

Celer⁽¹⁾: a fast Lasso solver with dual extrapolation

Joseph Salmon

Université de Montpellier

Joint work with:

Alexandre Gramfort (INRIA)

Mathurin Massias (INRIA)

⁽¹⁾Constraint Elimination for the Lasso with Extrapolated Residuals

Table of Contents

Lasso basics

Speeding up Lasso solvers

A new dual construction

The Lasso^{(2),(3)}: least squares and sparsity

$$\hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} \underbrace{\frac{1}{2} \|\mathbf{y} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1}_{\mathcal{P}(\mathbf{w})}$$

- ▶ $y \in \mathbb{R}^n$: observations
- ▶ $X = [\mathbf{x}_1, \dots, \mathbf{x}_p] \in \mathbb{R}^{n \times p}$: design matrix, p features
- ▶ $\lambda > 0$: trade-off parameter between data-fit and regularization
- ▶ sparsity: for λ large, $\|\hat{\mathbf{w}}\|_0 = \#\{j \in [p] : \hat{\mathbf{w}}_j \neq 0\} \ll p$

Rem: uniqueness is not guaranteed

⁽²⁾R. Tibshirani. "Regression Shrinkage and Selection via the Lasso". In: *J. R. Stat. Soc. Ser. B Stat. Methodol.* 58.1 (1996), pp. 267–288.

⁽³⁾S. S. Chen and D. L. Donoho. "Atomic decomposition by basis pursuit". In: *SPIE*. 1995.

Duality for the Lasso

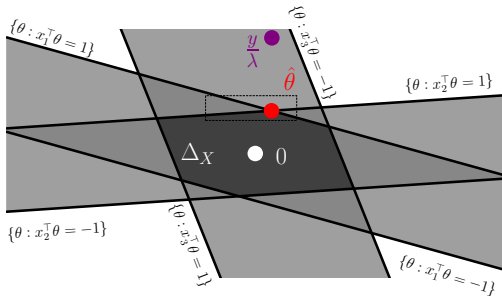
$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \Delta_X} \underbrace{\frac{1}{2} \|\mathbf{y}\|^2 - \frac{\lambda^2}{2} \|\mathbf{y}/\lambda - \boldsymbol{\theta}\|^2}_{\mathcal{D}(\boldsymbol{\theta})}$$

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

Duality for the Lasso

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

$\Delta_X = \{\theta \in \mathbb{R}^n : \forall j \in [p], |\mathbf{x}_j^\top \theta| \leq 1\}$: **dual feasible set**

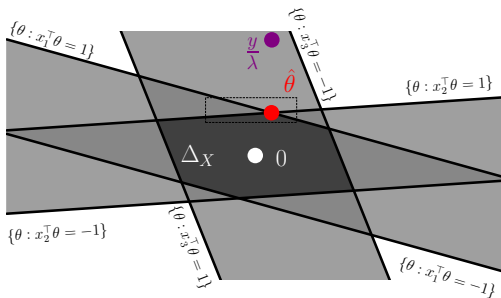


Toy visualization example: $n = 2, p = 3$

Duality for the Lasso

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

$\Delta_X = \{\theta \in \mathbb{R}^n : \forall j \in [p], |\mathbf{x}_j^\top \theta| \leq 1\}$: **dual feasible set**

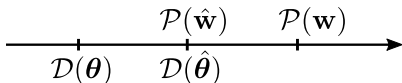


Projection problem: $\hat{\theta} = \Pi_{\Delta_X}(\mathbf{y}/\lambda)$

Duality gap and stopping criterion

For any primal-dual pair $(\mathbf{w}, \boldsymbol{\theta}) \in \mathbb{R}^p \times \Delta_X$:

$$\mathcal{P}(\mathbf{w}) \geq \mathcal{P}(\hat{\mathbf{w}}) = \mathcal{D}(\hat{\boldsymbol{\theta}}) \geq \mathcal{D}(\boldsymbol{\theta})$$



Duality gap : $\text{gap}(\mathbf{w}, \boldsymbol{\theta}) := \mathcal{P}(\mathbf{w}) - \mathcal{D}(\boldsymbol{\theta})$

upper bound on **suboptimality gap** : $\mathcal{P}(\mathbf{w}) - \mathcal{P}(\hat{\mathbf{w}})$

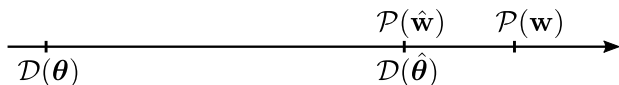
$$\forall \mathbf{w} \in \mathbb{R}^p, (\exists \boldsymbol{\theta} \in \Delta_X, \text{gap}(\mathbf{w}, \boldsymbol{\theta}) \leq \epsilon) \Rightarrow \mathcal{P}(\mathbf{w}) - \mathcal{P}(\hat{\mathbf{w}}) \leq \epsilon$$

i.e., \mathbf{w} is an ϵ -solution whenever $\text{gap}(\mathbf{w}, \boldsymbol{\theta}) \leq \epsilon$

Duality gap and stopping criterion

For any primal-dual pair $(\mathbf{w}, \boldsymbol{\theta}) \in \mathbb{R}^p \times \Delta_X$:

$$\mathcal{P}(\mathbf{w}) \geq \mathcal{P}(\hat{\mathbf{w}}) = \mathcal{D}(\hat{\boldsymbol{\theta}}) \geq \mathcal{D}(\boldsymbol{\theta})$$



Duality gap : $\text{gap}(\mathbf{w}, \boldsymbol{\theta}) := \mathcal{P}(\mathbf{w}) - \mathcal{D}(\boldsymbol{\theta})$

upper bound on **suboptimality gap** : $\mathcal{P}(\mathbf{w}) - \mathcal{P}(\hat{\mathbf{w}})$

$$\forall \mathbf{w} \in \mathbb{R}^p, (\exists \boldsymbol{\theta} \in \Delta_X, \text{gap}(\mathbf{w}, \boldsymbol{\theta}) \leq \epsilon) \Rightarrow \mathcal{P}(\mathbf{w}) - \mathcal{P}(\hat{\mathbf{w}}) \leq \epsilon$$

i.e., \mathbf{w} is an ϵ -solution whenever $\text{gap}(\mathbf{w}, \boldsymbol{\theta}) \leq \epsilon$

Solving the Lasso problem

So-called “*smooth + separable*” problem

- ▶ In signal processing: use ISTA/FISTA⁽⁴⁾ (proximal algorithms)
- ▶ In ML: state-of-the-art algorithm when X is not an implicit operator: coordinate descent (CD)^{(5), (6)}

⁽⁴⁾A. 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.

⁽⁵⁾J. Friedman et al. “Pathwise coordinate optimization”. In: *Ann. Appl. Stat.* 1.2 (2007), pp. 302–332.

⁽⁶⁾P. Tseng. “Convergence of a block coordinate descent method for nondifferentiable minimization”. In: *J. Optim. Theory Appl.* 109.3 (2001), pp. 475–494.

Solving the Lasso: cyclic CD

To minimize :
$$\mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

To minimize :
$$\mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

|

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

To minimize :
$$\mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

$$\mathbf{w}_1^t \leftarrow \arg \min_{\mathbf{w}_1 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1, \mathbf{w}_2^{t-1}, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

To minimize :
$$\mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

$$\mathbf{w}_1^t \leftarrow \arg \min_{\mathbf{w}_1 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1, \mathbf{w}_2^{t-1}, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_2^t \leftarrow \arg \min_{\mathbf{w}_2 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

$$\text{To minimize : } \mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

$$\mathbf{w}_1^t \leftarrow \arg \min_{\mathbf{w}_1 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1, \mathbf{w}_2^{t-1}, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_2^t \leftarrow \arg \min_{\mathbf{w}_2 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_3^t \leftarrow \arg \min_{\mathbf{w}_3 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2^t, \mathbf{w}_3, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

$$\text{To minimize : } \mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

$$\mathbf{w}_1^t \leftarrow \arg \min_{\mathbf{w}_1 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1, \mathbf{w}_2^{t-1}, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_2^t \leftarrow \arg \min_{\mathbf{w}_2 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_3^t \leftarrow \arg \min_{\mathbf{w}_3 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2^t, \mathbf{w}_3, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

\vdots

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

Solving the Lasso: cyclic CD

$$\text{To minimize : } \mathcal{P}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \sum_{j=1}^p \mathbf{x}_j \mathbf{w}_j\|^2 + \lambda \sum_{j=1}^p |\mathbf{w}_j|$$

Algorithm: Cyclic CD

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $t = 1, \dots, T$ **do**

$$\mathbf{w}_1^t \leftarrow \arg \min_{\mathbf{w}_1 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1, \mathbf{w}_2^{t-1}, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_2^t \leftarrow \arg \min_{\mathbf{w}_2 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2, \mathbf{w}_3^{t-1}, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

$$\mathbf{w}_3^t \leftarrow \arg \min_{\mathbf{w}_3 \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2^t, \mathbf{w}_3, \dots, \mathbf{w}_{p-1}^{t-1}, \mathbf{w}_p^{t-1})$$

\vdots

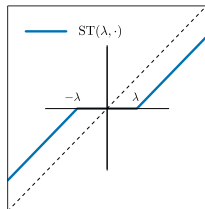
$$\mathbf{w}_p^t \leftarrow \arg \min_{\mathbf{w}_p \in \mathbb{R}} \mathcal{P}(\mathbf{w}_1^t, \mathbf{w}_2^t, \mathbf{w}_3^t, \dots, \mathbf{w}_{p-1}^t, \mathbf{w}_p)$$

cf. Tseng (2001), Friedman *et al.* (2007), Wu *et al.* (2008),
Nesterov (2012), Beck *et al.* (2013), ...

CD update: soft-thresholding

Coordinate-wise minimization is easy:

$$\mathbf{w}_j \leftarrow \text{ST} \left(\frac{\lambda}{\|\mathbf{x}_j\|^2}, \mathbf{w}_j + \frac{\mathbf{x}_j^\top (\mathbf{y} - X\mathbf{w})}{\|\mathbf{x}_j\|^2} \right)$$



► 1 update is $\mathcal{O}(n)$

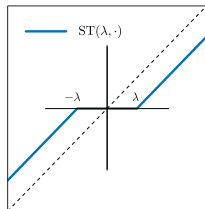
Variants: minimize *w.r.t.* \mathbf{w}_j with j chosen at random, or shuffle order every epoch (1 epoch = p updates)

⁽⁷⁾R. J. Tibshirani. "Dykstra's Algorithm, ADMM, and Coordinate Descent: Connections, Insights, and Extensions". In: *NIPS*. 2017, pp. 517–528.

CD update: soft-thresholding

Coordinate-wise minimization is easy:

$$\mathbf{w}_j \leftarrow \text{ST} \left(\frac{\lambda}{\|\mathbf{x}_j\|^2}, \mathbf{w}_j + \frac{\mathbf{x}_j^\top (\mathbf{y} - X\mathbf{w})}{\|\mathbf{x}_j\|^2} \right)$$



► 1 update is $\mathcal{O}(n)$

Variants: minimize *w.r.t.* \mathbf{w}_j with j chosen at random, or shuffle order every epoch (1 epoch = p updates)

Rem: equivalent to performing Dykstra Algorithm in the dual⁽⁷⁾

⁽⁷⁾R. J. Tibshirani. "Dykstra's Algorithm, ADMM, and Coordinate Descent: Connections, Insights, and Extensions". In: *NIPS*. 2017, pp. 517–528.

Choice of dual point

Primal-dual link at optimum:

$$\hat{\boldsymbol{\theta}} = (\mathbf{y} - X\hat{\mathbf{w}})/\lambda$$

⁽⁸⁾J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Choice of dual point

Primal-dual link at optimum:

$$\hat{\boldsymbol{\theta}} = (\mathbf{y} - X\hat{\mathbf{w}})/\lambda$$

Standard approach⁽⁸⁾: at epoch t , corresponding to primal \mathbf{w}^t and **residuals** $\mathbf{r}^t := \mathbf{y} - X\mathbf{w}^t$, choose

$$\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{res}}^t := \mathbf{r}^t/\lambda$$

⁽⁸⁾J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Choice of dual point

Primal-dual link at optimum:

$$\hat{\boldsymbol{\theta}} = (\mathbf{y} - X\hat{\mathbf{w}})/\lambda$$

Standard approach⁽⁸⁾: at epoch t , corresponding to primal \mathbf{w}^t and **residuals** $\mathbf{r}^t := \mathbf{y} - X\mathbf{w}^t$, choose

$$\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{res}}^t := \mathbf{r}^t/\lambda$$

Beware: might not be feasible!

⁽⁸⁾J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Choice of dual point

Primal-dual link at optimum:

$$\hat{\boldsymbol{\theta}} = (\mathbf{y} - X\hat{\mathbf{w}})/\lambda$$

Standard approach⁽⁸⁾: at epoch t , corresponding to primal \mathbf{w}^t and **residuals** $\mathbf{r}^t := \mathbf{y} - X\mathbf{w}^t$, choose

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

residuals rescaling

⁽⁸⁾J. Mairal. "Sparse coding for machine learning, image processing and computer vision". PhD thesis. École normale supérieure de Cachan, 2010.

Choice of dual point

Primal-dual link at optimum:

$$\hat{\boldsymbol{\theta}} = (\mathbf{y} - X\hat{\mathbf{w}})/\lambda$$

Standard approach⁽⁸⁾: at epoch t , corresponding to primal \mathbf{w}^t and **residuals** $\mathbf{r}^t := \mathbf{y} - X\mathbf{w}^t$, choose

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

residuals rescaling

► Convergence: $\lim_{t \rightarrow +\infty} \boldsymbol{\theta}_{\text{res}}^t = \hat{\boldsymbol{\theta}}$ provided $\lim_{t \rightarrow +\infty} \mathbf{w}^t = \mathbf{w}$

► $\mathcal{O}(np)$ to compute (= 1 epoch of CD)

→ rule of thumb: compute $\boldsymbol{\theta}_{\text{res}}^t$ and $\text{gap}(\mathbf{w}^t, \boldsymbol{\theta}_{\text{res}}^t)$ every 10 epochs

⁽⁸⁾J. 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 Lasso solvers

A new dual construction

Speeding up solvers

$$\hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - X\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1$$

Key property leveraged: we expect sparse solutions/small **supports**

$$\mathcal{S}_{\hat{\mathbf{w}}} := \{j \in [p] : \hat{\mathbf{w}}_j \neq 0\}$$

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

$$\hat{\mathbf{w}}_{\mathcal{S}_{\hat{\mathbf{w}}}} \in \arg \min_{w \in \mathbb{R}^{|\hat{\mathbf{w}}|_0}} \frac{1}{2} \|\mathbf{y} - X_{\mathcal{S}_{\hat{\mathbf{w}}}} w\|^2 + \lambda \|w\|_1$$

Usually $\|\hat{\mathbf{w}}\|_0 \ll p$; hence second problem much simpler

Technical details

- ▶ The primal solution/support might not be unique!
- ▶ For simplicity let us assume uniqueness, otherwise consider instead the **equicorrelation set**⁽⁹⁾:

$$E := \left\{ j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1 \right\} = \left\{ j \in [p] : \left| \mathbf{x}_j^\top \left(\frac{\mathbf{y} - X\hat{\mathbf{w}}}{\lambda} \right) \right| = 1 \right\}$$

⁽⁹⁾R. J. Tibshirani. "The lasso problem and uniqueness". In: *Electron. J. Stat.* 7 (2013), pp. 1456–1490.

Technical details

- ▶ The primal solution/support might not be unique!
- ▶ For simplicity let us assume uniqueness, otherwise consider instead the **equicorrelation set**⁽⁹⁾:

$$E := \left\{ j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1 \right\} = \left\{ j \in [p] : \left| \mathbf{x}_j^\top \left(\frac{\mathbf{y} - X\hat{\mathbf{w}}}{\lambda} \right) \right| = 1 \right\}$$

Grail of sparse solvers: identify $\mathcal{S}_{\hat{\mathbf{w}}}$, solve only on $\mathcal{S}_{\hat{\mathbf{w}}}$

Practical observation: generally $\#\mathcal{S}_{\hat{\mathbf{w}}} \ll p$

⁽⁹⁾R. J. Tibshirani. "The lasso problem and uniqueness". In: *Electron. J. Stat.* 7 (2013), pp. 1456–1490.

Speeding-up solvers

Two approaches:

- ▶ safe screening^{(10),(11)} (**backward approach**): remove feature j when it is certified that $j \notin \mathcal{S}_{\hat{w}}$
- ▶ working set⁽¹²⁾ (**forward approach**): focus on j 's very likely to be in $\mathcal{S}_{\hat{w}}$

Rem: hybrid approaches possible, e.g., strong rules⁽¹³⁾

⁽¹⁰⁾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.

⁽¹¹⁾A. Bonnefoy et al. "A dynamic screening principle for the lasso". In: *EUSIPCO. 2014*.

⁽¹²⁾T. B. Johnson and C. Guestrin. "Blitz: A Principled Meta-Algorithm for Scaling Sparse Optimization". In: *ICML. 2015*, pp. 1171–1179.

⁽¹³⁾R. Tibshirani et al. "Strong rules for discarding predictors in lasso-type problems". In: *J. R. Stat. Soc. Ser. B Stat. Methodol.* 74.2 (2012), pp. 245–266.

Duality comes into play: gap screening

We want to identify $E = \{j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1\}$...
... but we can't get it without $\hat{\mathbf{w}}$!

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

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

⁽¹⁴⁾E. Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: *J. Mach. Learn. Res.* 18.128 (2017), pp. 1–33.

Duality comes into play: gap screening

We want to identify $E = \{j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1\}$...
... but we can't get it without $\hat{\mathbf{w}}$!

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

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}} |\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 \Rightarrow |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| < 1 \Rightarrow j \notin E$$

⁽¹⁴⁾E. Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: *J. Mach. Learn. Res.* 18.128 (2017), pp. 1–33.

Duality comes into play: gap screening

We want to identify $E = \{j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1\}$...
... but we can't get it without $\hat{\mathbf{w}}$!

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

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}} |\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 \Rightarrow |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| < 1 \Rightarrow j \notin E \Rightarrow \hat{\mathbf{w}}_j = 0$$

⁽¹⁴⁾E. Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: *J. Mach. Learn. Res.* 18.128 (2017), pp. 1–33.

Duality comes into play: gap screening

We want to identify $E = \{j \in [p] : |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| = 1\}$...
... but we can't get it without $\hat{\mathbf{w}}$!

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

$$\sup_{\boldsymbol{\theta} \in \mathcal{C}} |\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 \Rightarrow |\mathbf{x}_j^\top \hat{\boldsymbol{\theta}}| < 1 \Rightarrow j \notin E \Rightarrow \hat{\mathbf{w}}_j = 0$$

Gap Safe screening rule⁽¹⁴⁾: \mathcal{C} is a ball of radius

$$\rho = \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})}$$
 centered at $\boldsymbol{\theta} \in \Delta_X$

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

⁽¹⁴⁾E. Ndiaye et al. "Gap Safe screening rules for sparsity enforcing penalties". In: *J. Mach. Learn. Res.* 18.128 (2017), pp. 1–33.

Table of Contents

Lasso basics

Speeding up Lasso solvers

A new dual construction

Back to dual choice

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

Two drawbacks of **residuals rescaling**:

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

Back to dual choice

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

Two drawbacks of **residuals rescaling**:

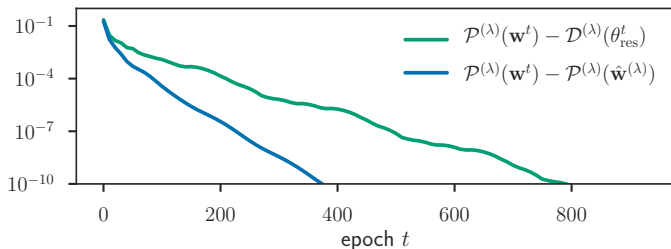
- ▶ ignores information from previous iterates
- ▶ workload "imbalanced": more efforts in primal than in dual

Back to dual choice

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

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 $\lambda = \lambda_{\max}/20$

Acceleration through residuals extrapolation⁽¹⁵⁾

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

⁽¹⁵⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

Acceleration through residuals extrapolation⁽¹⁵⁾

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

extrapolation!

→ use the same idea to infer $\lim_{t \rightarrow \infty} \mathbf{r}^t = \lambda \hat{\boldsymbol{\theta}}$

⁽¹⁵⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

Extrapolation justification

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

$$r_t = ar_{t-1} + b \quad (|a| < 1, b \in \mathbb{R}) \quad \text{with} \quad \lim_{t \rightarrow \infty} r_t = r^*$$

we have

$$r_t - r^* = a(r_{t-1} - r^*)$$

Aitken's Δ^2 : 2 unknowns, so 2 equations/3 points r_t, r_{t-1}, r_{t-2} are enough to find r^* !⁽¹⁶⁾

Rem: Aitken's rule replaces r_{n+1} by

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

⁽¹⁶⁾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

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

t	$\sum_{i=0}^t \frac{(-1)^i}{2i+1}$	Δ^2
0	1.0000	-
1	0.66667	-
2	0.86667	0.79167
3	0.72381	0.78333
4	0.83492	0.78631
5	0.74401	0.78492
6	0.82093	0.78568
7	0.75427	0.78522
8	0.81309	0.78552
9	0.76046	0.78531

Approximate Minimal Polynomial Extrapolation (AMPE)

Approximate Minimal Polynomial Extrapolation: generalization for vector autoregressive (VAR) process

$$\mathbf{r}_{k+1} - \mathbf{r}^* = A(\mathbf{r}_k - \mathbf{r}^*), \quad \text{where } A \text{ is a matrix}$$

This leads to:

$$\sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) = \sum_{k=1}^K c_k A^k (\mathbf{r}_0 - \mathbf{r}^*)$$

Approximate Minimal Polynomial Extrapolation (AMPE)

Approximate Minimal Polynomial Extrapolation: generalization for vector autoregressive (VAR) process

$$\mathbf{r}_{k+1} - \mathbf{r}^* = A(\mathbf{r}_k - \mathbf{r}^*), \quad \text{where } A \text{ is a matrix}$$

This leads to:

$$\sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) = \sum_{k=1}^K c_k A^k (\mathbf{r}_0 - \mathbf{r}^*)$$

Under the constraint: $\sum_{k=1}^K c_k = 1$, one has:

$$\sum_{k=1}^K c_k \mathbf{r}_k - \mathbf{r}^* = \left(\sum_{k=1}^K c_k A^k \right) (\mathbf{r}_0 - \mathbf{r}^*)$$

Approximate Minimal Polynomial Extrapolation (AMPE)

Approximate Minimal Polynomial Extrapolation: generalization for vector autoregressive (VAR) process

$$\mathbf{r}_{k+1} - \mathbf{r}^* = A(\mathbf{r}_k - \mathbf{r}^*), \quad \text{where } A \text{ is a matrix}$$

This leads to:

$$\sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) = \sum_{k=1}^K c_k A^k (\mathbf{r}_0 - \mathbf{r}^*)$$

Under the constraint: $\sum_{k=1}^K c_k = 1$, one has:

$$\sum_{k=1}^K c_k \mathbf{r}_k - \mathbf{r}^* = \left(\sum_{k=1}^K c_k A^k \right) (\mathbf{r}_0 - \mathbf{r}^*)$$

Consequence: approximate \mathbf{r}^* by a combination of \mathbf{r}_k 's

$$\min_{\mathbf{c}^\top \mathbf{1}_K = 1} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) \right\|, \quad \text{where } \mathbf{1}_K = (1, \dots, 1)^\top \in \mathbb{R}^K$$

(Continued)



$\min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) \right\|$ can not be solved, \mathbf{r}^* unknown!

► Note that

$$\mathbf{r}_k - \mathbf{r}_{k-1} = (\mathbf{r}_k - \mathbf{r}^*) - (\mathbf{r}_{k-1} - \mathbf{r}^*) = (A - \text{Id})A^{k-1}(\mathbf{r}_0 - \mathbf{r}^*)$$

(Continued)



$\min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) \right\|$ can not be solved, \mathbf{r}^* unknown!

► Note that


$$\mathbf{r}_k - \mathbf{r}_{k-1} = (\mathbf{r}_k - \mathbf{r}^*) - (\mathbf{r}_{k-1} - \mathbf{r}^*) = (A - \text{Id})A^{k-1}(\mathbf{r}_0 - \mathbf{r}^*)$$

► Hence, if $\text{Id} - A$ is **non singular** and $\sum_{k=1}^K c_k A^{k-1} = 0$, one must have $\sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) = 0$:

Realistic program:

$$\min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

(Continued)

 $\min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}^*) \right\|$ can not be solved, \mathbf{r}^* unknown!

► Note that

$$\mathbf{r}_k - \mathbf{r}_{k-1} = (\mathbf{r}_k - \mathbf{r}^*) - (\mathbf{r}_{k-1} - \mathbf{r}^*) = (A - \text{Id})A^{k-1}(\mathbf{r}_0 - \mathbf{r}^*)$$

► Hence, if $\text{Id} - A$ is **non singular** and $\sum_{k=1}^K c_k A^{k-1} = 0$, one must have $\sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) = 0$:

Realistic program:

$$\min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

Extrapolated dual point⁽¹⁷⁾

- ▶ Keep track of K past residuals $\mathbf{r}^t, \dots, \mathbf{r}^{t+1-K}$
- ▶ Solve (linear system resolution+normalization):

$$c^* = \arg \min_{c^\top \mathbf{1}_K=1} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

⁽¹⁷⁾M. Massias, A. Gramfort, and J. Salmon. "Celer: a Fast Solver for the Lasso with Dual Extrapolation". In: *ICML*. 2018.

Extrapolated dual point⁽¹⁷⁾

- ▶ Keep track of K past residuals $\mathbf{r}^t, \dots, \mathbf{r}^{t+1-K}$
- ▶ Solve (linear system resolution+normalization):

$$\mathbf{c}^* = \arg \min_{\mathbf{c}^\top \mathbf{1}_K=1} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

- ▶ Extrapolate:

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

⁽¹⁷⁾M. Massias, A. Gramfort, and J. Salmon. “Celer: a Fast Solver for the Lasso with Dual Extrapolation”. In: *ICML*. 2018.

Extrapolated dual point⁽¹⁷⁾

- ▶ Keep track of K past residuals $\mathbf{r}^t, \dots, \mathbf{r}^{t+1-K}$
- ▶ Solve (linear system resolution+normalization):

$$c^* = \arg \min_{c^\top \mathbf{1}_{K=1}} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

- ▶ Extrapolate:

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

- ▶ Get dual feasible point:

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

⁽¹⁷⁾M. Massias, A. Gramfort, and J. Salmon. “Celer: a Fast Solver for the Lasso with Dual Extrapolation”. In: *ICML*. 2018.

Extrapolated dual point⁽¹⁷⁾

- ▶ Keep track of K past residuals $\mathbf{r}^t, \dots, \mathbf{r}^{t+1-K}$
- ▶ Solve (linear system resolution+normalization):

$$c^* = \arg \min_{c^\top \mathbf{1}_K=1} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

- ▶ Extrapolate:

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

- ▶ Get dual feasible point:

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

⁽¹⁷⁾M. Massias, A. Gramfort, and J. Salmon. “Celer: a Fast Solver for the Lasso with Dual Extrapolation”. In: *ICML*. 2018.

Extrapolated dual point⁽¹⁷⁾

- ▶ Keep track of K past residuals $\mathbf{r}^t, \dots, \mathbf{r}^{t+1-K}$
- ▶ Solve (linear system resolution+normalization):

$$c^* = \arg \min_{c^\top \mathbf{1}_K=1} \left\| \sum_{k=1}^K c_k (\mathbf{r}_k - \mathbf{r}_{k-1}) \right\|$$

- ▶ Extrapolate:

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

- ▶ Get dual feasible point:

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

$K = 5$ is (already) enough in practice !

⁽¹⁷⁾M. Massias, A. Gramfort, and J. Salmon. "Celer: a Fast Solver for the Lasso with Dual Extrapolation". In: *ICML*. 2018.

Guarantees?

- ▶ Convergence of θ_{accel}^t ?
- ▶ Quadratic problem to solve?
Add Ridge/Tikhonov regularization if needed

Guarantees?

- ▶ Convergence of θ_{accel}^t ?
- ▶ Quadratic problem to solve?
Add Ridge/Tikhonov regularization if needed

Guarantees?

- ▶ Convergence of θ_{accel}^t ?
- ▶ Quadratic problem to solve?
Add Ridge/Tikhonov regularization if needed

Guarantees?

► Convergence of θ_{accel}^t ?

► Quadratic problem to solve?

Add Ridge/Tikhonov regularization if needed

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

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

Guarantees?

- ▶ Convergence of θ_{accel}^t ?
- ▶ Quadratic problem to solve?
Add Ridge/Tikhonov regularization if needed

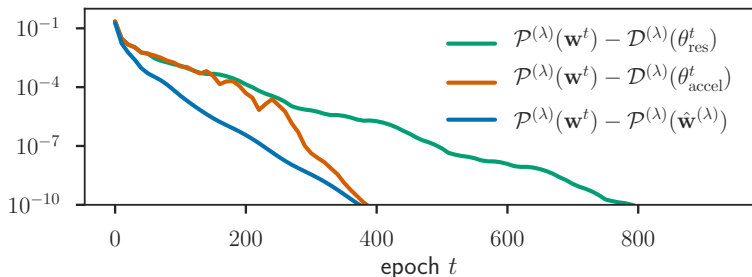
θ_{accel}^t is $\mathcal{O}(np + K^2n)$ to compute, so compute θ_{res}^t as well and pick the best, so use

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

Cost (including stopping criterion evaluation):

- ▶ classical: evaluate 1 dual point every 10 CD epoch $\approx 11np$
- ▶ new : evaluate 2 dual points every 10 CD epoch $\approx 12np$

Does it work for duality gap evaluation?



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

- ▶ θ_{res} is bad
- ▶ θ_{accel} gives a tighter bound

Which algorithm to produce w^t ?

Key assumption for extrapolation⁽¹⁸⁾: \mathbf{r}^t follows a VAR.

- ▶ True with ISTA for Lasso, once support is identified⁽¹⁹⁾ (but ISTA/FISTA slow on our statistical scenarios)

⁽¹⁸⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

⁽¹⁹⁾J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

Which algorithm to produce w^t ?

Key assumption for extrapolation⁽¹⁸⁾: \mathbf{r}^t follows a VAR.

- ▶ True with ISTA for Lasso, once support is identified⁽¹⁹⁾ (but ISTA/FISTA slow on our statistical scenarios)
- ▶ Idem for cyclic CD (though $\text{Id} - A$ is singular)

⁽¹⁸⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

⁽¹⁹⁾J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

Which algorithm to produce w^t ?

Key assumption for extrapolation⁽¹⁸⁾: \mathbf{r}^t follows a VAR.

- ▶ True with ISTA for Lasso, once support is identified⁽¹⁹⁾ (but ISTA/FISTA slow on our statistical scenarios)
- ▶ Idem for cyclic CD (though $\text{Id} - A$ is singular)

⁽¹⁸⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

⁽¹⁹⁾J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

Which algorithm to produce w^t ?

Key assumption for extrapolation⁽¹⁸⁾: \mathbf{r}^t follows a VAR.

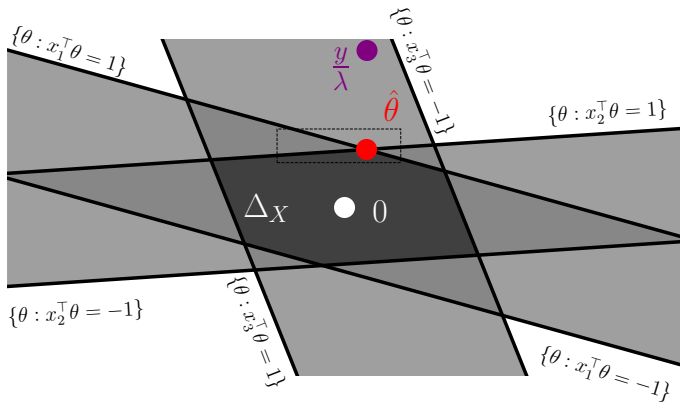
- ▶ True with ISTA for Lasso, once support is identified⁽¹⁹⁾ (but ISTA/FISTA slow on our statistical scenarios)
- ▶ Idem for cyclic CD (though $\text{Id} - A$ is singular)

Rem: Shuffle/Random CD breaks the VAR regularity

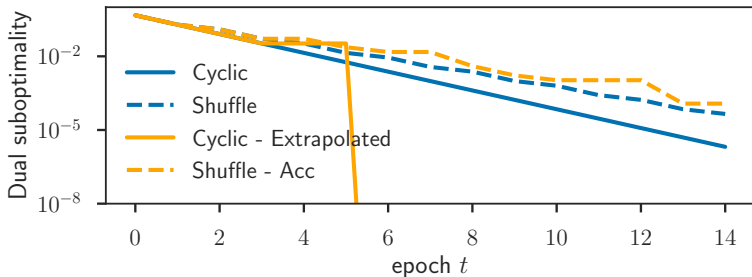
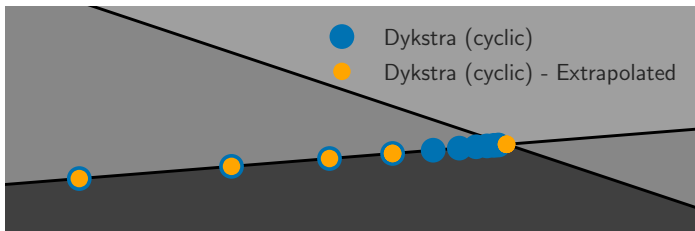
⁽¹⁸⁾D. Scieur, A. d'Aspremont, and F. Bach. "Regularized Nonlinear Acceleration". In: *NIPS*. 2016, pp. 712–720.

⁽¹⁹⁾J. Liang, J. Fadili, and G. Peyré. "Local Linear Convergence of Forward–Backward under Partial Smoothness". In: *NIPS*. 2014, pp. 1970–1978.

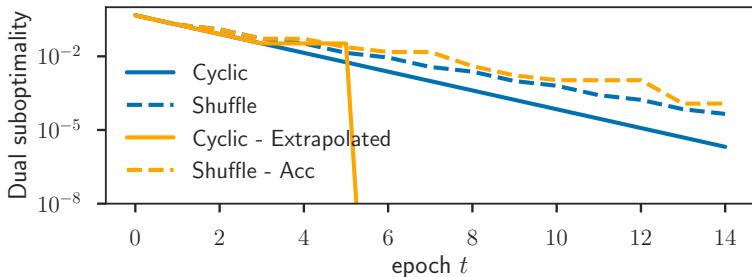
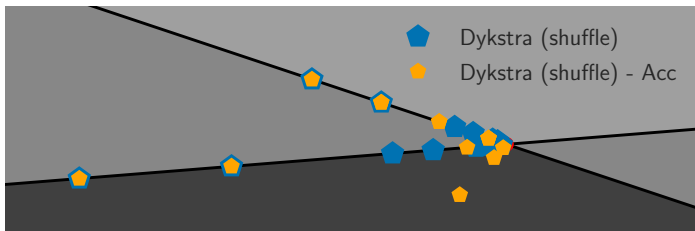
Back to toy example



Toy dual zoom: cyclic



Toy dual zoom: shuffle



Screening vs Working sets

$$|\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 - \|\mathbf{x}_j\| \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})} \Rightarrow \hat{\mathbf{w}}_j = 0$$

Screening vs Working sets

$$|\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 - \|\mathbf{x}_j\| \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})} \Rightarrow \hat{\mathbf{w}}_j = 0$$

$$\iff$$

$$d_j(\boldsymbol{\theta}) > \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})} \Rightarrow \hat{\mathbf{w}}_j = 0$$

$$\text{with } d_j(\boldsymbol{\theta}) := \frac{1 - |\mathbf{x}_j^\top \boldsymbol{\theta}|}{\|\mathbf{x}_j\|}$$

Interpretation: $d_j(\boldsymbol{\theta})$ larger than threshold \rightarrow exclude feature j

Screening vs Working sets

$$|\mathbf{x}_j^\top \boldsymbol{\theta}| < 1 - \|\mathbf{x}_j\| \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})} \Rightarrow \hat{\mathbf{w}}_j = 0$$

$$\iff$$

$$d_j(\boldsymbol{\theta}) > \sqrt{\frac{2}{\lambda^2} \text{gap}(\mathbf{w}, \boldsymbol{\theta})} \Rightarrow \hat{\mathbf{w}}_j = 0$$

$$\text{with } d_j(\boldsymbol{\theta}) := \frac{1 - |\mathbf{x}_j^\top \boldsymbol{\theta}|}{\|\mathbf{x}_j\|}$$

Interpretation: $d_j(\boldsymbol{\theta})$ larger than threshold \rightarrow exclude feature j

Alternative: Solve subproblem with small $d_j(\boldsymbol{\theta})$ only (WS)

Working/active set

Algorithm: Generic WS algorithm

Initialization: $\mathbf{w}^0 = 0 \in \mathbb{R}^p$

for $it = 1, \dots, it_{\max}$ **do**

 define working set $\mathcal{W}_{it} \subset [p]$

 approximately solve Lasso restricted to features in \mathcal{W}_{it}

 update $\mathbf{w}_{\mathcal{W}_{it}}$

3 questions for working sets

- ▶ How to prioritize features?

3 questions for working sets

- ▶ How to prioritize features? → use $d_j(\boldsymbol{\theta})$

3 questions for working sets

- ▶ How to prioritize features? \rightarrow use $d_j(\boldsymbol{\theta})$
- ▶ How many features in WS?

3 questions for working sets

- ▶ How to prioritize features? → use $d_j(\boldsymbol{\theta})$
- ▶ How many features in WS? → start small (say 100), double at each WS definition. Features cannot leave the WS

3 questions for working sets

- ▶ How to prioritize features? → use $d_j(\theta)$
- ▶ How many features in WS? → start small (say 100), double at each WS definition. Features cannot leave the WS
- ▶ What accuracy should be targeted to solve the subproblem?

3 questions for working sets

- ▶ How to prioritize features? → use $d_j(\theta)$
- ▶ How many features in WS? → start small (say 100), double at each WS definition. Features cannot leave the WS
- ▶ What accuracy should be targeted to solve the subproblem?
→ use same as required for whole problem

3 questions for working sets

- ▶ How to prioritize features? → use $d_j(\boldsymbol{\theta})$
- ▶ How many features in WS? → start small (say 100), double at each WS definition. Features cannot leave the WS
- ▶ What accuracy should be targeted to solve the subproblem?
→ use same as required for whole problem

Convergence Guaranteed!

3 questions for working sets

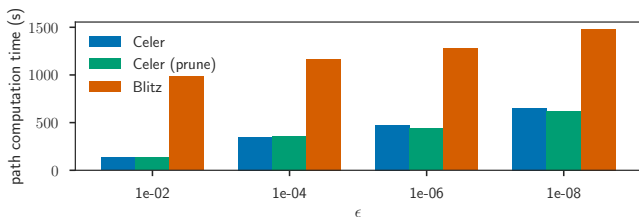
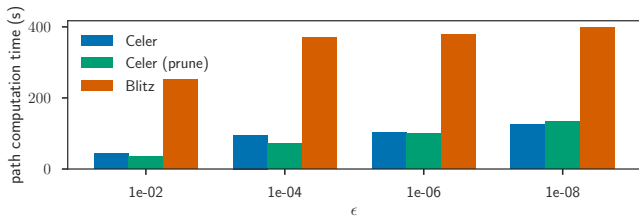
- ▶ How to prioritize features? → use $d_j(\theta)$
- ▶ How many features in WS? → start small (say 100), double at each WS definition. Features cannot leave the WS
- ▶ What accuracy should be targeted to solve the subproblem?
→ use same as required for whole problem

Convergence Guaranteed!

Rem: pruning variant also tested without much benefit (working set can decrease in size & features can leave the working set)

Comparison

State-of-the-art WS solver for sparse problems: Blitz⁽²⁰⁾



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

⁽²⁰⁾T. 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

 0 releases

 3 contributors

 BSD-3-Clause

Branch: **master** ▾

[New pull request](#)

[Create new file](#)

[Upload files](#)

[Find file](#)

[Clone or download](#) ▾



glemaitre and **mathurinm** [MRG] Make coverage great again (#21)

Latest commit cb5629e 8 days ago

 [celer](#)

Replacing nosetest with pytest (#13)

9 days ago

 [README.md](#)

celer

build passing codecov 92%

Fast algorithm to solve the Lasso with dual extrapolation

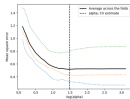
Documentation

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

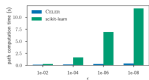
Examples gallery

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

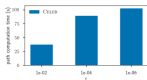
Examples Gallery ¶



Run LassoCV for cross-validation on Leukemia



Lasso path computation on Leukemia dataset



Lasso path computation on Finance/log1p

Drop-in sklearn replacement

-
- ~~from sklearn.linear_model import Lasso, LassoCV~~
 - from celer import Lasso, LassoCV
-

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)

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.

Fork me on GitHub

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=0, tol=1e-06, prune=0, fit_intercept=True)

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.

Fork me on GitHub

Conclusion

Duality matters at several levels for the Lasso:

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

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, etc.?
- ▶ Can we prove convergence of θ_{accel} ? rates?

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, etc.?
- ▶ Can we prove convergence of θ_{accel} ? rates?

Feedback welcome on the online code!



Powered with **MooseTeX**

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.
- ▶ Beck, A. and L. Tetruashvili. “On the convergence of block coordinate type methods”. In: *SIAM J. Imaging Sci.* 23.4 (2013), pp. 651–694.
- ▶ 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.

References II

- ▶ 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.
- ▶ 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. “Celer: a Fast Solver for the Lasso with Dual Extrapolation”. In: *ICML*. 2018.

References III

- ▶ Ndiaye, E. et al. “Gap Safe screening rules for sparsity enforcing penalties”. In: *J. Mach. Learn. Res.* 18.128 (2017), pp. 1–33.
- ▶ Nesterov, Y. “Efficiency of coordinate descent methods on huge-scale optimization problems”. In: *SIAM J. Optim.* 22.2 (2012), pp. 341–362.
- ▶ 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, pp. 4384–4394.
- ▶ 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. “Dykstra’s Algorithm, ADMM, and Coordinate Descent: Connections, Insights, and Extensions”. In: *NIPS*. 2017, pp. 517–528.

References IV

- ▶ 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.
- ▶ Wu, T. T. and K. Lange. “Coordinate descent algorithms for lasso penalized regression”. In: *Ann. Appl. Stat.* (2008), pp. 224–244.

Dykstra Algorithm

Goal: find the projection of z on the intersection of convex set C_1, \dots, C_p , providing the projections $\Pi_{C_1}, \dots, \Pi_{C_p}$ are available.

Algorithm: DYKSTRA'S ALTERNATING PROJECTION

input : $\Pi_{C_1}, \dots, \Pi_{C_p}, z$

init : $\theta = z, q_1 = 0, \dots, q_p = 0$

for $t = 1, \dots$ **do**

for $j = 1, \dots, p$ **do**

$\tilde{\theta} \leftarrow \theta + q_j$

$\theta \leftarrow \Pi_{C_j}(\tilde{\theta})$

$q_j \leftarrow \tilde{\theta} - \theta$

return θ

Similarities with correlation screening^{(21), (22)}

$$d_j(\boldsymbol{\theta}) := \frac{1 - |\mathbf{x}_j^\top \boldsymbol{\theta}|}{\|\mathbf{x}_j\|}$$

⁽²¹⁾ 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.

⁽²²⁾ S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: *NIPS*. 2017, pp. 4384–4394.

Similarities with correlation screening^{(21), (22)}

$$d_j(\boldsymbol{\theta}) := \frac{1 - |\mathbf{x}_j^\top \boldsymbol{\theta}|}{\|\mathbf{x}_j\|}$$

Lasso case with $\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{res}}$ and normalized \mathbf{x}_j 's:

$$1 - d_j(\boldsymbol{\theta}) \propto |\mathbf{x}_j^\top \mathbf{r}^t|$$

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

⁽²¹⁾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.

⁽²²⁾S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: *NIPS*. 2017, pp. 4384–4394.

Similarities with correlation screening^{(21), (22)}

$$d_j(\boldsymbol{\theta}) := \frac{1 - |\mathbf{x}_j^\top \boldsymbol{\theta}|}{\|\mathbf{x}_j\|}$$

Lasso case with $\boldsymbol{\theta} = \boldsymbol{\theta}_{\text{res}}$ and normalized \mathbf{x}_j 's:

$$1 - d_j(\boldsymbol{\theta}) \propto |\mathbf{x}_j^\top \mathbf{r}^t|$$

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

BUT our strength is that we can use any $\boldsymbol{\theta}$, in particular $\boldsymbol{\theta}_{\text{accel}}$

⁽²¹⁾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.

⁽²²⁾S. Stich, A. Raj, and M. Jaggi. "Safe Adaptive Importance Sampling". In: *NIPS*. 2017, pp. 4384–4394.