The smoothed multivariate square-root Lasso: an optimization lens on concomitant estimation

Joseph Salmon

http://josephsalmon.eu IMAG, Univ. Montpellier, CNRS

Series of works with: Quentin Bertrand (INRIA) Mathurin Massias (University of Genova) Olivier Fercoq (Institut Polytechnique de Paris) Alexandre Gramfort (INRIA)

Table of Contents

Neuroimaging The M/EEG problem

Estimation procedures

Sparsity and Multi-task approaches Smoothing interpretation of concomitant and $\sqrt{\rm Lasso}$ Optimization algorithm

The M/EEG inverse problem

observe magnetoelectric field outside the scalp (100 sensors)
 reconstruct cerebral activity inside the brain (10,000 locations)



 $n \ll p$: ill-posed problem

▶ Motivation: identify brain regions responsible for the signals

▶ Applications: epilepsy treatment, brain aging, anesthesia risks

M/EEG inverse problem for brain imaging

sensors: electric and magnetic fields during a cognitive task



MEG elements: magnometers and gradiometers







Device

Sensors

Detail of a sensor

M/EEG = MEG + EEG



Photo Credit: Stephen Whitmarsh

Table of Contents

Neuroimaging The M/EEG problem Stastistical model

Estimation procedures

Sparsity and Multi-task approaches Smoothing interpretation of concomitant and $\sqrt{\text{Lasso}}$ Optimization algorithm

Source modeling



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

Design matrix - Forward operator



Mathematical model: linear regression



Experiments repeated r **times**



M/EEG specifity #1: combined measurements







Device



Sensor detail

Structure of Y and X: $\begin{pmatrix}
X_{\text{EEG}} \\
X_{\text{grad}} \\
\vdots \\
X_{\text{mag}}
\end{pmatrix}
\begin{pmatrix}
Y_{\text{EEG}} \\
Y_{\text{grad}} \\
\vdots \\
Y_{\text{mag}}
\end{pmatrix}$

Sensor types & noise structure









M/EEG specificity #2: averaging repetitions of experiment



M/EEG specificity #2: averaging repetitions of experiment



M/EEG specificity #2: averaged signals

Averaging 5 repetitions (EEG only)



Limit on the repetitions: subject/patient fatigue

A multi-task framework

Multi-task regression notation:

- n observations (number of sensors)
- ▶ T tasks (temporal information)
- ▶ p features (spatial description)
- \blacktriangleright *r* number of repetitions for the experiment
- $Y^{(1)}, \ldots, Y^{(r)} \in \mathbb{R}^{n \times T}$ observation matrices; $\bar{Y} = \frac{1}{r} \sum_{l} Y^{(l)}$
- $X \in \mathbb{R}^{n \times p}$ forward matrix

$$Y^{(l)} = X \mathbf{B}^* + S_* \mathbf{E}^{(l)}, \quad \text{where}$$

B* ∈ ℝ^{p×T} : true source activity matrix (unknown)
 S_{*} ∈ Sⁿ₊₊ co-standard deviation matrix⁽¹⁾ (unknown)
 E⁽¹⁾,...,E^(r) ∈ ℝ^{n×T} : white noise (standard Gaussian)

Table of Contents

Neuroimaging The M/EEG problem Stastistical model

Estimation procedures

Sparsity and Multi-task approaches

Smoothing interpretation of concomitant and $\sqrt{ ext{Lasso}}$

Signals can often be represented combining few atoms/features:

Fourier decomposition for sounds



⁽²⁾I. Daubechies. Ten lectures on wavelets. SIAM, 1992.

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- Wavelets for images (1990's)⁽²⁾



⁽²⁾I. Daubechies. Ten lectures on wavelets. SIAM, 1992.

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- Wavelets for images (1990's)⁽²⁾
- Dictionary learning for images (2000's)⁽³⁾



⁽²⁾I. Daubechies. Ten lectures on wavelets. SIAM, 1992.

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- ▶ Wavelets for images (1990's)⁽²⁾
- Dictionary learning for images (2000's)⁽³⁾
- Neuroimaging: measurements assumed to be explained by a few active brain sources



⁽²⁾I. Daubechies. Ten lectures on wavelets. SIAM, 1992.

Justification for dipolarity assumption

Sparsity holds: dipolar patterns equivalent to focal sources

- short duration
- simple cognitive task
- repetitions of experiment average out other sources
- ICA recovers dipolar patterns,⁽⁴⁾ well modeled by focal sources:



⁽⁴⁾A. Delorme et al. "Independent EEG sources are dipolar". In: *PloS one* 7.2 (2012), e30135.

(Structured) Sparsity inducing penalties⁽⁵⁾

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



Sparse support: no structure X

Lasso penalty

$$\|\mathbf{B}\|_1 \triangleq \sum_{j=1}^p \sum_{t=1}^T |\mathbf{B}_{jt}|$$

⁽⁵⁾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.

(Structured) Sparsity inducing penalties⁽⁵⁾

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



Sparse support: group structure 🗸

Group-Lasso penalty

$$\|\mathbf{B}\|_{2,1} \triangleq \sum_{j=1}^p \|\mathbf{B}_{j:}\|_2$$

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

⁽⁵⁾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.

Data-fitting term and experiment repetitions

• Classical estimator: use averaged⁽⁶⁾ signal \bar{Y}

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

How to take advantage of the number of repetitions? Intuitive estimator:

$$\hat{\mathbf{B}}^{\mathsf{repet}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nTr} \sum_{l=1}^{r} \left\| Y^{(l)} - X\mathbf{B} \right\|_{F}^{2} + \lambda \|\mathbf{B}\|_{2,1} \right)$$

⁽⁶⁾& whitened, say using baseline data

Data-fitting term and experiment repetitions

• Classical estimator: use averaged⁽⁶⁾ signal \bar{Y}

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

How to take advantage of the number of repetitions? Intuitive estimator:

$$\hat{\mathbf{B}}^{\mathsf{repet}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nTr} \sum_{l=1}^{r} \left\| Y^{(l)} - X\mathbf{B} \right\|_{F}^{2} + \lambda \|\mathbf{B}\|_{2,1} \right)$$

▶ Fail: $\hat{B}^{\text{repet}} = \hat{B}$ (because of datafit $\|\cdot\|_F^2$)

Data-fitting term and experiment repetitions

• Classical estimator: use averaged⁽⁶⁾ signal \bar{Y}

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

How to take advantage of the number of repetitions? Intuitive estimator:

$$\hat{\mathbf{B}}^{\mathsf{repet}} \in \operatorname*{arg\,min}_{\mathbf{B} \in \mathbb{R}^{p \times T}} \left(\frac{1}{2nTr} \sum_{l=1}^{r} \left\| Y^{(l)} - X\mathbf{B} \right\|_{F}^{2} + \lambda \|\mathbf{B}\|_{2,1} \right)$$

▶ Fail: $\hat{B}^{\text{repet}} = \hat{B}$ (because of datafit $\|\cdot\|_F^2$)

 \hookrightarrow investigate other datafits

⁽⁶⁾& whitened, say using baseline data

Table of Contents

Neuroimaging The M/EEG problem Stastistical model

Estimation procedures

Sparsity and Multi-task approaches Smoothing interpretation of concomitant and $\sqrt{\text{Lasso}}$ Optimization algorithm

Lasso^{(7), (8)}: the "modern least-squares"⁽⁹⁾

$$\hat{\beta} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \frac{1}{2n} \left\| y - X\beta \right\|^2 + \lambda \left\| \beta \right\|_1$$

• $y \in \mathbb{R}^n$: observations

• $X \in \mathbb{R}^{n \times p}$: design matrix

• sparsity: for λ large enough, $\|\hat{\beta}\|_0 \ll p$

⁽⁷⁾ R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.

⁽⁸⁾S. S. Chen and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.

⁽⁹⁾E. J. Candès, M. B. Wakin, and S. P. Boyd. "Enhancing Sparsity by Reweighted l₁ Minimization". In: J. Fourier Anal. Applicat. 14.5-6 (2008), pp. 877–905.

Lasso and optimal $\lambda^{(10),(11)}$

For $y = X\beta^* + \sigma_*\varepsilon$, $\varepsilon \sim \mathcal{N}(0, \mathrm{Id}_n)$ and X satisfying the "Restricted Eigenvalue" property, if $\lambda = 2\sigma_*\sqrt{\frac{2\log(p/\delta)}{n}}$, then

$$\frac{1}{n} \left\| X\beta^* - X\hat{\beta} \right\|^2 \le \frac{18}{\kappa_{s^*}^2} \frac{\sigma_*^2 s^*}{n} \log\left(\frac{p}{\delta}\right)$$

with probability $1-\delta,$ where $\hat{\beta}$ is a Lasso solution

<u>Rem</u>: optimal rate in the minimax sense (up to constant/log term)

BUT σ_* is unknown in practice !

⁽¹⁰⁾ 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.

⁽¹¹⁾ A. S. Dalalyan, M. Hebiri, and J. Lederer. "On the Prediction Performance of the Lasso". In: *Bernoulli* 23.1 (2017), pp. 552–581.

Other datafit: the $\sqrt{Lasso}^{(12)}$

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

optimal $\lambda \propto \sigma_*$

Confirmed in practice:



(12) 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.

Other datafit: the $\sqrt{Lasso}^{(12)}$

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

optimal λ adaptive to σ_*

Confirmed in practice:



(12) 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.

Unhappy optimizer

 $\sqrt{\text{Lasso}}$: non-smooth+non-smooth \hookrightarrow use *Concomitant Lasso*⁽¹³⁾:

$$(\hat{\beta}, \hat{\sigma}) \in \underset{\beta \in \mathbb{R}^{p}, \sigma > 0}{\operatorname{arg\,min}} \quad \frac{\|y - X\beta\|^{2}}{2n\sigma} + \frac{\sigma}{2} + \lambda \|\beta\|_{1}$$

same solutions when $||y - X\hat{\beta}_{\sqrt{\text{Lasso}}}|| \neq 0$, but jointly convex, non smooth + separable: solvable by alternate min.⁽¹⁴⁾ in β and σ



(13) A. B. Owen. "A robust hybrid of lasso and ridge regression". In: Contemporary Mathematics 443 (2007), pp. 59–72.

⁽¹⁴⁾T. Sun and C.-H. Zhang. "Scaled sparse linear regression". In: *Biometrika* 99.4 (2012), pp. 879–898.

Unhappy optimizer

 $\sqrt{\text{Lasso}}$: non-smooth+non-smooth \hookrightarrow use *Concomitant Lasso*⁽¹³⁾:

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

same solutions when $||y - X\hat{\beta}_{\sqrt{\text{Lasso}}}|| \neq 0$, but jointly convex, smooth + separable: solvable by alternate min.⁽¹⁴⁾ in β and σ



(13) A. B. Owen. "A robust hybrid of lasso and ridge regression". In: Contemporary Mathematics 443 (2007), pp. 59–72.

⁽¹⁴⁾T. Sun and C.-H. Zhang. "Scaled sparse linear regression". In: Biometrika 99.4 (2012), pp. 879-898.

"Concomitant": smoothing the $\sqrt{\text{Lasso}}^{(17)}$

"Huberization":
replace
$$\frac{\|\cdot\|}{\sqrt{n}}$$
 by a smooth approximation
huber $_{\underline{\sigma}}(z) = \begin{cases} \frac{\|z\|^2}{2n\underline{\sigma}} + \frac{\underline{\sigma}}{2} & \text{if } \frac{\|z\|}{\sqrt{n}} \leq \underline{\sigma} \\ \frac{\|z\|}{\sqrt{n}} & \text{if } \frac{\|z\|}{\sqrt{n}} > \underline{\sigma} \end{cases}$
 $= \min_{\underline{\sigma} \geq \underline{\sigma}} \left(\frac{\|z\|^2}{2n\sigma} + \frac{\sigma}{2} \right) = \frac{1}{\sqrt{n}} \|\cdot\| \Box \left(\frac{1}{2n\underline{\sigma}} \|\cdot\|^2 + \frac{\underline{\sigma}}{2} \right)(z)$

Leads to the Smoothed^{(15),(16)} Concomitant Lasso formulation:

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

⁽¹⁵⁾A. Beck and M. Teboulle. "Smoothing and first order methods: A unified framework". In: SIAM J. Optim. 22.2 (2012), pp. 557–580.

⁽¹⁶⁾Y. Nesterov. "Smooth minimization of non-smooth functions". In: *M. Prog.* 103.1 (2005), pp. 127–152.

(17) E. Ndiaye et al. "Efficient Smoothed Concomitant Lasso Estimation for High Dimensional Regression". In: Journal of Physics: Conference Series 904.1 (2017), p. 012006.

Smoothing aparté^{(18),(19)}

Smoothing: for $\underline{\sigma} > 0$, a "smoothed" version of f is $f_{\underline{\sigma}}$

$$f_{\underline{\sigma}} = \underline{\sigma}\omega\left(\frac{\cdot}{\underline{\sigma}}\right)\Box f$$
, where $f\Box g(x) = \inf_{u}\{f(u) + g(x-u)\}$

• ω is a predefined smooth function (s.t. $\nabla \omega$ is Lipschitz)

	Fourier: $\mathcal{F}(f)$	Fenchel/Legendre: f^*
	convolution: *	inf-convolution:
Kernel smoothing analogy:	$\mathcal{F}(f \star g) = \mathcal{F}(f) \cdot \mathcal{F}(g)$	$(f \Box g)^* = f^* + g^*$
	$Gaussian:\mathcal{F}(g)=g$	$\omega = \frac{\ \cdot\ ^2}{2}: \omega^* = \omega$
	$f_h = \frac{1}{h}g\left(\frac{\cdot}{h}\right) \star f$	$f_{\underline{\sigma}} = \underline{\sigma} \omega \left(\frac{\cdot}{\underline{\sigma}} \right) \Box f$

(18) Y. Nesterov. "Smooth minimization of non-smooth functions". In: Math. Program. 103.1 (2005), pp. 127–152.

(19) A. Beck and M. Teboulle. "Smoothing and first order methods: A unified framework". In: SIAM J. Optim. 22.2 (2012), pp. 557–580.

















Smoothing other norms

Smoothing Frobenius norm yields a trivial gen. of conco Lasso

More interesting: S. van de Geer introduced the pivotal multivariate √Lasso,⁽²⁰⁾ using trace/nuclear norm for data-fitting

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

hard to solve, statistical analysis makes stringent assumptions

Smoothing the datafit makes optim. and stats easier!

⁽²⁰⁾S. van de Geer. Estimation and testing under sparsity. École d'Été de Probabilités de Saint-Flour. 2016.

Smoothing the nuclear norm⁽²¹⁾

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 γ_i 's are the singular values of Z

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

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

 $^{^{(21)}\}mathsf{Q}.$ Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: NeurIPS. 2019.

Smoothing of the multivariate \sqrt{Lasso}

Smoothed Generalized Concomitant Lasso (SGCL)⁽²²⁾:

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

Concomitant Lasso with Repetitions (CLaR)⁽²³⁾:

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

(22) M. Massias et al. "Generalized concomitant multi-task Lasso for sparse multimodal regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

⁽²³⁾Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NeurIPS*. 2019.

Simulations : row support identification

▶
$$n = 150$$
, $p = 500$, $T = 100$

► X Toeplitz-correlated

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



Table of Contents

Neuroimaging The M/EEG problem Stastistical model

Estimation procedures

Sparsity and Multi-task approaches Smoothing interpretation of concomitant and $\sqrt{\rm Lasso}$ Optimization algorithm

SGCL and CLaR: alternate updates

Alternate minimization converges

 \underline{B} update (S fixed): standard Multi-task Lasso optimization, off-the-shelf techniques and lots of refinements

S update (B fixed):

$$\underset{S \succeq \underline{\sigma}}{\operatorname{arg\,min}} \left(\frac{1}{2n} \operatorname{Tr}[Z^{\top} S^{-1} Z] + \frac{1}{2n} \operatorname{Tr}(S) \right)$$

closed-form solution : clipped sqrt of eigen value decomposition of $\frac{1}{T}(\bar{Y} - XB)(\bar{Y} - XB)^{\top} \text{ or } \frac{1}{rT}\sum_{l=1}^{r}(Y^{(l)} - XB)(Y^{(l)} - XB)^{\top}$

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

Algorithm: Concomitant Lasso w. Repetitions (CLaR)

input : $X \in \mathbb{R}^{n \times p}, Y^{(1)}, \dots, Y^{(r)} \in \mathbb{R}^{n \times T}, \sigma > 0, \lambda > 0$ init : $B = 0_{p,q}, R = Y$ for iter = $1, \ldots, do$ $S \leftarrow \mathsf{SpectralClipping}(\frac{1}{T_r}\sum_{l}^r (Y^{(l)} - XB)(Y^{(l)} - XB)^\top, \sigma)$ // closed-form sol. of min. in S: EVD + clipping sqrt of eigenvalues at level σ for j = 1, ..., p do $L_{i} = X_{i}^{\top} S^{-1} X_{i}$ // Lipschitz constants for j = 1, ..., p do $R \leftarrow R + X_{:i}B_{i:i}$ // partial residual update $\mathbf{B}_{j:} \leftarrow \mathrm{BST}\left(X_{:j}^{\top}S^{-1}R/L_j, \lambda nT/L_j\right)$ // coef. update $R \leftarrow R - X_{i}B_{i}$ // residual update return B.S.

 $\label{eq:complexity} \begin{array}{l} \underline{\text{Complexity?}} \\ \overline{\text{Fine, if we store } S^{-1}X \text{, and } S^{-1}R \text{ instead of } R. \\ \\ \text{Need eigenvalue decomposition though } \mathcal{O}(n^3) \text{ (here } n \approx 100) \end{array}$

Statistical properties for i.i.d. case⁽²⁴⁾



- i.i.d. Gaussian noise
- X satisfying the "mutual incoherence" property
 λ ∝ √log p/T√n (independent of σ*)
 c₁ σ ≤ σ* ≤ c₂ σ

 \implies with probability at least $1 - ne^{-cT/n}$

$$\frac{1}{T} \|\mathbf{B}^* - \hat{\mathbf{B}}\|_{2,\infty} \le C\boldsymbol{\sigma}_* \frac{1}{T} \sqrt{\frac{\log p}{n}}$$

(²⁴)M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal regression". In: AISTATS. 2020.

Real data experiments



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

Links

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

A. Gramfort

▶ Papers: arXiv / personal webpage^{(25), (26), (27)}

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

⁽²⁵⁾ M. Massias et al. "Generalized concomitant multi-task Lasso for sparse multimodal regression". In: AISTATS. vol. 84. 2018, pp. 998–1007.

 $^{^{(26)}}$ Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NeurIPS*. 2019.

⁽²⁷⁾M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal regression". In: AISTATS. 2020.

References I

- Beck, A. and M. Teboulle. "Smoothing and first order methods: A unified framework". In: SIAM J. Optim. 22.2 (2012), pp. 557–580.
- Belloni, A., V. Chernozhukov, and L. Wang. "Square-root Lasso: pivotal recovery of sparse signals via conic programming". In: *Biometrika* 98.4 (2011), pp. 791–806.
- Bertrand, Q. et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NeurIPS*. 2019.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov. "Simultaneous analysis of Lasso and Dantzig selector". In: Ann. Statist. 37.4 (2009), pp. 1705–1732.
- Candès, E. J., M. B. Wakin, and S. P. Boyd. "Enhancing Sparsity by Reweighted l₁ Minimization". In: *J. Fourier Anal. Applicat.* 14.5-6 (2008), pp. 877–905.

References II

- Chen, S. S. and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.
- Dalalyan, A. S., M. Hebiri, and J. Lederer. "On the Prediction Performance of the Lasso". In: *Bernoulli* 23.1 (2017), pp. 552–581.
- Daubechies, I. Ten lectures on wavelets. SIAM, 1992.
- Delorme, A. et al. "Independent EEG sources are dipolar". In: *PloS* one 7.2 (2012), e30135.
- Massias, M. et al. "Generalized concomitant multi-task Lasso for sparse multimodal regression". In: AISTATS. Vol. 84. 2018, pp. 998–1007.
- Massias, M. et al. "Support recovery and sup-norm convergence rates for sparse pivotal regression". In: AISTATS. 2020.
- Ndiaye, E. et al. "Efficient Smoothed Concomitant Lasso Estimation for High Dimensional Regression". In: *Journal of Physics: Conference Series* 904.1 (2017), p. 012006.

References III

- Nesterov, Y. "Smooth minimization of non-smooth functions". In: M. Prog. 103.1 (2005), pp. 127–152.
 - ."Smooth minimization of non-smooth functions". In: Math. Program. 103.1 (2005), pp. 127–152.
- Obozinski, G., 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.
- Olshausen, B. A. and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: Vision research (1997).
- Owen, A. B. "A robust hybrid of lasso and ridge regression". In: Contemporary Mathematics 443 (2007), pp. 59–72.
- Sun, T. and C.-H. Zhang. "Scaled sparse linear regression". In: Biometrika 99.4 (2012), pp. 879–898.

References IV

- Tibshirani, R. "Regression Shrinkage and Selection via the Lasso". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.
- van de Geer, S. Estimation and testing under sparsity. École d'Été de Probabilités de Saint-Flour. 2016.

Statistical assumptions

<u>Gaussian noise</u>: the entries $E_{i,j}$ are i.i.d. $\mathcal{N}(0, \sigma_*^2)$ random variables.

<u>Mutual incoherence</u>: The Gram matrix $\Psi \triangleq \frac{1}{n}X^{\top}X$ satisfies

$$\Psi_{jj}=1$$
 , and $\max_{j'
eq j} \left|\Psi_{jj'}
ight| \leq rac{1}{7lpha s},\, orall j\in [p]$,

for some integer $s \ge 1$ and some constant $\alpha > 1$.

<u>Residuals bound</u>: For the multivariate square-root Lasso, $\hat{E}^{\top}\hat{E}$ is invertible, and there exists η such that

$$\|(\frac{1}{T}\hat{\mathbf{E}}^{\top}\hat{\mathbf{E}})^{\frac{1}{2}}\|_{2} \le C\sigma^{*}$$

 $\frac{\text{Smoothing parameter value: } \underline{\sigma}, \ \bar{\sigma} \ \text{and} \ \eta \ \text{verify: } \underline{\sigma} \leq \frac{\sigma^*}{\sqrt{2}} \ \text{and} \ \bar{\sigma} = (2+\eta)\sigma^* \ \text{with} \ \eta \geq 1.$

Competitors

• (smoothed) $\ell_{2,1}$ -MLE

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

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

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

 $\underline{\text{Rem}}: \ell_{2,1}\text{-}\text{MLE}$ and $\ell_{2,1}\text{-}\text{MLER}$ are bi-convex but not jointly convex

• MRCER has an additional term $\mu \|\Sigma^{-1}\|$ w.r.t. $\ell_{2,1}$ -MLER