Celer: Fast solver for the Lasso with dual extrapolation

Mathurin Massias

https://mathurinm.github.io INRIA

Joint work with: Alexandre Gramfort (INRIA) Joseph Salmon (Télécom ParisTech)

To appear in ICML 2018

Table of Contents

Lasso basics

Speeding up solvers

A new dual construction

The Lasso^{1,2}

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

• $y \in \mathbb{R}^n$: observations

- $X = [X_1, \dots, X_p] \in \mathbb{R}^{n \times p}$: design matrix
- $\lambda > 0$: trade-off parameter between data-fit and regularization
- sparsity: for λ large enough, $\|\hat{\beta}\|_0 \ll p$

<u>Rem</u>: uniqueness is not guaranteed, more later

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

²S. S. Chen and D. L. Donoho. "Atomic decomposition by basis pursuit". In: SPIE. 1995.

Duality for the Lasso

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

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

Duality for the Lasso

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

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



Toy visualization example: n = p = 2

Duality for the Lasso

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

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



Projection problem: $\hat{\theta} = \prod_{\Delta_X} (y/\lambda)$

Solving the Lasso

So-called *smooth* + *separable* problem

▶ In signal processing: use ISTA/FISTA³

► In ML: state-of-the-art algorithm when X is not an implicit operator: coordinate descent (CD)^{4,5}

Iterative algorithm: minimize $\mathcal{P}(\beta) = \mathcal{P}(\beta_1, \dots, \beta_p)$ w.r.t. β_1 , then β_2 , etc.

³A. Beck and M. Teboulle. "A fast iterative shrinkage-thresholding algorithm for linear inverse problems". In: SIAM J. Imaging Sci. 2.1 (2009), pp. 183–202.

⁴J. Friedman et al. "Pathwise coordinate optimization". In: Ann. Appl. Stat. 1.2 (2007), pp. 302–332.

⁵P. Tseng. "Convergence of a block coordinate descent method for nondifferentiable minimization". In: J. Optim. Theory Appl. 109.3 (2001), pp. 475–494.

To minimize
$$\mathcal{P}(\beta) = \frac{1}{2} \|y - \sum_{j=1}^{p} X_j \beta_j\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
:

Initialisation: $\beta^0 = 0 \in \mathbb{R}^p$

To minimize
$$\mathcal{P}(eta) = rac{1}{2} \|y - \sum_{j=1}^p X_j eta_j\|^2 + \lambda \sum_{j=1}^p |eta_j|$$
:

Initialisation: $\beta^0 = 0 \in \mathbb{R}^p$ for $t = 1, \dots, T$ do

To minimize
$$\mathcal{P}(\beta) = \frac{1}{2} \|y - \sum_{j=1}^{p} X_j \beta_j\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
:

Initialisation:
$$\beta^0 = 0 \in \mathbb{R}^p$$

for $t = 1, ..., T$ do
 $\beta_1^t \leftarrow \arg \min_{\substack{\beta_1 \in \mathbb{R}}} \mathcal{P}(\beta_1, \beta_2^{t-1}, \beta_3^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_p^{t-1})$

To minimize
$$\mathcal{P}(\beta) = \frac{1}{2} \|y - \sum_{j=1}^{p} X_j \beta_j\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
:

Initialisation:
$$\beta^{0} = 0 \in \mathbb{R}^{p}$$

for $t = 1, ..., T$ do
 $\beta_{1}^{t} \leftarrow \underset{\beta_{1} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}, \beta_{2}^{t-1}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
 $\beta_{2}^{t} \leftarrow \underset{\beta_{2} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}^{t}, \beta_{2}^{-}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$

To minimize
$$\mathcal{P}(\beta) = \frac{1}{2} \|y - \sum_{j=1}^{p} X_j \beta_j\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
:

Initialisation:
$$\beta^{0} = 0 \in \mathbb{R}^{p}$$

for $t = 1, ..., T$ do
 $\beta_{1}^{t} \leftarrow \underset{\beta_{1} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}, \beta_{2}^{t-1}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
 $\beta_{2}^{t} \leftarrow \underset{\beta_{2} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}^{t}, \beta_{2}^{t}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
 $\beta_{3}^{t} \leftarrow \underset{\beta_{3} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}^{t}, \beta_{2}^{t}, \beta_{3}^{t}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$

To minimize
$$\mathcal{P}(eta) = rac{1}{2} \|y - \sum_{j=1}^p X_j eta_j\|^2 + \lambda \sum_{j=1}^p |eta_j|$$
:

Initialisation:
$$\beta^{0} = 0 \in \mathbb{R}^{p}$$

for $t = 1, ..., T$ do
 $\beta_{1}^{t} \leftarrow \underset{\beta_{1} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}, \beta_{2}^{t-1}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
 $\beta_{2}^{t} \leftarrow \underset{\beta_{2} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}^{t}, \beta_{2}^{t}, \beta_{3}^{t-1}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
 $\beta_{3}^{t} \leftarrow \underset{\beta_{3} \in \mathbb{R}}{\operatorname{arg\,min}} \mathcal{P}(\beta_{1}^{t}, \beta_{2}^{t}, \beta_{3}^{t}, ..., \beta_{p-1}^{t-1}, \beta_{p}^{t-1})$
:

To minimize
$$\mathcal{P}(\beta) = \frac{1}{2} \|y - \sum_{j=1}^{p} X_j \beta_j\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|$$
:

CD update: soft-thresholding

Coordinate-wise minimization is easy:

$$\beta_j \leftarrow \operatorname{ST}\left(\frac{\lambda}{\|X_j\|^2}, \beta_j + \frac{X_j^\top (y - X\beta)}{\|X_j\|^2}\right)$$



1 update is $\mathcal{O}(n)$

<u>Variants</u>: minimize *w.r.t.* β_j with j chosen at random, or **shuffle** order every epoch (1 epoch = p updates)

CD update: soft-thresholding

Coordinate-wise minimization is easy:

$$\beta_j \leftarrow \operatorname{ST}\left(\frac{\lambda}{\|X_j\|^2}, \beta_j + \frac{X_j^\top (y - X\beta)}{\|X_j\|^2}\right)$$



1 update is $\mathcal{O}(n)$

<u>Variants</u>: minimize *w.r.t.* β_j with j chosen at random, or **shuffle** order every epoch (1 epoch = p updates)

When do we stop?

Duality gap as a stopping criterion

For any primal-dual pair (β, θ) :

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

$$\begin{array}{c|c} \mathcal{P}(\hat{\beta}) & \mathcal{P}(\beta) \\ \hline \mathcal{D}(\theta) & \mathcal{D}(\hat{\theta}) \end{array}$$

The duality gap $\mathcal{P}(\beta) - \mathcal{D}(\theta) =: gap(\beta, \theta)$ is an upper bound of the suboptimality gap $\mathcal{P}(\beta) - \mathcal{P}(\hat{\beta})$:

$$\forall \beta, (\exists \theta \in \Delta_X, \operatorname{\mathsf{gap}}(\beta, \theta) \le \epsilon) \Rightarrow \mathcal{P}(\beta) - \mathcal{P}(\hat{\beta}) \le \epsilon$$

i.e., β is an ϵ -solution

Duality gap as a stopping criterion

For any primal-dual pair (β, θ) :

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



The duality gap $\mathcal{P}(\beta) - \mathcal{D}(\theta) =: gap(\beta, \theta)$ is an upper bound of the suboptimality gap $\mathcal{P}(\beta) - \mathcal{P}(\hat{\beta})$:

$$\forall \beta, (\exists \theta \in \Delta_X, \operatorname{\mathsf{gap}}(\beta, \theta) \le \epsilon) \Rightarrow \mathcal{P}(\beta) - \mathcal{P}(\hat{\beta}) \le \epsilon$$

i.e., β is an ϵ -solution

Primal-dual link at optimum:

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

⁶J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Primal-dual link at optimum:

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

Standard approach⁶: at epoch t, corresponding to iterate β^t and residuals $r^t := y - X\beta^t$, take

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

⁶J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Primal-dual link at optimum:

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

Standard approach⁶: at epoch t, corresponding to iterate β^t and residuals $r^t := y - X\beta^t$, take

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

It is not necessarily feasible!

⁶J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Primal-dual link at optimum:

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

Standard approach⁶: at epoch t, corresponding to iterate β^t and residuals $r^t := y - X\beta^t$, take

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

residuals rescaling

⁶J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Primal-dual link at optimum:

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

Standard approach⁶: at epoch t, corresponding to iterate β^t and residuals $r^t := y - X\beta^t$, take

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

residuals rescaling

- converges to $\hat{\theta}$ (provided β^t converges to $\hat{\beta}$)
- ► $\mathcal{O}(np)$ to compute (= 1 epoch of CD) → rule of thumb: compute θ_{res}^t and gap every f = 10 epochs

⁶J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Table of Contents

Lasso basics

Speeding up solvers

A new dual construction

Speeding up solvers

Lasso motivation: we expect sparse solutions/small supports

$$\mathcal{S}_{\hat{\beta}} := \left\{ j \in [p] : \hat{\beta}_j \neq 0 \right\}$$

"the solution restricted to its support solves the problem restricted to features in this support"

$$\hat{\beta}_{\mathcal{S}_{\hat{\beta}}} \in \operatorname*{arg\,min}_{b \in \mathbb{R}^{\|\hat{\beta}\|_{0}}} \frac{1}{2} \|y - X_{\mathcal{S}_{\hat{\beta}}}b\|^{2} + \lambda \|b\|_{1}$$

Usually $\|\hat{\beta}\|_0 \ll p$ so the second problem is much simpler.

Technical detail

- The primal solution/support might not be unique!
- But $\hat{\theta}$ is unique and so is the *equicorrelation set*⁷:

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

 $\blacktriangleright \ \, {\rm For \ any \ primal \ solution, \ } \mathcal{S}_{\hat{\beta}} \subset E$

Technical detail

- The primal solution/support might not be unique!
- But $\hat{\theta}$ is unique and so is the *equicorrelation set*⁷:

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

▶ For any primal solution, $S_{\hat{\beta}} \subset E$

Grail of sparse solvers: identify E, solve on E<u>Practical observation</u>: generally $\#E \ll p$

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

Speeding-up solvers

Two approaches:

- safe screening^{8,9} (backward approach): remove feature j when it is certified that j ∉ E
- ▶ working set¹⁰ (forward approach): focus on j's very likely to be in E

<u>Rem</u>: hybrid approach: strong rules¹¹

⁸L. El Ghaoui, V. Viallon, and T. Rabbani. "Safe feature elimination in sparse supervised learning". In: J. Pacific Optim. 8.4 (2012), pp. 667–698.

⁹A. Bonnefoy et al. "A dynamic screening principle for the lasso". In: EUSIPCO. 2014.

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

¹¹R. 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.

We want to identify $E = \left\{ j \in [p] : |X_j^\top \hat{\theta}| = 1 \right\}$... but we can't get it without $\hat{\beta}$

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

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}} |\boldsymbol{X}_j^\top \boldsymbol{\theta}| < 1 \Rightarrow |\boldsymbol{X}_j^\top \hat{\boldsymbol{\theta}}| < 1$$

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

We want to identify $E = \left\{ j \in [p] : |X_j^\top \hat{\theta}| = 1 \right\}$... but we can't get it without $\hat{\beta}$

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

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

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

We want to identify $E = \left\{ j \in [p] : |X_j^\top \hat{\theta}| = 1 \right\}$... but we can't get it without $\hat{\beta}$

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

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

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

We want to identify $E = \left\{ j \in [p] : |X_j^\top \hat{\theta}| = 1 \right\}$... but we can't get it without $\hat{\beta}$

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 rule¹²: C is a ball of radius $r = \sqrt{\frac{2}{\lambda^2}} gap(\beta, \theta)$

$$\forall (\beta, \theta), |X_j^\top \theta| < 1 - ||X_j|| r \Rightarrow \hat{\beta}_j = 0$$

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

Working/active set

Algorithm: Generic WS algorithm

Table of Contents

Lasso basics

Speeding up solvers

A new dual construction

Back to dual choice

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

Two drawbacks of residuals rescaling:

- ignores information from previous iterates
- workload "imbalanced": more efforts in primal than in dual

Back to dual choice

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

Two drawbacks of residuals rescaling:

- ignores information from previous iterates
- workload "imbalanced": more efforts in primal than in dual
Back to dual choice

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

Two drawbacks of residuals rescaling:

- ignores information from previous iterates
- workload "imbalanced": more efforts in primal than in dual



Leukemia dataset (p=7129, n=72), for $\lambda=\lambda_{\max}/20$

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

Acceleration through residuals extrapolation¹³

What is the limit of $(0, \frac{1}{2}, \frac{3}{4}, \frac{7}{8}, \frac{15}{16}, \ldots)$?

¹³D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: NIPS. 2016, pp. 712–720.

Acceleration through residuals extrapolation¹³

What is the limit of
$$(0, \frac{1}{2}, \frac{3}{4}, \frac{7}{8}, \frac{15}{16}, \ldots)$$
?

extrapolation!

$$ightarrow$$
 use the same idea to infer $\lim_{t
ightarrow\infty}r^t=\lambda\hat{ heta}$

¹³D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: NIPS. 2016, pp. 712–720.

Extrapolation justification

If $(x_t)_{t\in\mathbb{N}}$ follows a converging autoregressive process (AR):

$$x_t = ax_{t-1} - b$$
 $(|a| < 1)$ with $\lim_{t \to \infty} x_t = x^*$

we have

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

Aitken's Δ^2 : 2 unknowns, so 2 equations/3 points x_t, x_{t-1}, x_{t-2} are enough to find x^* !¹⁴

¹⁴A. 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	$\sum_{i=0}^{t} \frac{(-1)^i}{2i+1}$	Δ^2
0	1.0000	_
1	0.66667	-
2	0.86667	0.7 9167
3	0.7 2381	0.78333
4	0.83492	0.78631
5	0.7 4401	0.78492
6	0.82093	0.78568
7	0.75427	0.78522
8	0.81309	0.78552
9	0.7 6046	0.7853 1

Generalization to $x_t \in \mathbb{R}^n$

AMPE (Approximate Minimal Polynomial Extrapolation): applies to Vector autoregressive (VAR) process

i.e., $x_t \in \mathbb{R}^n$, and $a \in \mathbb{R}$ becomes $A \in \mathbb{R}^{n \times n}$

More difficult to eliminate A (unobserved), the idea is to approximate its minimal polynomial.

• fix K = 5 (small)

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

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

- fix K = 5 (small)
- keep track of K past residuals r^t, \ldots, r^{t+1-K}
- ▶ form $U^t = [r^{t+1-K} r^{t-K}, \dots, r^t r^{t-1}] \in \mathbb{R}^{n \times K}$

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

- fix K = 5 (small)
- keep track of K past residuals r^t, \ldots, r^{t+1-K}
- form $U^t = [r^{t+1-K} r^{t-K}, \dots, r^t r^{t-1}] \in \mathbb{R}^{n \times K}$
- solve $(U^t)^{\top} U^t z = \mathbf{1}_K$

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

- fix K = 5 (small)
- keep track of K past residuals r^t, \ldots, r^{t+1-K}
- form $U^t = [r^{t+1-K} r^{t-K}, \dots, r^t r^{t-1}] \in \mathbb{R}^{n \times K}$

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

► $c = \frac{z}{z^\top \mathbf{1}_K}$

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

- fix K = 5 (small)
- keep track of K past residuals r^t, \ldots, r^{t+1-K}
- form $U^t = [r^{t+1-K} r^{t-K}, \dots, r^t r^{t-1}] \in \mathbb{R}^{n \times K}$

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

► $c = \frac{z}{z^{\top} \mathbf{1}_K}$

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

▶ keep track of K past residuals
$$r^t, \ldots, r^{t+1-K}$$
▶ form $U^t = [r^{t+1-K} - r^{t-K}, \ldots, r^t - r^{t-1}] \in \mathbb{R}^{n \times K}$
▶ solve $(U^t)^\top U^t z = \mathbf{1}_K$
▶ $c = \frac{z}{z^\top \mathbf{1}_K}$
 $\left\{ r^t \quad \text{if } t \leq K \right\}$

• fix K = 5 (small)

$$r_{\text{accel}}^{t} = \begin{cases} r, & \text{if } t \leq K \\ \sum_{k=1}^{K} c_k r^{t+1-k}, & \text{if } t > K \end{cases}$$

(affine combination, goes outside convex hull of $\{r^t, \ldots, r^{t+1-K}\}$)

¹⁵M. Massias, A. Gramfort, and J. Salmon. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].

Extrapolated dual point

$$r_{\text{accel}}^{t} = \begin{cases} r^{t}, & \text{if } t \leq K \\ \sum_{k=1}^{K} c_{k} r^{t+1-k}, & \text{if } t > K \end{cases}$$
$$\overline{\theta_{\text{accel}}^{t} := r_{\text{accel}}^{t} / \max(\lambda, \|X^{\top} r_{\text{accel}}^{t}\|_{\infty})}$$

• convergence of θ_{accel}^t ?

• $(U^t)^\top U^t z = \mathbf{1}_K \rightarrow \text{linear system solving}$?

 $\theta_{\rm accel}$ is $\mathcal{O}(np+K^2n+K^3)$ to compute, so compute $\theta_{\rm res}$ as well and pick the best

use
$$\theta^t = \underset{\theta \in \{\theta_{\text{res}}^t, \theta_{\text{accel}}^t, \theta^{t-1}\}}{\arg \max} \mathcal{D}(\theta)$$

 $\theta_{\rm accel}$ is $\mathcal{O}(np+K^2n+K^3)$ to compute, so compute $\theta_{\rm res}$ as well and pick the best

use
$$\theta^t = \operatorname*{arg\,max}_{\theta \in \{\theta^t_{\mathrm{res}}, \theta^t_{\mathrm{accel}}, \theta^{t-1}\}} \mathcal{D}(\theta)$$

Final cost of 10 CD epochs + gap computation \approx 12 np vs 11 np in classical approach

Does it work?



Leukemia dataset (p = 7129, n = 72), for $\lambda = \lambda_{\text{max}}/20$ (consistent finding across datasets)

• $heta_{
m res}$ is bad

- $heta_{
 m accel}$ gives a tighter bound
- θ_{accel} does not behave erratically

Key assumption for extrapolation¹⁶: r^t follows a VAR.

► True for ISTA and the Lasso, once support is identified¹⁷ (but ISTA/FISTA slow on our statistical scenarios)

¹⁶D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: NIPS. 2016, pp. 712–720.
¹⁷J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: NIPS. 2014, pp. 1970–1978.

Key assumption for extrapolation¹⁶: r^t follows a VAR.

 True for ISTA and the Lasso, once support is identified¹⁷ (but ISTA/FISTA slow on our statistical scenarios)

Conjecture: it is also true for cyclic CD

 ¹⁶D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.
 ¹⁷J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

Key assumption for extrapolation¹⁶: r^t follows a VAR.

- True for ISTA and the Lasso, once support is identified¹⁷ (but ISTA/FISTA slow on our statistical scenarios)
- Conjecture: it is also true for cyclic CD

 ¹⁶D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.
 ¹⁷J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

Key assumption for extrapolation¹⁶: r^t follows a VAR.

- True for ISTA and the Lasso, once support is identified¹⁷ (but ISTA/FISTA slow on our statistical scenarios)
- Conjecture: it is also true for cyclic CD
- Rem: : Shuffle CD breaks the regularity

 ¹⁶D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.
 ¹⁷J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

2D example again



Toy dual zoom: cyclic





Toy dual zoom: shuffle





Better safe screening

Recall Gap Safe screening rule:

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

better dual point \Rightarrow better safe screening

Better safe screening

Recall Gap Safe screening rule:

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

better dual point \Rightarrow better safe screening



Finance dataset: $(p=1.5 imes10^6,n=1.5 imes10^4)$, $\lambda=\lambda_{
m max}/5$

Screening vs Working sets

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

Screening vs Working sets

$$\begin{split} X_j^\top \theta | < 1 - \|X_j\| \sqrt{\frac{2}{\lambda^2} \mathsf{gap}(\beta, \theta)} \Rightarrow \hat{\beta}_j &= 0 \\ \Leftrightarrow \\ d_j(\theta) > \sqrt{\frac{2}{\lambda^2} \mathsf{gap}(\beta, \theta)} \Rightarrow \hat{\beta}_j &= 0 \\ \text{with } d_j(\theta) &:= \frac{1 - |X_j^\top \theta|}{\|X_j\|} \end{split}$$

 $d_j(\theta)$ larger than threshold \rightarrow exclude feature j

Screening vs Working sets

$$\begin{split} X_j^\top \theta | < 1 - \|X_j\| \sqrt{\frac{2}{\lambda^2}} \mathsf{gap}(\beta, \theta) \Rightarrow \hat{\beta}_j &= 0 \\ \Leftrightarrow \\ d_j(\theta) > \sqrt{\frac{2}{\lambda^2}} \mathsf{gap}(\beta, \theta) \Rightarrow \hat{\beta}_j &= 0 \\ \text{with } d_j(\theta) &:= \frac{1 - |X_j^\top \theta|}{\|X_j\|} \end{split}$$

 $d_j(\theta)$ larger than threshold \rightarrow exclude feature j

Alternative: Solve subproblem with small $d_i(\theta)$ only (WS)

how to prioritize features?

• how to prioritize features? \rightarrow use $d_j(\theta)$

- how to prioritize features? \rightarrow use $d_j(\theta)$
- how many features in WS?

- how to prioritize features? \rightarrow use $d_j(\theta)$
- ▶ how many features in WS? → start at 100, double at each WS definition. Features cannot leave the WS

- how to prioritize features? \rightarrow use $d_j(\theta)$
- \blacktriangleright how many features in WS? \rightarrow start at 100, double at each WS definition. Features cannot leave the WS
- solve the subproblem with which precision?
3 questions for working sets

- how to prioritize features? \rightarrow use $d_j(\theta)$
- ▶ how many features in WS? → start at 100, double at each WS definition. Features cannot leave the WS
- \blacktriangleright solve the subproblem with which precision? \rightarrow use same as required for whole problem

3 questions for working sets

- how to prioritize features? \rightarrow use $d_j(\theta)$
- ▶ how many features in WS? → start at 100, double at each WS definition. Features cannot leave the WS
- \blacktriangleright solve the subproblem with which precision? \rightarrow use same as required for whole problem

Not so fancy, but guarantees convergence.

Pruning variant: working set can decrease in size & features can leave the working set

Similarities^{18,19}

$$d_j(\theta) := \frac{1 - |X_j^\top \theta|}{\|X_j\|}$$

¹⁹S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: NIPS. 2017.

¹⁸J. Fan and J. Lv. "Sure independence screening for ultrahigh dimensional feature space". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 70.5 (2008), pp. 849–911.

Similarities^{18,19}

$$d_j(\theta) := \frac{1 - |X_j^\top \theta|}{\|X_j\|}$$

Lasso case with $\theta = \theta_{res}$ and normalized X_j 's:

$$1 - d_j(\theta) \propto |X_j^\top r^t|$$

small $d_j(\theta) = \text{high correlation with residuals/high norm of partial gradient of data-fitting term...}$

¹⁸J. Fan and J. Lv. "Sure independence screening for ultrahigh dimensional feature space". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 70.5 (2008), pp. 849–911.

¹⁹S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: NIPS. 2017.

Similarities^{18,19}

$$d_j(\theta) := \frac{1 - |X_j^\top \theta|}{\|X_j\|}$$

Lasso case with $\theta = \theta_{res}$ and normalized X_j 's:

$$1 - d_j(\theta) \propto |X_j^\top r^t|$$

small $d_j(\theta) = \text{high correlation with residuals/high norm of partial gradient of data-fitting term...}$

BUT our strength is that we can use any θ , in particular θ_{accel}

¹⁸ J. Fan and J. Lv. "Sure independence screening for ultrahigh dimensional feature space". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 70.5 (2008), pp. 849–911.

¹⁹S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: NIPS. 2017.

Comparison

State-of-the-art WS solver for sparse problems: Blitz²⁰



Finance dataset, Lasso path of 10 (top) or 100 (bottom) λ 's from $\lambda_{\rm max}$ to $\lambda_{\rm max}/100$

²⁰T. B. Johnson and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: ICML. 2015, pp. 1171–1179.

Reusable science

https://github.com/mathurinm/celer: code with continuous integration, code coverage, bug tracker

mathurinm / celer Fast solver for the Lasso https://mathurinm.github.io/celer/ Edit Add topics 75 commits ₽ 6 branches O releases 3 contributors ർ BSD-3-Clause Branch: master + New pull request Create new file Unload files Find file Clone or download 😥 🖞 alemaitre and mathurinm [MRG] Make coverage great again (#21) Latest commit cb5629e 8 days ago celer Replacing nosetests with pytest (#13) 9 days ago FREADME.md

celer



Fast algorithm to solve the Lasso with dual extrapolation

Documentation

Please visit https://mathurinm.github.io/celer/ for the latest version of the documentation.

Examples gallery

https://mathurinm.github.io/celer: documentation
(examples, API)

Examples Gallery¶



validation on Leukemia

Run LassoCV for cross-



Lasso path computation on Leukemia dataset



Lasso path computation on Finance/log1p



Drop-in sklearn replacement

1 from sklearn.linear_model import Lasso, LassoCV

2 from celer import Lasso, LassoCV



The optimization objective for Lasso is:

(1 / (2 * n_samples)) * ||y - X beta||^2_2 + alpha * ||beta||_1

Parameters: alpha : float, optional

Constant that multiplies the L1 term. Defaults to 1.0. alpha = θ is equivalent to an ordinary least square. For numerical reasons, using alpha = θ with the Lasso object is not advised.

max_iter : int, optional

The maximum number of iterations (subproblem definitions)

gap_freq : int

Number of coordinate descent epochs between each duality gap computations.

Drop-in sklearn replacement

1 from sklearn.linear_model import Lasso, LassoCV

2 from celer import Lasso, LassoCV

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

celer.Lasso

class celer. Lasso (alpha=1.0, max_iter=100, gap_freq=10, max_epochs=50000, p0=10, verbose= tol=1e-06, prune=0, fit_intercept=True)

Fort me on Cit

Lasso scikit-learn estimator based on Celer solver

The optimization objective for Lasso is:

(1 / (2 * n_samples)) * ||y - X beta||^2_2 + alpha * ||beta||_1

Parameters: alpha : float, optional

Constant that multiplies the L1 term. Defaults to 1.0. alpha = 0 is equivalent to an ordinary least square. For numerical reasons, using alpha = 0 with the Lasso object is not advised.

max_iter : int, optional

The maximum number of iterations (subproblem definitions)

gap_freq : int

Number of coordinate descent epochs between each duality gap computations.

Conclusion

Duality matters at several levels for the Lasso:

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

Key improvement: residuals rescaling \rightarrow residuals extrapolation

Future works:

- Can it work for sparse logreg, group Lasso?
- Can we prove convergence of θ_{accel} and give rates?

Feedback welcome on the online code!

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.
- Beck, A. and M. Teboulle. "A fast iterative shrinkage-thresholding algorithm for linear inverse problems". In: SIAM J. Imaging Sci. 2.1 (2009), pp. 183–202.
- 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.
- Fan, J. and J. Lv. "Sure independence screening for ultrahigh dimensional feature space". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 70.5 (2008), pp. 849–911.

References II

- Friedman, J. et al. "Pathwise coordinate optimization". In: Ann. Appl. Stat. 1.2 (2007), pp. 302–332.
- Johnson, T. B. and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: ICML. 2015, pp. 1171–1179.
- Liang, J., J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: NIPS. 2014, pp. 1970–1978.
- 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. "Dual Extrapolation for Faster Lasso Solvers". In: ArXiv e-prints (2018). arXiv: 1802.07481 [stat.ML].
- Ndiaye, E. et al. "Gap Safe screening rules for sparsity enforcing penalties". In: J. Mach. Learn. Res. 18.128 (2017), pp. 1–33.

References III

- Scieur, D., A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.
- Stich, S., A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: NIPS. 2017.
- Tibshirani, R. "Regression Shrinkage and Selection via the Lasso". In: J. R. Stat. Soc. Ser. B Stat. Methodol. 58.1 (1996), pp. 267–288.
- Tibshirani, R. J. "The lasso problem and uniqueness". In: Electron. J. Stat. 7 (2013), pp. 1456–1490.
- Tibshirani, R. 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.
- Tseng, P. "Convergence of a block coordinate descent method for nondifferentiable minimization". In: J. Optim. Theory Appl. 109.3 (2001), pp. 475–494.