} + n_k + n_l)$  covariance matrix.

$$\hat{\mathbf{Y}}_{k0} = \mathbf{X}_{k0} \hat{\boldsymbol{\beta}} + \mathbf{A}_0^{-1} \mathbf{b}_0$$

$$\mathbf{A}_0 = \mathbf{A}_0(\hat{\boldsymbol{\theta}}), \mathbf{b}_0 = \mathbf{b}_0(\hat{\boldsymbol{\theta}}; X_{kl} \text{ pairs}, Y_{kl} \text{ pairs})$$

## Prediction properties

$$\hat{Y}_{k0} - Y_{k0} \sim N(0, G_0^{-1})$$

$$G_0 = H_0 J_0^{-1} H_0,$$

$$H_0 = -E \left( \frac{d^2 l_{CL}^k}{dY_{k0}^2} \right), \quad J_0 = \text{Var} \left( \frac{dl_{CL}^k}{dY_{k0}} \right)$$



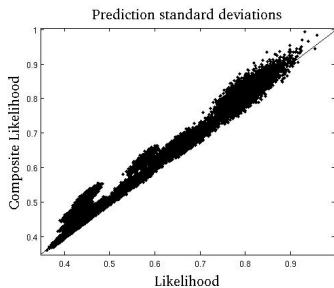
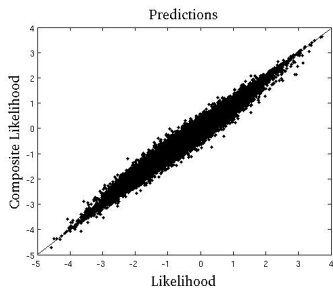
# Synthetic data

Unit square domain.  $n = 2000$ . 2000 replicates of data. 200 prediction sites.

Covariance  $\Sigma(h) = \tau^2 I(h = 0) + \sigma^2(1 + \phi h) \exp(-\phi h)$ ,  $h = |\mathbf{s}_i - \mathbf{s}_j|$ .

Effective range 1/4th of domain.

# Likelihood vs composite with $M = 100$ blocks.

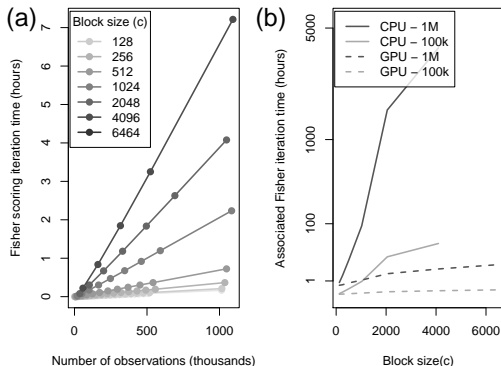


Results for block size  $3^2$  and  $10^2$ 

|                                  | L      | CL, $3 \times 3$ | CL, $10 \times 10$ |
|----------------------------------|--------|------------------|--------------------|
| MSE $\hat{\beta}_1$              | 0.05   | 0.05             | 0.06               |
| MSE $\hat{\beta}_2$              | 0.17   | 0.21             | 0.23               |
| MSE $\hat{\theta}_1$             | 0.014  | 0.018            | 0.018              |
| MSE $\hat{\theta}_2$             | 0.0036 | 0.0043           | 0.0054             |
| MSE $\hat{\theta}_3$             | 0.0007 | 0.0008           | 0.0008             |
| Coverage (0.95) $\hat{\beta}_1$  | 0.96   | 0.95             | 0.96               |
| Coverage (0.95) $\hat{\beta}_2$  | 0.93   | 0.92             | 0.92               |
| Coverage (0.95) $\hat{\theta}_1$ | 0.93   | 0.92             | 0.91               |
| Coverage (0.95) $\hat{\theta}_2$ | 0.94   | 0.94             | 0.91               |
| Coverage (0.95) $\hat{\theta}_3$ | 0.95   | 0.95             | 0.94               |
| MSPE                             | 193    | 195              | 204                |
| Mean pred. coverage (0.95)       | 0.95   | 0.95             | 0.95               |
| Computing time (sec)             | 76     | 39               | 12                 |

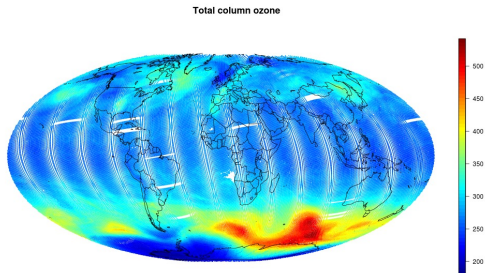


# Parallelization



Matrix factorization on GPU. For-loop harder - require splitting the data on resource units.

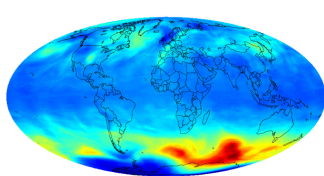
## Satellite dataset: $n \sim 200.000$



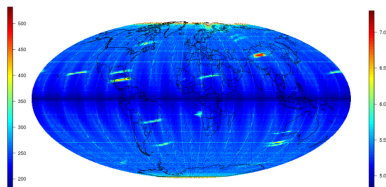
$$Y(\mathbf{s}) = \beta_0 + w(\mathbf{s}) + \epsilon(\mathbf{s}).$$

Estimate parameters. Use plug in for prediction.

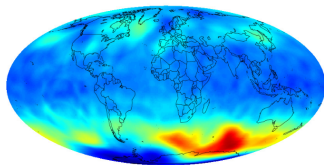
# Satellite dataset: Fixed rank kriging and CL



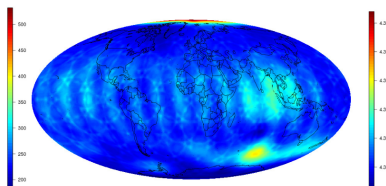
(a) CL predictions



(b) CL standard errors



(c) FRK predictions



(d) FRK standard errors

## Satellite dataset: Prediction results for a hold-out set

|                 | FRK, 4 | FRK, 3 | CL, reg 15 | CL, fast 30 |
|-----------------|--------|--------|------------|-------------|
| MSPE            | 44.3   | 88.1   | 25.7       | 26.0        |
| Pred cov (0.95) | 0.81   | 0.71   | 0.96       | 0.96        |
| Timing (min)    | 6      | 2      | 40         | 4           |

## Local composite spatial modeling

- ▶ Split parameter estimation and spatial prediction
- ▶ Just predicting in (pairs of) blocks... overlapping blocks...

Unifying these elements instead:

- ▶ Form a *process* on a local graph.
- ▶ Ease hierarchical Bayesian modeling.

## Predictions on local domains, building a process

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \dots p(x_n|x_{n-1}, \dots, x_1) \\ &= p(x_1) \prod_{i=2}^n p(x_i|x_1, \dots, x_{i-1}) \end{aligned}$$

Assuming a neighborhood. Some of the conditioning can be ignored.

$$p(x_1, x_2, \dots, x_n) = p(x_1) \prod_{i=2}^n p(x_i|x_j, j \in N_i^{\rightarrow})$$

Where  $N_i^{\rightarrow}$  represents the neighborhood of  $i$  for the indexes below  $i$ . (Similar to a Cholesky factorization.)

## Datta et al. (2016) formed a neighborhood graph

$$p(\mathbf{w}_S) = p(\mathbf{w}(s_1)) p(\mathbf{w}(s_2) | \mathbf{w}(s_1)) \\ \dots p(\mathbf{w}(s_k) | \mathbf{w}(s_{k-1}), \dots, \mathbf{w}(s_1)) ,$$

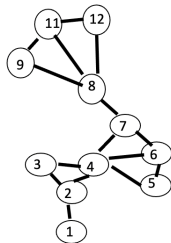
$$\tilde{p}(\mathbf{w}_S) = \prod_{i=1}^k p(\mathbf{w}(s_i) | \mathbf{w}_{N(s_i)})$$

$$\tilde{p}(\mathbf{w}_S) = \prod_{i=1}^k N(\mathbf{w}(s_i) | \mathbf{B}_{s_i} \mathbf{w}_{N(s_i)}, \mathbf{F}_{s_i})$$

$$\mathbf{B}_{s_i} = \mathbf{C}_{s_i, N(s_i)} \mathbf{C}_{N(s_i)}^{-1}$$

$$\mathbf{F}_{s_i} = \mathbf{C}(s_i, s_i) - \mathbf{C}_{s_i, N(s_i)} \mathbf{C}_{N(s_i)}^{-1} \mathbf{C}_{N(s_i), s_i}$$

Neighborhood structure is defined by edges between nodes. And by order of variables.



# Nearest neighbor Gaussian process

$$\tilde{p}(\mathbf{w}_{\mathcal{U}} | \mathbf{w}_{\mathcal{S}}) = \prod_{i=1}^r p(\mathbf{w}(\mathbf{u}_i) | \mathbf{w}_{N(\mathbf{u}_i)})$$

$$\tilde{p}(\mathbf{w}_{\mathcal{U}} | \mathbf{w}_{\mathcal{S}}) = \prod_{i=1}^r N(\mathbf{w}(\mathbf{u}_i) | \mathbf{B}_{\mathbf{u}_i} \mathbf{w}_{N(\mathbf{u}_i)}, \mathbf{F}_{\mathbf{u}_i}) = N(\mathbf{B}_{\mathcal{U}} \mathbf{w}_{\mathcal{S}}, \mathbf{F}_{\mathcal{U}})$$

$$\tilde{p}(\mathbf{w}_{\mathcal{V}}) = \int \tilde{p}(\mathbf{w}_{\mathcal{U}} | \mathbf{w}_{\mathcal{S}}) \tilde{p}(\mathbf{w}_{\mathcal{S}}) \prod_{\{\mathbf{s}_i \in \mathcal{S} \setminus \mathcal{V}\}} d(\mathbf{w}(\mathbf{s}_i))$$

where  $\mathcal{U} = \mathcal{V} \setminus \mathcal{S}$ .

Node set:

$$\mathcal{S} = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_k\}$$

Outside node set:

$$\mathcal{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$$

$$\tilde{\mathbf{C}}(\mathbf{v}_1, \mathbf{v}_2; \boldsymbol{\theta})$$

$$= \begin{cases} \tilde{\mathbf{C}}_{\mathbf{s}_i, \mathbf{s}_j} & \text{if } \mathbf{v}_1 = \mathbf{s}_i \text{ and } \mathbf{v}_2 = \mathbf{s}_j \text{ are both in } \mathcal{S}, \\ \mathbf{B}_{\mathbf{v}_1} \tilde{\mathbf{C}}_{N(\mathbf{v}_1), \mathbf{s}_j} & \text{if } \mathbf{v}_1 \notin \mathcal{S} \\ \mathbf{B}_{\mathbf{v}_1} \tilde{\mathbf{C}}_{N(\mathbf{v}_1), N(\mathbf{v}_2)} \mathbf{B}'_{\mathbf{v}_2} + \delta_{\mathbf{v}_1} & \frac{1}{\prod_{i=1}^k \sqrt{\det(\mathbf{F}_{\mathbf{s}_i})}} \exp\left(-\frac{1}{2} \sum_{i=1}^k (\mathbf{w}(\mathbf{s}_i) - \mathbf{B}_{\mathbf{s}_i} \mathbf{w}_{N(\mathbf{s}_i)})' \mathbf{F}_{\mathbf{s}_i}^{-1} (\mathbf{w}(\mathbf{s}_i) - \mathbf{B}_{\mathbf{s}_i} \mathbf{w}_{N(\mathbf{s}_i)})\right) \\ & \text{are not in } \mathcal{S} \end{cases}$$

Sparse:

$$\tilde{\mathbf{C}}_{\mathcal{S}}^{-1}$$



## Nearest neighbor Gaussian process

$$\frac{1}{\prod_{i=1}^k \sqrt{\det(\mathbf{F}_{s_i})}} \exp\left(-\frac{1}{2} \sum_{i=1}^k (\mathbf{w}(s_i) - \mathbf{B}_{s_i} \mathbf{w}_{N(s_i)})' \mathbf{F}_{s_i}^{-1} (\mathbf{w}(s_i) - \mathbf{B}_{s_i} \mathbf{w}_{N(s_i)})\right)$$

$$(\tilde{\mathbf{C}}_{\mathcal{S}})^{-1} = \mathbf{B}'_{\mathcal{S}} \mathbf{F}_{\mathcal{S}}^{-1} \mathbf{B}_{\mathcal{S}} \quad \det((\mathbf{B}'_{\mathcal{S}} \mathbf{F}_{\mathcal{S}}^{-1} \mathbf{B}_{\mathcal{S}})^{-1}) = \prod \det(\mathbf{F}_{s_i})$$

Hierarchical model evaluation:  $p(\boldsymbol{\theta}) \times \prod_{j=1}^l IG(\tau_j^2 \mid a_{\tau_j}, b_{\tau_j}) \times N(\boldsymbol{\beta} \mid \boldsymbol{\mu}_{\beta}, \mathbf{V}_{\beta}) \times N(\mathbf{w}_{\mathcal{U}} \mid \mathbf{B}_{\mathcal{U}} \mathbf{w}_{\mathcal{S}}, \mathbf{F}_{\mathcal{U}})$

$$\times N(\mathbf{w}_{\mathcal{S}} \mid \mathbf{0}, \tilde{\mathbf{C}}_{\mathcal{S}}) \times \prod_{i=1}^n N(y(t_i) \mid \mathbf{X}(t_i)' \boldsymbol{\beta} + \mathbf{Z}(t_i)' \mathbf{w}(t_i), \mathbf{D}),$$

|            | True | Full<br>Gaussian Process | Order by<br>$y$ -coordinates | NNGP ( $\mathcal{S} = \mathcal{T}$ )<br>Order by<br>$x$ -coordinates |
|------------|------|--------------------------|------------------------------|--|
| $\sigma^2$ | 1    | 0.640 (0.414, 1.297)     | 0.712 (0.449, 1.530)         | 0.757 (0.479, 1.501)   |
| $\tau^2$   | 0.1  | 0.107 (0.098, 0.117)     | 0.106 (0.097, 0.114)         | 0.107 (0.099, 0.117)   |
| $\phi$     | 6    | 8.257 (4.056, 13.408)    | 8.294 (3.564, 12.884)        | 7.130 (3.405, 11.273)  |