The smoothed multivariate square-root Lasso (optimizational and statistical handling of correlated noise)

Mathurin Massias

https://mathurinm.github.io

Series of works with: Quentin Bertrand (INRIA) Olivier Fercoq (Institut Polytechnique de Paris) Alexandre Gramfort (INRIA) Joseph Salmon (IMAG, Univ. Montpellier, CNRS)

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

Signals can often be represented combining few atoms/features:

Fourier decomposition for sounds



⁽¹⁾I. Daubechies. Ten lectures on wavelets. SIAM, 1992.

 $^{(2)}B.$ A. Olshausen and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: Vision research (1997).

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- Wavelets for images (1990's)⁽¹⁾



 (1) I. Daubechies. Ten lectures on wavelets. SIAM, 1992.
 (2) B. A. Olshausen and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: Vision research (1997).

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- ▶ Wavelets for images (1990's)⁽¹⁾
- Dictionary learning for images (2000's)⁽²⁾



 (1) I. Daubechies. Ten lectures on wavelets. SIAM, 1992.
 (2) B. A. Olshausen and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: Vision research (1997).

Signals can often be represented combining few atoms/features:

- Fourier decomposition for sounds
- ▶ Wavelets for images (1990's)⁽¹⁾
- Dictionary learning for images (2000's)⁽²⁾
- Here we assume that measurements are explained by a few active brain sources



 (1) L. Daubechies. Ten lectures on wavelets. SIAM, 1992.
 (2) B. A. Olshausen and D. J. Field. "Sparse coding with an overcomplete basis set: A strategy employed by V1?" In: Vision research (1997).

Justification for dipolarity assumption

- short duration
- simple cognitive task
- repetitions of experiment average out other sources
- ICA recovers dipolar patterns,⁽³⁾ well modelled by focal sources:



⁽³⁾A. Delorme et al. "Independent EEG sources are dipolar". In: *PloS one* 7.2 (2012), e30135.

Mathematical model: linear regression



Lasso^{(4), (5)}: the "modern least-squares"⁽⁶⁾

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

- $y \in \mathbb{R}^n$: observations
- $X \in \mathbb{R}^{n \times p}$: design matrix
- sparsity: for λ large enough, $\|\hat{\beta}\|_0 \ll p$

⁽⁴⁾ R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.

⁽⁵⁾S. S. Chen and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.

⁽⁶⁾E. J. Candès, M. B. Wakin, and S. P. Boyd. "Enhancing Sparsity by Reweighted *l*₁ Minimization". In: *J. Fourier Anal. Applicat.* 14.5-6 (2008), pp. 877–905.

Sparsity inducing penalties⁽⁷⁾

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



Sparse support: no structure X

Lasso penalty

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

⁽⁷⁾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.

Sparsity inducing penalties⁽⁷⁾

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



Sparse support: group structure 🗸

Group-Lasso penalty

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

where $B_{j:} = j$ -th row of B

⁽⁷⁾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.

M/EEG specifity #1: combined measurements







Device

Sensors

Sensor detail

Structure of Y and X:

$$\begin{pmatrix} Y_{\text{EEG}} \\ \vdots \\ Y_{\text{grad}} \\ \vdots \\ Y_{\text{mag}} \end{pmatrix} \qquad \begin{pmatrix} X_{\text{EEG}} \\ \vdots \\ Z_{\text{grad}} \\ \vdots \\ Z_{\text{mag}} \end{pmatrix}$$



9 / 23

M/EEG specificity #2: averaging repetitions of experiment



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; $\overline{Y} = \frac{1}{r} \sum_{l} Y^{(l)}$
- $X \in \mathbb{R}^{n imes p}$ forward matrix

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

•
$$\mathbf{B}^* \in \mathbb{R}^{p \times T}$$
 : true source activity matrix (unknown)

S_{*} ∈ Sⁿ₊₊ co-standard deviation matrix⁽⁸⁾ (unknown)
 E⁽¹⁾,...,E^(r) ∈ ℝ^{n×T} : white Gaussian noise

Data-fitting term

• Classical estimator: use averaged $^{(9)}$ signal \bar{Y}

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

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

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

Data-fitting term

• Classical estimator: use averaged $^{(9)}$ signal \bar{Y}

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

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

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

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

Data-fitting term

• Classical estimator: use averaged $^{(9)}$ signal \bar{Y}

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

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

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

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

 \hookrightarrow investigate other datafits

⁽⁹⁾& whitened, say using baseline data

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

Theorem

For $y = X\beta^* + \sigma_*\varepsilon$, ε standard Gaussian 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 $\sqrt{Lasso}^{(12)}$

Same guarantees, but optimal λ independent of σ_* (*pivotality*):

$$\widehat{\beta}_{\sqrt{\text{Lasso}}} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \left(\frac{1}{\sqrt{n}} \left\| y - X\beta \right\| + \lambda \left\| \beta \right\|_1 \right)$$

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 $\sqrt{Lasso}^{(12)}$

Same guarantees, but optimal λ independent of σ_* (*pivotality*):

$$\widehat{\beta}_{\sqrt{\text{Lasso}}} \in \operatorname*{arg\,min}_{\beta \in \mathbb{R}^p} \left(\frac{1}{\sqrt{n}} \left\| y - X\beta \right\| + \lambda \left\| \beta \right\|_1 \right)$$

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

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

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

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





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

Unhappy optimizer

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

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

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





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

Concomitant origin: smoothing the $\sqrt{\mathrm{Lasso}}^{(16)}$

"Huberization": replace $\|\cdot\|$ by a smooth approximation:

$$\begin{aligned} \mathsf{huber}_{\underline{\sigma}}(\|z\|) &= \begin{cases} \frac{\|z\|^2}{2\underline{\sigma}} + \frac{\underline{\sigma}}{2} & \text{if } \|z\| \leq \underline{\sigma} \\ \|z\| & \text{if } \|z\| > \underline{\sigma} \\ &= \min_{\underline{\sigma} \geq \underline{\sigma}} \left(\frac{\|z\|^2}{2\sigma} + \frac{\sigma}{2} \right) = \|\cdot\| \square(\frac{1}{2\underline{\sigma}}\|\cdot\|^2 + \frac{\underline{\sigma}}{2}) \end{aligned}$$

Leads to the Smoothed (14), (15) Concomitant Lasso formulation:

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

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

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

⁽¹⁶⁾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 other norms

Smoothing Frobenius norm yields a trivial gen. of conco Lasso

More interesting: S. van de Geer introduced the pivotal multivariate $\sqrt{Lasso}^{(17)}$:

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

hard to solve, statistical analysis makes stringent assumptions

Smoothing the datafit makes optim. and stats easier!

⁽¹⁷⁾S. van de Geer. Estimation and testing under sparsity. École d'Été de Probabilités de Saint-Flour. 2016.

Smoothing the nuclear norm⁽¹⁸⁾

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

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

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

 $^{^{(18)}}$ 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)⁽¹⁹⁾:

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

Concomitant Lasso with Repetitions (CLaR)⁽²⁰⁾:

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

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

⁽²⁰⁾Q. Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: *NeurIPS*. 2019.

SGCL and CLaR: alternate updates

Alternate minimization converges

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

S update (B fixed):

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

closed-form solution involving clipped EVD of:

$$\frac{1}{T}(\bar{Y} - XB)(\bar{Y} - XB)^{\top} \text{ or } \frac{1}{rT} \sum_{l=1}^{r} (Y^{(l)} - XB)(Y^{(l)} - XB)^{\top}$$

Support recovery guarantees⁽²¹⁾

For the multivariate $\sqrt{\text{Lasso}/\text{its}}$ smoothed version, with Gaussian noise, mutual incoherence (α) and bounded residuals/correct value of smoothing parameter (η).

Let $C = (1 + \frac{16}{7(\alpha - 1)})$, $A \ge \sqrt{2}$, and $\lambda = \frac{2\sqrt{2}}{\sqrt{nq}}(1 + A\sqrt{(\log p)/q})$. There exists $c \ge 1/64$ s.t. with proba $\ge 1 - p^{1 - A^2/2} - 2ne^{-cq/n}$,

$$\frac{1}{q} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,\infty} \le C(3+\eta)\lambda\sigma^*$$

Moreover if

$$\min_{j \in \mathcal{S}^*} \frac{1}{q} \| \mathbf{B}_{j:}^* \|_2 > 2C(3+\eta)\lambda\sigma^*$$

then with the same probability:

$$\mathcal{S}^* = \hat{\mathcal{S}} \triangleq \{ j \in [p] : \frac{1}{q} \| \hat{\mathbf{B}}_{j:} \|_2 > C(3+\eta)\lambda\sigma^* \}$$

⁽²¹⁾ M. Massias et al. "Support recovery and sup-norm convergence rates for sparse pivotal regression". In: AISTATS. 2020.

Real data experiments



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

Figure: Left auditory stimulations (n = 102, p = 7498, T = 76, r = 63) Sources found in the left and right hemispheres

- 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

- ▶ Papers: arXiv / personal webpage^{(22), (23), (24)}
- Python code online for CLaR https://github.com/QB3/CLaR

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

 $^{^{(23)}\}mathsf{Q}.$ Bertrand et al. "Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso". In: NeurIPS. 2019.

^{(&}lt;sup>24</sup>) 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.
- Tibshirani, R. "Regression Shrinkage and Selection via the Lasso".
 In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.

References IV

 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 of E_1, \ldots, E_n 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 $\max_{j'
eq j} \left| \Psi_{jj'}
ight| \leq rac{1}{7 lpha s}, \, orall j \in [p]$,

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

$$\|(\frac{1}{q}\hat{\mathbf{E}}^{\top}\hat{\mathbf{E}})^{\frac{1}{2}}\|_{2} \le (2+\eta)\sigma^{*}$$

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

Alternate minimization



Complexity? OK if we store $S^{-1}X$, and $S^{-1}R$ instead of R.

Smoothing aparté^{(25), (26)}

<u>Motivation</u>: smooth a non-smooth function f to ease optimization Smoothing: for $\mu > 0$, a "smoothed" version of f is f_{μ}

$$f_{\mu} = \mu \omega \left(\frac{\cdot}{\mu}\right) \Box f$$
, 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:	$\mathcal{F}(f\star g)=\mathcal{F}(f)\cdot\mathcal{F}(g)$	$(f \Box g)^* = f^* + g^*$
	$Gaussian:\mathcal{F}(g)=g$	$\omega = \frac{\ \cdot\ ^2}{2}: \omega^* = \omega$
	$f_h = \frac{1}{h}g\left(\frac{\cdot}{h}\right) \star f$	$f_{\mu} = \mu \omega \left(rac{\cdot}{\mu} ight) \Box f$

 $^{(25)}$ Y. Nesterov. "Smooth minimization of non-smooth functions". In: Math. Program. 103.1 (2005), pp. 127–152.

⁽²⁶⁾A. Beck and M. Teboulle. "Smoothing and first order methods: A unified framework". In: SIAM J. Optim. 22.2 (2012), pp. 557–580.

Competitors

• (smoothed) $\ell_{2,1}$ -MLE

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

► and its repetitions version (ℓ_{2,1}-MLER):

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

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