# Scalable Gaussian Processes on Physically Constrained Domains

**Spatial Modelling with Intrinsic Geometry** 

#### Bora Jin, Amy H. Herring, David Dunson Duke University





• Groundwater



- Groundwater
- Neighborhood/wildlife disconnected by barriers



- Groundwater
- Neighborhood/wildlife disconnected by barriers
- Coastlines



https://coastwatch.pfeg.noaa.gov/erddap/griddap/erdMH1chlamday.graph

## **BORA-GP**

Background Goals Methods Analysis Conclusions



### Background

- Many applications with unique geometry including boundaries or barriers.
- Traditional spatial Gaussian process (GP) models ignore the unique geometry of the domain.
  - Inappropriate smoothing over physical barriers
  - Likely produce sub-optimal results
- Traditional GP has high computational cost  $\mathcal{O}(n^3)$ .
  - Prohibitive already if n > 50,000.
  - Recent & active development of scalable GP models whose computational cost grows linearly with n.







- Aim to construct a new <u>scalable</u> GP model that incorporates <u>constrained</u> <u>domains</u>.
- In particular, we want the new GP to
  - Produce physically sensible kriging,
  - Induce physically sensible covariance behavior,
  - Mimic geodesic distance-based covariance without loss of scalability.
    - Typical geodesic distance estimation is also an expensive operation,  $O(n^3)$ .



)11ke

Bayesian hierarchical spatial regression

$$y = X\beta + w + \epsilon, \quad \epsilon \sim N(0, \tau^2 I_n)$$
$$w \sim GP(0, C(\cdot, \cdot | \theta))$$

• Computational bottleneck in computing the joint density of  $w_{obs} = (w_1, \dots, w_n)^T$   $p(w_1, \dots, w_n | \theta) = N(0, C_{\theta}) \propto |C_{\theta}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}w_{obs}^T C_{\theta}^{-1} w_{obs}\right)$ Costly!  $O(n^3)$  Costly!  $O(n^3)$ 

$$\circ \quad w_i = w(s_i)$$

 $\circ$  *n* is the number of spatial locations.

- Vecchia (1988)
  - Approximate the joint density with lower-dimensional conditional densities to ease computational burden.

$$p(w_1, \cdots, w_n) = \prod_{i=1}^n p(w_i | w_1, \cdots, w_{i-1})$$

- Lots of redundant information in the conditioning set  $\{w_1, \dots, w_{i-1}\}$  for large *i*'s.
- $p(w_i|w_1, \dots, w_{i-1}) \approx p(w_i|w_{im})$  where  $w_{im}$  is a conditioning set of <u>at most</u>  $m \ (\ll n)$  elements from  $\{w_1, \dots, w_{i-1}\}$ .

$$p(w_1, \cdots, w_n) \approx \tilde{p}(w_1, \cdots, w_n) = \prod_{i=1}^n p(w_i | w_{im})$$



- Vecchia (1988) to Nearest Neighbor Gaussian Process (NNGP) (2016)
  - One possible choice of *w*<sub>im</sub>



- $w_{im} = \{ \mathbf{X}_{i}, \mathbf{X}_{2}, w_{3}, \cdots, \mathbf{X}_{i-2}, w_{i-1} \}$

conditional independence

zeros in a precision matrix – sparsity



- Vecchia (1988) to Nearest Neighbor Gaussian Process (NNGP) (2016)
  - One possible choice of w<sub>im</sub>



- NNGP extends finite-dimensional sparsity to sparsity-inducing spatial processes.
- Constructs a scalable and valid spatial process based on nearest neighbors.
- Empirically shows that the above-mentioned w<sub>im</sub> performs better than other alternatives in a very wide range of scenarios.



• Visualization of sparsity using directed acyclic graphs (DAGs)





Duke

• What if we know that our measurements lie in a constrained domain?



Barrier Overlap-Removal Acyclic Directed Graph Gaussian Process (BORA-



- BORA-GP
  - **1.** Specify a multivariate normal distribution over a fixed finite set  $S = \{s_1, \dots, s_k\} \subset D$

$$w_{S} = (w(s_{1}), \cdots, w(s_{k}))^{T} \sim N(0, \tilde{C}_{S}) = \prod_{i=1}^{k} N(w(s_{i}); H_{s_{i}}w_{N(s_{i})}, R_{s_{i}})$$

• 
$$H_{s_i} = C_{s_i, N(s_i)} C_{N(s_i)}^{-1}, R_{s_i} = C_{s_i} - C_{s_i, N(s_i)} C_{N(s_i)}^{-1} C_{N(s_i), s_i}$$

- *C*: base covariance function
- N(s<sub>i</sub>): set of spatial locations in S whose straight line to s<sub>i</sub> does not overlap barriers (physically sensible neighbors of s<sub>i</sub> of size min(i, m))

• BORA-GP

2. Extend it to the whole domain  $\ensuremath{\mathcal{D}}$ 

 $\forall s \in \mathcal{D} \setminus S, \qquad w(s) \sim N(H_s w_{N(s)}, R_s)$ 

•  $N(s) \subset S$  physically sensible neighbors of size m

■ 1 & 2 define a valid scalable process  $w \sim \text{BORA-GP}(0, \tilde{C}(\cdot, \cdot | \theta))$ 

•  $\tilde{C}$  is non-stationary.

#### **Analysis**

Duke

• Univariate Bayesian spatial regression

```
y = w + \epsilon, w \sim \text{BORA-GP}(0, \tilde{C}(\cdot, \cdot | \theta))
```

y: log chlorophyll-a level (mg/m3) of March 2021







		BORA-GP	Full GP	NNGP
	RMSPE*	0.165	0.172	0.178
	95% CI coverage	98.57%	98.42%	98.42%
Duke	Mean 95% CI width	1.110	1.113	1.172

\*RMSPE: Root Mean Squared Prediction Error











Analysis		40%		
Delawa	re Bay	BORA-GP	Full GP	NNGP
RMS	SPE	0.257	0.346	0.356
95% CI c	overage	92.98%	85.96%	87.72%
Mean 95%	6 CI width	1.171	1.150	1.210

	Chesapeake Bay	BORA-GP	Full GP	NNGP
	RMSPE	0.214	0.296	0.300
	95% CI coverage	98.18%	95.45%	95.45%
	Mean 95% CI width	1.157	1.131	1.192
D	Juke			











#### Estimated stationary covariance from full GP



#### Estimated non-stationary covariance from BORA-GP





#### Estimated stationary covariance from full GP



Estimated non-stationary covariance from BORA-GP



#### **Conclusions**

• Does BORA-GP achieve all the goals?

✓ Scalable

 $m \ll n$  neighbors only (Sparsity-inducing DAG)

- ✓ Account for the domain
  - Physically sensible kriging
    Using physically sensible neighbors
  - Physically sensible covariance behavior

The implied non-stationary covariance does not go beyond barriers

 Mimicking geodesic distance-based covariance without loss of scalability

> No need to estimate geodesic distance Covariance resembles water flow



#### **References & Acknowledgement**

- Datta, A., Banerjee, S., Finley, A.O., and Gelfand, A.E. (2016). Hierarchical nearestneighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, *111*(514), 800-812.
- Vecchia, A.V. (1988). Estimation and model identification for continuous spatial processes. Journal of the Royal Statistical Society: Series B (Methodological), 50(2), 297-312.
- BORA-GP soon to be available in arXiv.
- Thanks, coauthors: Amy H. Herring and David Dunson!
- Thanks for the support from NIEHS R01ES027498 and R01ES028804!