# Exploiting structure in sparse GLMs for fast & safe support identification

Mathurin Massias (University of Genoa)

Joint work with:

Alexandre Gramfort (INRIA) Joseph Salmon (Université de Montpellier) Samuel Vaiter (CNRS, UMB)

# **Table of Contents**

The Lasso

Exploiting regularity

Sparse GLMs

More solvers speed-up

# The Lasso<sup>1,2</sup>

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

• 
$$y \in \mathbb{R}^n$$
: observations

- $X = [X_1| \dots |X_p] \in \mathbb{R}^{n \times p}$ : design matrix
- sparsity: for  $\lambda$  large enough,  $\|\hat{\beta}\|_0 \ll p$
- popular solver in ML: coordinate descent (used here)

<sup>&</sup>lt;sup>1</sup>R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.

<sup>&</sup>lt;sup>2</sup>S. S. Chen and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.

#### Duality for the Lasso

$$\hat{\theta} = \underset{\theta \in \Delta_X}{\operatorname{arg\,max}} \underbrace{\frac{1}{2} \|y\|^2 - \frac{\lambda^2}{2} \|y/\lambda - \theta\|^2}_{\mathcal{D}(\theta)}$$

 $\Delta_X = \left\{ \theta \in \mathbb{R}^n \, : \, \forall j \in [p], \; |X_j^\top \theta| \leq 1 \right\}: \text{ dual feasible set}$ 



Example: n = 2, p = 3

#### **Primal-dual links**

$$\hat{\theta} = (y - X\hat{\beta})/\lambda$$

$$\mathcal{P}(\beta) \ge \mathcal{P}(\hat{\beta}) = \mathcal{D}(\hat{\theta}) \ge \mathcal{D}(\theta)$$



 $\forall \beta, (\exists \theta \in \Delta_X, \operatorname{dgap}(\beta, \theta) \le \epsilon) \Rightarrow \mathcal{P}(\beta) - \mathcal{P}(\hat{\beta}) \le \epsilon$  $\beta \text{ is an } \epsilon \text{-solution whenever } \operatorname{dgap}(\beta, \theta) \le \epsilon$ 

# Choice of dual point

$$\hat{\theta} = (y - X\hat{\beta})/\lambda$$

Standard approach<sup>3</sup>: at epoch *t*, corresponding to primal  $\beta^{(t)}$  and residuals  $r^{(t)} := y - X\beta^{(t)}$ , take

$$\theta = \theta_{\rm res}^{(t)} := r^{(t)} / \lambda$$

<sup>&</sup>lt;sup>3</sup>J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

# Choice of dual point

$$\hat{\theta} = (y - X\hat{\beta})/\lambda$$

Standard approach<sup>3</sup>: at epoch t, corresponding to primal  $\beta^{(t)}$  and residuals  $r^{(t)} := y - X\beta^{(t)}$ , take

$$\theta = \theta_{\text{res}}^{(t)} := r^{(t)} / \max(\lambda, \|X^{\top} r^{(t)}\|_{\infty})$$

#### residuals rescaling

- converges to  $\hat{\theta}$
- ►  $\mathcal{O}(np)$  to compute (= 1 epoch of CD)  $\hookrightarrow$  rule of thumb: compute  $\theta_{\text{res}}^{(t)}$  and dgap every 10 epochs

<sup>&</sup>lt;sup>3</sup>J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

# Slower convergence of dual

$$\theta_{\text{res}}^{(t)} = r^{(t)} / \max(\lambda, \|X^{\top} r^{(t)}\|_{\infty})$$



$$\lambda_{\max} = \| X^\top y \|_\infty$$
 is the smallest  $\lambda$  giving  $\hat{\beta} = 0$ 

# Slower convergence of dual

Unconstrained dual of Elastic net:



# **Table of Contents**

The Lasso

Exploiting regularity

Sparse GLMs

More solvers speed-up





























# VAR regularity in residuals

#### **Theorem**<sup>4</sup>

Under uniqueness assumption, ISTA/CD achieves sign id.:  $\operatorname{sign} \beta_j^{(t)} = \operatorname{sign} \hat{\beta}_j$ . Then, Lasso residuals are Vector AutoRegressive (VAR):

$$r^{(t+1)} = Ar^{(t)} + b$$

 $\hookrightarrow$  we could fit a VAR to infer  $\lim_{t\to\infty} r^{(t)} = \lambda \hat{\theta}$ 

We do not know when the sign is identified, n points is high Need a cheaper solution  $\hookrightarrow$  extrapolation

<sup>&</sup>lt;sup>4</sup>M. Massias, A. Gramfort, and J. Salmon. "Celer: a fast solver for the Lasso with dual extrapolation". In: *ICML*. 2018, pp. 3321–3330.

# Simple example: extrapolation in 1D

1D autoregressive process:

$$x^{(t)} = ax^{(t-1)} + b \underset{t \to \infty}{\to} x^*$$

we have

$$x^{(t)} - x^* = a(x^{(t-1)} - x^*)$$
$$x^{(t-1)} - x^* = a(x^{(t-2)} - x^*)$$

"Aitken's  $\Delta^2$  ": 2 unknowns, so 2 eqs or 3 points  $x^{(t)}, x^{(t-1)}, x^{(t-2)}$  are enough to find  $x^*!^5$ 

<sup>&</sup>lt;sup>5</sup>A. Aitken. "On Bernoulli's numerical solution of algebraic equations". In: *Proceedings of the Royal Society of Edinburgh* 46 (1926), pp. 289–305.

# Aitken application

$$\lim_{t \to \infty} \sum_{i=0}^{t} \frac{(-1)^i}{2i+1} = \frac{\pi}{4} = 0.785398...$$

t	$\sum_{i=0}^{t} \frac{(-1)^i}{2i+1}$	$\Delta^2$
0	1.0000	_
1	0.66667	-
2	0.86667	<b>0.7</b> 9167
3	<b>0.7</b> 2381	0.78333
4	0.83492	0.78631
5	<b>0.7</b> 4401	0.78492
6	0.82093	0.78568
7	0.75427	0.78522
8	0.81309	0.78552
9	<b>0.7</b> 6046	<b>0.7853</b> 1

# Generalization<sup>6</sup> to VAR $r^{(t)} \in \mathbb{R}^n$

• fix 
$$K = 5$$
 (small)

• keep track of K past residuals  $r^{(t)}, \ldots, r^{(t+1-K)}$ 

► 
$$U^{(t)} = [r^{(t+1-K)} - r^{(t-K)}, \dots, r^{(t)} - r^{(t-1)}] \in \mathbb{R}^{n \times K}$$

► solve 
$$(U^{(t)})^{\top}U^{(t)}z = \mathbf{1}_K$$
  
►  $c = z/z^{\top}\mathbf{1}_K$ 

$$r_{\text{accel}}^{(t)} \triangleq \sum_{k=1}^{K} c_k r^{(t+1-k)}$$

$$\theta_{\text{accel}}^{(t)} \triangleq r_{\text{accel}}^{(t)} / \max(\lambda, \|X^{\top} r_{\text{accel}}^{(t)}\|_{\infty})$$

Cost:  $\mathcal{O}(K^3 + K^2n + np)$ 

<sup>&</sup>lt;sup>6</sup>D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: NIPS. 2016, pp. 712–720.

#### Dual extrapolation for the Lasso



Leukemia dataset: p = 7129, n = 72,  $\lambda = \lambda_{\max}/10$ 

# VAR after sign identification

Wlog, support of  $\hat{\beta}$  :  $\{1, \dots, S\}$  (other coordinates stay at 0) Consider 1 epoch of CD:

$$\beta^{(t)} \to \beta^{(t+1)}$$

Decomposed into non-zero coordinate updates

$$\beta^{(t)} = \tilde{\beta}^{(0)} \xrightarrow{1} \tilde{\beta}^{(1)} \xrightarrow{2} \dots \xrightarrow{S} \tilde{\beta}^{(S)} = \beta^{(t+1)}$$



 $\tilde{\beta}^{(s)} = \tilde{\beta}^{(s-1)}$  except at coordinate s:

$$\begin{split} \tilde{\beta}_{s}^{(s)} &= \mathrm{ST}\left(\tilde{\beta}_{s}^{(s-1)} + \frac{1}{\|X_{s}\|^{2}}X_{s}^{\top}(y - X\tilde{\beta}^{(s-1)}), \frac{\lambda}{\|X_{s}\|^{2}}\right) \\ &= \tilde{\beta}_{s}^{(s-1)} + \frac{1}{\|X_{s}\|^{2}}X_{s}^{\top}(y - X\tilde{\beta}^{(s-1)}) - \frac{\lambda\operatorname{sign}(\hat{\beta}_{s})}{\|X_{s}\|^{2}} \end{split}$$

#### VAR after sign identification

$$X\tilde{\beta}^{(s)} = \underbrace{\left(\mathrm{Id}_n - \frac{1}{\|X_s\|^2} X_s X_s^\top\right)}_{A_s \in \mathbb{R}^{n \times n}} X\tilde{\beta}^{(s-1)} + \underbrace{\frac{X_s^\top y - \lambda \operatorname{sign}(\hat{\beta}_s)}{\|X_s\|^2} X_s}_{b_s \in \mathbb{R}^n}$$

So for the full epoch  $t \rightarrow t + 1$ :

$$\begin{split} X\tilde{\beta}^{(S)} &= A_S X\tilde{\beta}^{(S-1)} + b_S \\ &= A_S A_{S-1} X\tilde{\beta}^{(S-2)} + A_S b_{S-1} + b_S \\ &= \underbrace{A_S \dots A_1}_A X\tilde{\beta}^{(0)} + \underbrace{A_S \dots A_2 b_1 + \dots + A_S b_{S-1} + b_S}_b \\ & \boxed{X\beta^{(t+1)} = A X\beta^{(t)} + b} \end{split}$$

# **Table of Contents**

The Lasso

Exploiting regularity

Sparse GLMs

More solvers speed-up

# VAR for other GLMs

sparse Log. reg. 
$$\underset{\beta \in \mathbb{R}^p}{\operatorname{arg\,min}} \sum_{i=1}^n \log \left(1 + \exp(-y_i \beta^\top x_i)\right) + \lambda \|\beta\|_1$$
  
Multi-task Lasso  $\underset{B \in \mathbb{R}^{p \times T}}{\operatorname{arg\,min}} \frac{1}{2} \|Y - XB\|^2 + \lambda \|B\|_{2,1}$ 

Log. reg. CD update after sign ID:

$$\tilde{\beta}_s^{(s)} = \tilde{\beta}_s^{(s-1)} - \frac{\gamma}{\|X_s\|^2} X_s^\top \nabla F(X \tilde{\beta}^{(s-1)}) - \frac{\gamma}{\|X_s\|^2} \lambda \operatorname{sign}(\hat{\beta}_s)$$

 $\nabla F$  not linear in  $X\beta$  if F is not a quadratic

# Structure for other GLMs

Solution: linearization of  $\nabla F$  around optimum

$$\nabla F(X\beta) = \nabla F(X\hat{\beta}) + \underbrace{D}_{n \times n, \text{ diagonal}} (X\beta - X\hat{\beta}) + o(X\beta - X\hat{\beta})$$

Leads to asymptotic VAR sequence:

$$D^{\frac{1}{2}}X\tilde{\beta}^{(s)} = \underbrace{\left(\mathrm{Id}_{n} - \frac{\gamma}{\|X_{s}\|^{2}}D^{\frac{1}{2}}X_{s}X_{s}^{\top}D^{-\frac{1}{2}}\right)}_{A_{s}}D^{\frac{1}{2}}X\tilde{\beta}^{(s-1)} + b_{s} + o$$

$$X\beta^{(t+1)} = AX\beta^{(t)} + b + o(X\beta - X\hat{\beta})$$

# Applicability to other models

Asymptotic VAR structure is still exploitable<sup>7</sup>



<sup>7</sup>M. Massias et al. "Dual extrapolation for sparse Generalized Linear Models". In: submission to JMLR (2019).

# **Table of Contents**

The Lasso

Exploiting regularity

Sparse GLMs

More solvers speed-up

# Speeding-up solvers

Two approaches:

- safe screening<sup>8,9</sup> (backward approach): remove feature j when it is certified that β̂<sub>j</sub> = 0
- ▶ working set<sup>10</sup> (forward approach): focus on j's for which it is very likely that  $\hat{\beta}_j \neq 0$ .

<sup>&</sup>lt;sup>8</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>9</sup>A. Bonnefoy et al. "A dynamic screening principle for the lasso". In: EUSIPCO. 2014.

<sup>&</sup>lt;sup>10</sup>T. B. Johnson and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: ICML. 2015, pp. 1171–1179.

# Key to identifying features

#### Equicorrelation set<sup>11</sup>

$$E := \left\{ j \in [p] \, : \, |X_j^\top \hat{\theta}| = 1 \right\} \stackrel{\mathsf{lasso}}{=} \left\{ j \in [p] \, : \, |X_j^\top (y - X \hat{\beta})| = \lambda \right\}$$

• For any primal solution,  $j \notin E \implies \hat{\beta}_j = 0$ 

Idea for speed-up: identify  ${\cal E},$  solve only on  ${\cal E}$ 

<u>Practical observation</u>: generally  $\#E \ll p$ 

<sup>11</sup>R. J. Tibshirani. "The lasso problem and uniqueness". In: Electron. J. Stat. 7 (2013), pp. 1456–1490.

# Duality again: gap screening

Cannot know in advance 
$$E = \left\{ j \in [p] \, : \, |X_j^ op \hat{ heta}| = 1 
ight\}$$

Good proxy: find a region  $\mathcal{C} \subset \mathbb{R}^n$  containing  $\hat{\theta}$ 

$$\sup_{\theta \in \mathcal{C}} |X_j^\top \theta| < 1 \Rightarrow |X_j^\top \hat{\theta}| < 1 \Rightarrow j \notin E \Rightarrow \hat{\beta}_j = 0$$

 $<sup>^{12}\</sup>mathsf{E}.$  Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: J. Mach. Learn. Res. 18.128 (2017), pp. 1–33.

# Duality again: gap screening

Cannot know in advance 
$$E = \left\{ j \in [p] \, : \, |X_j^ op \hat{ heta}| = 1 
ight\}$$

Good proxy: find a region  $\mathcal{C} \subset \mathbb{R}^n$  containing  $\hat{\theta}$ 

$$\sup_{\theta \in \mathcal{C}} |X_j^\top \theta| < 1 \Rightarrow |X_j^\top \hat{\theta}| < 1 \Rightarrow j \notin E \Rightarrow \hat{\beta}_j = 0$$

Gap Safe screening rule<sup>12</sup>: C is a ball of radius  $\rho = \sqrt{\frac{2}{\lambda^2}} dgap(\beta, \theta)$  centered at  $\theta \in \Delta_X$ 

$$\forall (\beta, \theta) \in \mathbb{R}^p \times \Delta_X, \quad |X_j^\top \theta| < 1 - ||X_j|| \rho \Rightarrow \hat{\beta}_j = 0$$

<sup>&</sup>lt;sup>12</sup>E. Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: J. Mach. Learn. Res. 18.128 (2017), pp. 1–33.

# Better Gap Safe screening<sup>13</sup>

$$\forall \theta \in \Delta_X, |X_j^\top \theta| < 1 - \|X_j\| \sqrt{\frac{2}{\lambda^2}} \mathsf{dgap}(\beta, \theta) \Rightarrow \hat{\beta}_j = 0$$

#### better dual point $\Rightarrow$ better safe screening



Finance dataset:  $p=1.5 imes 10^6, n=1.5 imes 10^4$ ,  $\lambda=\lambda_{
m max}/5$ 

<sup>&</sup>lt;sup>13</sup>O. Fercoq, A. Gramfort, and J. Salmon. "Mind the duality gap: safer rules for the lasso". In: *ICML*. 2015, pp. 333–342.

# Working sets

# Screening can be used aggressively to define WS, therefore a **better dual point also helps**:



#### news20 dataset, coarse and fine Lasso paths computation

# **Online code**

Fast & pip-installable Cython code, continuous integration, bug tracker, code coverage

Figures & doc at https://mathurinm.github.io/celer



2 from celer import Lasso, LassoCV

From 10,000 s to 50 s for cross-validation on Finance

# Conclusion

Duality matters at several levels for sparse GLMs:

- stopping criterion
- safe feature identification (screening or working set)

Lasso: Exploiting the VAR structure of  $X\beta^{(t)} \hookrightarrow$  better dual

Generalization

- ▶ any twice differentiable separable (samples) data-fitting term
- group penalties (multitask Lasso)

with proof of asymptotic VAR structure & extrapolation is useful

Code: https://github.com/mathurinm/celer Paper: https://arxiv.org/abs/1907.05830

# **References** I

- Aitken, A. "On Bernoulli's numerical solution of algebraic equations". In: *Proceedings of the Royal Society of Edinburgh* 46 (1926), pp. 289–305.
- Bonnefoy, A. et al. "A dynamic screening principle for the lasso". In: EUSIPCO. 2014.
- Chen, S. S. and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.
- El Ghaoui, L., V. Viallon, and T. Rabbani. "Safe feature elimination in sparse supervised learning". In: *J. Pacific Optim.* 8.4 (2012), pp. 667–698.
- Fercoq, O., A. Gramfort, and J. Salmon. "Mind the duality gap: safer rules for the lasso". In: *ICML*. 2015, pp. 333–342.
- Johnson, T. B. and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: ICML. 2015, pp. 1171–1179.

# **References II**

- Mairal, J. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.
- Massias, M., A. Gramfort, and J. Salmon. "Celer: a fast solver for the Lasso with dual extrapolation". In: *ICML*. 2018, pp. 3321–3330.
- Massias, M. et al. "Dual extrapolation for sparse Generalized Linear Models". In: *submission to JMLR* (2019).
- Ndiaye, E. et al. "Gap Safe screening rules for sparsity enforcing penalties". In: J. Mach. Learn. Res. 18.128 (2017), pp. 1–33.
- Scieur, D., A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: NIPS. 2016, pp. 712–720.
- 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 III**

Tibshirani, R. J. "The lasso problem and uniqueness". In: Electron. J. Stat. 7 (2013), pp. 1456–1490.

#### Intuition for extrapolation

Take a VAR sequence  $x^{(t)} \rightarrow x^*$ :

$$x^{(t+1)} - x^* = A(x^{(t)} - x^*)$$

Cayley-Hamilton: find coefficients<sup>14</sup> s.t.  $\sum_{k=0}^{n} a_k A^k = 0$ 

$$\sum_{k=0}^{n} a_k (x^{(t+k+1)} - x^*) = \sum_{k=0}^{n} a_k A^k (x^{(t)} - x^*) = 0$$
$$\hookrightarrow x^* \in \text{Span} (x^{(t)}, \dots, x^{(t+n+1)})$$

 $\hookrightarrow$  approximate  $x^*$  under the form  $x_{\mathrm{acc}} = \sum_{k=1}^{K} c_k x^{(t+k)}$ 

minimizing  $\|x_{\mathrm{acc}} - (Ax_{\mathrm{acc}} + b)\|$  leads to the previous formulas

$$^{14}\mathrm{wlog,}$$
 assume  $\sum\nolimits_{0}^{n}a_{k}=1$  if  $\|A\|<1$ 

# Aitken's rule

For a converging sequence  $(r_n)_{n\in\mathbb{N}}$ , Aitken's rule replaces  $r_{n+1}$  by

$$\Delta^2 = r_n + \frac{1}{\frac{1}{r_{n+1} - r_n} - \frac{1}{r_n - r_{n-1}}}$$