# Faster and still safe: Combining screening techniques and structured dictionaries to accelerate the Lasso

Cássio F. DANTAS,

cassio.fraga-dantas@inria.fr,

Rémi GRIBONVAL

remi.gribonval@inria.fr

**Accelerate the Lasso optimization**
by combining two strategies :

1) Safe Screening Rules

2) Fast Structured Dictionaries

*Inria*
inventors for the digital world

# Contents

*Inria*
inventors for the digital world

# 01

## Context

# Lasso problem

The l1-regularized least squares.

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \underbrace{\frac{1}{2}\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1}_{P(\boldsymbol{\beta})}$$

Denoting :

- $\mathbf{y} \in \mathbb{R}^N$        the observation vector;

- $\mathrm{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_K] \in \mathbb{R}^{N \times K}$    the design matrix (or dictionary);

- $\boldsymbol{\beta} \in \mathbb{R}^{\boldsymbol{K}}$        the sparse representation vector;

- $\lambda > 0$        parameter controlling the sparsity of the solution.

*Inria*
inventors for the digital world

# Dual Lasso

Dual formulation of the Lasso problem :

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Delta_{\mathbf{x}}}{\operatorname{argmax}} \quad \underbrace{\frac{1}{2}\|\mathbf{y}\|_2^2 - \frac{\lambda^2}{2}\left\|\boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda}\right\|_2^2}_{D(\boldsymbol{\theta})}$$

Denoting :

- $\boldsymbol{\theta} \in \mathbb{R}^N$        the dual variable;

- $\Delta_{\mathbf{x}} = \{\boldsymbol{\theta} \in \mathbb{R}^N : \|\mathbf{X}^T\boldsymbol{\theta}\|_\infty \leq 1\}$    the feasible set;

*Inria*
inventors for the digital world

# Motivation

- **Iterative algorithms** are often used to solve the Lasso problem.

- Exemple : ISTA (Iterative Shrinkage-Thresholding Algorithm)

$$\textbf{while} \text{ not converged } \textbf{do}$$
$$\boldsymbol{\beta}_{t+1} \leftarrow \text{ST}_{\frac{\lambda}{L_t}}\left(\boldsymbol{\beta}_t + \frac{1}{L_t}\mathbf{X}_t^T(\mathbf{y} - \mathbf{X}_t\boldsymbol{\beta}_t)\right)$$

- Two **matrix-vector multiplications** at each iteration.

Quadratic complexity!

Can it be reduced?

*Innía*
inventors for the digital world

# 02

Fast Structured Dictionaries

inventors for the digital world

# Structure ⇒ Acceleration

**Accelerate matrix-vector multiplications**

Constrain the dictionary matrix to have a certain type of structure.

Examples :

- Kronecker product

- Sparse factors

- Circulant factors

- (…)

# Structured Approximation

If the dictionary matrix $\mathbf{X}$ is not structured, find a structured approximation $\tilde{\mathbf{X}}$.

$$\tilde{\mathbf{X}} = \mathbf{X} + \mathbf{E},$$

where $\mathbf{E}$ is the approximation error matrix and $\mathbf{e}_j$ is its $j$-th column.

Inría
inventors for the digital world

# Algorithm (high level view)

1) Start Lasso optimization by using the structured $\tilde{\mathbf{X}}$, to take advantage of its reduced multiplication cost.

2) As the algorithm approaches the solution, switch back to the original dictionary $\mathbf{X}$.

**while** switching criterion not met **do**

$$\boldsymbol{\beta}_{t+1} \leftarrow \text{ST}_{\frac{\lambda}{L_t}}\left(\boldsymbol{\beta}_t + \frac{1}{L_t}\tilde{\mathbf{X}}_t^T(\mathbf{y} - \tilde{\mathbf{X}}_t\boldsymbol{\beta}_t)\right)$$

**while** not converged **do**

$$\boldsymbol{\beta}_{t+1} \leftarrow \text{ST}_{\frac{\lambda}{L_t}}\left(\boldsymbol{\beta}_t + \frac{1}{L_t}\mathbf{X}_t^T(\mathbf{y} - \mathbf{X}_t\boldsymbol{\beta}_t)\right)$$

*Inría*
inventors for the digital world

# Algorithm (high level view)

1) Start Lasso optimization by using the structured $\tilde{\mathbf{X}}$ , to take advantage of its reduced multiplication cost.

2) As the algorithm approaches the solution, switch back to the original dictionary $\mathbf{X}$.

**while** switching criterion not met **do**

$$\boldsymbol{\beta}_{t+1} \leftarrow \text{ST}_{\frac{\lambda}{L_t}}\left( \boldsymbol{\beta}_t + \frac{1}{L_t}\tilde{\mathbf{X}}_t^T(\mathbf{y} - \tilde{\mathbf{X}}_t\boldsymbol{\beta}_t) \right)$$

**while** not converged **do**

$$\boldsymbol{\beta}_{t+1} \leftarrow \text{ST}_{\frac{\lambda}{L_t}}\left( \boldsymbol{\beta}_t + \frac{1}{L_t}\mathbf{X}_t^T(\mathbf{y} - \mathbf{X}_t\boldsymbol{\beta}_t) \right)$$

Inría
inventors for the digital world

# 03

## Safe Screening Rules

*Inria*
inventors for the digital world
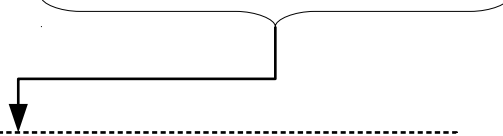
# Safe Screening

- Rules for identifying inactive dictionary atoms, before solving the problem.
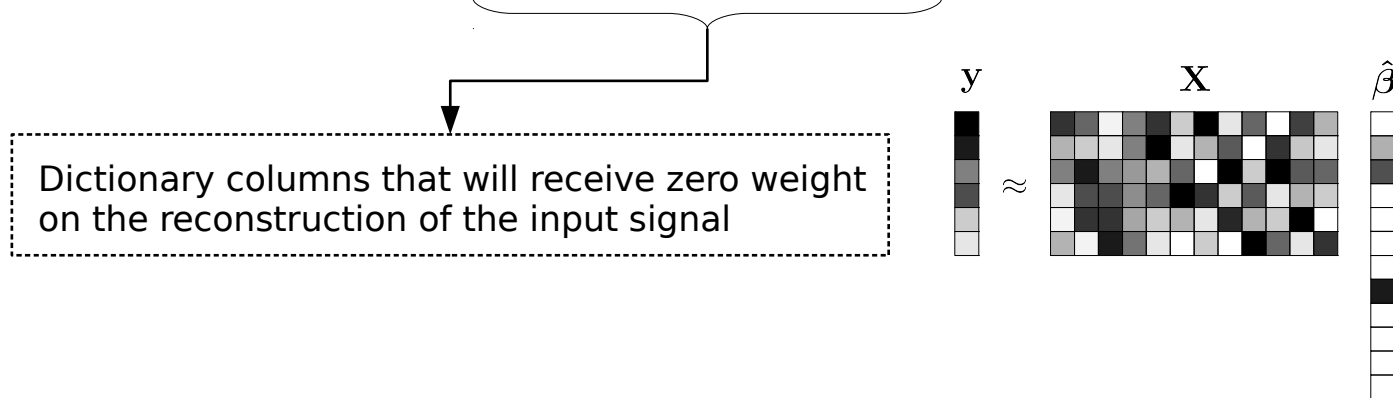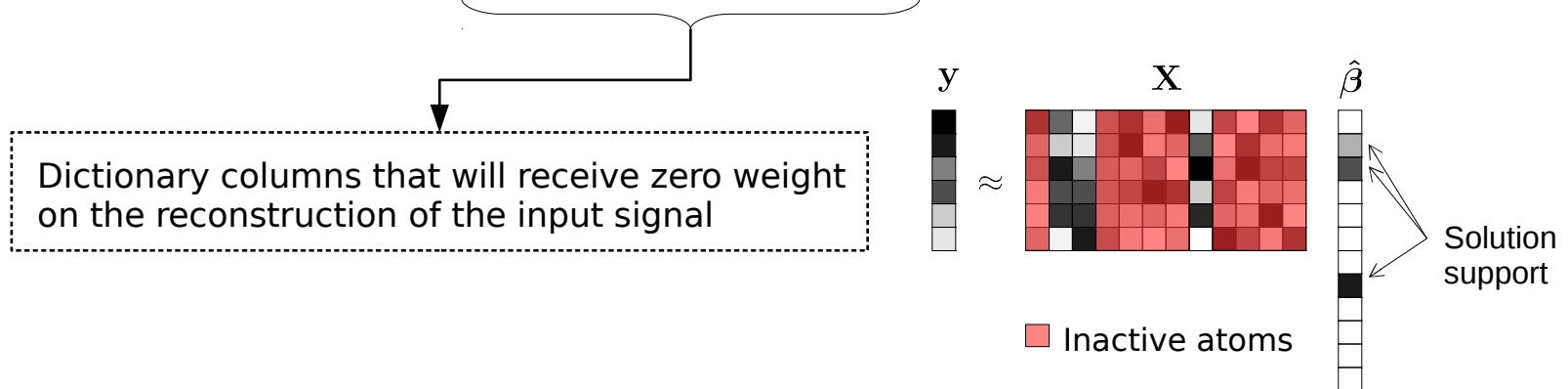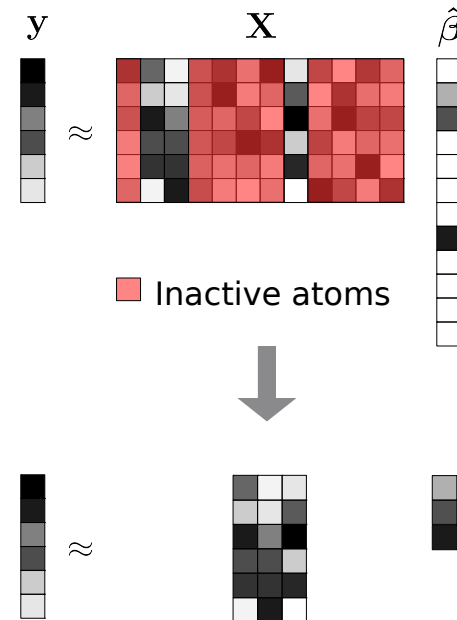
# Safe Screening

- Rules for identifying inactive dictionary atoms, before solving the problem.

Dictionary columns that will receive zero weight on the reconstruction of the input signal

# Safe Screening

- Rules for identifying inactive dictionary atoms, before solving the problem.

Dictionary columns that will receive zero weight on the reconstruction of the input signal



$$\mathbf{y} \approx \mathbf{X} \hat{\boldsymbol{\beta}}$$

# Safe Screening

- Rules for identifying inactive dictionary atoms, before solving the problem.

Dictionary columns that will receive zero weight on the reconstruction of the input signal



$y \approx X \hat{\beta}$

Solution support

⬛ Inactive atoms

# Safe Screening

- Rules for identifying inactive dictionary atoms, before solving the problem.

Dictionary columns that will receive zero weight on the reconstruction of the input signal



$\mathbf{y} \qquad \mathbf{X} \qquad \hat{\boldsymbol{\beta}}$

■ Inactive atoms

- We can eliminate such atoms.

- Zero risk of false eliminations!

Inria
inventors for the digital world

# Screening Test

- Function $\mu\left(\mathbf{x}_j\right)$ of the atom $\mathbf{x}_j$

$$\mu(\mathbf{x}_j) < 1 \quad \implies \quad \mathbf{x}_j \quad \text{is surely inactive.}$$

# Screening Test

- Function $\mu(\mathbf{x}_j)$ of the atom $\mathbf{x}_j$

$$\mu(\mathbf{x}_j) < 1 \quad \implies \quad \mathbf{x}_j \quad \text{is surely inactive.}$$

$$\mathbf{X} \longrightarrow \boxed{\begin{array}{c} \mu(\mathbf{x}_j) \\ \forall j \end{array}} \longrightarrow$$
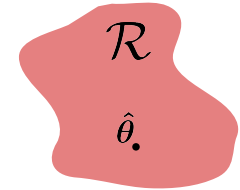
Rejection set:

$$\mathcal{A}^{\mathrm{c}} = \{j \in \{1, \dots, K\} : \mu(\mathbf{x}_j) < 1\}$$

Preserved set:

$$\mathcal{A} = \{j \in \{1, \dots, K\} : \mu(\mathbf{x}_j) \geq 1\}$$

# Screening Test – In practice

Given a region $\mathcal{R}$ (**safe region**) which contains $\hat{\boldsymbol{\theta}}$.
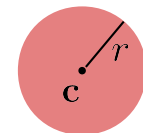
$$\mu_{\mathcal{R}}(\mathbf{x}_j) = \max_{\boldsymbol{\theta} \in \mathcal{R}} |\mathbf{x}_j^T \boldsymbol{\theta}|$$

**Sphere test**

Safe region is a closed l2-ball with center **c** and radius $r$ : $\quad \mathcal{R} = B(\mathbf{c}, r)$
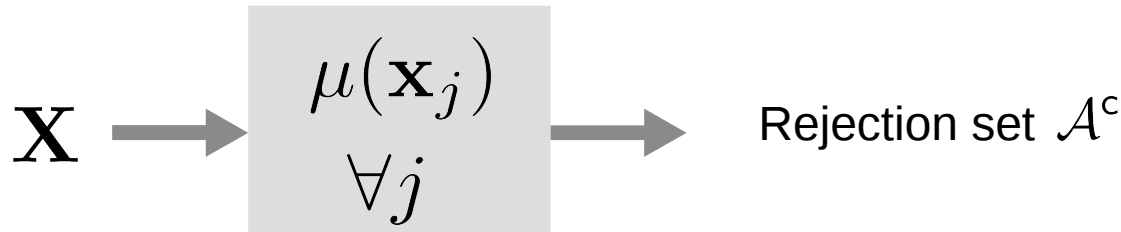
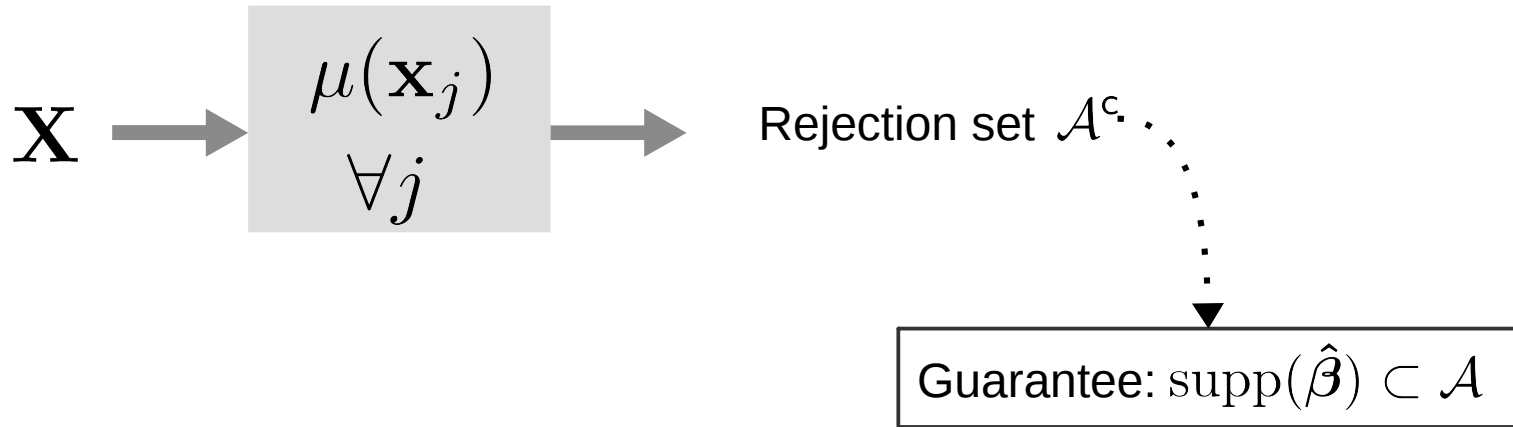$$\mu_{B(\mathbf{c}, r)}(\mathbf{x}_j) = |\mathbf{x}_j^T \mathbf{c}| + r\|\mathbf{x}_j\|_2.$$

*Inria*
inventors for the digital world

# 04

## Screening Rules with Approximate Dictionaries

*Inria*
inventors for the digital world

$$\hat{\boldsymbol{\beta}} = \operatorname*{argmin}_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$$

$$\mathbf{X} \longrightarrow \boxed{\begin{array}{c} \mu(\mathbf{x}_j) \\ \forall j \end{array}} \longrightarrow \text{Rejection set } \mathcal{A}^{\mathsf{c}}$$

$$\boxed{\text{Guarantee: } \operatorname{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}}$$

*Innia* — inventors for the digital world

$$\hat{\boldsymbol{\beta}} = \operatorname*{argmin}_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

$\mathbf{X}$ $\longrightarrow$ $\boxed{\begin{array}{c} \mu(\mathbf{x}_j) \\ \forall j \end{array}}$ $\longrightarrow$ Rejection set $\mathcal{A}^{\mathsf{c}}$

Guarantee: $\operatorname{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}$

Inría

inventors for the digital world

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\arg\min} \frac{1}{2}\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$$

$\mathbf{X}$ ⟶ $\mu(\mathbf{x}_j)$ $\forall j$ ⟶ Rejection set $\mathcal{A}^{\mathsf{c}}$

Guarantee: $\mathrm{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}$

$$\tilde{\mathbf{x}}_j = \mathbf{x}_j + \mathbf{e}_j$$

$\tilde{\mathbf{X}}$ ⟶ $\tilde{\mu}(\tilde{\mathbf{x}}_j)$ $\forall j$ ⟶ Rejection set $\mathcal{A}^{\mathsf{c}}$

*Inría*
inventors for the digital world

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\mathrm{argmin}} \frac{1}{2} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

$\mathbf{X} \longrightarrow$ $\mu(\mathbf{x}_j)$ $\forall j$ $\longrightarrow$ Rejection set $\mathcal{A}^{\mathsf{c}}$

Guarantee: $\mathrm{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}$

$$\tilde{\mathbf{x}}_j = \mathbf{x}_j + \mathbf{e}_j$$

$\tilde{\mathbf{X}} \longrightarrow$ $\tilde{\mu}(\tilde{\mathbf{x}}_j)$ $\forall j$ $\longrightarrow$ Rejection set $\mathcal{A}^{\mathsf{c}}$

*Inría*
inventors for the digital world

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2}\|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$$

$\mathbf{X}$ $\longrightarrow$ $\begin{array}{c} \mu(\mathbf{x}_j) \\ \forall j \end{array}$ $\longrightarrow$ Rejection set $\mathcal{A}^{\mathrm{c}}$

Guarantee: $\operatorname{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}$

$$\tilde{\mathbf{x}}_j = \mathbf{x}_j + \mathbf{e}_j$$

$\tilde{\mathbf{X}}$ $\longrightarrow$ $\begin{array}{c} \tilde{\mu}(\tilde{\mathbf{x}}_j) \\ \forall j \end{array}$ $\longrightarrow$ Rejection set $\mathcal{A}^{\mathrm{c}}$

$$\|\mathbf{e}_j\|_2 \quad \forall j$$

*Inria*
inventors for the digital world

# Extending sphere tests

Suppose a safe sphere $B(\mathbf{c}, r)$ given.

**Sphere test :** $\qquad \mu_{B(\mathbf{c}, r)}(\mathbf{x}_j) = |\mathbf{x}_j^T \mathbf{c}| + r\|\mathbf{x}_j\|_2.$

A certain *« security margin »* must be added to account for the atom approximation error.

*Inria*
inventors for the digital world

# Extending sphere tests

Suppose a safe sphere $B(\mathbf{c}, r)$ given.

**Sphere test :**
$$\mu_{B(\mathbf{c},r)}(\mathbf{x}_j) = \left|\mathbf{x}_j^T \mathbf{c}\right| + r\|\mathbf{x}_j\|_2.$$

A certain *« security margin »* must be added to account for the atom approximation error.

**Sphere test with approximate dictionary :**

$$\tilde{\mu}_{B(\mathbf{c},r)}(\tilde{\mathbf{x}}_j) = \left|\tilde{\mathbf{x}}_j^T \mathbf{c}\right| + \|\mathbf{e}_j\|_2\|\mathbf{c}\|_2 + r\|\mathbf{x}_j\|_2.$$

*Inria*
inventors for the digital world

# Obtaining a safe sphere

**GAP safe sphere :**

Given a primal-dual estimation $(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ at iteration $t$.

$$\mathbf{c} = \boldsymbol{\theta}_t$$

$$r = \frac{1}{\lambda} \sqrt{2G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)}$$

$t = 1$

with $G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t) = P(\boldsymbol{\beta}_t) - D(\boldsymbol{\theta}_t)$ the duality gap at iteration $t$.

$P(\boldsymbol{\beta}_t)$

$D(\boldsymbol{\theta}_t)$

$t$

*Inria*
inventors for the digital world

# Obtaining a safe sphere

**GAP safe sphere :**

Given a primal-dual estimation $(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ at iteration $t$.

$$\mathbf{c} = \boldsymbol{\theta}_t$$

$$r = \frac{1}{\lambda}\sqrt{2G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)}$$

$t = 2$



with $G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t) = P(\boldsymbol{\beta}_t) - D(\boldsymbol{\theta}_t)$ the duality gap at iteration $t$.

*Inría*
inventors for the digital world

# Obtaining a safe sphere

**GAP safe sphere :**

Given a primal-dual estimation $(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ at iteration $t$.

$$\mathbf{c} = \boldsymbol{\theta}_t \qquad\qquad t = 6$$

$$r = \frac{1}{\lambda}\sqrt{2G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)}$$

with $G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t) = P(\boldsymbol{\beta}_t) - D(\boldsymbol{\theta}_t)$ the duality gap at iteration $t$.

inventors for the digital world

# Obtaining a safe sphere

**GAP safe sphere :**

Given a primal-dual estimation $(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ at iteration $t$.

$$\mathbf{c} = \boldsymbol{\theta}_t \qquad\qquad t = 10$$

$$r = \frac{1}{\lambda}\sqrt{2G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)}$$

with $G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t) = P(\boldsymbol{\beta}_t) - D(\boldsymbol{\theta}_t)$ the duality gap at iteration $t$.

# Obtaining a safe sphere

**GAP safe sphere:**

Given a primal-dual estimation $(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ at iteration $t$ .

$$\mathbf{c} = \boldsymbol{\theta}_t$$

$$r = \frac{1}{\lambda}\sqrt{2G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)}$$

**GAP safe sphere with approximate dictionary:**

$G(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ cannot be calculated, since $P(\boldsymbol{\beta}_t) = \frac{1}{2}\|\mathbf{X}\boldsymbol{\beta}_t - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}_t\|_1$ depends on $\mathbf{X}$ .

Instead, we use a modified primal $\quad \tilde{P}(\boldsymbol{\beta}_t) = \frac{1}{2}\|\tilde{\mathbf{X}}\boldsymbol{\beta}_t - \mathbf{y}\|_2^2 + \lambda\|\boldsymbol{\beta}_t\|_1$

$$\tilde{\mathbf{c}} = \tilde{\boldsymbol{\theta}}_t$$

$$\tilde{r} = \frac{1}{\lambda}\sqrt{2\tilde{G}(\boldsymbol{\beta}_t, \tilde{\boldsymbol{\theta}}_t) + \textcolor{red}{2\delta'}}$$

with $\textcolor{red}{\delta' = \|\tilde{\boldsymbol{\rho}}_t\|_2\|\mathbf{E}\|\|\boldsymbol{\beta}_t\|_2 + \frac{1}{2}\|\mathbf{E}\|^2\|\boldsymbol{\beta}_t\|_2^2}$

*Inría*
inventors for the digital world

# Dynamic screening



Guarantee 1: $\hat{\boldsymbol{\theta}} \in \mathcal{R}_t$

$(\boldsymbol{\beta}_t, \boldsymbol{\theta}_t)$ → GAP sphere → Safe region $\mathcal{R}_t$

$\mathbf{X}$ → $\mu(\mathbf{x}_j)$ $\forall j$ → Rejection set $\mathcal{A}^{\mathsf{c}}$

Guarantee 2: $\mathrm{supp}(\hat{\boldsymbol{\beta}}) \subset \mathcal{A}$

# Extended dynamic screening

We now have safe screening rules that manipulate an approximate dictionary.

But, what's the impact of the numerous security margins? Is it still worth it?

# 05

## Results

*Inria*
inventors for the digital world

2500 x 10000 dictionary

# Running times per iteration

- **Less inactive atoms are identified by the extended screening.**

- **BUT, structured dictionary makes the initial iterations much faster.**

*Inría*
inventors for the digital world

# Running times per iteration

- **Less inactive atoms are identified by the extended screening.**

- **BUT, structured dictionary makes the initial iterations much faster.**



A-GAP: $\|\mathbf{e}_j\| = 10^{-2}$

$\tilde{\mathbf{X}}$

$\mathbf{X}$

Iteration

$\log_{10}(\lambda/\lambda_{\max})$

Inría
inventors for the digital world

# Running times per iteration

# 06

## Conclusion

Inria
inventors for the digital world

# Conclusion

- The proposed approach combines screening rules and fast approximate dictionaries.

- It reduces even further the execution time w.r.t screening rules alone.

**Potential extensions**

- Other region types (e.g. domes)

- Other problems than Lasso (e.g. Group-Lasso, Regularized Logistic Regression)

*Inria*
inventors for the digital world

# Thank you!

Questions?

Contact me: cassio.fraga-dantas@inria.fr

*Inría*
inventors for the digital world

# Screening test – Details

Dual formulation of the Lasso problem :

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\mathrm{argmax}} \quad \frac{1}{2}\|\mathbf{y}\|_2^2 - \frac{\lambda^2}{2}\left\|\boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda}\right\|_2^2$$

$$\mathrm{s.t.} \qquad \|\mathbf{X}^T\boldsymbol{\theta}\|_\infty \leq 1.$$



Projection problem !

At the dual solution $\hat{\boldsymbol{\theta}}$ :

- ➢ Constraints on $\mathbf{x}_1$ and $\mathbf{x}_2$ are active, i.e.

$$|\mathbf{x}_1^T\hat{\boldsymbol{\theta}}| = 1, \quad |\mathbf{x}_2^T\hat{\boldsymbol{\theta}}| = 1$$

- ➢ Constraints on $\mathbf{x}_3$ is inactive, i.e.

$$|\mathbf{x}_3^T\hat{\boldsymbol{\theta}}| < 1$$

inventors for the digital world

# Screening test – Details

- Every dictionary atom for which $|\mathbf{x}_j^T \hat{\boldsymbol{\theta}}| < 1$ is **inactive**.

- Then, simply calculate $|\mathbf{x}_j^T \hat{\boldsymbol{\theta}}|$ for all $j$ and discard all atoms for which the result is smaller than 1.

⚠ Dual solution $\hat{\boldsymbol{\theta}}$ is not known.

✅ Identify a region $\mathcal{R}$ (**safe region**) which contains $\hat{\boldsymbol{\theta}}$ .

$$\mathcal{R}$$
$$\hat{\theta}$$

Sufficient condition :    $\forall \quad \boldsymbol{\theta} \in \mathcal{R}, \quad |\mathbf{x}_j^T \boldsymbol{\theta}| < 1 \quad \Longrightarrow \quad \mathbf{x}_j$ is inactive

*Inría*
inventors for the digital world

# Swithing criterion

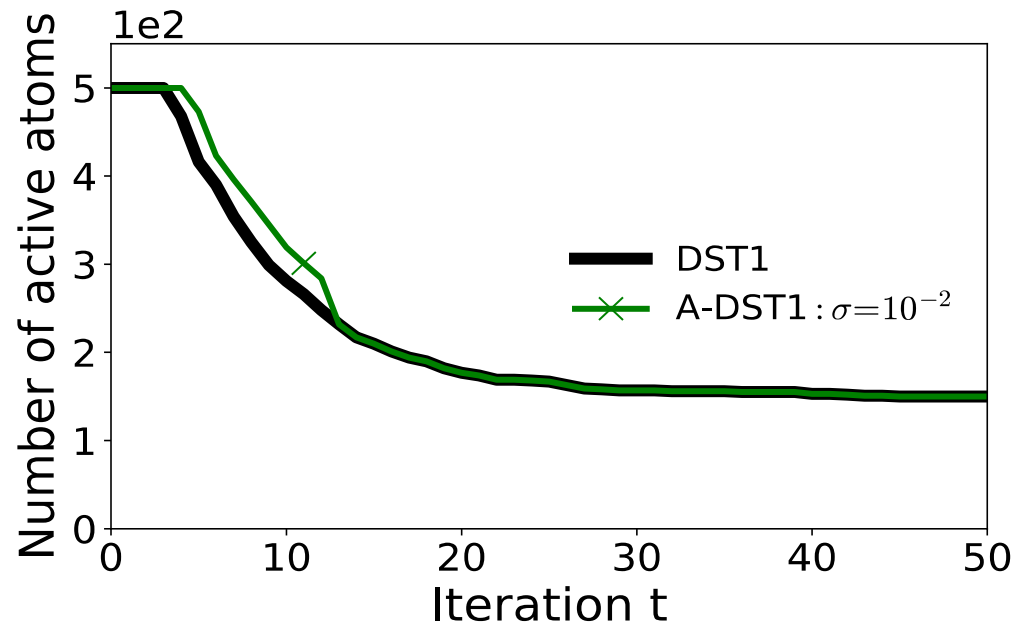Reasons to switch back from $\tilde{X}$ to $X$ :

- **Convergence**: to avoid converging to the solution of the approximate problem.

  - The higher the approximation error, the sooner we need to switch.

- **Screening ratio**: the number of active atoms may become so small that the use of $\tilde{X}$ does not pay off anymore.
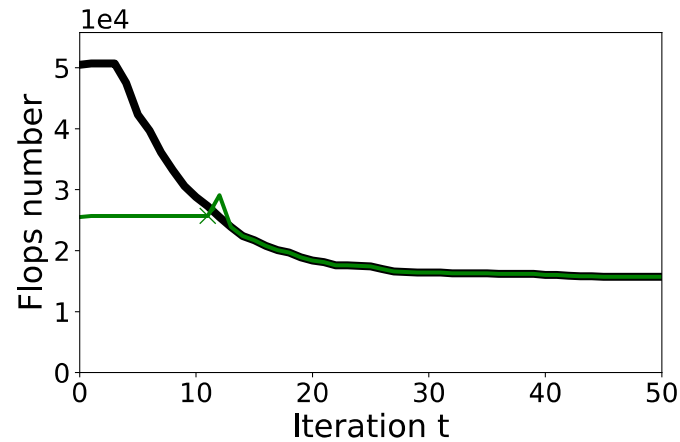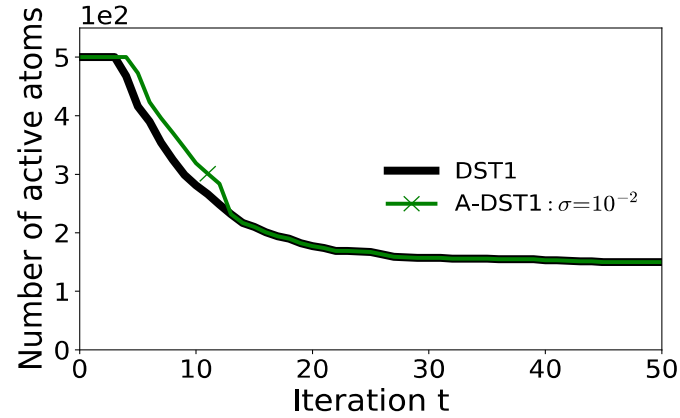
# Comparison

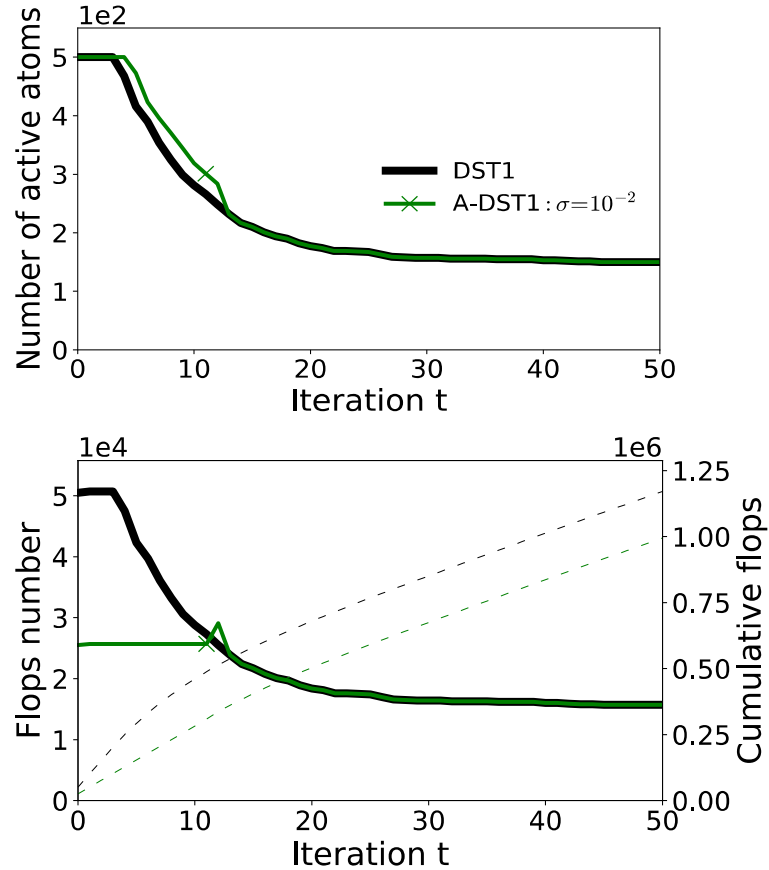**Less inactive atoms are identified by the extended screening.**

# Swithing criterion

# Complexity reduction

# Complexity reduction

# Impact of the Approximation Error