# Conjugate Bayesian Modeling and Inference In High-dimensional Spatial Statistics: Conquering New Challenges

Sudipto Banerjee

SIAM, February 24th, 2022
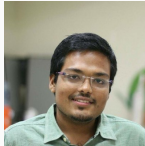
# Collaborators


Abhirup Datta (JHU)


Debangan Dey (JHU)


Barbara Engelhardt (Princeton)


Andrew Finley (MSU)

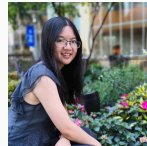
Rajarshi Guhaniyogi (TAMU)
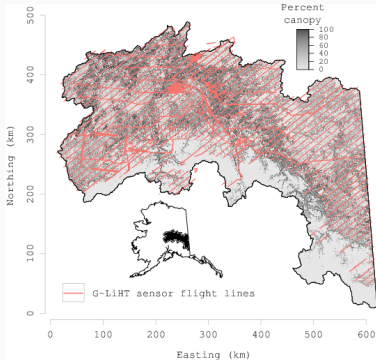

Andrew Jones (Princeton)


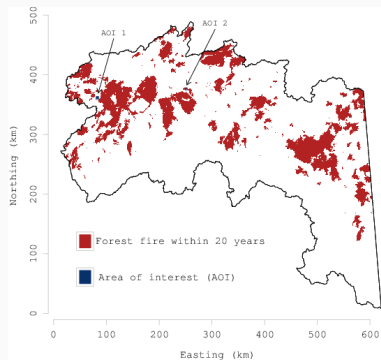Didong Li (Princeton/UCLA)


Michele Peruzzi (Duke)


Lu Zhang (Columbia)

1

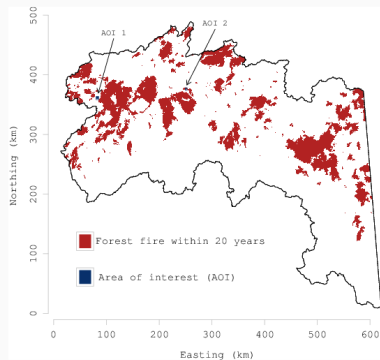# Example: Alaska Tanana Valley Forest Height Dataset (FD-CMAB, *JCGS*, 2019)



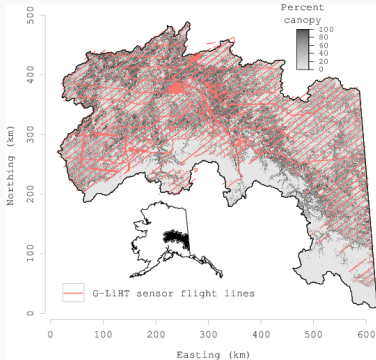Forest height and tree cover

Forest fire history

- Forest height (red lines) data from LiDAR at $10 \times 10^6$ locations
- Knowledge of forest height is important for biomass assessment, carbon management etc

Forest height and tree cover



Forest fire history

- Goal: High-resolution domain-wide prediction maps of forest height
- Covariates: Domain-wide tree cover (grey) and forest fire history (red patches) in the last 20 years

## Analyzing the data

Models used:

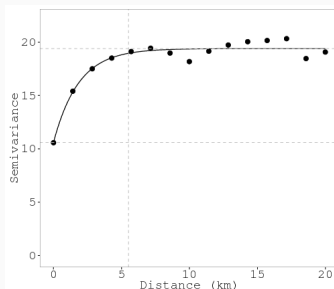- Non-spatial regression: $y_{FH} = \beta_0 + \beta_{tree}x_{tree} + \beta_{fire}x_{fire} + \epsilon$



**Figure:** Variogram (defined as $\text{var}\{Z(\ell + h) - Z(\ell)\}$) of the residuals from non-spatial regression indicates strong spatial pattern

## Bayesian regression for BIG DATA

- Conjugate Bayesian hierarchical linear model:

$$y_i \mid \beta, \sigma^2 \stackrel{ind}{\sim} N(x_i^\top \beta, \sigma^2) \,, \; i = 1, 2, \dots, n \,;$$
$$\beta \mid \sigma^2 \sim N(\mu_\beta, \sigma^2 V_\beta) \,; \quad \sigma^2 \sim IG(a, b) \,.$$

- Exact Bayesian inference:

$$\sigma^2 \mid y \sim IG(a^*, b^*) \quad \beta \mid \sigma^2, y \sim N(Mm, \sigma^2 M) \,, \quad \text{where}$$
$$m = V_\beta^{-1} \mu_\beta + X^\top y \,, \quad M^{-1} = V_\beta^{-1} + X^\top X \,,$$
$$a^* = a + n/2 \,, \quad b^* = \mu_\beta^\top V_\beta^{-1} \mu_\beta + y^\top y - m^\top M^{-1} m \,.$$

- What if the data cannot be stored/loaded into available workspace?

- HADOOP: Map-Reduce framework (Divide & Conquer) with cloud computing.

## Bayesian regression on HADOOP

- Partition data as $D_k = \{y_k, X_k\}$, $k = 1, 2, \ldots, K$, where each $y_k$ is $n_k \times 1$, $X_k$ is $n_k \times p$ and $N = \sum_{k=1}^{K} n_k$.

- Sequential ("streaming") updates:

$$p(\beta, \sigma^2 \,|\, D_1, \ldots, D_{k+1}) \propto p(\beta, \sigma^2 \,|\, D_1, \ldots, D_k) \times p(D_{k+1} \,|\, \beta, \sigma^2)$$

- Parallel architecture: compute simultaneously...

$$m_k = V_\beta^{-1} + X_k^\top y_k \text{ and } M_k^{-1} = V_\beta^{-1} + X_k^\top X_k \,;$$
$$m = \sum_{k=1}^{K} (m_k - (1 - 1/K)\mu_\beta) \text{ and } M^{-1} = \sum_{k=1}^{K} (M_k^{-1} - (1 - 1/K)V_\beta^{-1}) \,.$$

- Depends (crucially) on independence across subsets; not suitable for spatial random fields.

## Geostatistical models for parallel architectures

- $y_{FH}(\ell) = \beta_0 + \beta_{tree} x_{tree}(\ell) + \beta_{fire} x_{fire}(\ell) + w(\ell) + \epsilon(\ell)$

- $w(\ell) \sim GP(0, C(\cdot, \cdot \,|\, \sigma^2, \phi))$

- $y_{FH} \sim N(X\beta, K_\theta)$ where $K_\theta$ is the spatial covariance matrix:
$$K_\theta = C_{(\sigma, \phi)} + \tau^2 I \,, \quad \text{where} \ \ \theta = \{\sigma, \phi, \tau\}$$
where $C_{(\sigma^2, \phi)}$ is the GP covariance matrix derived from $C(\cdot, \cdot \,|\, \sigma^2, \phi)$.

- Massive data: divide and conquer?

- Bayesian model averaging? Predictive stacking? Exchangability?

- Meta-Kriging (GB, *Technometrics* 2018): find convex combination of subset-posteriors closest to the full posterior.

- Analyze "compressed data": Compressive sensing; Data sketching.

## Bayesian Hierarchical Models

[data | process, parameters] $\times$ [process | parameters] $\times$ [parameters]

- Construct a joint model...

$$p(\theta, \tau, \beta) \times p(w \,|\, \theta) \times p(\tilde{w} \,|\, w, \theta) \times p(y \,|\, \beta, w, \tau) \times p(\tilde{Y} \,|\, \tilde{w}, \theta, \tau)$$

- Posterior inference for parameters and the process:

$$p(\theta, \tau, \beta, w, \tilde{w}, \tilde{Y} \,|\, y) \propto p(\theta, \tau, \beta, w \,|\, y) \times p(\tilde{w} \,|\, w, \theta) \times p(\tilde{Y} \,|\, \tilde{w}, \theta, \tau)$$

- Multivariate example with $Y = \{Y_j(s_i)\}$ for $j = 1, 2, \ldots, m$ variables:

$$MN(Y \,|\, XB, K_{\theta, \tau}, \Sigma) \times MN(B \,|\, \mu_B, V_B, \Sigma) \times IW(\Sigma \,|\, a, S) \times p(\theta, \tau) \,.$$

# Constructing GPs from Graphs

spNNGP

meshed

## Sparse precision matrices (e.g., Vecchia's approximation; NNGP)

$$N(w \mid 0, K_\theta) \approx N(w \mid 0, \tilde{K}_\theta) \; ; \; \tilde{K}_\theta^{-1} = (I - A)^\top D^{-1}(I - A)$$



n=1000, m=10,
Sparsity: 97%

$I - A$             $D^{-1}$             $\tilde{K}_\theta^{-1}$

- $\det(\tilde{K}_\theta^{-1}) = \prod_{i=1}^n D_{ii}^{-1}$, $\tilde{K}_\theta^{-1}$ is sparse with $O(nm^2)$ entries

- Computing $A$ and $D$

```
for(i in 1:(n-1)) {
  Pa = N[i+1] # neighbors of i+1
  a[i+1,Pa] = solve(K[Pa,Pa], K[i+1, Pa])
  d[i+1,i+1] = K[i+1,i+1] - dot(K[i+1, Pa],a[i+1,Pa])
}
```

- We need to solve $n-1$ linear systems of size at most $m \times m$ in parallel.
- Quadratic form:

```
qf(u,v,A,D) = u[1] * v[1] / D[1,1]
for(i in 2:n) {
    qf(u,v,A,D) = qf(u,v,A,D) + (u[i] - dot(A[i,N(i)], u[N(i)]))
        *(v[i] - dot(A[i,N(i)], v[N(i)]))/D[i,i]
    }
```

- Determinant: $\det(\tilde{K}_\theta) = \prod_{i=1}^{n} \texttt{d[i,i]}$

# Alaska Tanana Valley data (Finley et al., *JCGS*, 2019)

|  | Conjugate NNGP | Collapsed NNGP | Response NNGP |
|---|---|---|---|
| $\beta_0$ | 2.51 | 2.41 (2.35, 2.47) | 2.37 (2.31, 2.42) |
| $\beta_{TC}$ | 0.02 | 0.02 (0.02, 0.02) | 0.02 (0.02, 0.02) |
| $\beta_{Fire}$ | 0.35 | 0.39 (0.34, 0.43) | 0.43 (0.39, 0.48) |
| $\sigma^2$ | 23.21 | 18.67 (18.50, 18.81) | 17.29 (17.13, 17.41) |
| $\tau^2$ | 1.21 | 1.56 (1.55, 1.56) | 1.55 (1.54, 1.55) |
| $\phi$ | 3.83 | 3.73 (3.70, 3.77) | 4.15 (4.13, 4.19) |
| CRPS | 0.84 | 0.86 | 0.86 |
| RMSPE | 1.71 | 1.73 | 1.72 |
| time (hrs.) | 0.002 | 319 | 38 |

**Table:** Parameter estimates and model comparison metrics for the Tanana valley dataset

- Conjugate model produces estimates and model comparison numbers very similar to the MCMC based NNGP models
- For $5 \times 10^6$ locations, conjugate model takes 7 seconds

- Complex dependencies are often modeled using CI graphs (Cox & Wermuth, 1996)



- But what about complex dependencies among processes (Each node is $\{w_i(s) : s \in \mathbb{R}^d\}$)? And a very large number of nodes, too?

- What does this mean : $w_i(\cdot) \perp w_j(\cdot) \,|\, \{w_{-(ij)}(\cdot)\}$? Dalhaus (2000):

$$\mathrm{cov}\,(z_i(s), z_j(s)) = 0 \quad \text{for all } s, s' \in \mathcal{D},$$

  where $z_i(s) = w_i(s) - \mathbb{E}[w_i(s) \,|\, \sigma(\{w_k(\cdot) : k \in \mathcal{V} \setminus \{i,j\}\})]$.

- Graphical GP (GGP): $\{w_i(\cdot) : i = 1, 2, \ldots, q\} \sim GGP_{\mathcal{G}}$ if $w_i(\cdot) \perp w_j(\cdot) \,|\, \{w_{-(i,j)}(\cdot)\}$ according to CI graph $\mathcal{G}$.

- Given a CI graph $\mathcal{G}$ and any cross-covariance function, there exists a unique (and optimal) $GGP_{\mathcal{G}}$ whose cross-covariance agrees with the given cross-covariance for all adjacent pairs in the graph.

- Constructing a GGP from a given $C(\mathcal{S})$ over a fixed finite set $\mathcal{S}$:
  1. Form an extended graph over $\mathcal{V} \times \mathcal{S}$ using *strong product* adjacency rules (to allow "stitching" across random fields);
  2. Use Dempster's (1972) covariance selection to specify $w(\mathcal{S}) \sim N(0, M(\mathcal{S}))$,
     2.1 $M_{ii}(\mathcal{S}) = C_{ii}(\mathcal{S})$ for each node $i$;
     2.2 Zeroes in $M(\mathcal{S})^{-1}$ correspond to CI relations in $\mathcal{G}$;
     2.3 $M_{ij}(\mathcal{S}) = C_{ij}(\mathcal{S})$ for all adjacent pairs in $\mathcal{G}$.
  3. Extend from finite set $\mathcal{S}$ to entire domain using predictive process with $\mathcal{S}$ as knots (Banerjee et al., 2008).

- DDB, 2021 also implement Bayesian inference for an unknown $\mathcal{G}$ using RJMCMC for (embeddable) decomposable graphs (Green & Thomas, 2013).

# Parallelizable Stitching of Gaussian Processes



**Figure:** Stitching Gaussian Processes. Left: Realizations of 4 univariate GPs. Right: Realization of a multivariate (4-dimensional) GGP created by stitching together the 4 univariate GPs from the left figure using the strong product graph over the 4 variables and 3 locations.

**Figure:** Chromatic sampling for GGP with a gem graph between 5 variables: Left: Gem graph and coloring used for chromatic sampling of the variable-specific parameters. Right: Coloring of the corresponding edge graph $\mathcal{G}_E(\mathcal{G}_V)$ used for chromatic sampling of the cross-covariance parameters. In chromatic sampling, we can use this coloring to sample nodes belonging to same color in parallel bringing down the complexity by significant amount.

Prediction performance for full analysis

Estimates of time-specific cross-correlations

# Burgeoning literature on DAG-based spatial models...

- Vecchia, A.V. (1988). Estimation and model identification for continuous spatial processes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 50, 297–312. DOI: https://doi.org/10.1111/j.2517-6161.1988.tb01729.x

- Stein, M.L., Chi, Z. and Welty, L.J. (2004), Approximating likelihoods for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66, 275–296. DOI: https://doi.org/10.1046/j.1369-7412.2003.05512.x

- Datta, A., Banerjee, S., Finley, A. O., and Gelfand, A. E. (2016). Hierarchical Nearest-Neighbor Gaussian Process models for large geostatistical datasets. *Journal of the American Statistical Association*, 111, 800–812. DOI: https://doi.org/10.1080/01621459.2015.1044091.

- Datta, A., Banerjee, S., Finley, A. O., Hamm, N. A. S., and Schaap, M. (2016b). Non-separable dynamic Nearest-Neighbour Gaussian Process models for large spatio-temporal data with an application to particulate matter analysis. *Annals of Applied Statistics*, 10, 1286–1316. DOI: https://doi.org/10.1214/16-AOAS931

- Zhang, L., Datta, A. and Banerjee, S. (2018). Practical Bayesian modelling and inference for massive spatial datasets on modest computing environments. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, **12**, 197–209. DOI: https://doi.org/10.1002/sam.11413

- Taylor-Rodriguez, D., Finley, A.O., Datta, A., Babcock, C., Andersen, H.E., Cook, B.C., Morton, D.C. and **Banerjee, S.** (2019). Spatial factor models for high-dimensional and large spatial data: An application in forest variable mapping. *Statistica Sinica*, 29, 1155–1180. DOI: https://doi.org/10.5705/ss.202018.0005.

- Katzfuss, M. and Guinness, J. (2021). A general framework for Vecchia approximations of gaussian processes. *Statistical Science*, 36, 124–141. DOI: https://doi.org/10.1214/19-STS755

- Peruzzi, M., Banerjee, S. and Finley, A.O. (in press). Highly scalable Bayesian geostatistical modeling via meshed Gaussian processes on partitioned domains. *Journal of the American Statistical Association*, DOI: https://doi.org/10.1080/01621459.2020.1833889

Thank You!