# Screen & Relax

Accélérer la résolution de l'Elastic-Net par identification du support de la solution

Théo Guyard[1,2], Cédric Herzet[2], Clément Elvira[3]

[1] INSA Rennes
[2] INRIA Rennes Bretagne Atlantique
[3] IETR CentraleSupélec

**GRETSI – 8 Septembre 2022**

# General context

# Sparse problem

**Ingredients of the problem**

- A target y

**Ingredients of the problem**

- A target y
- A dictionary $A = \{a_i\}_{i \in \mathcal{I}}$ made of atoms

**Ingredients of the problem**

- A target y
- A dictionary $A = \{a_i\}_{i \in \mathcal{I}}$ made of atoms

**Objective**

- Find a sparse linear combination of atoms that well approximates the target through a given model

# Sparse problem

**Ingredients of the problem**

- A target y
- A dictionary $A = \{a_i\}_{i \in \mathcal{I}}$ made of atoms

**Objective**

- Find a sparse linear combination of atoms that well approximates the target through a given model

**Rough formulation**

> **Problem**
> Find x sparse such that $y \simeq \text{Model}(Ax)$

# Sparse problem

**Ingredients of the problem**

- A target $y$
- A dictionary $A = \{a_i\}_{i \in \mathcal{I}}$ made of atoms

**Objective**

- Find a sparse linear combination of atoms that well approximates the target through a given model

**Rough formulation**

> **Problem**
> Find $x$ sparse such that $y \simeq \text{Model}(Ax)$

Remark : Entries of $x$ weight each atom in the linear combination.

# The Elastic-Net problem

## Formulation and properties

**Target problem**

Solve

$$x^\star = \arg\min_x \left\{ \mathrm{P}(x) = \underbrace{\tfrac{1}{2}\|y - Ax\|_2^2}_{f(Ax)} + \underbrace{\lambda(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_2)}_{\lambda g(x)} \right\} \qquad (P)$$

where $\lambda > 0$ and $\sigma \in ]0, 1[$ are tuning hyperparameters.

**Target problem**

Solve

$$x^\star = \arg\min_x \left\{ P(x) = \underbrace{\tfrac{1}{2}\|y - Ax\|_2^2}_{f(Ax)} + \underbrace{\lambda(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_2)}_{\lambda g(x)} \right\} \qquad (P)$$

where $\lambda > 0$ and $\sigma \in\; ]0, 1[$ are tuning hyperparameters.

**Properties of** $(P)$

- Convex non-smooth problem
- Least-squares : Ensures a good reconstruction of the target
- $\ell_1$-norm : Enforces sparsity
- $\ell_2$-norm : Promotes desirable properties

## Formulation and properties

**Target problem**

Solve

$$x^\star = \arg\min_x \left\{ P(x) = \underbrace{\tfrac{1}{2}\|y - Ax\|_2^2}_{f(Ax)} + \underbrace{\lambda\left(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_2\right)}_{\lambda g(x)} \right\} \qquad (P)$$

where $\lambda > 0$ and $\sigma \in ]0, 1[$ are tuning hyperparameters.

**Properties of** $(P)$

- Convex non-smooth problem
- Least-squares : Ensures a good reconstruction of the target
- $\ell_1$-norm : Enforces sparsity
- $\ell_2$-norm : Promotes desirable properties

**Solving** $(P)$

- Broad class of solution methods (gradient-based, pivot-based, ...)
- Acceleration strategies (backtracking, screening tests, ...)

# Screening and Relaxing tests

# Main idea

**In the context of a sparse problem**

## Main idea

**In the context of a sparse problem**

- Where are zero and non-zero entries in $x^\star$ ?

## Main idea

**In the context of a sparse problem**

- Where are zero and non-zero entries in $x^\star$ ?
- Can we accelerate a given solution method if we knew some ?

# Main idea

**In the context of a sparse problem**

- Where are zero and non-zero entries in $x^\star$ ?
- Can we accelerate a given solution method if we knew some ?
- → Spoiler alert : Yes and yes ! We can leverage duality.

# Main idea

**In the context of a sparse problem**

- Where are zero and non-zero entries in $x^\star$ ?
- Can we accelerate a given solution method if we knew some ?

$\rightarrow$ Spoiler alert : Yes and yes ! We can leverage duality.

| **Primal** problem |
| --- |
| $x^\star = \arg\min_x P(x)$   $(P)$ |

$\equiv$

| **Dual** problem |
| --- |
| $u^\star = \arg\max_u D(u)$   $(D)$ |

with optimality conditions linking $x^\star$ and $u^\star$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$



$$|a_i^T u^\star| \leq \lambda \sigma \quad \Longleftrightarrow \quad x_i^\star = 0$$
$$|a_i^T u^\star| > \lambda \sigma \quad \Longleftrightarrow \quad x_i^\star \neq 0$$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$



$$|a_i^T u^\star| \leq \lambda\sigma \iff x_i^\star = 0$$
$$|a_i^T u^\star| > \lambda\sigma \iff x_i^\star \neq 0$$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$



$$|a_i^T u^\star| \leq \lambda \sigma \quad \Longleftrightarrow \quad x_i^\star = 0$$
$$|a_i^T u^\star| > \lambda \sigma \quad \Longleftrightarrow \quad x_i^\star \neq 0$$



We can identify zeros and non-zeros in $x^\star$

$$a_i^T u^\star \in \lambda \partial g(x_i^\star)$$



$\partial g(x_i) < -\sigma$    $\partial g(x_i) > \sigma$

$\partial g(x_i) = [-\sigma, \sigma]$

$$|a_i^T u^\star| \leq \lambda \sigma \iff x_i^\star = 0$$
$$|a_i^T u^\star| > \lambda \sigma \iff x_i^\star \neq 0$$



$\bullet\ u^\star$

$a_6$
$a_5$
$a_4$
$a_3$
$a_2$
$a_1$

We can identify zeros and non-zeros in $x^\star$
... but we need $u^\star$ !

# Relaxed optimality condition

Let $u^\star \in \mathbb{S}(c, r)$, then

$$|a_i^T c| + r \leq \lambda\sigma \quad \implies \quad x_i^\star = 0 \qquad \text{(screening test)}$$

Let $u^\star \in \mathbb{S}(\mathsf{c}, r)$, then

$$
\begin{aligned}
|\mathsf{a}_i^T \mathsf{c}| + r \leq \lambda \sigma &\implies x_i^\star = 0 &&\text{(screening test)} \\
|\mathsf{a}_i^T \mathsf{c}| - r > \lambda \sigma &\implies x_i^\star \neq 0 &&\text{(relaxing test)}
\end{aligned}
$$

Let $u^\star \in \mathbb{S}(c, r)$, then

$$
\begin{aligned}
|a_i^T c| + r \leq \lambda\sigma &\implies x_i^\star = 0 \qquad \text{(screening test)} \\
|a_i^T c| - r > \lambda\sigma &\implies x_i^\star \neq 0 \qquad \text{(relaxing test)}
\end{aligned}
$$

$\rightarrow$ No $u^\star$ needed anymore, but only a "safe region" containing it !

# Relaxed optimality condition

Let $u^\star \in \mathbb{S}(c, r)$, then

$$|a_i^T c| + r \leq \lambda\sigma \quad \implies \quad x_i^\star = 0 \qquad \text{(screening test)}$$
$$|a_i^T c| - r > \lambda\sigma \quad \implies \quad x_i^\star \neq 0 \qquad \text{(relaxing test)}$$

$\rightarrow$ No $u^\star$ needed anymore, but only a "safe region" containing it !

# Dimensionality reduction

**With screening test**

Zero entries identified in $x^\star$ can be discarded from the problem without changing the objective value.

**With screening test**

Zero entries identified in $x^\star$ can be discarded from the problem without changing the objective value.

**With relaxing test**

Non-zero entries identified in $x^\star$ can be expressed as a linear combination of all the other entries.

**Notations** $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ : Set of zero/non-zero/unclassified entries in $x^\star$

**Notations** $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ : Set of zero/non-zero/unclassified entries in $x^\star$

### Initial problem

$$x^\star = \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ P(x) = \tfrac{1}{2}\|y - Ax\|_2^2 + \lambda(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_2^2) \right\}$$

$n$ dimensional problem

**Notations** $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ : Set of zero/non-zero/unclassified entries in $x^\star$

### Initial problem

$$x^\star = \text{argmin}_{x \in \mathbb{R}^n} \left\{ P(x) = \tfrac{1}{2} \|y - Ax\|_2^2 + \lambda(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_2^2) \right\}$$

$n$ dimensional problem

$$\lessgtr$$

### Reduced problem

$$\left\{ \begin{array}{ll} x^\star_{\bar{\mathcal{S}}} & = \text{argmin}_{x \in \mathbb{R}^{|\bar{\mathcal{S}}|}} \left\{ \tilde{P}(x) = \tfrac{1}{2} \|\tilde{y} - \tilde{A}x\|_2^2 + \lambda(\sigma\|x\|_1 + \tfrac{1-\sigma}{2}\|x\|_M^2) \right\} \\ x^\star_{\mathcal{S}_1} & = Bx^\star_{\bar{\mathcal{S}}} + b \\ x^\star_{\mathcal{S}_0} & = 0 \end{array} \right.$$

# Problem reformulation

**Notations** $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ : Set of zero/non-zero/unclassified entries in $x^\star$

**Initial problem**

$$x^\star = \mathrm{argmin}_{x \in \mathbb{R}^n} \left\{ P(x) = \tfrac{1}{2} \| y - Ax \|_2^2 + \lambda(\sigma \| x \|_1 + \tfrac{1-\sigma}{2} \| x \|_2^2) \right\}$$

$n$ dimensional problem

$\wr$

**Reduced problem**

$$\begin{cases} x_{\bar{\mathcal{S}}}^\star & = \mathrm{argmin}_{x \in \mathbb{R}^{|\bar{\mathcal{S}}|}} \left\{ \tilde{P}(x) = \tfrac{1}{2} \| \tilde{y} - \tilde{A}x \|_2^2 + \lambda(\sigma \| x \|_1 + \tfrac{1-\sigma}{2} \| x \|_M^2) \right\} \\ x_{\mathcal{S}_1}^\star & = Bx_{\bar{\mathcal{S}}}^\star + b \\ x_{\mathcal{S}_0}^\star & = 0 \end{cases}$$

$n - |\mathcal{S}_0| - |\mathcal{S}_1|$ dimensional problem (similar structure)

Some linear algebra operations (negligible computational cost)

**Algorithm 1:** "Screen & Relax" solving procedure

**Input:** A, y, $\lambda$, $\sigma$, $x^{(0)}$

1 $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}}) \leftarrow (\emptyset, \emptyset, \{1, \ldots, n\})$

2 **while** *convergence criterion is not met* **do**

**Algorithm 2:** "Screen & Relax" solving procedure

**Input:** A, y, $\lambda$, $\sigma$, $x^{(0)}$

1 $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}}) \leftarrow (\emptyset, \emptyset, \{1, \ldots, n\})$

2 **while** *convergence criterion is not met* **do**

3 | Update the current iterate

**Algorithm 3:** "Screen & Relax" solving procedure

**Input:** A, y, $\lambda$, $\sigma$, $x^{(0)}$

1   $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}}) \leftarrow (\emptyset, \emptyset, \{1, \ldots, n\})$

2 **while** *convergence criterion is not met* **do**

3      Update the current iterate

4      Update $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ with screening and relaxing tests

**Algorithm 4:** "Screen & Relax" solving procedure

**Input:** A, y, $\lambda$, $\sigma$, $x^{(0)}$

1  $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}}) \leftarrow (\emptyset, \emptyset, \{1, \ldots, n\})$

2  **while** *convergence criterion is not met* **do**

3      Update the current iterate

4      Update $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ with screening and relaxing tests

5      Reduce the problem

# Dynamic Screen & Relax principle

---
**Algorithm 5:** "Screen & Relax" solving procedure

**Input:** A, y, $\lambda$, $\sigma$, $x^{(0)}$

---

1   $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}}) \leftarrow (\emptyset, \emptyset, \{1, \ldots, n\})$

2   **while** *convergence criterion is not met* **do**

3      Update the current iterate

4      Update $(\mathcal{S}_0, \mathcal{S}_1, \bar{\mathcal{S}})$ with screening and relaxing tests

5      Reduce the problem

6      **if** $\bar{\mathcal{S}} = \emptyset$ **then**

7          The solution is available in closed form

8      **end**

9   **end**

---

# Some numerical results

**Synthetic data generation**

- Generate the dictionary A randomly
- Generate a k-sparse vector $x^\dagger$
- Set $y = Ax^\star + \text{noise}$
- Solve $(P)$ with a tailored method

**Synthetic data generation**

- Generate the dictionary A randomly
- Generate a k-sparse vector $x^\dagger$
- Set $y = Ax^\star + \text{noise}$
- Solve ($P$) with a tailored method

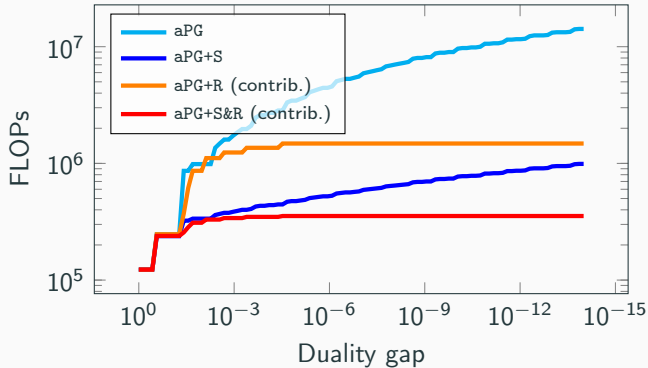**Concurrent methods**

- Accelerated proximal-gradient algorithm
- Same with screening tests
- Same with relaxing tests
- Same with screening and relaxing tests

# Experimental setup

**Synthetic data generation**

- Generate the dictionary A randomly
- Generate a k-sparse vector $x^\dagger$
- Set $y = Ax^\star + \text{noise}$
- Solve $(P)$ with a tailored method

**Concurrent methods**

- Accelerated proximal-gradient algorithm
- Same with screening tests
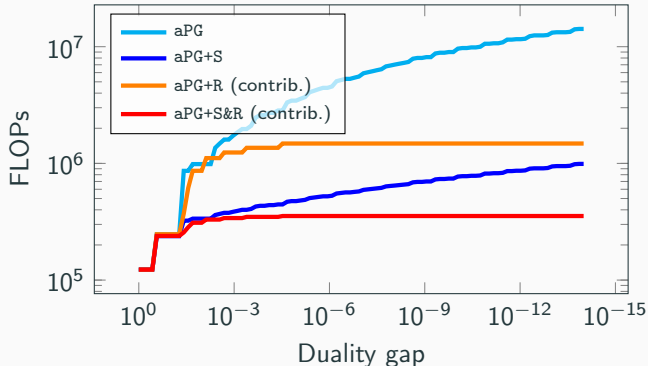- Same with relaxing tests
- Same with screening and relaxing tests

**Metrics**

- Duality gap : How close is the objective from its optimal value
- FLOPs : Number of linear algebra operations performed
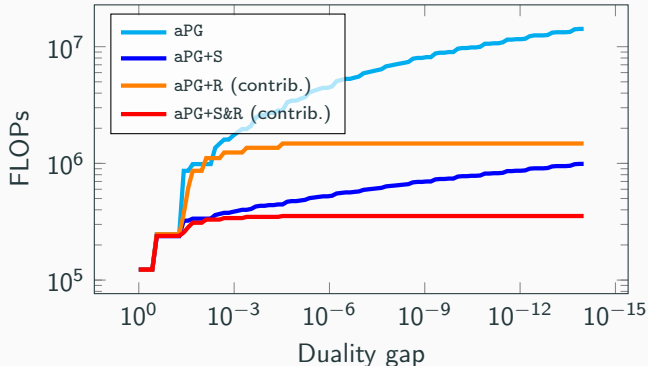
# Classical convergence scheme

- Complexity reduction with screening and/or relaxing
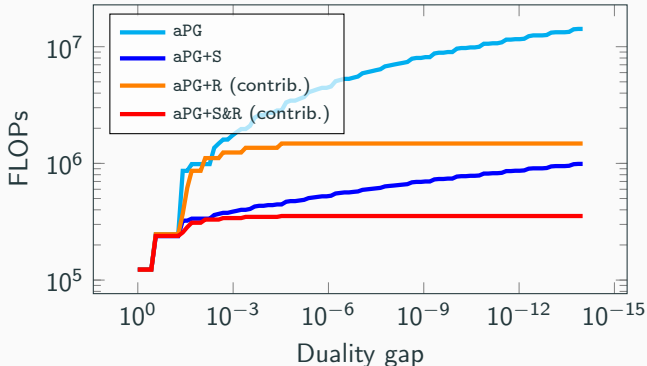
- Complexity reduction with screening and/or relaxing
- Convergence to machine precision at some point

# Classical convergence scheme



- Complexity reduction with screening and/or relaxing
- Convergence to machine precision at some point
- Gains depend on the sparsity in $x^\star$
  - Too few non-zeros : relaxing has little impact
  - Too many non-zeros : problem updates become binding