# Support recovery and sup-norm convergence rates for sparse pivotal estimation

**Quentin Bertrand** 

Joint work with: Mathurin Massias (University of Genova) Alexandre Gramfort (INRIA) Joseph Salmon (IMAG, Univ Montpellier, CNRS)

## **Table of Contents**

#### Motivation

Calibrating  $\lambda$  for the Lasso

Generalization to other estimators

Experiments

## M/EEG inverse problem for brain imaging

sensors: electric and magnetic fields during a cognitive task
 goal: which parts of the brain are responsible for the signals?
 applications: epilepsy treatment, brain aging, anesthesia risks



## M/EEG data



## Source modeling (discretization with voxels)



 $\mathbf{B}^* \in \mathbb{R}^{p \times T}$ 

## The M/EEG inverse problem: modeling



n << p 6/41

## Very noisy data: must repeat recordings



## Very noisy data: must repeat recordings

average of 5 (top) / 10 (middle) / 50 (bottom) repetitions



### Noise covariance for each type of sensor

$$Y = XB + E, \quad E \sim \mathcal{N}(0, \Sigma) \tag{1}$$



• 3 different sensors  $\implies$  3 different noise structures

## A Multi-Task framework

Multi-Task regression notation:

- n observations (e.g., number of sensors)
- ▶ T tasks (e.g., temporal information)
- ▶ p features
- $\blacktriangleright$  *r* number of repetitions
- $Y^{(1)}, \ldots, Y^{(r)} \in \mathbb{R}^{n \times T}$  observation matrices;  $\overline{Y} = \frac{1}{r} \sum_{l} Y^{(l)}$
- $X \in \mathbb{R}^{n imes p}$  design matrix (known)

$$Y^{(l)} = X\mathbf{B}^* + S\mathbf{E}^{(l)}$$

where

B\* ∈ ℝ<sup>p×T</sup> : true source activity matrix (unknown)
 S ∈ S<sup>n</sup><sub>++</sub> co-standard deviation matrix (unknown)
 E<sup>(1)</sup>,...,E<sup>(r)</sup> ∈ ℝ<sup>n×T</sup> : white Gaussian noise

## Multi-Task penalties<sup>(1)</sup>

Popular convex penalties considered:

$$\hat{\mathbf{B}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left( \frac{1}{2nT} \| \bar{Y} - X\mathbf{B} \|_{F}^{2} + \lambda \Omega(\mathbf{B}) \right)$$



Sparse support: no structure

Penalty: Lasso type

$$\Omega(\mathbf{B}) = \|\mathbf{B}\|_{1} = \sum_{j=1}^{p} \sum_{k=1}^{T} |\mathbf{B}_{j,k}|$$

Parameter 
$$\hat{\mathbf{B}} \in \mathbb{R}^{p \times T}$$

<sup>&</sup>lt;sup>(1)</sup>G. Obozinski, B. Taskar, and M. I. Jordan. "Joint covariate selection and joint subspace selection for multiple classification problems". In: *Statistics and Computing* 20.2 (2010), pp. 231–252.

## Multi-Task penalties<sup>(1)</sup>

Popular convex penalties considered: Multi-Task Lasso (MTL)

$$\hat{\mathbf{B}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left( \frac{1}{2nT} \| \bar{Y} - X\mathbf{B} \|_{F}^{2} + \lambda \Omega(\mathbf{B}) \right)$$



Sparse support: group structure 🗸

Penalty: Group-Lasso type

$$\Omega(\mathbf{B}) = \|\mathbf{B}\|_{2,1} = \sum_{j=1}^{p} \|\mathbf{B}_{j,:}\|_{2}$$

where  $B_{j,:}$  the *j*-th row of B

<sup>(1)</sup>G. Obozinski, B. Taskar, and M. I. Jordan. "Joint covariate selection and joint subspace selection for multiple classification problems". In: *Statistics and Computing* 20.2 (2010), pp. 231–252.

## Summary



Figure: What you have.



Figure: What you want.

This is typically done using optimization based estimators:

$$\hat{\mathbf{B}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left( \frac{1}{2nT} \| \bar{Y} - X\mathbf{B} \|_{F}^{2} + \lambda \Omega(\mathbf{B}) \right)$$

## **Table of Contents**

Motivation

Calibrating  $\lambda$  for the Lasso

Generalization to other estimators

Experiments

Reminder on the Lasso theory<sup>(2)</sup> (i.i.d. case, Single-Task,  $y = X\beta + \sigma^* \varepsilon$ )

$$\hat{\boldsymbol{\beta}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \frac{1}{2n} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\|_{2}^{2} + \lambda \left\| \boldsymbol{\beta} \right\|_{1}$$

Theorem

i.i.d. Gaussian noise

► + X satisfying the "Restricted Eigenvalue" property ► +  $\lambda = 2\sigma_* \sqrt{\frac{2\log(p/\delta)}{n}}$ ► ⇒ with probability  $1 - \delta$ :  $\frac{1}{n} \| X \beta^* - X \hat{\beta}^{(\lambda)} \|^2 \le \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$ 

#### **BUT** $\sigma_*$ is unknown in practice !

(2) P. J. Bickel, Y. Ritov, and A. B. Tsybakov. "Simultaneous analysis of Lasso and Dantzig selector". In: Ann. Statist. 37.4 (2009), pp. 1705–1732. Reminder on the Lasso theory<sup>(2)</sup> (i.i.d. case, Single-Task,  $y = X\beta + \sigma^* \varepsilon$ )

$$\hat{\boldsymbol{\beta}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \frac{1}{2n} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\|_{2}^{2} + \lambda \left\| \boldsymbol{\beta} \right\|_{1}$$

Theorem

i.i.d. Gaussian noise

► + X satisfying the "Restricted Eigenvalue" property ► +  $\lambda = 2\sigma_* \sqrt{\frac{2\log(p/\delta)}{n}}$ ► ⇒ with probability  $1 - \delta$ :  $\frac{1}{n} \| X \beta^* - X \hat{\beta}^{(\lambda)} \|^2 \le \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$ 

#### **BUT** $\sigma_*$ is unknown in practice !

<sup>(2)</sup>P. J. Bickel, Y. Ritov, and A. B. Tsybakov. "Simultaneous analysis of Lasso and Dantzig selector". In: Ann. Statist. 37.4 (2009), pp. 1705–1732.

## Reminder on the square-root Lasso<sup>(3)(4)(5)</sup> (i.i.d. case, Single-Task)

$$\hat{\boldsymbol{\beta}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{\sqrt{n}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|_2 + \lambda \, \|\boldsymbol{\beta}\|_1$$
**Theorem**

i.i.d. Gaussian noise
+ X satisfying the "Restricted Eigenvalue" property
+ λ = 2√(2 log(p/\delta))/n
⇒ with high probability:  $\frac{1}{n} ||X\beta^* - X\hat{\beta}^{(\lambda)}||^2 \le \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$ 

#### $\lambda$ does not depend on $\sigma_*$ anymore!

<sup>(5)</sup>C. Giraud. Introduction to high-dimensional statistics. Vol. 138. CRC Press, 2014.

<sup>&</sup>lt;sup>(3)</sup>A. Belloni, V. Chernozhukov, and L. Wang. "Square-root Lasso: pivotal recovery of sparse signals via conic programming". In: *Biometrika* 98.4 (2011), pp. 791–806.

<sup>&</sup>lt;sup>(4)</sup>T. Sun and C.-H. Zhang. "Scaled sparse linear regression". In: *Biometrika* 99.4 (2012), pp. 879–898.

## Example

$$\hat{\boldsymbol{\beta}}_{\text{Lasso}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \left( \frac{1}{2n} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\|^{2} + \lambda \left\| \boldsymbol{\beta} \right\|_{1} \right)$$

Confirmed in practice:



## Example

$$\hat{\boldsymbol{\beta}}_{\sqrt{\text{Lasso}}} \in \operatorname*{arg\,min}_{\boldsymbol{\beta} \in \mathbb{R}^{p}} \left( \frac{1}{\sqrt{n}} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\| + \lambda \left\| \boldsymbol{\beta} \right\|_{1} \right)$$

Confirmed in practice:



## The Smoothed Concomitant Lasso<sup>(6)</sup> (i.i.d. case, Single-Task)

$$\hat{\boldsymbol{\beta}}^{(\lambda)} \in \underset{\boldsymbol{\beta} \in \mathbb{R}^{p}}{\arg\min} \frac{1}{\sqrt{n}} \underbrace{\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|_{2}}_{\text{non-smooth}} + \lambda \underbrace{\|\boldsymbol{\beta}\|_{1}}_{\text{non-smooth}}$$

$$\text{Idea: replacing } \|\cdot\|_{2} \text{ by } \underbrace{\|\cdot\|_{2} \Box \underline{\sigma} \, \omega\left(\frac{\cdot}{\underline{\sigma}}\right)}_{\text{smooth}} (z) = \underset{\boldsymbol{\sigma} \geq \underline{\sigma}}{\min} \left(\frac{\|\boldsymbol{z}\|_{2}^{2}}{2\sigma} + \frac{\sigma}{2}\right)$$

$$(\hat{\boldsymbol{\beta}}^{(\lambda)}, \hat{\boldsymbol{\sigma}}^{(\lambda)}) \in \underset{\boldsymbol{\beta} \in \mathbb{R}^{p}, \boldsymbol{\sigma} \geq \underline{\sigma}}{\arg\min} \frac{\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\boldsymbol{\beta}\|_{1}$$

jointly convex: alternate minimization

<sup>&</sup>lt;sup>(6)</sup>E. Ndiaye et al. "Efficient Smoothed Concomitant Lasso Estimation for High Dimensional Regression". In: Journal of Physics: Conference Series 904.1 (2017), p. 012006.

















## Solving the Smooth Concomitant Lasso

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \sigma} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

**Jointly convex** formulation : can be optimized by alternate minimization w.r.t.  $\beta$  and  $\sigma$  (gradient Lipschitz)

#### Alternate iteratively:

Fix  $\sigma$ : do one epoch of Lasso coordinate descent on  $\beta_j$ s: For  $j \in 1...p$  $\hat{\beta}_j \leftarrow \operatorname*{arg\,min}_{\beta_j \in \mathbb{R}} \frac{\|y - X\beta\|^2}{2n} + \lambda \sigma |\beta_j|$  (Lasso step)

## Solving the Smooth Concomitant Lasso

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \sigma} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

**Jointly convex** formulation : can be optimized by alternate minimization w.r.t.  $\beta$  and  $\sigma$  (gradient Lipschitz)

#### Alternate iteratively:

 Fix σ: do one epoch of Lasso coordinate descent on β<sub>j</sub>s: For j ∈ 1...p β<sub>j</sub> ← arg min μ||y - Xβ||<sup>2</sup>/2n + λσ|β<sub>j</sub>| (Lasso step)

 Fix β: closed form solution to update σ ∂ ← max (μ||y - Xβ||/√n, g) (Noise estimation step)

## Solving the Smooth Concomitant Lasso

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \underline{\sigma}} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

**Jointly convex** formulation : can be optimized by alternate minimization w.r.t.  $\beta$  and  $\sigma$  (gradient Lipschitz)

#### Alternate iteratively:

 Fix σ: do one epoch of Lasso coordinate descent on β<sub>j</sub>s: For j ∈ 1...p β<sub>j</sub> ← arg min μ|y - Xβ||<sup>2</sup>/2n + λσ|β<sub>j</sub>| (Lasso step)
 Fix β: closed form solution to update σ δ ← max (μ|y - Xβ||/√n, σ) (Noise estimation step)

## The Smoothed Concomitant Lasso (i.i.d. case, Single-Task)

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \underline{\sigma}} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

#### Questions:

- What is the impact of the smoothing on statistical properties? (sup-norm convergence? support recovery?)
- How to choose  $\underline{\sigma}$ ?
- Can this estimator (with unknown noise level σ\*) generalize for correlated noise (with unknown covariance matrix Σ\*)?

## Statistical properties<sup>(7)</sup>



- i.i.d. Gaussian noise
- ▶ + X satisfying the "mutual incoherence" property ▶ +  $\lambda \sim \frac{\sqrt{\log p}}{\sqrt{n}}$
- $\blacktriangleright + \underline{\sigma} \le \sigma_* / \sqrt{2}$

 $\blacktriangleright \implies$  with high probability:

$$\|\boldsymbol{\beta}^* - \hat{\boldsymbol{\beta}}\|_{\infty} \le C\boldsymbol{\sigma}_* \sqrt{\frac{\log p}{n}}$$

<sup>&</sup>lt;sup>(7)</sup>M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal estimation". In: AISTATS (2020).

## SCL support recovery perforance as a function of $\lambda$ and $\underline{\sigma}$



Figure: (Synthetic data, n = 50, p = 1000) Hard recovery loss for different values of SNR for the SCL.

## The Smoothed Concomitant Lasso (i.i.d. case, Single-Task)

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \sigma} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

#### Questions:

- ► What is the impact of the smoothing on statistical properties? (sup-norm convergence? support recovery?) ✓
- How to choose  $\underline{\sigma}$ ?  $\checkmark$
- Can this estimator (with unknown noise level σ\*) generalize for correlated noise (with unknown covariance matrix Σ\*)?

## **Table of Contents**

Motivation

Calibrating  $\lambda$  for the Lasso

Generalization to other estimators

Experiments

## Smoothing other norms

More interesting: S. van de Geer introduced the pivotal *multivariate*  $\sqrt{Lasso^{(8)}}$ :

$$\underset{\mathbf{B}\in\mathbb{R}^{p\times T}}{\arg\min}\frac{1}{\sqrt{nq}}\|\bar{Y}-X\mathbf{B}\|_{*}+\lambda\|\mathbf{B}\|_{2,1}$$

hard to solve ,<sup>(9)</sup> statistical analysis makes stringent assumptions

Smoothing the datafit makes optim. and stats easier!

<sup>&</sup>lt;sup>(8)</sup>S. van de Geer and B. Stucky. "X 2-confidence sets in high-dimensional regression". In: Statistical analysis for high-dimensional data. Springer, 2016, pp. 279–306.

<sup>&</sup>lt;sup>(9)</sup>A. J. Molstad. "Insights and algorithms for the multivariate square-root lasso". In: *arXiv preprint arXiv:1909.05041* (2019).

## Smoothing the nuclear norm<sup>(10)</sup>

Nuclear norm (Schatten-1 norm, or trace norm):  $Z \in \mathbb{R}^{n \times T}$ 

$$\left\|Z\right\|_* = \sum_{i=1}^{n \wedge T} \gamma_i$$

where the  $\gamma_i$ 's are the singular values of Z

$$\begin{aligned} \|\cdot\|_* \Box \left(\frac{1}{2\underline{\sigma}} \|\cdot\|^2 + \frac{n}{2}\right)(Z) &= \sum_i \mathsf{huber}_{\underline{\sigma}}\left(\gamma_i\right) \\ &= \min_{S \succeq \underline{\sigma}} \left(\frac{1}{2} \|Z\|_{S^{-1}}^2 + \frac{1}{2}\operatorname{Tr}(S)\right) \end{aligned}$$

where  $||Z||_{S^{-1}}^2 \triangleq \operatorname{Tr}(Z^{\top}S^{-1}Z)$ 

 $<sup>^{(10)}</sup>$  Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: NIPS (2019).

## Generalization ? Yes ! (correlated Gaussian noise, Multi-Task)



#### Benefits

▶ jointly convex formulation

Drawbacks:

Statistically: estimate a matrix of size  $n^2$ .  $\overline{\mathcal{O}(n^2)}$  parameters to estimate for S only nT observations

(11) M. Massias et al. "Generalized Concomitant Multi-Task Lasso for Sparse Multimodal Regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

## Generalization ? Yes ! (correlated Gaussian noise, Multi-Task)



#### Benefits

jointly convex formulation

#### Drawbacks:

Statistically: estimate a matrix of size  $n^2$ .  $\overline{\mathcal{O}(n^2)}$  parameters to estimate for <u>S</u> only nT observations

<sup>(11)</sup> M. Massias et al. "Generalized Concomitant Multi-Task Lasso for Sparse Multimodal Regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

## Statistical properties<sup>(12)</sup>



i.i.d. Gaussian noise

► + X satisfying the "mutual incoherence" property ► +  $\lambda \sim \frac{\sqrt{\log p}}{T\sqrt{n}}$ ► +  $\underline{\sigma} \leq \sigma_*/\sqrt{2}$ ► ⇒ with high probability: at least  $1 - \dots - 2ne^{-cT/n}$  $\frac{1}{T} \|\mathbf{B}^* - \hat{\mathbf{B}}\|_{2,\infty} \leq C\sigma_* \frac{1}{T} \sqrt{\frac{\log p}{n}}$ 

<sup>(12)</sup>M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal estimation". In: AISTATS (2020).

## SGCL support recovery perforance as a function of $\lambda$ and $\underline{\sigma}$



Figure: (Synthetic data, n = 50, p = 1000, T = 20) Hard recovery loss for different values of SNR for the SGCL.

## Can take advantage of the repetitions? Yes!

#### Smoothed Generalized Concomitant Lasso (SGCL)<sup>(13)</sup>:



**Concomitant Lasso with Repetitions** (CLaR)<sup>(14)</sup>:

$$(\hat{\mathbf{B}}^{\text{CLaR}}, \hat{S}^{\text{CLaR}}) \in \underset{\substack{\mathbf{B} \in \mathbb{R}^{p \times T} \\ S \in \mathbb{S}^{n}_{++}, S \succeq \underline{\sigma}}}{\operatorname{arg\,min}} \quad \frac{\sum_{l=1}^{r} \left\| Y^{(l)} - X\mathbf{B} \right\|_{S^{-1}}^{2}}{2nTr} + \frac{\operatorname{Tr}(S)}{2n} + \lambda \left\| \mathbf{B} \right\|_{2,1}$$

(13) M. Massias et al. "Generalized Concomitant Multi-Task Lasso for Sparse Multimodal Regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

 $^{(14)}$ Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NIPS* (2019).

## SGCL and CLaR computations: ${\rm B}$ update

Jointly convex: alternate minimization converges

**B** Update (S fixed): "smooth + non-smooth" optimization

(Block) Coordinate Descent (Soft-Threshold.) : update  ${\rm B}$  row-wise

#### Possible refinements:

- ▶ (Gap) safe screening rules<sup>(15)</sup>
- ► Strong rules<sup>(16)</sup>
- ► Active sets methods<sup>(17)</sup> etc.

<sup>&</sup>lt;sup>(15)</sup>L. El Ghaoui, V. Viallon, and T. Rabbani. "Safe feature elimination in sparse supervised learning". In: J. Pacific Optim. 8.4 (2012), pp. 667–698.

<sup>&</sup>lt;sup>(16)</sup>R. Tibshirani et al. "Strong rules for discarding predictors in lasso-type problems". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 74.2 (2012), pp. 245–266.

<sup>&</sup>lt;sup>(17)</sup>T. B. Johnson and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: ICML. 2015, pp. 1171–1179.

## SGCL and CLaR computations: ${\rm B}$ update

Jointly convex: alternate minimization converges

**B** Update (S fixed): "smooth + non-smooth" optimization

(Block) Coordinate Descent (Soft-Threshold.) : update B row-wise

Possible refinements:

- ▶ (Gap) safe screening rules<sup>(15)</sup>
- ▶ Strong rules<sup>(16)</sup>
- Active sets methods<sup>(17)</sup> etc.

<sup>(15)</sup> L. El Ghaoui, V. Viallon, and T. Rabbani. "Safe feature elimination in sparse supervised learning". In: J. Pacific Optim. 8.4 (2012), pp. 667–698.

<sup>(16)</sup> R. Tibshirani et al. "Strong rules for discarding predictors in lasso-type problems". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 74.2 (2012), pp. 245–266.

<sup>&</sup>lt;sup>(17)</sup>T. B. Johnson and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: *ICML*. 2015, pp. 1171–1179.

## SGCL and CLaR computations: S update

#### **S** Update (B fixed):

For SGCL and CLaR the problem can be reformulated as

$$\hat{\boldsymbol{S}} = \operatorname*{arg\,min}_{\boldsymbol{S} \in \mathbb{S}^n_{++}, \boldsymbol{S} \succeq \underline{\sigma}} \left( \frac{1}{2n} \underbrace{\operatorname{Tr}[\boldsymbol{Z}^\top \boldsymbol{S}^{-1} \boldsymbol{Z}]}_{\|\boldsymbol{Z}\|_{\boldsymbol{S}^{-1}}^2} + \frac{1}{2n} \operatorname{Tr}(\boldsymbol{S}) \right)$$

## SGCL and CLaR computations: S update

#### **S** Update (B fixed):

For SGCL and CLaR the problem can be reformulated as

$$\hat{\boldsymbol{S}} = \operatorname*{arg\,min}_{\boldsymbol{S} \in \mathbb{S}^n_{++}, \boldsymbol{S} \succeq \underline{\sigma}} \left( \frac{1}{2n} \underbrace{\operatorname{Tr}[\boldsymbol{Z}^\top \boldsymbol{S}^{-1} \boldsymbol{Z}]}_{\|\boldsymbol{Z}\|_{\boldsymbol{S}^{-1}}^2} + \frac{1}{2n} \operatorname{Tr}(\boldsymbol{S}) \right)$$

<u>Closed-form solution</u> (Spectral clipping):

if  $U^{\top} \operatorname{diag}(s_1, \ldots, s_n) U$  is the spectral decomposition of  $ZZ^{\top}$ :

$$\hat{S} = U^{\top} \operatorname{diag}(\max(\underline{\sigma}, \sqrt{s_1}), \dots, \max(\underline{\sigma}, \sqrt{s_n}))U$$

 $\underline{Rem}:$  as in the classical concomitant Lasso, at each step CLaR and SGCL estimate alternatively B and S

## SGCL and CLaR computations: S update

#### **S** Update (B fixed):

For SGCL and CLaR the problem can be reformulated as

$$\hat{\boldsymbol{S}} = \operatorname*{arg\,min}_{\boldsymbol{S} \in \mathbb{S}^n_{++}, \boldsymbol{S} \succeq \underline{\sigma}} \left( \frac{1}{2n} \underbrace{\operatorname{Tr}[\boldsymbol{Z}^\top \boldsymbol{S}^{-1} \boldsymbol{Z}]}_{\|\boldsymbol{Z}\|_{\boldsymbol{S}^{-1}}^2} + \frac{1}{2n} \operatorname{Tr}(\boldsymbol{S}) \right)$$

<u>Closed-form solution</u> (Spectral clipping):

if  $U^{\top} \operatorname{diag}(s_1, \ldots, s_n) U$  is the spectral decomposition of  $ZZ^{\top}$ :

$$\hat{S} = U^{\top} \operatorname{diag}(\max(\underline{\sigma}, \sqrt{s_1}), \dots, \max(\underline{\sigma}, \sqrt{s_n}))U$$

 $\underline{\rm Rem}:$  as in the classical concomitant Lasso, at each step CLaR and SGCL estimate alternatively  $\rm B$  and S

## The Smoothed Concomitant Lasso (i.i.d. case, Single-Task)

$$(\hat{\beta}^{(\lambda)}, \hat{\sigma}^{(\lambda)}) \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^{p}, \sigma \geq \sigma} \frac{\|y - X\beta\|_{2}^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

#### Questions:

- ► What is the impact of the smoothing on statistical properties? (sup-norm convergence? support recovery?) ✓
- How to choose  $\underline{\sigma}$ ?  $\checkmark$
- Can this estimator (with unknown noise level σ<sup>\*</sup>) generalize for correlated noise (with unknown covariance matrix Σ<sup>\*</sup>)? ✓

## **Table of Contents**

Motivation

Calibrating  $\lambda$  for the Lasso

Generalization to other estimators

Experiments

## Simulated scenarios

▶ 
$$n = 150$$
,  $p = 500$ ,  $q = 100$ 

► X Toeplitz-correlated

•  $S^*$  Toeplitz matrix:  $S^*_{i,j} = \rho_{S^*}^{|i-j|}$ ,  $\rho_{S^*} \in ]0,1[$ 



## **Real data**



(a) ours (b) MTL

Figure: Real data, left auditory stimulations (n = 102, p = 7498, q = 76, r = 63) Sources found in the left hemisphere (top) and the right hemisphere (bottom) after left auditory stimulations.

- expected: 2 sources (one in each auditory cortex)
- $\lambda$  chosen such that  $\|\hat{B}\|_{2,0} = 2$

## **Real data**



(a) ours (b) SGCL (c) MLER (d) MLE (e) MRCER (f) MTL

Figure: Real data, left auditory stimulations (n = 102, p = 7498, q = 76, r = 63) Sources found in the left hemisphere (top) and the right hemisphere (bottom) after left auditory stimulations.

- expected: 2 sources (one in each auditory cortex)
- $\lambda$  chosen such that  $\|\hat{\mathbf{B}}\|_{2,0} = 2$
- deep sources for SGCL and  $\ell_{2,1}$ -MRCER (not visible)

Statistical analysis of the smoothing for the square-root Lasso and the multivariate square-root Lasso

• Guidelines on how to set the smoothing parameter  $\underline{\sigma}$ 

- Statistical analysis of the smoothing for the square-root Lasso and the multivariate square-root Lasso
- Guidelines on how to set the smoothing parameter  $\underline{\sigma}$
- New estimator to handle correlated noise and repetitions in Multi-Task

- Statistical analysis of the smoothing for the square-root Lasso and the multivariate square-root Lasso
- Guidelines on how to set the smoothing parameter  $\underline{\sigma}$
- New estimator to handle correlated noise and repetitions in Multi-Task

Improved support identification

- Statistical analysis of the smoothing for the square-root Lasso and the multivariate square-root Lasso
- Guidelines on how to set the smoothing parameter  $\underline{\sigma}$
- New estimator to handle correlated noise and repetitions in Multi-Task
- Improved support identification
- Future work: study the influence of <u>σ</u> from the optimization point of view, see if the analysis holds for the LAD-Lasso

- Statistical analysis of the smoothing for the square-root Lasso and the multivariate square-root Lasso
- Guidelines on how to set the smoothing parameter  $\underline{\sigma}$
- New estimator to handle correlated noise and repetitions in Multi-Task
- Improved support identification
- Future work: study the influence of <u>σ</u> from the optimization point of view, see if the analysis holds for the LAD-Lasso

## Merci!

"All models are wrong but some come with good open source implementation and good documentation to use these."

A. Gramfort

Python code online https://github.com/QB3/CLaR

Papers: arXiv<sup>(18), (19)(20)</sup>

<sup>(18)</sup> M. Massias et al. "Generalized Concomitant Multi-Task Lasso for Sparse Multimodal Regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

 $<sup>^{(19)}</sup>$  Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: NIPS (2019).

<sup>&</sup>lt;sup>(20)</sup>M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal estimation". In: AISTATS (2020).

## Competitors

• (smoothed) 
$$\ell_{2,1}$$
-MLE

$$(\hat{\mathbf{B}}, \hat{\boldsymbol{\Sigma}}) \in \underset{\substack{\mathbf{B} \in \mathbb{R}^{p \times q} \\ \boldsymbol{\Sigma} \succeq \underline{\sigma}^2/r^2}}{\operatorname{arg\,min}} \left\| \bar{Y} - X\mathbf{B} \right\|_{\boldsymbol{\Sigma}^{-1}}^2 - \log \det(\boldsymbol{\Sigma}^{-1}) + \lambda \left\| \mathbf{B} \right\|_{2,1} ,$$

• and its repetitions version ( $\ell_{2,1}$ -MLER):

$$(\hat{\mathbf{B}}, \hat{\boldsymbol{\Sigma}}) \in \underset{\substack{\mathbf{B} \in \mathbb{R}^{p \times q} \\ \boldsymbol{\Sigma} \succeq \boldsymbol{\sigma}^2}}{\operatorname{arg\,min}} \sum_{1}^{r} \left\| \boldsymbol{Y}^{(l)} - \boldsymbol{X} \mathbf{B} \right\|_{\boldsymbol{\Sigma}^{-1}}^{2} - \log \det(\boldsymbol{\Sigma}^{-1}) + \lambda \left\| \mathbf{B} \right\|_{2,1} .$$

▶  $\ell_{2,1}$ -MLE and  $\ell_{2,1}$ -MLER are bi-convex but not jointly convex

## Smoothing of matrix norm

#### Huber-like formula for the Frobenius norm

$$\begin{aligned} \|\cdot\|_{F} \Box \underline{\sigma} \omega \left(\frac{\cdot}{\underline{\sigma}}\right) (Z) &= \begin{cases} \frac{\|Z\|_{F}^{2}}{2\underline{\sigma}} + \frac{\underline{\sigma}}{2}, & \text{if } \|Z\|_{F} \leq \underline{\sigma} \\ \|Z\|_{F}, & \text{if } \|Z\|_{F} > \underline{\sigma} \end{cases} \\ &= \min_{\underline{\sigma} \geq \underline{\sigma}} \left(\frac{\|Z\|_{F}^{2}}{2\sigma} + \frac{\sigma}{2}\right) \end{aligned}$$

What about other norms ?

## Smoothing of matrix norm

Huber-like formula for the Frobenius norm

$$\|\cdot\|_{F} \Box \underline{\sigma} \omega \left(\frac{\cdot}{\underline{\sigma}}\right) (Z) = \begin{cases} \frac{\|Z\|_{F}^{2}}{2\underline{\sigma}} + \frac{\sigma}{2}, & \text{if } \|Z\|_{F} \leq \underline{\sigma} \\ \|Z\|_{F}, & \text{if } \|Z\|_{F} > \underline{\sigma} \\ \end{cases} \\ = \min_{\underline{\sigma} \geq \underline{\sigma}} \left(\frac{\|Z\|_{F}^{2}}{2\sigma} + \frac{\sigma}{2}\right) \end{cases}$$

#### What about other norms ?

Huber-like formula for the nuclear/trace norm

$$\|\cdot\|_{s,1} \Box \omega_{\underline{\sigma}}(Z) = \begin{cases} \frac{\|Z\|_F^2}{2\underline{\sigma}} + \frac{\underline{\sigma}}{2}n \wedge q, & \text{if } \|Z\|_{\infty} \leq \underline{\sigma} \\ \frac{1}{2\underline{\sigma}} \sum_i \gamma_i^2 - (\gamma_i \wedge \underline{\sigma} - \gamma_i)^2, & \text{if } \|Z\|_{\infty} > \underline{\sigma} \\ = \min_{S \succeq \underline{\sigma}} \frac{1}{2} \|Z\|_{S^{-1}}^2 + \frac{1}{2} \operatorname{Tr}(S) \end{cases}$$

 $\gamma_i$ : singular values of Z $\|Z\|_{S^{-1}}^2 := \operatorname{Tr}(Z^{ op}S^{-1}Z)$  Mahalanobis distance

## Smoothing of matrix norm

Huber-like formula for the Frobenius norm

$$\|\cdot\|_{F} \Box \underline{\sigma} \omega \left(\frac{\cdot}{\underline{\sigma}}\right) (Z) = \begin{cases} \frac{\|Z\|_{F}^{2}}{2\underline{\sigma}} + \frac{\underline{\sigma}}{2}, & \text{if } \|Z\|_{F} \leq \underline{\sigma} \\ \|Z\|_{F}, & \text{if } \|Z\|_{F} > \underline{\sigma} \\ \end{cases} \\ = \min_{\underline{\sigma} \geq \underline{\sigma}} \left(\frac{\|Z\|_{F}^{2}}{2\sigma} + \frac{\sigma}{2}\right) \end{cases}$$

#### What about other norms ? Huber-like formula for the nuclear/trace norm

$$\|\cdot\|_{s,1} \square \omega_{\underline{\sigma}}(Z) = \begin{cases} \frac{\|Z\|_F^2}{2\underline{\sigma}} + \frac{\underline{\sigma}}{2}n \wedge q, & \text{if } \|Z\|_{\infty} \leq \underline{\sigma}\\ \frac{1}{2\underline{\sigma}} \sum_i \gamma_i^2 - (\gamma_i \wedge \underline{\sigma} - \gamma_i)^2, & \text{if } \|Z\|_{\infty} > \underline{\sigma}\\ = \min_{S \succeq \underline{\sigma}} \frac{1}{2} \|Z\|_{S^{-1}}^2 + \frac{1}{2} \operatorname{Tr}(S) \end{cases}$$

 $\gamma_i:$  singular values of Z  $\|Z\|_{S^{-1}}^2:=\mathrm{Tr}(Z^\top S^{-1}Z)$  Mahalanobis distance