

# Efficient Hierarchical Bayesian Inference for Spatio-temporal Regression Models in Neuroimaging

Ali Hashemi

Technische Universität Berlin

joint work with Yijing Gao, Chang Cai, Sanjay Ghosh,  
Klaus-Robert Müller, Srikantan S. Nagarajan, and Stefan Haufe

35<sup>th</sup> Conference on Neural Information Processing Systems (NeurIPS 2021)



# Multi-task Linear Regression

$$\mathbf{Y}_g = \mathbf{L}\mathbf{X}_g + \mathbf{E}_g$$

$$\mathbf{Y}_g \in \mathbb{R}^{M \times T}$$

$$\mathbf{X}_g \in \mathbb{R}^{N \times T}$$

$$\mathbf{E}_g \in \mathbb{R}^{N \times T}$$

$$\mathbf{L} \in \mathbb{R}^{M \times N}$$

**Spatio-temporal generative model**

for  $g = 1, \dots, G$ ,  $G$ : #sample blocks or tasks

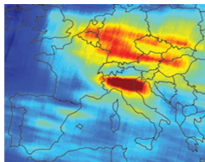
$M$ : #measurements or observations,  $T$ : #Samples,

$N$ : #coefficients or source components,

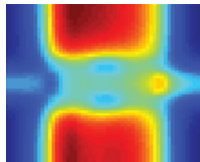
forward matrix (**known**): maps  $\mathbf{X}_g$  to  $\mathbf{Y}_g$

**Goal:** Estimate  $\{\mathbf{X}_g\}_{g=1}^G$  given  $\mathbf{L}$  and  $\{\mathbf{Y}_g\}_{g=1}^G$ :

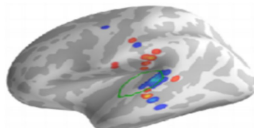
- ▶ Inverse problem in physics
- ▶ Multiple measurement vector (MMV) recovery problem in signal processing



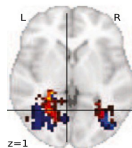
Temperature monitoring of climate [S. Beirle et al. 2003]



Temperature monitoring of CPU/GPU [J. Ranieri et al. 2012]



EEG/MEG Source Localization [H. Janati et al. 2020]

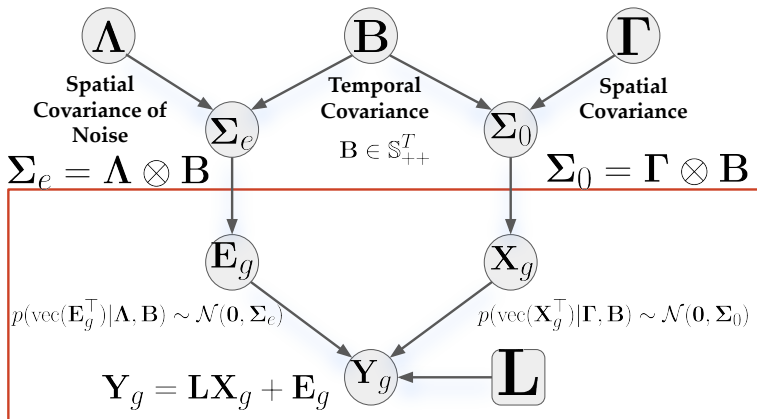


fMRI data analysis [M. B. Cai, et al. 2020]

# Hierarchical Bayesian Learning

Spatio-temporal dynamics of model parameters and noise are modeled to have **Kronecker product covariance structure**.

Probabilistic graphical model:



Posterior source distribution:  $p(\text{vec}(\mathbf{X}_g^\top) | \text{vec}(\mathbf{Y}_g^\top), \mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) \sim \mathcal{N}(\bar{\mathbf{x}}_g, \mathbf{\Sigma}_x)$   
with

$$\bar{\mathbf{x}}_g = \text{vec}(\bar{\mathbf{X}}_g^\top) = \mathbf{\Sigma}_0 \mathbf{D}^\top \tilde{\mathbf{\Sigma}}_y^{-1} \mathbf{y}_g$$

$$\mathbf{\Sigma}_x = \mathbf{\Sigma}_0 - \mathbf{\Sigma}_0 \mathbf{D}^\top \tilde{\mathbf{\Sigma}}_y^{-1} \mathbf{D} \mathbf{\Sigma}_0$$

$$\tilde{\mathbf{\Sigma}}_y = \mathbf{\Sigma}_y \otimes \mathbf{B}$$

$$\mathbf{\Sigma}_y = \mathbf{L} \mathbf{\Gamma} \mathbf{L}^\top + \mathbf{\Lambda},$$

where  $\mathbf{D} = \mathbf{L} \otimes \mathbf{I}_T$ .

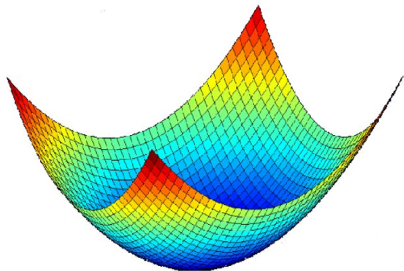
$\mathbf{\Gamma}$ ,  $\mathbf{\Lambda}$ ,  $\mathbf{B}$  are learned by minimizing the negative log marginal likelihood (Type-II) loss,  $-\log p(\mathbf{Y} | \mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ .

$$\text{Type - II Loss : } \mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_y| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^\top)$$

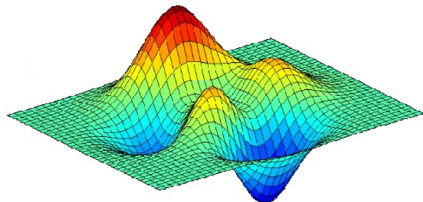
$$\text{Type - II Loss : } \mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_y| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^T)$$

$$\text{Type - II Loss : } \mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_y| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^T)$$

- ① **Non-convex** Type-II ML loss function: non-trivial to solve.



**convex function**



**non-convex function**

$$\text{Type - II Loss : } \mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_y| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^T)$$

- 1 **Non-convex** Type-II ML loss function: non-trivial to solve.
- 2 Most contributions in the literature *neglect the temporal structure* and are based on **MAP (Type-I)** estimation.
- 3 A few works that model the temporal dynamics often involve a *computationally demanding inference* scheme mostly based on expectation-maximization (EM).

- ▶ Derive novel Type-II algorithms that automatically learn the temporal structure
  - ① Exploit the intrinsic Riemannian geometry of temporal autocovariance matrices.
  - ② For stationary dynamics described by Toeplitz matrices, we employ the theory of circulant embeddings.
- ▶ Devise an efficient inference based on majorization-minimization optimization with guaranteed convergence properties.

To this end, we present a series of theorems resulting in a novel and efficient hierarchical Bayesian inference for spatio-temporal multi-task regression models.



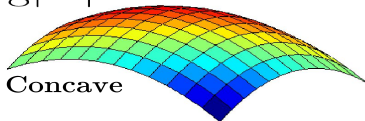
## Theorem (Majorizing function for temporal covariance update)

Optimizing  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{B}$  is equivalent to optimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

$$\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) = \text{tr}((\mathbf{B}^k)^{-1} \mathbf{B}) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}),$$

where  $\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^{\top} (\mathbf{\Sigma}_y^k)^{-1} \mathbf{Y}_g$ .

$\log |\mathbf{B}|$



$$\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_y| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top})$$

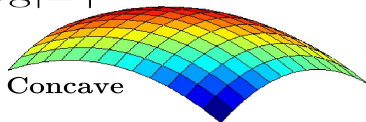
## Theorem (Majorizing function for temporal covariance update)

Optimizing  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{B}$  is equivalent to optimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

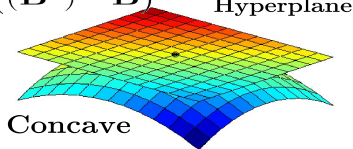
$$\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) = \text{tr}((\mathbf{B}^k)^{-1}\mathbf{B}) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}),$$

where  $\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{Y}_g$ .

$\log|\mathbf{B}|$



$\text{tr}((\mathbf{B}^k)^{-1}\mathbf{B})$



$$\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log|\mathbf{\Sigma}_y| + M \log|\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_y^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^\top)$$

# Convex Majorizing Functions

## Theorem (Majorizing function for temporal covariance update)

Optimizing  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{B}$  is equivalent to optimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

$$\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) = \text{tr}((\mathbf{B}^k)^{-1}\mathbf{B}) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}),$$

where  $\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{Y}_g$ .

## Theorem (Majorizing function for spatial covariance update)

Let  $\mathbf{H} = \text{diag}(\mathbf{h})$ ,  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ ,  $\mathbf{\Phi} := [\mathbf{L}, \mathbf{I}]$ , and  $\mathbf{\Sigma}_y = \mathbf{\Phi} \mathbf{H} \mathbf{\Phi}^\top$ . Then, optimizing  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{H}$  is equivalent to minimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

$$\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}^k) = \mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k) = \text{tr}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{\Phi} \mathbf{H}) + \text{tr}(\mathbf{M}_{\text{SN}}^k \mathbf{H}^{-1}),$$

where  $\mathbf{M}_{\text{SN}}^k := \mathbf{H}^k \mathbf{\Phi}^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{M}_{\text{space}}^k (\mathbf{\Sigma}_y^k)^{-1} \mathbf{\Phi} \mathbf{H}^k$ ,

$\mathbf{M}_{\text{space}}^k := \frac{1}{TG} \sum_{g=1}^G \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^\top$ .

# Convex Majorizing Functions

## Theorem (Majorizing function for temporal covariance update)

Optimizing  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{B}$  is equivalent to optimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

$$\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) = \text{tr}((\mathbf{B}^k)^{-1}\mathbf{B}) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}),$$

where  $\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{Y}_g$ .

## Theorem (Majorizing function for spatial covariance update)

Let  $\mathbf{H} = \text{diag}(\mathbf{h})$ ,  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ ,  $\mathbf{\Phi} := [\mathbf{L}, \mathbf{I}]$ , and  $\mathbf{\Sigma}_y = \mathbf{\Phi} \mathbf{H} \mathbf{\Phi}^\top$ . Then, optimizing  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  with respect to  $\mathbf{H}$  is equivalent to minimizing the following convex surrogate function, which majorizes  $\mathcal{L}_{kron}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$ :

$$\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}^k) = \mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k) = \text{tr}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{\Phi} \mathbf{H}) + \text{tr}(\mathbf{M}_{\text{SN}}^k \mathbf{H}^{-1}),$$

where  $\mathbf{M}_{\text{SN}}^k := \mathbf{H}^k \mathbf{\Phi}^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{M}_{\text{space}}^k (\mathbf{\Sigma}_y^k)^{-1} \mathbf{\Phi} \mathbf{H}^k$ ,

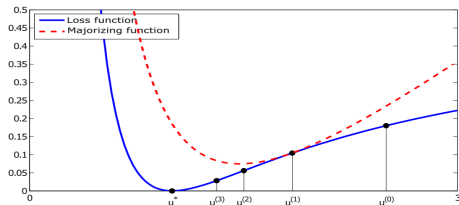
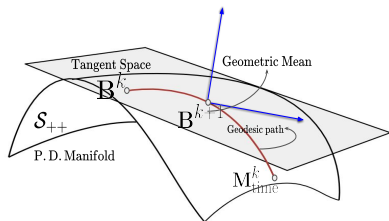
$\mathbf{M}_{\text{space}}^k := \frac{1}{TG} \sum_{g=1}^G \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^\top$ .

## Theorem (Geometric mean)

The cost function  $\mathcal{L}_{\text{CONV}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$  is strictly geodesically convex with respect to the P.D. manifold and its minimum with respect to  $\mathbf{B}$  can be attained by iterating the following update rule until convergence:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2},$$

which leads to a majorization-minimization (MM) algorithm with convergence guarantees  $\rightsquigarrow$  **Full Dugh**

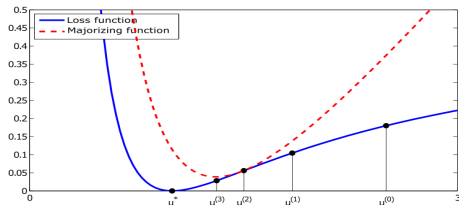
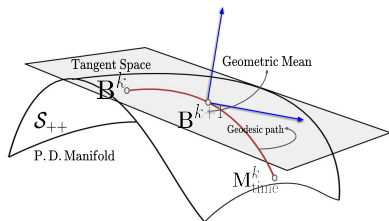


## Theorem (Geometric mean)

The cost function  $\mathcal{L}_{\text{CONV}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$  is strictly geodesically convex with respect to the P.D. manifold and its minimum with respect to  $\mathbf{B}$  can be attained by iterating the following update rule until convergence:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2},$$

which leads to a majorization-minimization (MM) algorithm with convergence guarantees  $\rightsquigarrow$  **Full Dugh**

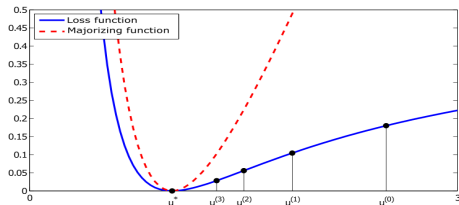
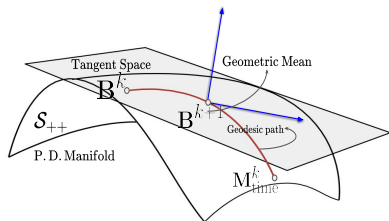


## Theorem (Geometric mean)

The cost function  $\mathcal{L}_{\text{CONV}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$  is strictly geodesically convex with respect to the P.D. manifold and its minimum with respect to  $\mathbf{B}$  can be attained by iterating the following update rule until convergence:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2},$$

which leads to a majorization-minimization (MM) algorithm with convergence guarantees  $\rightsquigarrow$  **Full Dugh**

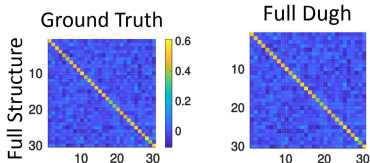


## Theorem (Geometric mean)

The cost function  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$  is strictly geodesically convex with respect to the P.D. manifold and its minimum with respect to  $\mathbf{B}$  can be attained by iterating the following update rule until convergence:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2},$$

which leads to a majorization-minimization (MM) algorithm with convergence guarantees  $\rightsquigarrow$  **Full Dugh**



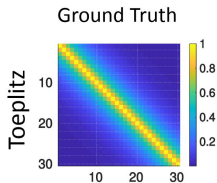


## Theorem (Geometric mean)

The cost function  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$  is strictly geodesically convex with respect to the P.D. manifold and its minimum with respect to  $\mathbf{B}$  can be attained by iterating the following update rule until convergence:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2},$$

which leads to a majorization-minimization (MM) algorithm with convergence guarantees  $\rightsquigarrow$  **Full Dugh**



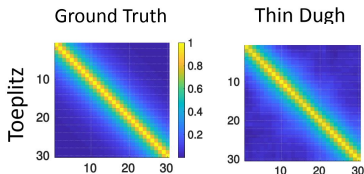
## Theorem (Temporal covariance update using circulant embedding)

Let  $\mathcal{L}_{\text{conv}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$  is constrained to the set of real-valued positive-definite Toeplitz matrices,  $\mathbf{B} \in \mathcal{B}^L : \mathbf{B} = \mathbf{Q}\mathbf{P}\mathbf{Q}^H$ , where  $\mathbf{P} = \text{diag}(\mathbf{p}) \in \mathbb{R}^{L \times L}$  with  $L > T$  be the circulant embedding of  $\mathbf{B}$ . Then the resulting constrained loss function is convex in  $\mathbf{p}$ , and its minimum with respect to  $\mathbf{p}$  can be obtained by iterating the following closed-form update rule until convergence:

$$p_l^{k+1} \leftarrow \sqrt{\frac{\hat{g}_l^k}{\hat{z}_l^k}} \text{ for } l = 1, \dots, L, \text{ where}$$

$$\hat{\mathbf{g}} := \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k)$$

$$\hat{\mathbf{z}} := \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q})$$



## Theorem (Temporal covariance update using circulant embedding)

$$p_i^{k+1} \leftarrow \sqrt{\frac{\hat{g}_i^k}{\hat{z}_i^k}} \text{ for } i = 1, \dots, L, \text{ where}$$

$$\hat{\mathbf{g}} := \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k)$$

$$\hat{\mathbf{z}} := \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q})$$

## Theorem (Spatial covariance with diagonal structure)

The cost function  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$  is convex in  $\mathbf{h}$ , and its minimum with respect to  $\mathbf{h}$  can be obtained according to the following closed-form update rule:

$$h_i^{k+1} \leftarrow \sqrt{\frac{g_i^k}{z_i^k}} \text{ for } i = 1, \dots, N + M, \text{ where}$$

$$\mathbf{g} := \text{diag}(\mathbf{M}_{\text{SN}}^k)$$

$$\mathbf{z} := \text{diag}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_y^k)^{-1} \mathbf{\Phi})$$

# Full and Thin Dugh

Combining this theoretical work, we developed a novel algorithm called “Dugh” for joint estimation of **spatial and temporal** covariances of **source and noise**.

## Algorithm 1: Full Dugh

**Input:** The lead field matrix  $\mathbf{L} \in \mathbb{R}^{M \times N}$  and  $G$  trials of measurement vectors  $\{\mathbf{Y}_g\}_{g=1}^G$ , where  $\mathbf{Y}_g \in \mathbb{R}^{M \times T}$ .

**Result:** Estimates of the source and noise variances  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , the temporal covariance  $\mathbf{B}$ , and the posterior mean  $\{\bar{\mathbf{x}}_g\}_{g=1}^G$  and covariance  $\Sigma_{\mathbf{x}}$  of the sources.

- 1 Choose a random initial value for  $\mathbf{B}$  as well as  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , and construct

$$\mathbf{H} = \text{diag}(\mathbf{h}) \text{ and } \mathbf{\Gamma} = \text{diag}([\gamma_1, \dots, \gamma_N]^\top).$$

- 2 Construct the augmented lead field  $\Phi = [\mathbf{L}, \mathbf{I}_M]$ .

- 3 Calculate the lead field  $\mathbf{D} = \mathbf{L} \otimes \mathbf{I}_T$  for vectorized sources.

- 4 Calculate the prior spatio-temporal covariance for the sources as  $\Sigma_0 = \mathbf{\Gamma} \otimes \mathbf{B}$ .

- 5 Calculate the spatial statistical covariance  $\Sigma_{\mathbf{y}} = \Phi \mathbf{H} \Phi^\top$ .

- 6 Calculate the spatio-temporal statistical covariance  $\tilde{\Sigma}_{\mathbf{y}} = \Sigma_{\mathbf{y}} \otimes \mathbf{B}$ .

- 7 Initialize  $k \leftarrow 1$

**repeat**

- 8 Calculate the posterior mean as  $\bar{\mathbf{x}}_g = \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{y}_g$ , for  $g = 1, \dots, G$ , where

$$\mathbf{y}_g = \text{vec}(\mathbf{Y}_g^\top) \in \mathbb{R}^{MT \times 1}.$$

- 9 Calculate  $\mathbf{M}_{\text{time}}^k$ , and update  $\mathbf{B}$  based on **Riemannian update on the manifold of PD. matrices**.

- 10 Calculate  $\mathbf{M}_{\text{SN}}^k$ , and update  $\mathbf{H}$ .

- 11  $k \leftarrow k + 1$

**until stopping condition is satisfied:**  $\|\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k\|_2 \leq \epsilon$  or  $k = k_{\text{max}}$ ;

- 12 Calculate the posterior covariance as  $\Sigma_{\mathbf{x}} = \Sigma_0 - \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{D} \Sigma_0$ .

## Algorithm 2: Thin Dugh

**Input:** The lead field matrix  $\mathbf{L} \in \mathbb{R}^{M \times N}$ , and  $G$  trials of measurement vectors  $\{\mathbf{Y}_g\}_{g=1}^G$ , where  $\mathbf{Y}_g \in \mathbb{R}^{M \times T}$ .

**Result:** Estimates of the source and noise variances  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , the temporal covariance  $\mathbf{B}$ , and the posterior mean  $\{\bar{\mathbf{x}}_g\}_{g=1}^G$ .

- 1 Choose a random initial value for  $\mathbf{p}$  as well as  $\mathbf{h}$ , and construct  $\mathbf{H} = \text{diag}(\mathbf{h})$  and  $\mathbf{P} = \text{diag}(\mathbf{p})$ .

- 2 Construct  $\mathbf{B} = \mathbf{Q} \mathbf{P} \mathbf{Q}^H$ , where  $\mathbf{Q} = [\mathbf{I}_M, \mathbf{0}] \mathbf{F}_L$  with  $L = 2T + 1$  and  $\mathbf{F}_L$  as DFT.

- 3 Construct the augmented lead field  $\Phi := [\mathbf{L}, \mathbf{I}_M]$ .

- 4 Calculate the prior spatio-temporal covariance for the sources as  $\Sigma_0 = \mathbf{\Gamma} \otimes \mathbf{B}$ .

- 5 Calculate the statistical covariance  $\Sigma_{\mathbf{y}} = \Phi \mathbf{H} \Phi^\top$ .

- 6 Calculate the spatio-temporal statistical covariance  $\tilde{\Sigma}_{\mathbf{y}} = \Sigma_{\mathbf{y}} \otimes \mathbf{B}$ .

- 7 Initialize  $k \leftarrow 1$

**repeat**

- 8 Calculate the posterior mean efficiently as  $\bar{\mathbf{x}}_g = \text{tr}(\mathbf{Q} \mathbf{P} (\mathbf{\Pi} \odot \mathbf{Q}^H \mathbf{Y}_g^\top \mathbf{U}_{\mathbf{x}}) (\mathbf{U}_{\mathbf{x}}^\top \mathbf{L} \mathbf{\Pi}^\top))$ , where  $\mathbf{L} \mathbf{\Pi} \mathbf{L}^\top = \mathbf{U}_{\mathbf{x}} \mathbf{D}_{\mathbf{x}} \mathbf{U}_{\mathbf{x}}^\top$  and  $[\mathbf{\Pi}]_{l,m} = \frac{1}{\sigma_{\mathbf{x}^k + p_l d_m}^2}$  for  $l = 1, \dots, L$  and  $m = 1, \dots, M$ .

- 9 Calculate  $\mathbf{M}_{\text{time}}^k$ , and update  $\mathbf{B}$  based on **Riemannian update for Toeplitz matrices using circulant embedding**.

- 10 Calculate  $\mathbf{M}_{\text{SN}}^k$ , and update  $\mathbf{H}$ .

- 11  $k \leftarrow k + 1$

**until stopping condition is satisfied:**  $\|\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k\|_2 \leq \epsilon$  or  $k = k_{\text{max}}$ ;

- 12 Calculate the posterior covariance as  $\Sigma_{\mathbf{x}} = \Sigma_0 - \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{D} \Sigma_0$ .

## Full Dugh: Temporal Covariance Update

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{1/2} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{1/2} (\mathbf{B}^k)^{1/2}$$

$$\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^\top (\Sigma_{\mathbf{y}}^k)^{-1} \mathbf{Y}_g$$

## Thin Dugh: Temporal Covariance Update

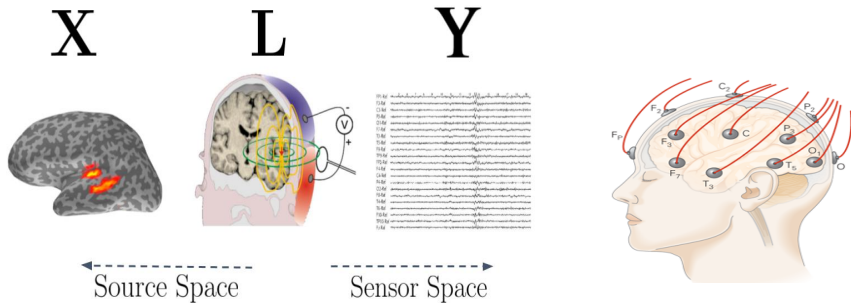
$$\mathbf{B} = \mathbf{Q} \mathbf{P} \mathbf{Q}^H, p_l^{k+1} \leftarrow \sqrt{\frac{\hat{g}_l^k}{\hat{z}_l^k}} \text{ for } l = 1, \dots, L$$

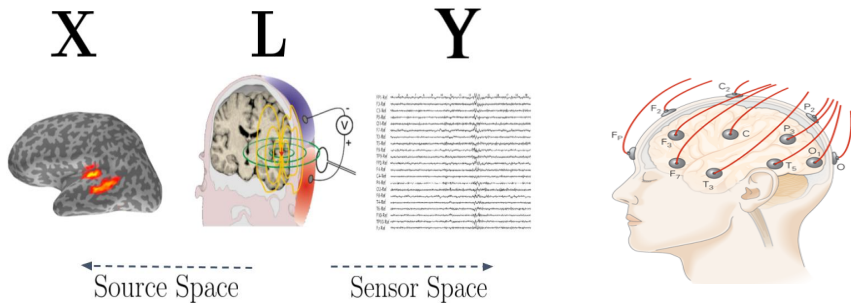
$$\hat{\mathbf{g}} := \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k)$$

$$\hat{\mathbf{z}} := \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q})$$

# Electromagnetic Brain Source Imaging (BSI)

**Electro-/Magnetoencephalography (E/MEG):** A non-invasive brain imaging technique with high temporal resolution (order of ms).



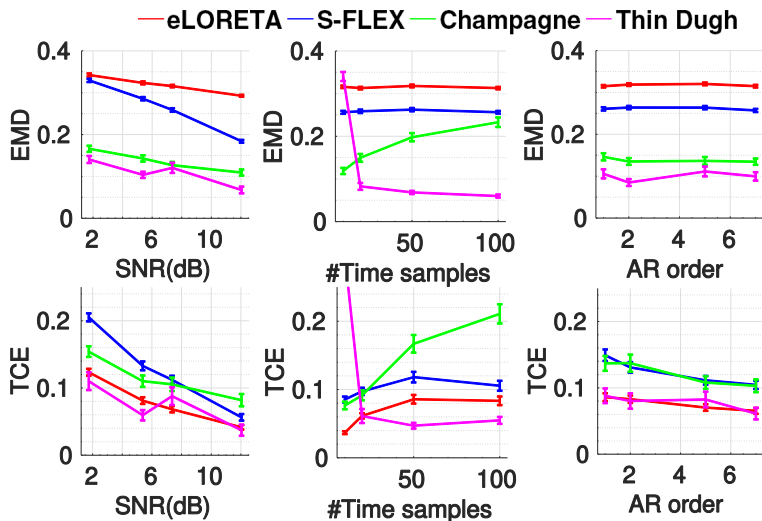


Ill-posed inverse problem: (#Sensors= 32 ~ 256 vs #Sources=  $10^3 \sim 10^4$ )

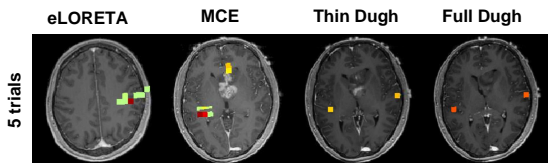
$$\mathbf{X}^* = \underset{\mathbf{X}}{\operatorname{argmin}} \underbrace{\|\mathbf{Y} - \mathbf{L}\mathbf{X}\|_F^2}_{\text{Likelihood: } p(\mathbf{Y}|\mathbf{X})} + \lambda \underbrace{\mathcal{R}(\mathbf{X})}_{\text{Prior: } p(\mathbf{X})}$$

- 1 Type-I MAP methods:  $\ell_1$ ,  $\ell_2$ ,  $\ell_{1,2}$ -norms, sparsity in transformed domains (Gabor).  
 [Pascual-Marqui et al., '07][Haufe et al, '08, '11][Gramfort et al., '12, '13][Castaño-Candamil et al., '15]
- 2 Type-II ML approaches: different sparse Bayesian learning (SBL) variants ignoring the temporal dynamics.  
 [Wipf et al., '09, '10, '11][Owen et al, '12][Cai et al., '17, '21]

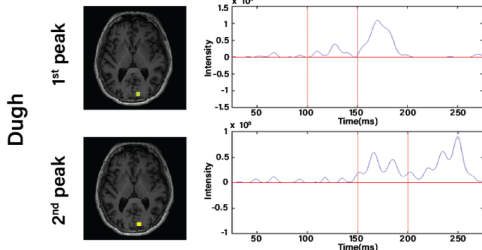
**Conclusion I:** Dugh consistently outperforms benchmark methods in the BSI literature according to all evaluation metrics.



**Conclusion II:** Dugh can provide accurate reconstruction even under extreme SNR conditions - superior to benchmarks.



## VEF





# Thank you for your attention!



## References



A. Hashemi, C. Cai, K.-R. Müller, S. S. Nagarajan and S. Haufe  
Joint Hierarchical Bayesian Learning of Full-structure Noise for Brain Source Imaging.  
*Medical Imaging meets NeurIPS (Med-NeurIPS) Workshop, 2020.*



A. Hashemi, Y. Gao, C. Cai, S. Ghosh, K.-R. Müller, S. S. Nagarajan and S. Haufe  
Joint Learning of Full-structure Noise in Hierarchical Bayesian Regression Models.  
*Preprint, 2021.* Draft is available on bioRxiv.



A. Hashemi, C. Cai, G. Kutyniok, K.-R. Müller, S. S. Nagarajan and S. Haufe  
Unification of sparse Bayesian learning algorithms for electromagnetic brain imaging with the majorization minimization framework.  
*NeuroImage 239, 2021.*



C. Cai, A. Hashemi, M. Diwakar, S. Haufe, K. Sekihara, S. S. Nagarajan  
Robust estimation of noise for electromagnetic brain imaging with the Champagne algorithm.  
*NeuroImage 225, 2021.*



A. Hashemi and S. Haufe  
Improving EEG Source Localization Through Spatio-Temporal Sparse Bayesian Learning.  
*26th IEEE European Signal Processing Conference (EUSIPCO), 2018.*



K. Sekihara and S. S. Nagarajan  
Electromagnetic Brain Imaging: A Bayesian Perspective.  
*Springer, 2015.*