

# Bayesian High-Dimensional Regression with Tensors and Distributed Computation with Space-Time Data

**Rajarshi Guhaniyogi, Ph.D**

Department of Statistics, Texas A & M University  
**Partially Supported by**



March 13, 2022

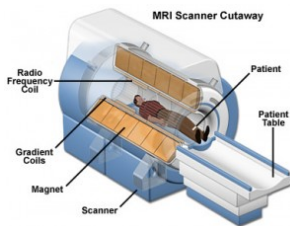


# Outline

- 1 Bayesian tensor regression
- 2 Bayesian tensor response regression
- 3 Bayesian symmetric tensor response regression
- 4 Distributed computation with space-time data

# Why Tensor Regression? Application in Primary Progressive Aphasia

Primary Progressive Aphasia is manifested in terms of language loss and indicates early stage of Alzheimers.



(a) MRI scan



(b) Atlas



(c) GM

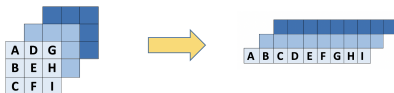
**Tensor predictor:** Structural MRI for 142 patients of language loss.

**scalar predictors:** gender, age.

**Response:** Language score representing degree of language loss.

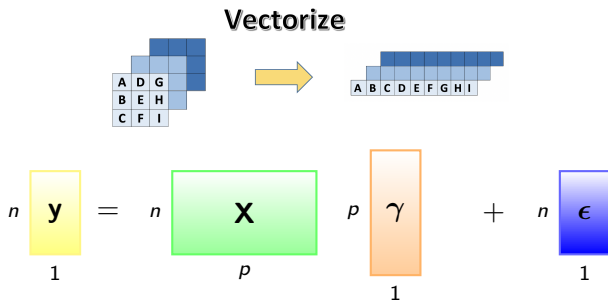
# Penalized Optimization: Unsatisfactory Predictive Performance

Vectorize



$$\begin{matrix} n \\ \mathbf{y} \\ 1 \end{matrix} = \begin{matrix} n \\ \mathbf{X} \\ p \end{matrix} \begin{matrix} p \\ \boldsymbol{\gamma} \\ 1 \end{matrix} + \begin{matrix} n \\ \boldsymbol{\epsilon} \\ 1 \end{matrix}$$

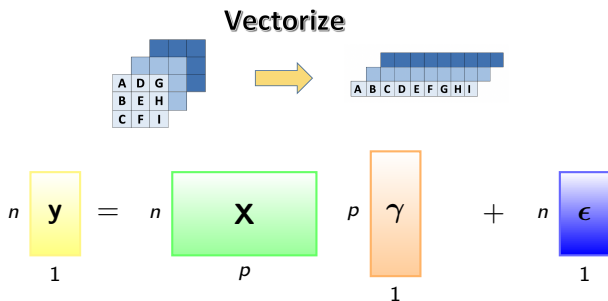
# Penalized Optimization: Unsatisfactory Predictive Performance



$\psi(\cdot)$  = convex penalty function,  $\zeta$  = tuning parameter

$\arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\gamma}\|^2 + \zeta \sum_{j=1}^p \psi(\gamma_j) \rightarrow$  **Penalized Opt.**

# Penalized Optimization: Unsatisfactory Predictive Performance

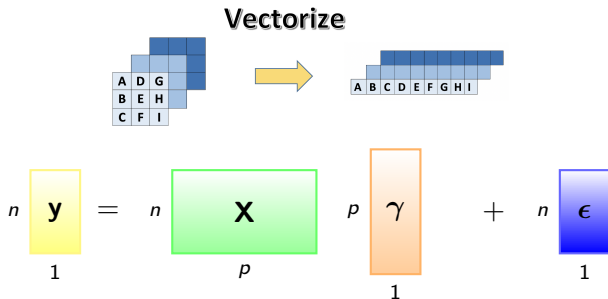


$\psi(\cdot)$  = convex penalty function,  $\zeta$  = tuning parameter

$\arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\gamma}\|^2 + \zeta \sum_{j=1}^p \psi(\gamma_j) \rightarrow$  **Penalized Opt.**

- LASSO (Tibshirani, 1996), Elastic Net (Zhou et al., 2005), tons of other variants.

# Penalized Optimization: Unsatisfactory Predictive Performance



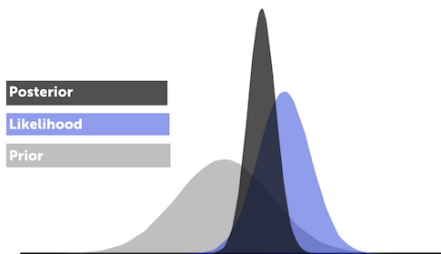
$\psi(\cdot)$  = convex penalty function,  $\zeta$  = tuning parameter

$\arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\gamma}\|^2 + \zeta \sum_{j=1}^p \psi(\gamma_j) \rightarrow$  **Penalized Opt.**

- LASSO (Tibshirani, 1996), Elastic Net (Zhou et al., 2005), tons of other variants.
- Unsatisfactory predictive uncertainty.

# Bayesian Inference

- Start with a prior distribution of  $\gamma$ .
- “Combine” data likelihood and prior distribution to obtain posterior distribution of  $\gamma$ .



- **point estimation** → mean of the posterior, **uncertainty** → 95% credible interval from the posterior.
- Markov Chain Monte Carlo (MCMC) and its variants exist to empirically estimate the posterior distribution of  $\gamma$ .



# Bayesian High Dim. Reg.: Unsuitable in High Dimension

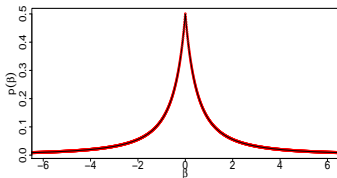
- Bayesians choose sparsity-favoring priors on  $\gamma \in \mathbb{R}^P$  which will set components of  $\gamma$  to be 0.

# Bayesian High Dim. Reg.: Unsuitable in High Dimension

- Bayesians choose sparsity-favoring priors on  $\gamma \in \mathbb{R}^P$  which will set components of  $\gamma$  to be 0.

## Spike & Slab Prior (Computationally Inefficient)

$$\gamma_j \sim \pi\delta_0 + (1 - \pi)g, \quad g \text{ is a cont. density.}$$



## Bayesian Shrinkage Prior (Statistically Inefficient)

$$\gamma_j \sim N(0, \zeta_j\tau), \quad \zeta_j \sim f_1, \quad \tau \sim f_2.$$

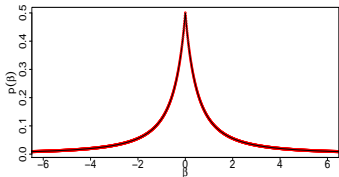
Marginally,  $\gamma_j$  has a heavy tailed density

# Bayesian High Dim. Reg.: Unsuitable in High Dimension

- Bayesians choose sparsity-favoring priors on  $\gamma \in \mathbb{R}^P$  which will set components of  $\gamma$  to be 0.

## Spike & Slab Prior (Computationally Inefficient)

$$\gamma_j \sim \pi \delta_0 + (1 - \pi)g, \quad g \text{ is a cont. density.}$$



## Bayesian Shrinkage Prior (Statistically Inefficient)

$$\gamma_j \sim N(0, \zeta_j \tau), \quad \zeta_j \sim f_1, \quad \tau \sim f_2.$$

Marginally,  $\gamma_j$  has a heavy tailed density

- Important shrinkage priors, Bayesian Lasso (Park et al., 2008; Hans, 2009), Horseshoe (Carvalho et al., 2009), Generalized Double Pareto (Armagan et al., 2013).

- Bayesians choose sparsity-favoring priors on  $\gamma \in \mathbb{R}^P$  which will set components of  $\gamma$  to be 0.

## Spike & Slab Prior (Computationally Inefficient)

$$\gamma_j \sim \pi \delta_0 + (1 - \pi)g, \quad g \text{ is a cont. density.}$$

## Serious Drawbacks of Penalization and Shrinkage

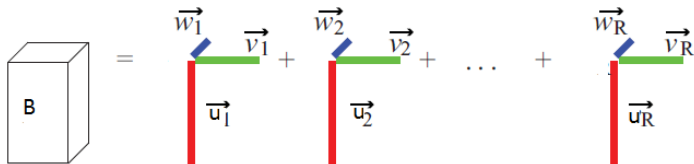
- $p = p_1 \times p_2 \times p_3$ , each  $p_i = 64$  typically, implies massive dimensional regression with close to half a million predictors  $\Rightarrow$  Infeasibility
- Misses out on wealth of information that the tensor valued images carry.
- Important shrinkage priors, Bayesian Lasso (Park et al., 2008; Hans, 2009), Horseshoe (Carvalho et al., 2009), Generalized Double Pareto (Armagan et al., 2013).

# Tensor Regression Model with PARAFAC Decomposition

## Data Model

$$y = \langle \mathbf{X}, \mathbf{B} \rangle + \mathbf{z}'\boldsymbol{\gamma} + \epsilon, \epsilon \sim N(0, \sigma^2)$$

rank-R PARAFAC decomposition of  $\mathbf{B}$  for dimension reduction



For  $D > 3$ , need a better notation  $\Rightarrow \mathbf{B} = \sum_{r=1}^R \boldsymbol{\beta}_1^{(r)} \circ \dots \circ \boldsymbol{\beta}_D^{(r)}$   
 $\boldsymbol{\beta}_j^{(r)} \in \mathbb{R}^{p_j}$ ,  $\circ$  denotes *outer product* between vectors.

## Data Model

$$y = \langle \mathbf{X}, \mathbf{B} \rangle + \mathbf{z}'\boldsymbol{\gamma} + \epsilon, \epsilon \sim N(0, \sigma^2)$$

## rank- $R$ PARAFAC decomposition of $\mathbf{B}$ for dimension reduction

### Advantages

- Number of parameters needed to model is  $R \sum_{j=1}^D p_j$  as opposed to  $\prod_{j=1}^D p_j \Rightarrow$  Dimension Reduction.
- Exploits neighborhood structure of  $\mathbf{X} \Rightarrow$  potentially better inference.

For  $D > 3$ , need a better notation  $\Rightarrow \mathbf{B} = \sum_{r=1}^R \boldsymbol{\beta}_1^{(r)} \circ \dots \circ \boldsymbol{\beta}_D^{(r)}$   
 $\boldsymbol{\beta}_j^{(r)} \in \mathbb{R}^{p_j}$ ,  $\circ$  denotes *outer product* between vectors.

# Need for a Multiway Shrinkage Prior on $B$

$$B = \begin{matrix} \vec{w}_1 & \vec{v}_1 \\ \hline \vec{u}_1 \end{matrix} + \begin{matrix} \vec{w}_2 & \vec{v}_2 \\ \hline \vec{u}_2 \end{matrix} + \dots + \begin{matrix} \vec{w}_R & \vec{v}_R \\ \hline \vec{u}_R \end{matrix}$$

Exchangable shrinkage across  $r=1, \dots, R$

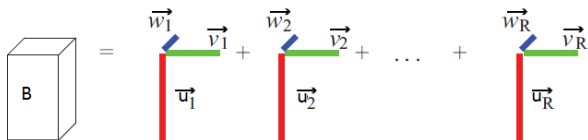
$$B = \begin{matrix} \vec{w}_1 & \vec{v}_1 \\ \hline \vec{u}_1 \end{matrix} + \begin{matrix} \vec{w}_2 & \vec{v}_2 \\ \hline \vec{u}_2 \end{matrix} + \dots + \begin{matrix} \vec{w}_R & \vec{v}_R \\ \hline \vec{u}_R \end{matrix}$$

Shrinkage within every rank

- Protects from overfitting due to a higher rank ( $R$ ) than needed.

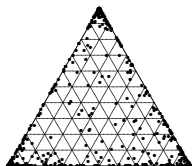
- Estimate tensor margins with an approximate sparsity.

# Multiway Shrinkage Prior for $B$ (G. et al. 2017, JMLR)

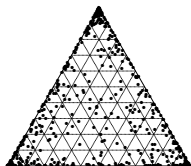


Exchangable shrinkage across  $r=1, \dots, R$

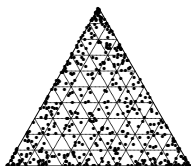
$\beta_j^{(r)} \sim N(\mathbf{0}, \text{diag}(w_{jr,1}, \dots, w_{jr,p_j}) \tau \phi_r)$ ,  $\phi_r$ 's rank specific parameters.  
Shrinkage across ranks:  $(\phi_1, \dots, \phi_R) \sim \text{Dirichlet}(\alpha, \dots, \alpha)$ ,  $\alpha > 0$ .



(a)  $\alpha = 0.2$



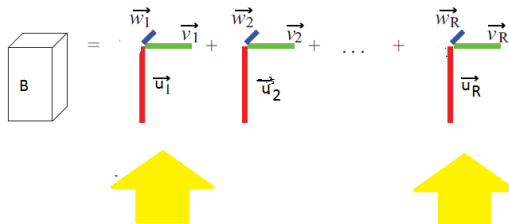
(b)  $\alpha = 0.3$



(c)  $\alpha = 0.5$



# Multiway Dirichlet Generalized Double Pareto Prior (M-DGDP)

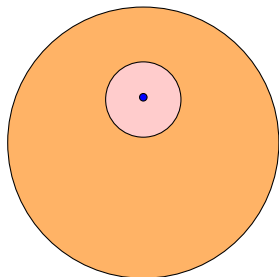


Shrinkage within every rank

$$w_{jr,k} \sim \text{Exp}(\lambda_{jr}^2/2), \quad \lambda_{jr} \sim \text{Ga}(a_\lambda, b_\lambda), \quad \tau \sim \text{IG}(a_\tau, b_\tau)$$

Integrating out  $\mathbf{W}_{jr}$

$\beta_{j,k}^{(r)} \mid \phi_r, \tau$  marginally follows GDP shrinkage prior.

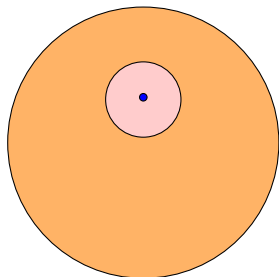


● True Model

$$(f(y|\mathbf{B}_n^0) = \mathcal{N}(\langle \mathbf{X}, \mathbf{B}_n^0 \rangle, \sigma^2))$$

Class of tensor reg. models fitted to the data

KL metric ball of radius  $\epsilon$  around the truth



● True Model  
 $(f(y|\mathbf{B}_n^0) = \mathcal{N}(\langle \mathbf{X}, \mathbf{B}_n^0 \rangle, \sigma^2))$

Class of tensor reg. models fitted to the data

KL metric ball of radius  $\epsilon$  around the truth

$$\mathcal{B}_n = \{ \mathbf{B}_n : \frac{1}{n} \sum_{i=1}^n \text{KL}(f(y_i|\mathbf{B}_n^0), f(y_i|\mathbf{B}_n)) < \epsilon \} \Rightarrow \textit{Neighborhood}$$

## Posterior Consistency

$$\Pi_n(\mathcal{B}_n^c) \rightarrow 0 \quad \text{under } \mathbf{B}_n^0 \quad \text{a.s. as } n \rightarrow \infty. \quad (1)$$

$\Pi_n$  posterior distribution given  $y_1, \dots, y_n$ .

## Theorem

The posterior is consistent under the following assumptions.

## Theorem

The posterior is consistent under the following assumptions.

- 1**  $\mathbf{B}_n^0 = \sum_{r=1}^{R_0} \beta_{1,n}^{0(r)} \circ \cdots \circ \beta_{D,n}^{0(r)}$  follows rank- $R_0$  decomposition,  $R > R_0$ . (Structure on the true coefficients)

## Theorem

The posterior is consistent under the following assumptions.

- 1**  $\mathbf{B}_n^0 = \sum_{r=1}^{R_0} \beta_{1,n}^{0(r)} \circ \cdots \circ \beta_{D,n}^{0(r)}$  follows rank- $R_0$  decomposition,  $R > R_0$ . (Structure on the true coefficients)
- 2**  $\sup_{l=1, \dots, p_{j,n}} |\beta_{j,n,l}^{0(r)}| < \infty$ , for all  $j = 1, \dots, D$ ;  $r = 1, \dots, R_0$ . (Structure on the true coefficients)

## Theorem

The posterior is consistent under the following assumptions.

- 1**  $\mathbf{B}_n^0 = \sum_{r=1}^{R_0} \beta_{1,n}^{0(r)} \circ \cdots \circ \beta_{D,n}^{0(r)}$  follows rank- $R_0$  decomposition,  $R > R_0$ . (Structure on the true coefficients)
- 2**  $\sup_{l=1, \dots, p_{j,n}} |\beta_{j,n,l}^{0(r)}| < \infty$ , for all  $j = 1, \dots, D$ ;  $r = 1, \dots, R_0$ . (Structure on the true coefficients)
- 3**  $\sum_{j=1}^D p_{j,n} \log(p_{j,n}) = o(n)$ . (Dimensions of tensor margins)

## Theorem

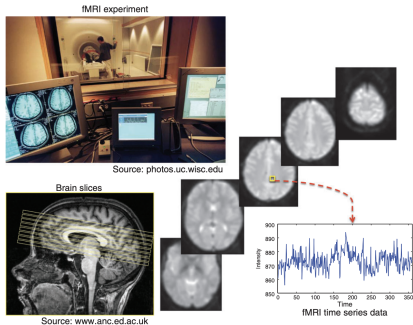
The posterior is consistent under the following assumptions.

- 1**  $\mathbf{B}_n^0 = \sum_{r=1}^{R_0} \beta_{1,n}^{0(r)} \circ \cdots \circ \beta_{D,n}^{0(r)}$  follows rank- $R_0$  decomposition,  $R > R_0$ . (Structure on the true coefficients)
- 2**  $\sup_{l=1, \dots, p_{j,n}} |\beta_{j,n,l}^{0(r)}| < \infty$ , for all  $j = 1, \dots, D$ ;  $r = 1, \dots, R_0$ . (Structure on the true coefficients)
- 3**  $\sum_{j=1}^D p_{j,n} \log(p_{j,n}) = o(n)$ . (Dimensions of tensor margins)

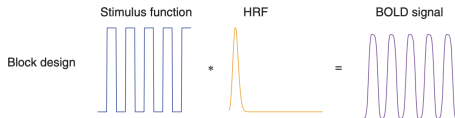
More involved convergence results can be found in G. (2017), JMVA.



# Motivation: Brain Activation Study



- $\mathbf{Y}_t$  is  $64 \times 64 \times 64$  dimensional tensor response at time  $t = 1, \dots, T$ .



- Identify brain voxels activated by an external stimulus.

# Tensor Response Regression Model

## Data Model

$$\mathbf{Y}_t = \mathbf{B}_1 x_{1t} + \cdots + \mathbf{B}_m x_{mt} + \mathbf{E}_t$$

- $\mathbf{Y}_t$  is a  $p_1 \times \cdots \times p_D$  dimensional tensor response,  $x_{1,t}, \dots, x_{m,t}$  are  $m$  predictors.
- $\mathbf{B}_1, \dots, \mathbf{B}_m$  are  $p_1 \times \cdots \times p_D$  dim. tensor coefficients.
- $\text{vec}(\mathbf{E}_t) \sim$  Stationary AR(1) process with the lag parameter  $\phi$ .

# Tensor Response Regression Model

## Data Model

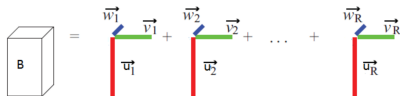
$$\mathbf{Y}_t = \mathbf{B}_1 x_{1t} + \cdots + \mathbf{B}_m x_{mt} + \mathbf{E}_t$$

- $\mathbf{Y}_t$  is a  $p_1 \times \cdots \times p_D$  dimensional tensor response,  $x_{1,t}, \dots, x_{m,t}$  are  $m$  predictors.
- $\mathbf{B}_1, \dots, \mathbf{B}_m$  are  $p_1 \times \cdots \times p_D$  dim. tensor coefficients.
- $\text{vec}(\mathbf{E}_t) \sim$  Stationary AR(1) process with the lag parameter  $\phi$ .

$$\begin{matrix} p_1 \\ \boxed{\mathbf{Y}_t} \\ p_2 \end{matrix} = \sum_{j=1}^m \begin{matrix} p_1 \\ \boxed{\mathbf{B}_j} \\ p_2 \end{matrix} \begin{matrix} 1 \\ \boxed{x_{jt}} \\ 1 \end{matrix} + \begin{matrix} p_1 \\ \boxed{\mathbf{E}_t} \\ p_2 \end{matrix}$$

- $(k, l)$ -th entry of  $\mathbf{B}_j$  determines the effect of  $j$ -th predictor on the  $(k, l)$ -th cell of the response tensor.

# Multiway Stick Breaking Prior for $B_j$ (Spencer et al., Psychometrika (2020); G. & Spencer, Bayesian Analysis (2021))



$$\bullet B_j = \sum_{r=1}^R \beta_{j,1}^{(r)} \circ \dots \circ \beta_{j,D}^{(r)}$$

Increasing Shrinkage across ranks  $r = 1, \dots, R$

- Shrinkage within every rank through generalized double pareto shrinkage prior.
- Shrinkage prior involves rank specific parameters  $\phi_{j,r}$ ,  $r = 1, \dots, R$ .
- They assume a stick breaking construction  $\phi_{j,1} = \xi_{j,1}, \phi_{j,2} = \xi_{j,2}(1 - \xi_{j,1}), \dots, \phi_{j,R} = \prod_{r=1}^{R-1} (1 - \xi_{j,r})$ .

# Theoretical Study: Bayesian Tensor Response Regression

$\mathbf{B}$  (tensor of dimensions  $m \times p_1 \times \dots \times p_D$ ): stacking tensor coefficients  $\mathbf{B}_1, \dots, \mathbf{B}_m$  together.

$$\mathcal{A}_T = \{ \mathbf{B} : \|\mathbf{B} - \mathbf{B}_0\|_2 < \epsilon \}, \quad \mathbf{B}_0 = \text{true value of } \mathbf{B}.$$

$\Pi_T(\cdot)$  is the posterior distribution of  $\mathbf{B}$  with  $T$  observations.

## Notion of Posterior Consistency

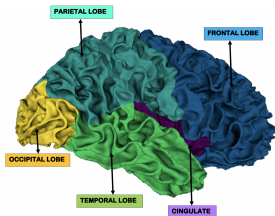
$$\Pi_T(\mathcal{A}_T^c) \rightarrow 0, \text{ a.s., when } T \rightarrow \infty.$$

Posterior consistency holds under the following conditions:

- 1  $\mathbf{B}_{0,j}$  assumes rank  $R_{0,j}$  PARAFAC decomposition with  $R > \max_j R_{0,j}$ .
- 2  $m \sum_{d=1}^D p_d \log(p_d) = o(T)$ ,  $s \log(m \prod_{d=1}^D p_d) = o(T)$ , where  $s$  is the number of nonzero entries in  $\mathbf{B}_0$ .
- 3 Covariate matrix has bounded singular values.

# Motivation: Brain Connectomes with Phenotypes

- **Data:** Brain connectome network ( $\mathbf{Y}_i$ ), creative achievement ( $x_i$ ) for subjects.
- $(k, l)$ -th entry of  $\mathbf{Y}_i$  represents “association” between  $k$ th and  $l$ th brain regions.

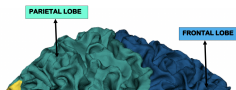


		Nodes				
		1	2	3	4	5
Nodes	1	0	a	b	c	d
	2		0	e	f	g
	3			0	h	i
	4				0	j
	5					0

- 68 Regions of interest (ROI), 34 in each hemisphere.
- 12 Lobes, 6 in each hemisphere.

# Motivation: Brain Connectomes with Phenotypes

- **Data:** Brain connectome network ( $\mathbf{Y}_i$ ), creative achievement ( $x_i$ ) for subjects.
- $(k, l)$ -th entry of  $\mathbf{Y}_i$  represents “association” between  $k$ th and  $l$ th brain regions.



		Nodes				
		1	2	3	4	5
Nodes	1	0	a	b	c	d
	2		0	e	f	g
	3			0	h	i
	4				0	j
	5					0

- 68 Regions of interest (ROI), 34 in each

## Inferential Goal

Develop regression of  $\mathbf{Y}_i$  on  $x_i$ , identify important network nodes related to creativity.

# Bayesian Symmetric Tensor on Vector Regression (Guha & G., Technometrics, 2021)

## Data Model

$$\mathbf{Y}_i = \mathbf{B}_1 x_{1i} + \cdots + \mathbf{B}_m x_{mi} + \mathbf{E}_i$$

- $\mathbf{Y}_i$  is a  $p \times \cdots \times p$  dimensional *symmetric* tensor response,  $x_{1i}, \dots, x_{mi}$  are  $m$  predictors.
- $\mathbf{B}_1, \dots, \mathbf{B}_m$  are  $p \times \cdots \times p$  dim. symmetric tensor coefficients.
- $\mathbf{B}_j$  follows a symmetric rank- $R$  PARAFAC decomposition  
$$\mathbf{B}_j = \sum_{r=1}^R \lambda_{j,r} \beta_j^{(r)} \circ \cdots \circ \beta_j^{(r)}.$$

$$\mathbf{B} = \lambda_1 \begin{matrix} \vec{u}_1 \\ \vec{u}_1 \\ \vec{u}_1 \end{matrix} + \dots + \lambda_R \begin{matrix} \vec{u}_R \\ \vec{u}_R \\ \vec{u}_R \end{matrix}$$

- $\beta_j^{(r)} = (\beta_{j,1}^{(r)}, \dots, \beta_{j,p}^{(r)})' \in \mathbb{R}^p$ ,  $\lambda_{j,r} \in \{0, 1\}$ .



# Influential Node Identification

- $\mathbf{u}_{j,k} = (\beta_{j,k}^{(1)}, \dots, \beta_{j,k}^{(R)})' = \mathbf{0}$  implies  $k$ -th node is unrelated to the  $j$ th predictor.
- Variable selection prior to identify important nodes,

$$\mathbf{u}_{j,k} \sim \begin{cases} N(\mathbf{0}, \mathbf{M}), & \text{if } \xi_{j,k} = 1 \\ \delta_{\mathbf{0}}, & \text{if } \xi_{j,k} = 0 \end{cases}, \quad \xi_{j,k} \sim \text{Ber}(\Delta),$$

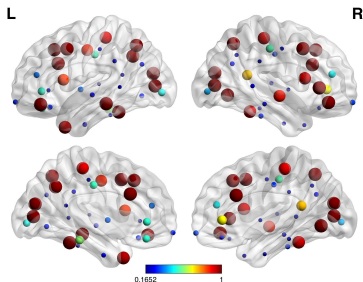
where  $\delta_{\mathbf{0}}$  is the Dirac-delta function at  $\mathbf{0}$ ,  $\mathbf{M}$  is a covariance matrix of order  $R \times R$ .

## Near Optimal Estimation of Predictive Density (Guha and G., 2021)

The predictive density of the proposed model can be estimated at a rate close to  $n^{-1/2}$  upto a  $\log(n)$  factor.

# Inference on Significant Nodes

- Compute posterior probability of  $\{\mathbf{u}_{j,k} = \mathbf{0}\}$  empirically from MCMC samples.
- $k$ th network node is related to the  $j$ th predictor if this probability is less than 0.5.



**13** frontal, **6** temporal ROIs are significant among **34** significant ROIs. (More than half of the identified ROIs).

- We offer posterior probabilities of each node being related to a predictor, which quantifies the statistical uncertainty.
- Posterior prob. close to 0 or 1 means less uncertainty with the decision.

## Space-Time Varying Coefficient Model

$$y(\mathbf{s}_i, t_i) = \mathbf{x}(\mathbf{s}_i, t_i)' \boldsymbol{\beta} + \mathbf{z}(\mathbf{s}_i, t_i)' \boldsymbol{\gamma}(\mathbf{s}_i, t_i) + \epsilon(\mathbf{s}_i, t_i)$$

- $p \times 1$  Fixed Effect
- $m \times 1$  Space-Time Varying Coefficients
- Non-Spatial Error following i.i.d.  $N(0, \tau^2)$

## Some Observations

- Only  $m$  of the  $p$  predictors have varying coefficients,  $m \leq p$ .
- $\mathbf{z}(\mathbf{s}_i, t_i) = \mathbf{1} \Rightarrow$  spatio-temporal geo-statistical model.
- $\mathbf{x}(\mathbf{s}_i, t_i) = \mathbf{z}(\mathbf{s}_i, t_i) \Rightarrow$  all predictor coefficients are space-time varying.

$$\{\gamma(\mathbf{s}, t) : (\mathbf{s}, t) \in \mathcal{D} \times \mathcal{T}\} \sim GP(\mathbf{0}, \mathbf{C}_\theta(\cdot, \cdot))$$

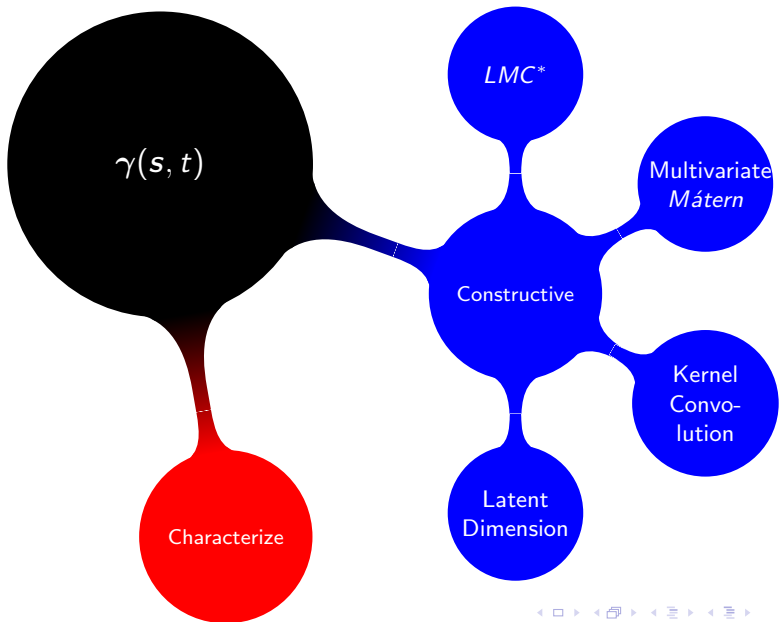
$$(\gamma(\mathbf{s}_1, t_1), \dots, \gamma(\mathbf{s}_n, t_n))' \sim N(\mathbf{0}, \mathbf{C}_\theta)$$

for any finite set of location-time tuples  $(\mathbf{s}_1, t_1), \dots, (\mathbf{s}_n, t_n)$ .

## Cross Covariance Kernel Matrix

- $\mathbf{C}_\theta(\cdot, \cdot)$  is the  $m \times m$  cross-covariance kernel matrix.
- $\mathbf{C}_\theta$  is the  $nm \times nm$  covariance matrix with the  $(i, j)$ -th block given by the  $m \times m$  matrix  $\mathbf{C}_\theta((\mathbf{s}_i, t_i), (\mathbf{s}_j, t_j))$ .

# Modeling Cross Covariance: Popular Approaches



# Full Likelihood from Gaussian Process (GP) Model

$\mathbf{y} = (y(\mathbf{s}_1, t_1), \dots, y(\mathbf{s}_n, t_n))'$ , ( $n \times 1$  vector)

$\mathbf{X} = [\mathbf{x}(\mathbf{s}_1, t_1) : \dots : \mathbf{x}(\mathbf{s}_n, t_n)]'$ , ( $n \times p$  matrix)

$\mathbf{Z} = \text{Block-diag}(\mathbf{z}(\mathbf{s}_1, t_1)', \dots, \mathbf{z}(\mathbf{s}_n, t_n)')$ , ( $n \times nm$  matrix)

- Model:  $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\mathbf{C}_\theta\mathbf{Z}' + \tau^2\mathbf{I})$ .
- Estimating parameters  $\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2$  from the likelihood

$$-\frac{\log(\det(\mathbf{Z}\mathbf{C}_\theta\mathbf{Z}' + \tau^2\mathbf{I}))}{2} - \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Z}\mathbf{C}_\theta\mathbf{Z}' + \tau^2\mathbf{I})^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2}$$

## Challenges

- Store  $\mathbf{Z}\mathbf{C}_\theta\mathbf{Z}' + \tau^2\mathbf{I}$
- Compute  $\text{Chol}(\mathbf{Z}\mathbf{C}_\theta\mathbf{Z}' + \tau^2\mathbf{I}) = \mathbf{L}\mathbf{L}'$ .

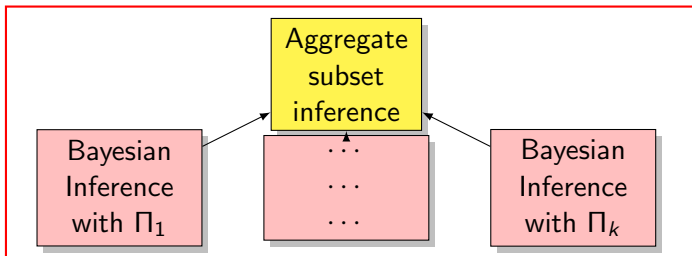
$n^3$  floating point operations per MCMC iteration  $\rightarrow$  Big-n problem

# Literature on Spatial/Spatio-Temporal Big Data

- **Low rank model** (Wahba, 1990; Higdon, 2001; Kamman & Wand, 2003; Paciorek, 2007; Lemos and Sanso, 2006; Banerjee et al., 2008; Cressie & Johannesson, 2008; Finley et al., 2009; Gramacy and Lee, 2008; **Guhaniyogi et al., 2011 & 2013**; Sang et al. 2012; Katzfuss, 2016).
- **Multiscale approaches** (Nychka, 2002; Johannesson et al., 2007; Tzeng and Huang, 2015; Nychka et al., 2015; Katzfuss, 2016; Katzfuss & Guinness, 2021, **Guhaniyogi & Sanso, 2017**).
- **Spectral approximations and composite likelihoods** (Fuentes, 2007; Eidvisk, 2016).
- **Sparsity**: Covariance tapering (Kaufman et al., 2008; Du et al., 2009; Sang et al., 2012; **Guhaniyogi, 2017**), INLA (Rue et al., 2009; Lindgren et al., 2011), 1agp (Gramacy and Apley, 2015), nearest neighbor processes (Stein et al., 2004; Stroud et al., 2014; Datta et al., 2016).

# Divide-and-Conquer Inference with Big Data

- Split the data  $\mathcal{S} = \{(\mathbf{s}_1, t_1), \dots, (\mathbf{s}_n, t_n)\}$ ,  $\mathcal{Y} = \{y(\mathbf{s}_1, t_1), \dots, y(\mathbf{s}_n, t_n)\}$ ,  $\mathcal{X} = \{\mathbf{x}(\mathbf{s}_1, t_1), \dots, \mathbf{x}(\mathbf{s}_n, t_n)\}$ ,  $\mathcal{Z} = \{\mathbf{z}(\mathbf{s}_1, t_1), \dots, \mathbf{z}(\mathbf{s}_n, t_n)\}$  into  $k$  exhaustive subsets  $\mathcal{S}_j, \mathcal{Y}_j, \mathcal{X}_j, \mathcal{Z}_j, j = 1, \dots, k$ .
- The  $j$ th subset posterior  $\Pi_j$  computed from the  $j$ -th subset containing  $M_j$  data points drawn randomly from the entire domain,  $M_1 + \dots + M_k \geq n$ .





# Construction of Subset Posteriors

## Subset Posterior: “Weak Learner” of Full Posterior

$$\Pi_j(\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2 | \mathcal{Y}_j, \mathcal{X}_j, \mathcal{Z}_j) \propto [ p(\mathcal{Y}_j | \mathcal{X}_j, \mathcal{Z}_j, \boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2) ]^{n/M_j} p(\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2)$$

- Likelihood:  $N(\mathbf{y}_j | \mathbf{X}_j \boldsymbol{\beta}, \mathbf{Z}_j \mathbf{C}_{\boldsymbol{\theta}, j} \mathbf{Z}_j' + \tau^2 \mathbf{I}_{M_j})$
- Prior Distribution
- These are “Stochastic Approximations” of the full posterior  $\Pi(\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2 | \mathcal{Y}, \mathcal{X}, \mathcal{Z})$ .
- For easier implementation you may get rid of the power  $n/M_j \Rightarrow$  satisfactory point estimation, wider confidence intervals.

How to combine  $\Pi_j$ 's optimally?

# Combine Subset Posteriors Marginally: Li et al. (2017), Biometrika; G. et al. (2022), Stat. Sci.

- Compute **Wasserstein mean**  $\bar{\Pi}$  of  $\Pi_1, \dots, \Pi_k$ .

Each  $\Pi_j$  multivariate normal  $\implies \bar{\Pi}$  multivariate normal

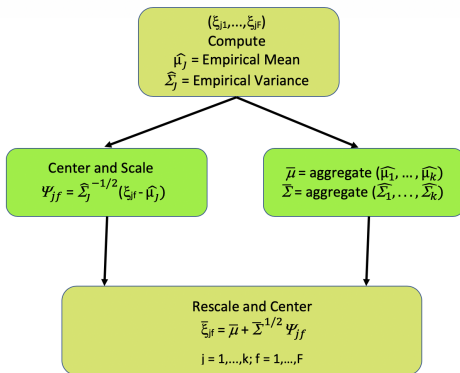
- $h(\boldsymbol{\theta}, \boldsymbol{\beta}) \in \mathbb{R}$ , a  $1D$  parametric function of  $(\boldsymbol{\theta}, \boldsymbol{\beta})$ .
- $\Pi_j^{-1}(u)$ :  $u$ th quantile of the  $j$ th subset posterior distribution of  $h(\boldsymbol{\theta}, \boldsymbol{\beta})$ ,  $u \in (0, 1)$ .
- If  $\bar{\Pi}^{-1}(u)$  is the  $u$ th quantile of the Wasserstein mean, then

$$\text{PIE Combination: } \bar{\Pi}^{-1}(u) = \frac{1}{k} \sum_{j=1}^k \Pi_j^{-1}(u), \forall u \in (0, 1)$$

- Combines marginals of posterior distribution separately.

# General Subset Posterior Aggregation Approach

- $\xi_{jf}$ :  $f$ -th post burn-in iterate of model parameters from the  $j$ -th subset.



Aggregated Monte Carlo (AMC) (G. et al., 2022, JMLR)

$$\bar{\mu} = \frac{1}{k} \sum_{j=1}^k \hat{\mu}_j, \quad \bar{\Sigma} = \text{AM}\{\hat{\Sigma}_1, \dots, \hat{\Sigma}_k\}$$

Wasserstein Posterior (WASP) (G. et al., 2022, JMLR)

$$\bar{\mu} = \frac{1}{k} \sum_{j=1}^k \hat{\mu}_j, \quad \bar{\Sigma} = \text{GM}\{\hat{\Sigma}_1, \dots, \hat{\Sigma}_k\}$$

- WASP uses geometric mean for aggregating  $\hat{\Sigma}_j$ 's.
- Geometric mean requires computing an iterative algorithm.

# Novelty vis-a-vis Existing Divide & Conquer Techniques

- **Aggregation of Subset Posteriors through Median** (Minsker et al., 2017)
- **Computation of Meta Posterior** (Guhaniyogi and Banerjee, 2017)
- **Consensus Monte Carlo (CMC)** (Scott et al., 2016), **Semiparametric Density Product (SDP)** (Neiswenger et al., 2014).
- **ADVI** (Kucukelbir et al., 2017).
- Both theory and practice for uncertainty quantification of parameters are unavailable for spatial/spatio-temporal process models.

For notational simplicity, we denote  $\mathbf{u} = (\mathbf{s}, t)$ , and assume  $M_1 = \dots = M_k = M = n/k$ .

## True Model ( $\mathcal{M}_0$ ) and Fitted Model ( $\mathcal{M}$ )

$$\mathcal{M} : y(\mathbf{u}) = z(\mathbf{u})' \gamma(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma(\mathbf{u}) = (\gamma_1(\mathbf{u}), \dots, \gamma_m(\mathbf{u}))'$$

$$\mathcal{M}_0 : y(\mathbf{u}) = z(\mathbf{u})' \gamma_0(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma_0(\mathbf{u}) = (\gamma_{0,1}(\mathbf{u}), \dots, \gamma_{0,m}(\mathbf{u}))'$$

- 1 The cross-covariance function for  $\gamma(\mathbf{u})$  has bounded eigenfunctions and polynomially decaying eigenvalues.

For notational simplicity, we denote  $\mathbf{u} = (\mathbf{s}, t)$ , and assume  $M_1 = \dots = M_k = M = n/k$ .

## True Model ( $\mathcal{M}_0$ ) and Fitted Model ( $\mathcal{M}$ )

$$\mathcal{M} : y(\mathbf{u}) = z(\mathbf{u})' \gamma(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma(\mathbf{u}) = (\gamma_1(\mathbf{u}), \dots, \gamma_m(\mathbf{u}))'$$

$$\mathcal{M}_0 : y(\mathbf{u}) = z(\mathbf{u})' \gamma_0(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma_0(\mathbf{u}) = (\gamma_{0,1}(\mathbf{u}), \dots, \gamma_{0,m}(\mathbf{u}))'$$

- 1** The cross-covariance function for  $\gamma(\mathbf{u})$  has bounded eigenfunctions and polynomially decaying eigenvalues.
- 2** The function  $\gamma_{0,g}(\mathbf{u})$  has  $\nu$  degrees of smoothness.
- 3** subsets are disjoint and subset size  $M$  must be greater than a certain fraction of  $n$  depending on  $\nu$ .

For notational simplicity, we denote  $\mathbf{u} = (\mathbf{s}, t)$ , and assume  $M_1 = \dots = M_k = M = n/k$ .

## True Model ( $\mathcal{M}_0$ ) and Fitted Model ( $\mathcal{M}$ )

$$\mathcal{M} : y(\mathbf{u}) = z(\mathbf{u})' \gamma(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma(\mathbf{u}) = (\gamma_1(\mathbf{u}), \dots, \gamma_m(\mathbf{u}))'$$

$$\mathcal{M}_0 : y(\mathbf{u}) = z(\mathbf{u})' \gamma_0(\mathbf{u}) + \epsilon(\mathbf{u}), \quad \gamma_0(\mathbf{u}) = (\gamma_{0,1}(\mathbf{u}), \dots, \gamma_{0,m}(\mathbf{u}))'$$

- 1 The cross-covariance function for  $\gamma(\mathbf{u})$  has bounded eigenfunctions and polynomially decaying eigenvalues.
- 2 The function  $\gamma_{0,g}(\mathbf{u})$  has  $\nu$  degrees of smoothness.
- 3 subsets are disjoint and subset size  $M$  must be greater than a certain fraction of  $n$  depending on  $\nu$ .
- 4 Subset posterior aggregation schemes follow a general rule satisfied by PIE, AMC and WASP.

Then  $\gamma_{0,g}$  is estimated close to the optimal rate of  $n^{-2\nu/(2\nu+3)}$ .

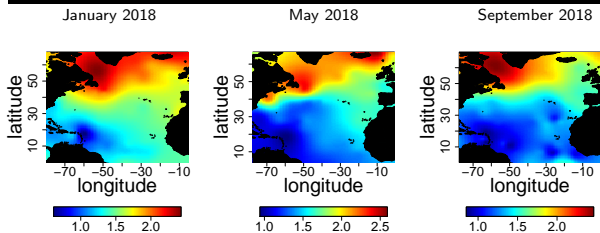


# Sea Surface Temperature (SST) and Sea Surface Salinity (SSS)

- Data from Hadley Center of the MET office in UK.
- Data on SST and SSS between  $0^{\circ} - 70^{\circ}$  N. latitude and  $0^{\circ} - 80^{\circ}$  W. longitude.
- We consider  $\sim 110K$  space-time observations on SST and SSS over the 12 months in 2018.

$$\text{Fit: } \text{SST}(\mathbf{s}, t) = \gamma_0(\mathbf{s}, t) + \gamma_1(\mathbf{s}, t) \text{SSS}(\mathbf{s}, t) + \epsilon(\mathbf{s}, t)$$

## Estimated Map of $\gamma_1(\mathbf{s}, t)$ in the North Atlantic



- 1 Divide data into 400 subsets.
- 2 An overall positive association between SSS and SST from equator to the pole.
- 3 In lower latitude, due to the pronounced salt accumulation as a result of excess heating and oceanic currents, SSS surges.
- 4 SSS decreases in comparison with SST during winter, except for the Brazilian coast due to the strong North Brazil Current.

# Predictive Inferential Accuracy

Predictive Inference on 600 hold out observations

	Coverage	MSPE	95% PI Length	Efficiency = $\frac{\log_2(\text{Eff. Sample Size})}{\text{Comp. Time}}$
AMC	0.98	2.92	6.60	9.99
PIE	0.97	2.93	5.93	-
CMC	0.90	74.95	24.71	1.06
WASP	0.98	2.92	5.66	9.99

## Some Important Findings

- 1 All divide-and-conquer schemes with the theoretical backing perform similarly.
- 2 Popular ML aggregation scheme CMC offers suboptimal inference.

# References

- Guhaniyogi, R. and Qamar, S. and Dunson, D.B. (2017). Bayesian tensor regression. *The Journal of Machine Learning Research*, 18(1), 2733–2763.
- Guhaniyogi, R. (2017). Convergence rate of Bayesian supervised tensor modeling with multiway shrinkage priors. *Journal of Multivariate Analysis*, 160, 157–168.
- Guhaniyogi, R. and Spencer, D. (2021). Bayesian tensor response regression with an application to brain activation studies. *Bayesian Analysis*, 16(4), 1221–1249.
- Spencer, D. and Guhaniyogi, R. and Prado, R. (2020). Joint Bayesian estimation of voxel activation and inter-regional connectivity in fMRI experiments. *Psychometrika*, 85(4), 845–869.
- Guha, S. and Guhaniyogi, R. (2021). Bayesian generalized sparse symmetric tensor-on-vector regression. *Technometrics*, 63(2), 160–170.
- Guhaniyogi, R. and Banerjee, S. (2018). Meta-kriging: Scalable Bayesian modeling and inference for massive spatial datasets. *Technometrics*, 60(4), 430–444.
- Guhaniyogi, R. and Li, C. and Savitsky, T. D. and Srivastava, S. (2022+). A divide-and-conquer Bayesian approach to large-scale kriging. *arXiv preprint arXiv:1712.09767*.
- Guhaniyogi, R. and Li, C. and Savitsky, T. D. and Srivastava, S. (2022). Distributed Bayesian varying coefficient modeling using a Gaussian process prior. Accepted, *Journal of Machine Learning Research (arXiv preprint arXiv:2006.00783)*.

# References

- Scott, S. L., Blocker, A. W., Bonassi, F. V., Chipman, H. A., George, E. I. and McCulloch, R. E. (2016). Bayes and big data: The consensus Monte Carlo algorithm. *International Journal of Management Science and Engineering Management*, 11(2), 78–88.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A. and Blei, D. M. (2017). Automatic differentiation variational inference. *Journal of Machine Learning Research*.
- Neiswanger, W. and Xing, E. (2017). Post-inference prior swapping. *International Conference on Machine Learning*, 2594–2602.
- Minsker, S., Srivastava, S., Lin, L. and Dunson, D. B. (2017). Robust and scalable Bayes via a median of subset posterior measures. *Journal of Machine Learning Research*, 18(1), 4488–4527.
- Li, C., Srivastava, S. and Dunson, D. B. (2017). Simple, scalable and accurate posterior interval estimation. *Biometrika*, 104(3), 665–680.
- Armagan, A. and Dunson, D. B. and Lee, J. (2013). Generalized double Pareto shrinkage. *Statistica Sinica*, 23(1), 119–143.
- Kolda, T. G. and Bader, B. W. (2009). Tensor decompositions and applications. *SIAM review*, 51(3), 455–500.
- Guhaniyogi, R. (2020). Bayesian methods for tensor regression. *Wiley StatsRef: Statistical References Online*, 1–18.