

# Screen & Relax for Sparse Support Identification

---

**Théo Guyard**<sup>\*</sup>, Cédric Herzet<sup>◇</sup>, Clément Elvira<sup>†</sup>

<sup>\*</sup>Inria, Rennes, France

<sup>◇</sup>Ensaï, Rennes, France

<sup>†</sup>CentraleSupélec, Rennes, France

SMAI-MODE days

INSA Lyon, France

March 28th, 2024

## General Context

---

## Sparse optimization

- Minimize a loss with a sparse optimizer
- Applications in signal processing, machine learning, statistics, etc...

## Sparse optimization

- **Minimize** a loss with a **sparse** optimizer
- Applications in signal processing, machine learning, statistics, etc...

### Problem of interest

$$x^* \in \operatorname{argmin}_x f(Ax) + \underbrace{\lambda \|x\|_1 + h(x)}_{g(x)}$$

## Sparse optimization

- **Minimize** a loss with a **sparse** optimizer
- Applications in signal processing, machine learning, statistics, etc...

### Problem of interest

$$x^* \in \operatorname{argmin}_x f(Ax) + \underbrace{\lambda \|x\|_1 + h(x)}_{g(x)}$$

## Working hypotheses

- $f(\cdot)$  and  $h(\cdot)$  are **proper**, **closed** and **convex** functions
- $f(\cdot)$  and  $h(\cdot)$  are **differentiable** with **Lipschitz-continuous** gradient
- $h(\cdot)$  is **separable**
  - With additional mild assumptions ensuring non-degeneracy of the problem

$$x^* \in \operatorname{argmin}_x f(Ax) + g(x)$$

## Solution methods

- Composite objective: smooth + non-smooth-separable
- **First-order** methods accessing  $\nabla f(\cdot)$ ,  $\partial g(\cdot)$ ,  $\operatorname{prox}_{f/g}(\cdot)$ , ...
  - Proximal gradient descent
  - Coordinate descent
  - Alternating direction method of multipliers
  - ...

# Solving sparse problems

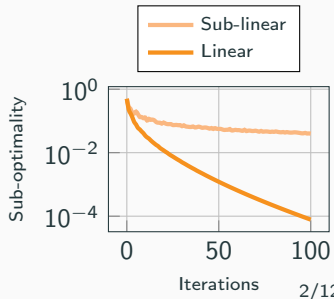
$$x^* \in \operatorname{argmin}_x f(Ax) + g(x)$$

## Solution methods

- Composite objective: smooth + non-smooth-separable
- **First-order** methods accessing  $\nabla f(\cdot)$ ,  $\partial g(\cdot)$ ,  $\operatorname{prox}_{f/g}(\cdot)$ , ...
  - Proximal gradient descent
  - Coordinate descent
  - Alternating direction method of multipliers
  - ...

## Convergence rates

- Sub-linear:  $P(x^{(k)}) - P(x^*) \leq C/k^\gamma$
- Linear:  $P(x^{(k)}) - P(x^*) \leq Ce^{-\gamma k}$ 
  - Asymptotically (Peyré et al., 2015)
  - Strong convexity (Aujol et al., 2023)



# Sparse structure exploitation

Variable with many useless and few informative entries



# Sparse structure exploitation

Variable with many useless and few informative entries

## Screening tests

- Identify zeros in  $x^*$
- Dimensionality shrinking
- Computational savings

# Sparse structure exploitation

Variable with many useless and few informative entries

## Screening tests

- Identify zeros in  $x^*$
- Dimensionality shrinking
- Computational savings

## Relaxing tests

- Identify non-zeros in  $x^*$
- Objective smoothing
- Super-linear convergence

# Sparse structure exploitation

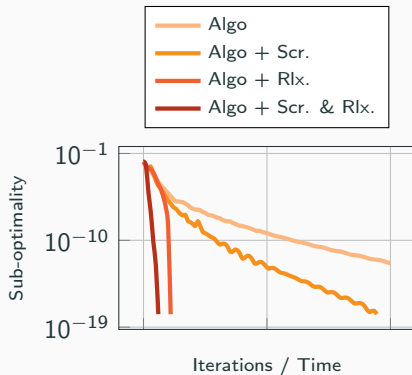
Variable with many useless and few informative entries

## Screening tests

- Identify zeros in  $x^*$
- Dimensionality shrinking
- Computational savings

## Relaxing tests

- Identify non-zeros in  $x^*$
- Objective smoothing
- Super-linear convergence

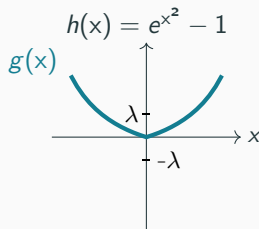
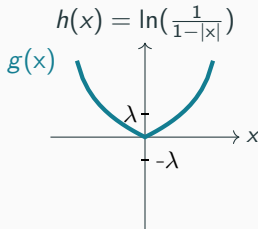
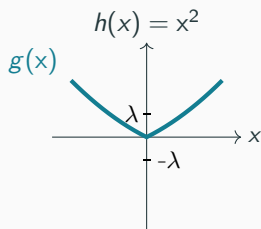


# Screening and Relaxing Tests

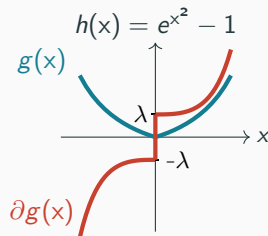
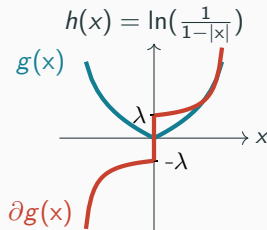
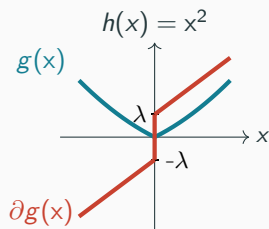
---

Characterize the nullity in  $x^*$  from  $\partial g(x^*)$

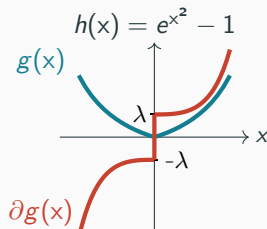
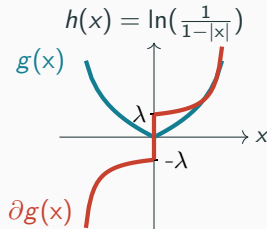
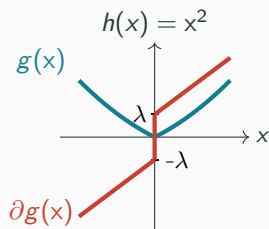
Characterize the nullity in  $x^*$  from  $\partial g(x^*)$



Characterize the nullity in  $x^*$  from  $\partial g(x^*)$



Characterize the nullity in  $x^*$  from  $\partial g(x^*)$



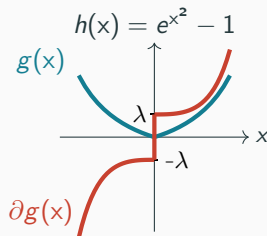
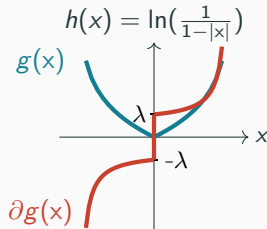
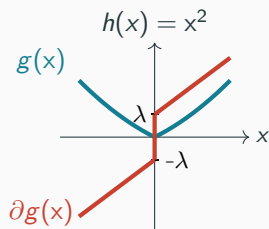
**Geometrical screening and relaxing test**

$$\partial_i g(x^*) \subset [-\lambda, \lambda] \iff x_i^* = 0$$

$$\partial_i g(x^*) \not\subset [-\lambda, \lambda] \iff x_i^* \neq 0$$



Characterize the nullity in  $x^*$  from  $\partial g(x^*)$



**Geometrical screening and relaxing test**

$$\partial_i g(x^*) \subset [-\lambda, \lambda] \iff x_i^* = 0$$

$$\partial_i g(x^*) \not\subset [-\lambda, \lambda] \iff x_i^* \neq 0$$

$\partial g(x^*)$  is not available

# Duality to rescue

Characterize the nullity in  $x^*$  from the dual problem

# Duality to rescue

Characterize the nullity in  $x^*$  from the dual problem

**Dual problem**

$$u^* \in \operatorname{argmax}_u -f^*(-u) - g^*(A^T u)$$

$$a_i^T u^* \in \partial_i g(x^*)$$

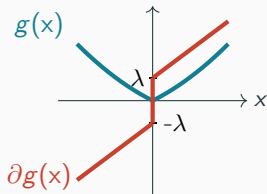
# Duality to rescue

Characterize the nullity in  $x^*$  from the dual problem

**Dual problem**

$$u^* \in \operatorname{argmax}_u -f^*(-u) - g^*(A^T u)$$

$$a_i^T u^* \in \partial_i g(x^*)$$



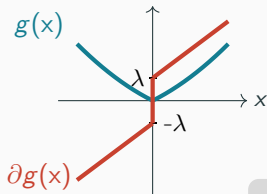
# Duality to rescue

Characterize the nullity in  $x^*$  from the dual problem

## Dual problem

$$u^* \in \operatorname{argmax}_u -f^*(-u) - g^*(A^T u)$$

$$a_i^T u^* \in \partial_i g(x^*)$$



## Screening and relaxing test

$$|a_i^T u^*| \leq \lambda \iff x_i^* = 0$$

$$|a_i^T u^*| > \lambda \iff x_i^* \neq 0$$

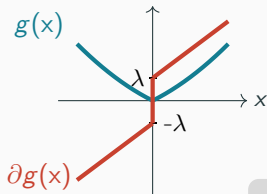
# Duality to rescue

Characterize the nullity in  $x^*$  from the dual problem

## Dual problem

$$u^* \in \operatorname{argmax}_u -f^*(-u) - g^*(A^T u)$$

$$a_i^T u^* \in \partial_i g(x^*)$$



## Screening and relaxing test

$$|a_i^T u^*| \leq \lambda \iff x_i^* = 0$$

$$|a_i^T u^*| > \lambda \iff x_i^* \neq 0$$

$u^*$  is not available

# Safe regions

Characterize the nullity in  $x^*$  from a safe region

# Safe regions

Characterize the nullity in  $x^*$  from a safe region

Safe region:  $u^* \in \mathcal{R}$

**Safe screening and relaxing test**

$$\max_{u \in \mathcal{R}} |a_i^T u| \leq \lambda \implies x_i^* = 0$$

$$\min_{u \in \mathcal{R}} |a_i^T u| > \lambda \implies x_i^* \neq 0$$



# Safe regions

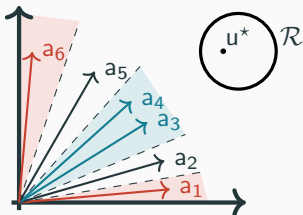
Characterize the nullity in  $x^*$  from a safe region

Safe region:  $u^* \in \mathcal{R}$

Safe screening and relaxing test

$$\max_{u \in \mathcal{R}} |a_i^T u| \leq \lambda \implies x_i^* = 0$$

$$\min_{u \in \mathcal{R}} |a_i^T u| > \lambda \implies x_i^* \neq 0$$



# Safe regions

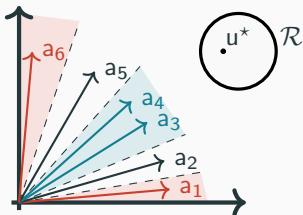
Characterize the nullity in  $x^*$  from a safe region

Safe region:  $u^* \in \mathcal{R}$

Safe screening and relaxing test

$$\max_{u \in \mathcal{R}} |a_i^T u| \leq \lambda \implies x_i^* = 0$$

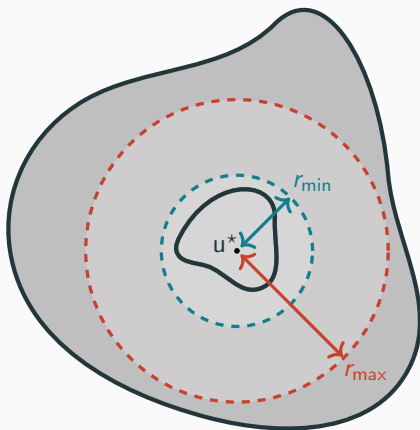
$$\min_{u \in \mathcal{R}} |a_i^T u| > \lambda \implies x_i^* \neq 0$$



Lasso problem

Limit case with  $h(x) = 0$

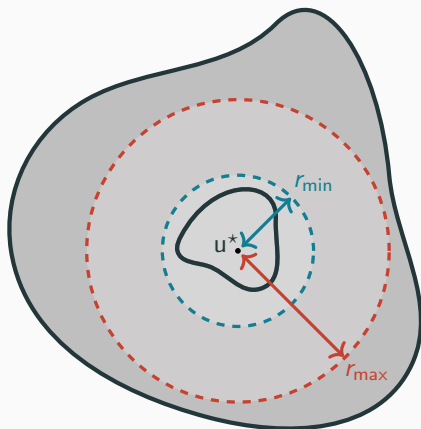
# Working regimes



$\mathcal{R} \subset \mathcal{S}(u^*, r_{\min}) \implies$  all tests passed

$\mathcal{R} \supset \mathcal{S}(u^*, r_{\max}) \implies$  no tests passed

# Working regimes



$\mathcal{R} \subset \mathcal{S}(u^*, r_{\min}) \implies$  all tests passed

$\mathcal{R} \supset \mathcal{S}(u^*, r_{\max}) \implies$  no tests passed

$$r_{\min} > 0$$

We know how to construct safe regions with a radius proportional to the (square root of) the duality gap



Guaranty to identify all zeros and non-zeros in  $x^*$  in finite time

# Implementation

---

## Initial problem

$$\min_x f(Ax) + \lambda \|x\|_1 + h(x)$$

# Reformulation

## Initial problem

$$\min_x f(Ax) + \lambda \|x\|_1 + h(x)$$

Perfect identif.

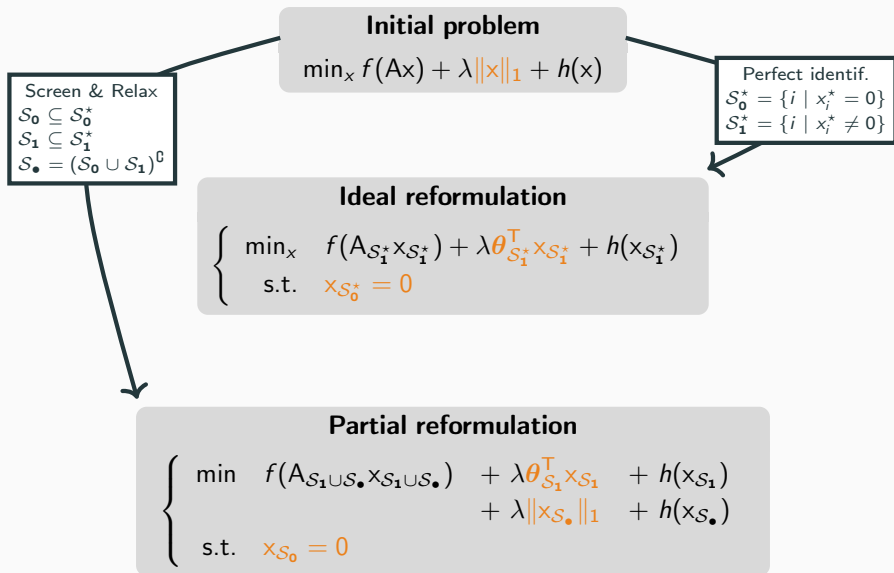
$$S_0^* = \{i \mid x_i^* = 0\}$$

$$S_1^* = \{i \mid x_i^* \neq 0\}$$

## Ideal reformulation

$$\begin{cases} \min_x & f(A_{S_1^*} x_{S_1^*}) + \lambda \theta_{S_1^*}^T x_{S_1^*} + h(x_{S_1^*}) \\ \text{s.t.} & x_{S_0^*} = 0 \end{cases}$$

# Reformulation





---

**Algorithm 1** Screen & Relax

---

```
1: initialize  $(\mathcal{S}_0, \mathcal{S}_1, \mathcal{S}_\bullet) = (\emptyset, \emptyset, \llbracket 1, n \rrbracket)$ 
2: repeat
3:   // Iterate update
4:    $x_{\mathcal{S}_\bullet}^k \leftarrow 1^{st}OrderIteration(x^{k-1})$ 
5:    $x_{\mathcal{S}_1}^k \leftarrow 2^{nd}OrderIteration(x^{k-1})$ 
6:    $x_{\mathcal{S}_0}^k \leftarrow 0$ 
7:   // Problem update
8:   Construct a new safe region  $\mathcal{R}^k$  from  $x^k$ 
9:   Update  $(\mathcal{S}_0, \mathcal{S}_1, \mathcal{S}_\bullet)$  using  $\mathcal{R}^k$ 
10: until convergence criterion is met
```

---

Iteration cost reduction:  $n \rightarrow n - |\mathcal{S}_0|$

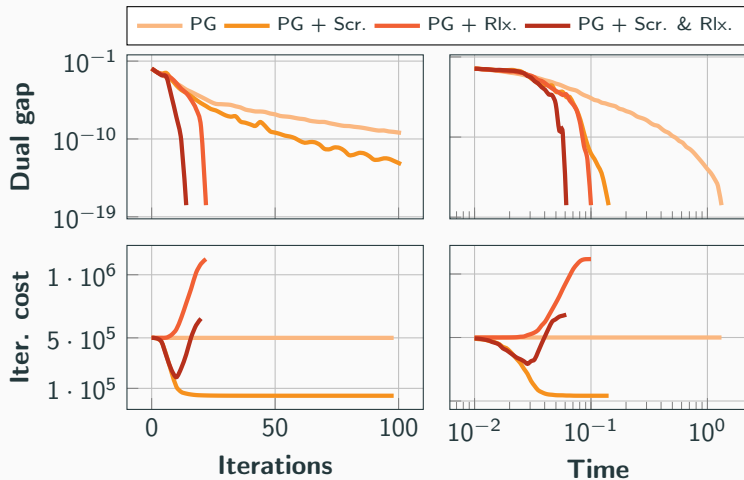
Faster convergence rate: (sub-)linear  $\rightarrow$  super-linear

# Numerics

---

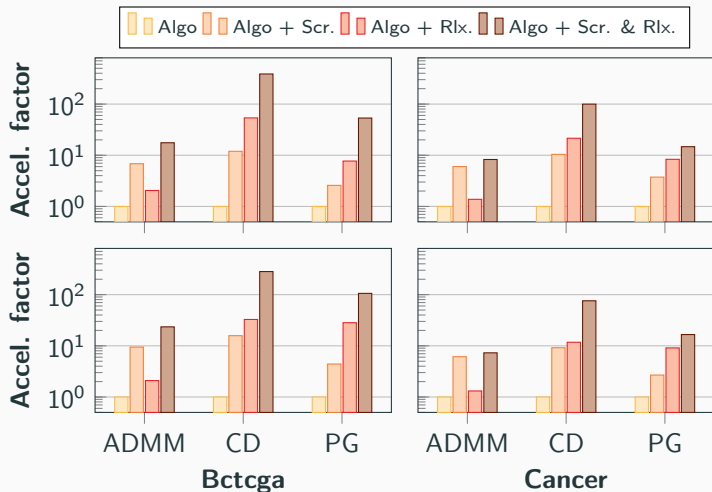
# Screen & Relax effects

$$\min_x f(Ax) + \lambda \|x\|_1 + h(x)$$



# Regularization path fitting

$$\min_x f(Ax) + \lambda \|x\|_1 + h(x)$$



$$h(x) = \|x\|_2^2$$

$$h(x) = \mathbf{1}^T \log\left(\frac{\mathbf{1}}{1 - |x|}\right)$$

Regression - 356 x 17,322

Classification - 100 x 12,600

## Take-home message

- Structure matters in sparse problems
- Identification of **zeros**
  - Screening tests
  - Allows for dimensionality reduction
  - Computational savings
- Identification of **non-zeros**
  - Relaxing tests
  - Allows for objective smoothing
  - Faster convergence rate
- **Screen & Relax** strategy to benefit from
  - Computational savings
  - Accelerated convergence

# Question time

