Sparse linear regression with spatially and temporally correlated support

Michael Riis Andersen

Technical University of Denmark

Joint work with

Aki Vehtari, Ole Winther (supervisor) & Lars K. Hansen (supervisor)

October 29, 2014

Sparse Linear Regression and Structured Sparsity

Goal

Given $\boldsymbol{y} \in \mathbb{R}^N$ and $\boldsymbol{A} \in \mathbb{R}^{N \times D}$ for N < D, find *structured* sparse $\boldsymbol{x} \in \mathbb{R}^D$ s.t.

$$y = Ax + noise$$



Sparse Linear Regression and Structured Sparsity

Goal

Given $\boldsymbol{y} \in \mathbb{R}^N$ and $\boldsymbol{A} \in \mathbb{R}^{N \times D}$ for N < D, find *structured* sparse $\boldsymbol{x} \in \mathbb{R}^D$ s.t.

$$y = Ax + noise$$



Multiple sparse regression problems

- In many applications we observe a series of problems of the $y_t = Ax_t + e_t$ for different points in time t = 1, 2.., T
- By defining as $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1 & \boldsymbol{x}_2 & \dots & \boldsymbol{x}_T \end{bmatrix} \in \mathbb{R}^{D \times T}$ and similar for $\boldsymbol{Y}, \boldsymbol{E}$, we can write

$$Y = AX + E$$
.

Multiple sparse regression problems

- In many applications we observe a series of problems of the $y_t = Ax_t + e_t$ for different points in time t = 1, 2.., T
- By defining as $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1 & \boldsymbol{x}_2 & \dots & \boldsymbol{x}_T \end{bmatrix} \in \mathbb{R}^{D \times T}$ and similar for $\boldsymbol{Y}, \boldsymbol{E}$, we can write

$$Y = AX + E$$
.

I.I.D. vs. spatial vs. spatio-temporal sparsity patterns for X



Application: EEG Source Localization

• We observe a multivariate time series $\boldsymbol{Y} \in \mathbb{R}^{N \times T}$ and want to infer the underlying sources $\boldsymbol{X} \in \mathbb{R}^{D \times T}$ given a forward $\boldsymbol{A} \in \mathbb{R}^{N \times D}$.



• The brain is modelled using a set of discrete current dipoles

Application: EEG Source Localization

• We observe a multivariate time series $\boldsymbol{Y} \in \mathbb{R}^{N \times T}$ and want to infer the underlying sources $\boldsymbol{X} \in \mathbb{R}^{D \times T}$ given a forward $\boldsymbol{A} \in \mathbb{R}^{N \times D}$.



• The brain is modelled using a set of discrete current dipoles

Spike and slab priors for promoting sparsity

• Assume x_i is composed of two variables

 $x_i = s_i \cdot c_i, \quad s_i \in \{0,1\}, \quad c_i \in \mathbb{R}$



Spike and slab priors for promoting sparsity

• Assume x_i is composed of two variables

 $x_i = s_i \cdot c_i, \quad s_i \in \{0,1\}, \quad c_i \in \mathbb{R}$



• In terms of probability distributions,

$$p(s_i) = \text{Ber}(p)$$

$$p(x_i|s_i) = (1 - s_i)\delta(x_i) + s_i\mathcal{N}(x_i|0,\tau)$$

Spike and slab priors for promoting sparsity

• Assume x_i is composed of two variables

 $x_i = s_i \cdot c_i, \quad s_i \in \{0,1\}, \quad c_i \in \mathbb{R}$



• In terms of probability distributions,

$$p(s_i) = \text{Ber}(p)$$

$$p(x_i|s_i) = (1 - s_i)\delta(x_i) + s_i\mathcal{N}(x_i|0,\tau)$$

• Marginalizing out s_i

$$p(x_i) = (1 - p)\delta(x_i) + p\mathcal{N}(x_i|0, \tau)$$

 \uparrow
Spike Slab



The structured spike and slab prior

• We want to build a prior distribution for *x* s.t. the binary variables *s* are spatially correlated



• Idea: Generate a set of correlated random variables

$$p(oldsymbol{\gamma}) = \mathcal{N}\left(oldsymbol{\gamma} \Big| oldsymbol{\mu}, oldsymbol{\Sigma}
ight)$$

• and transform them into probabilities using a map $\phi : \mathbb{R} \to (0, 1)$ $p(s_i | \gamma_i) = \text{Ber} (s_i | \phi (\gamma_i))$ $p(x_i | s_i) = (1 - s_i) \, \delta (x_i) + s_i \mathcal{N} (x_i | 0, \tau)$

• Σ now determines the correlation structure of the support of x and

$$p(s_i = 1) = \phi\left(\frac{\mu_i}{\sqrt{1 + \Sigma_{ii}}}\right)$$

$$p(\boldsymbol{\gamma}) = \mathcal{N}\left(\boldsymbol{\gamma} | \boldsymbol{\mu}, \boldsymbol{\Sigma}\right),$$

 $p(s_i | \gamma_i) = \text{Ber}\left(s_i | \phi(\gamma_i)\right)$



$$p(oldsymbol{\gamma}) = \mathcal{N}\left(oldsymbol{\gamma} ig | oldsymbol{\mu}, oldsymbol{\Sigma}
ight), \ p(s_i ig | \gamma_i) = ext{Ber}\left(s_i ig | \phi\left(\gamma_i
ight)
ight)$$



$$p(oldsymbol{\gamma}) = \mathcal{N}\left(oldsymbol{\gamma} ig | oldsymbol{\mu}, oldsymbol{\Sigma}
ight), \ p(s_i ig | \gamma_i) = ext{Ber}\left(s_i ig | \phi\left(\gamma_i
ight)
ight)$$



1) $oldsymbol{\gamma} \sim \mathcal{N}\left(oldsymbol{\mu}, oldsymbol{\Sigma}
ight)$

$$p(oldsymbol{\gamma}) = \mathcal{N}\left(oldsymbol{\gamma} ig | oldsymbol{\mu}, oldsymbol{\Sigma}
ight), \ p(s_i ig \gamma_i) = ext{Ber}\left(s_i ig \phi\left(\gamma_i
ight)
ight)$$



1) $oldsymbol{\gamma} \sim \mathcal{N}\left(oldsymbol{\mu}, oldsymbol{\Sigma}
ight)$ 2)

2) Probabilities $\phi(oldsymbol{\gamma})$

$$p(\boldsymbol{\gamma}) = \mathcal{N}\left(\boldsymbol{\gamma} \middle| \boldsymbol{\mu}, \boldsymbol{\Sigma}\right),$$
$$p(s_i \middle| \gamma_i) = \text{Ber}\left(s_i \middle| \phi\left(\gamma_i\right)\right)$$



$$p(\boldsymbol{\gamma}) = \mathcal{N}\left(\boldsymbol{\gamma} \middle| \boldsymbol{\mu}, \boldsymbol{\Sigma}\right),$$
$$p(s_i \middle| \gamma_i) = \text{Ber}\left(s_i \middle| \phi\left(\gamma_i\right)\right)$$



Compare to sample from standard i.i.d. spike and slab prior



7 / 27

Spatiotemporal spike and slab priors

• Extending the prior to the spatio-temporal case is straightforward

• Let
$$m{\Gamma} = egin{bmatrix} m{\gamma}_1 & m{\gamma}_2 \dots m{\gamma}_T \end{bmatrix}$$
 and $m{S} = egin{bmatrix} m{S}_1 & m{S}_2 & \dots & m{S}_T \end{bmatrix}$

• Imposing a GP on Γ

$$p(\boldsymbol{\Gamma}) = \mathcal{N}(\boldsymbol{\Gamma}|\boldsymbol{\mu}, \boldsymbol{\Sigma}_{0})$$

$$p(\boldsymbol{X}|\boldsymbol{S}) = \prod_{t=1}^{T} \prod_{i=1}^{D} \left[(1 - s_{t,i})\delta(x_{t,i}) + s_{t,i}\mathcal{N}(x_{t,i}|0, \tau) \right]$$

$$p(\boldsymbol{S}|\boldsymbol{\Gamma}) = \prod_{t=1}^{T} \prod_{i=1}^{D} \operatorname{Ber}\left(z_{t,i}|\phi(\gamma_{t,i})\right)$$

 \bullet Assuming regular sampling in time: $\Sigma_0 = \Sigma_{\mathsf{temporal}} \otimes \Sigma_{\mathsf{spatial}}$

Inference using the spatio-temporal spike and slab prior

- Recall the model $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{E}$
- Assuming isotropic Gaussian noise leads to a posterior distribution of the form



$$p(\boldsymbol{X}, \boldsymbol{S}, \boldsymbol{\Gamma} | \boldsymbol{Y}) \propto \prod_{t=1}^{T} \mathcal{N} \left(\boldsymbol{y}_{t} | \boldsymbol{A} \boldsymbol{x}_{t}, \sigma_{0}^{2} \boldsymbol{I} \right) \prod_{t=1}^{T} \prod_{i=1}^{D} \left[(1 - \boldsymbol{s}_{t,i}) \delta \left(\boldsymbol{x}_{t,i} \right) + \boldsymbol{s}_{t,i} \mathcal{N} \left(\boldsymbol{x}_{i} | \boldsymbol{0}, \tau_{0} \right) \right] \\ \prod_{t=1}^{T} \prod_{i=1}^{D} \operatorname{Ber} \left(\boldsymbol{s}_{t,i} | \phi \left(\gamma_{t,i} \right) \right) \mathcal{N} \left(\boldsymbol{\Gamma} | \boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0} \right)$$

- Intractable due to the product over mixtures
- We use (parallel) expectation propagation for approximate inference

Approximate inference using EP

• Factorization of the true posterior

$$p(\boldsymbol{X}, \boldsymbol{S}, \boldsymbol{\Gamma} | \boldsymbol{Y}) \propto \prod_{t=1}^{T} \underbrace{\mathcal{N}\left(\boldsymbol{y}_{t} | \boldsymbol{A}\boldsymbol{x}_{t}, \sigma_{0}^{2}\boldsymbol{I}\right)}_{f_{1,t}(\boldsymbol{x}_{t})} \prod_{t=1}^{T} \prod_{i=1}^{D} \underbrace{\left[(1 - s_{i,t})\delta\left(\boldsymbol{x}_{i,t}\right) + s_{i,t}\mathcal{N}\left(\boldsymbol{x}_{i} | \boldsymbol{0}, \tau_{0}\right)\right]}_{f_{2,t,i}(\boldsymbol{x}_{i,t}, s_{i,t})}$$
$$\prod_{t=1}^{T} \prod_{i=1}^{D} \underbrace{\operatorname{Ber}\left(s_{i,t} | \phi\left(\gamma_{i,t}\right)\right)}_{f_{3,t,i}(s_{i,t}, \gamma_{i,t})} \underbrace{\mathcal{N}\left(\boldsymbol{\Gamma} | \boldsymbol{\mu}_{0}, \boldsymbol{\Sigma}_{0}\right)}_{f_{4}(\boldsymbol{\Gamma})}$$

• Site approximations

$$f_{2,t,i}(x_{i,t}, s_{i,t}) \approx \tilde{f}_{2,t,i}(x_{i,t}, s_{i,t}) = \mathcal{N}(x_{i,t} | \tilde{m}_{2,t,i}, \tilde{v}_{2,t,i}) \operatorname{Ber}(s_{i,t} | \phi(\tilde{\gamma}_{2,t,i}))$$

$$f_{3,t,i}(s_{i,t}, \gamma_{i,t}) \approx \tilde{f}_{3,t,i}(s_{i,t}, \gamma_{i,t}) = \operatorname{Ber}(s_{i,t} | \phi(\tilde{\gamma}_{3,t,i})) \mathcal{N}(\gamma_{i,t} | \tilde{\mu}_{3,t,i}, \tilde{\Sigma}_{3,t,i})$$

• Approximate posterior distribution has the form

$$Q(\boldsymbol{X}, \boldsymbol{S}, \boldsymbol{\Gamma}) = \prod_{t=1}^{T} \mathcal{N}\left(\boldsymbol{x}_{t} | \hat{\boldsymbol{m}}_{t}, \hat{\boldsymbol{V}}_{t}\right) \prod_{t=1}^{T} \operatorname{Ber}\left(\boldsymbol{s}_{t} | \phi\left(\hat{\boldsymbol{\gamma}}_{t}\right)\right) \mathcal{N}\left(\boldsymbol{\Gamma} | \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}\right).$$

Toy example: Recovery of Shepp-Logan Phantom I

Synthetic data generated using

 $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}_0 + \boldsymbol{e},$

- Spatial example, i.e. T = 1
- ullet Image size is 100 imes 100, i.e. $oldsymbol{x}_0\in\mathbb{R}^{10^4}$
- Number of non-zero pixels is K = 1723
- Number of measurements is $1.5 \mathcal{K} pprox 2585$, i.e. $oldsymbol{y} \in \mathbb{R}^{2585}$
- SNR is 20dB
- Matrix **A** is Gaussian I.I.D., i.e. $A_{ij} \sim \mathcal{N}(0, 1)$

ullet $\Sigma=$ squared exponential covariance defined on the image grid



Toy example: Recovery of Shepp-Logan Phantom II



Toy example: Recovery of Shepp-Logan Phantom II



The computational bottlenecks

• The computational bottlenecks are the updates of the posterior covariance matrices for $\boldsymbol{X} \in \mathbb{R}^{D \times D}$ and $\boldsymbol{\Gamma} \in \mathbb{R}^{DT \times DT}$

• For
$$\mathbf{x}_t$$
 for $t = 1, ..., T$:

$$\hat{\boldsymbol{V}}_t^* = \left[rac{1}{\sigma^2} \boldsymbol{A}^T \boldsymbol{A} + \boldsymbol{V}_{2,t}^{-1}
ight]^{-1}$$

For Γ:

$$\hat{\Sigma}^* = \left[\Sigma_0^{-1} + \Sigma_{3,t}^{-1}
ight]^{-1}$$

• Direct implementation $\mathcal{O}(TD^3)$ and $\mathcal{O}(T^3D^3)$, respectively.

Updating posterior covariance for X

- Direct implementation: $\mathcal{O}(T^3D^3)$
- For **X** we can use the assumptions $N \ll D$ and $\mathbf{V}_{1,t}^{-1} = \frac{1}{\sigma^2} \mathbf{A}^T \mathbf{A}$ and apply the matrix inversion lemma, so we can update each $\mathbf{V}_{1,t}$ as

$$ilde{oldsymbol{V}}_t = ilde{oldsymbol{V}}_{2,t} - ilde{oldsymbol{V}}_{2,t} oldsymbol{A}^T \left(\sigma_0^2 oldsymbol{I} + oldsymbol{A} ilde{oldsymbol{V}}_{2,t} oldsymbol{A}^T
ight)^{-1} oldsymbol{A} ilde{oldsymbol{V}}_{2,t}.$$

- We only need diagonal of V_t , i.e. $\mathcal{O}(N^2D)$.
- Total cost per iteration: $\mathcal{O}(TN^2D)$.
- Unfortunately, we cannot do the same for the posterior covariance matrix for $\Gamma \in \mathbb{R}^{DT \times DT}$

To reduce the computational complexity, we consider three different approximation schemes for the posterior covariance matrix for Γ

Approximation		Complexity	Storage
Standard EP	(EP)	$\mathcal{O}\left(T^{3}D^{3} ight)$	$\mathcal{O}\left(T^2 D^2\right)$
Low rank	(LR)	$\mathcal{O}(TDK^2)$	O (TDK)
Grouping/subsampling	(G)	$\mathcal{O}\left(T_g^3 D_g^3\right)$	$\mathcal{O}\left(T_{g}^{2}D_{g}^{2}\right)$
Common site precision	(CP)	$\mathcal{O}\left(T^2D+TD^2\right)$	$\mathcal{O}\left(T^2 + D^2 + TD\right)$

- ullet The costs above is only for updating posterior for Γ
- The cost for updating **X** is $\mathcal{O}(N^2D)$

The low rank approximation

• Use low rank approximation for prior covariance matrix for Γ

$$\mathbf{\Sigma}_0 = \mathbf{\Sigma}_{\mathsf{spatial}} \otimes \mathbf{\Sigma}_{\mathsf{temporal}} pprox oldsymbol{USU}^{\mathcal{T}} + oldsymbol{\Lambda}$$

where $\boldsymbol{U} \in \mathbb{R}^{DT \times K}$, $\boldsymbol{S} \in \mathbb{R}^{K \times K}$ and $\boldsymbol{\Lambda} \in \mathbb{R}^{DT \times DT}$ is a diagonal matrix

- We can obtain eigenvectors of Σ_0 from eigendecompositions of Σ_{spatial} and Σ_{temporal}
- We can now apply matrix inversion lemma as before to get $\mathcal{O}(TDK^2)$
- In practice: K can increases with both D and T!

Grouping/subsampling

• If the length-scales are large relative to the grid size, we can group neighbouring spatial locations and assign a single $\overline{\gamma}$ variable for each group and thus reducing the dimension of Γ

$$p(\boldsymbol{S}|\overline{\Gamma}) = \prod_{i,t} \operatorname{Ber}\left(s_{i,t} | \phi\left(\overline{\gamma}_{g(t,i)}\right)\right), \quad p(\overline{\Gamma}) = \mathcal{N}\left(\Gamma | \mu_0, \Sigma_0\right)$$

where $g: \mathcal{N} \times \mathcal{N} \to \mathcal{N}$ is function that assign each pair (t, i) to a group.

• Two variables in the same group do *not* share the same binary $z_{i,t}$ variable, but only the same probability of being active.



The common site precision approximation

$$f_{3,t,i}(s_{i,t},\gamma_{i,t}) \approx \tilde{f}_{3,t,i}(s_{i,t},\gamma_{i,t}) = \mathsf{Ber}\left(s_{i,t}|\phi\left(\tilde{\gamma}_{3,t,i}\right)\right) \mathcal{N}\left(\gamma_{i,t}|\tilde{\mu}_{3,t,i},\tilde{\Sigma}_{3,t,i}\right)$$

 $\bullet\,$ The idea is to approximate the site precisions/variances for $\Gamma\,$ with their common mean precision/variance

$$\begin{split} \left(\boldsymbol{\Sigma}_{s} \otimes \boldsymbol{\Sigma}_{t} + \boldsymbol{\tilde{\Sigma}}_{3}\right)^{-1} &\approx \left(\boldsymbol{\Sigma}_{s} \otimes \boldsymbol{\Sigma}_{t} + \boldsymbol{\bar{\Sigma}}_{3}\boldsymbol{I}\right)^{-1} \\ &= \left[\left(\boldsymbol{U}_{s} \otimes \boldsymbol{U}_{t}\right)\left(\boldsymbol{S}_{s} \otimes \boldsymbol{S}_{t}\right)\left(\boldsymbol{U}_{s}^{\mathsf{T}} \otimes \boldsymbol{U}_{t}^{\mathsf{T}}\right) + \boldsymbol{\bar{\Sigma}}_{3}\boldsymbol{I}\right]^{-1} \\ &= \left(\boldsymbol{U}_{s} \otimes \boldsymbol{U}_{t}\right)\left(\boldsymbol{S}_{s} \otimes \boldsymbol{S}_{t} + \boldsymbol{\bar{\Sigma}}_{3}\boldsymbol{I}\right)^{-1}\left(\boldsymbol{U}_{s}^{\mathsf{T}} \otimes \boldsymbol{U}_{t}^{\mathsf{T}}\right), \end{split}$$

- Using properties of Kronecker products, we can compute the diagonal of the covariance matrix in $O(T^2D + D^2T)$
- If the spatial positions are on a grid: $\Sigma_{s}=\Sigma_{x}\otimes\Sigma_{y}\otimes\Sigma_{z}$
- This approximation is very suitable for grid data
 Michael Riis Andersen

Comparing the approximation schemes for 1D problem



Figure: Data is generated from model y = Ax + e.

Michael Riis Andersen

Evaluation of methods

• Repeat for 100 different samples



Michael Riis Andersen

Spatiotemporal example

- We illustrate the spatiotemporal prior with a small simulated example
- We create a test signal $\boldsymbol{X} \in \mathbb{R}^{D \times T}$ with non-stationary sparsity pattern



Figure: Sparsity pattern for test signal

- We generate observations from the model Y = AX + E with D = 100, T = 30, $N = \frac{1}{3}D$ and the EP scheme to reconstruct to true signal
- $\Sigma_0 = \Sigma_{\textit{spatial}} \otimes \Sigma_{\textit{temporal}}$ with sq. exponential covariance functions
- Hyperparameters are estimated by maximizing the approximate marginal likelihood $p(\mathbf{Y})$



Michael Riis Andersen

- Subjects are presented with pictures of faces and "scrambled" faces
- Data is preprocessed and averaged over trials and subjects
- Number of sources D = 5124, number of EEG sensors N = 128 and number of time points T = 161, i.e. $X \in \mathbb{R}^{5124 \times 161}$ and $Y \in \mathbb{R}^{128 \times 161}$.



Face



EEG example: Face vs. scrambled face visual stimuli

- Subjects are presented with pictures of faces and "scrambled" faces
- Data is preprocessed and averaged over trials and subjects
- Number of sources D = 5124, number of EEG sensors N = 128 and number of time points T = 161, i.e. $X \in \mathbb{R}^{5124 \times 161}$ and $Y \in \mathbb{R}^{128 \times 161}$.



EEG example: Face vs. scrambled face visual stimuli

- Common precision approximation
- $\bullet\,$ Sq. exp. kernel for $\Gamma,\,1$ temporal length scale and 1 spatial length scale
- $N_{\text{train}} = 110$ for training and $N_{\text{test}} = 18$ for validation



Sensors are correlated, so this is only a sanity check at best

Michael Riis Andersen

Active sources

• Number of active sources as a fnc. of time, i.e. $\sum_{i} p(s_{i,t}|\mathbf{Y})$ vs. t.



Active sources

• Number of active sources as a fnc. of time, i.e. $\sum_{i} p(s_{i,t}|\mathbf{Y})$ vs. t.



0.9

0.8

0.6

0.5

0.4

0.2

0.1

• Active sources at t = 0.2s





Thank you for listening!



Any questions?

Number of variables proportional to

$$D \cdot T = 5124 \cdot 161 \approx 0.8M$$

$\%$ of variance explained in Σ_0	K	Size of $\boldsymbol{U} \in \mathbb{R}^{DT imes K}$
0.95	3961	24.3GB
0.90	2962	18.2GB
T 11 1		

Table: Low rank approximation for EEG example