
Méthodes de screening pour les moindres carrés non-négatifs

MERJANI AYMANE
MOUQED MOHAMED REDA
AHMED EL BAJDALI
BADR SAISSI
SAAD AAFFOUTE

Encadrants :
CLEMENT ELVIRA

Résumé

Ce projet s'intéresse principalement à l'étude d'une méthode d'accélération proposée par El Ghaoui pour la relaxation convexe de problèmes de représentations parcimonieuses, baptisé safe screening. Cette procédure repose sur deux éléments centraux. Le premier repose sur le fait que les solutions de problèmes convexes contiennent un très grand nombre de zéros. Le deuxième élément est la connaissance de la position des zéros de la solution va permettre de transformer le problème en un deuxième plus simple et de plus petite dimension. L'étude de ce projet permettra de réduire le temps de calcul des algorithmes pour résoudre ce type de problèmes, la mémoire utilisée, la puissance calculatoire et l'empreinte carbone de la résolution du problème.

Table des matières

1	Introduction	4
1.1	Contexte	4
1.2	Notations	4
1.2.1	Applications vectorielles	4
1.2.2	Modélisation	5
1.3	Nécessité d'un terme de régularisation	6
1.3.1	Overfitting et underfitting	6
1.3.2	Pénalisation ou régularisation	6
2	Estimateur Lasso	7
2.1	Convexité et non-différentiabilité	7
2.2	Dualité de Fenchel	9
2.2.1	Problème dual	9
2.2.2	Relation entre $\hat{\phi}$ et \hat{x}	10
2.3	Règles de screening	11
2.3.1	Sous-différentiel	11
2.3.2	Sous-différentiel de la fonction f_λ	12
2.3.3	Règles de screening	13
2.4	Expression de λ_{max}	13
2.5	Choix du pas de descente	14
2.6	Descente du gradient proximal	16
2.6.1	Opérateur proximal	16
2.6.2	Algorithme de descente de gradient proximal	16
2.6.3	Résultats intermédiaires	17
2.7	Étude du Gap	18
3	Problème des moindres carrés non-négatifs	20
3.1	Utilité de ce problème	20
3.2	Convexité et non-différentiabilité	20
3.3	Dualité de Fenchel	20
3.4	Sous-différentiel	21
3.5	Test de screening	21
3.6	Sphère Gap pour les moindres carrés non-négatifs	22
4	Implémentation	23
4.1	Pas de descente	23
4.1.1	Théorie autour du pas	23
4.1.2	Actualisation des pas	24
4.2	Une région safe - Sphère Gap	25
4.3	Gradient Proximal avec screening	26
4.4	Résultats de l'implémentation	27

4.4.1	Nombre d'it�rations	28
4.4.2	Temps d'ex�cution	28
4.4.3	Empreinte carbone	29
4.5	Impl�mentation Moindres Carr�s non-n�gatifs	30
5	Conclusion	32

1 Introduction

1.1 Contexte

Dans un contexte où la gestion de données devient centrale, il est important de trouver des solutions afin de minimiser le stockage et de réduire les calculs. C'est donc dans cette optique que la méthode de screening est mise au point, afin d'identifier et de sélectionner rapidement les données pertinentes nécessaires. L'objectif est de réduire le set de données, sans en réduire la qualité, et ainsi économiser de l'espace, du temps et de l'énergie.

1.2 Notations

1.2.1 Applications vectorielles

Dans ce document, on utilisera les applications vectorielles suivantes (toutes ces applications sont des normes vectorielles sauf la première) :

$$\begin{aligned} \forall x \in \mathbb{R}^n, n \in \mathbb{N}_+^* \\ \|x\|_0 &= \sum_{i=1}^n 1_{\{x_i \neq 0\}} \\ \|x\|_1 &= \sum_{i=1}^n |x_i| \\ \|x\|_2 &= \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}} \\ \|x\|_\infty &= \max_{1 \leq i \leq n} |x_i| \end{aligned}$$

1.2.2 Modélisation

- $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$ est un vecteur de dimension m ,

- $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$ est le vecteur de dimension n des paramètres inconnus du modèle,

- $A = \begin{pmatrix} a_{11} & a_{12} & \dots & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & \dots & a_{2n} \\ \vdots & & & & \\ \vdots & & & & \\ & & & & \vdots \\ & & & & \vdots \\ a_{m1} & \dots & \dots & a_{mn-1} & a_{mn} \end{pmatrix}$ est une matrice de taille $m \times n$.

1.3 N cessit  d'un terme de r gularisation

1.3.1 Overfitting et underfitting

Les deux probl mes les plus courants en Machine Learning sont : l'overfitting et l'underfitting. L'overfitting se produit lorsque le mod le est trop complexe et s'ajuste trop  troitement aux donn es d'entra nement, conduisant   une mauvaise g n ralisation sur de nouvelles donn es. En revanche, l'underfitting se produit lorsque le mod le est trop simple et ne parvient pas   saisir les motifs des donn es d'entra nement, entra nant  galement une mauvaise performance. Trouver le bon  quilibre entre la complexit  du mod le et la quantit  de donn es d'entra nement est essentiel pour obtenir de bons r sultats en machine learning.

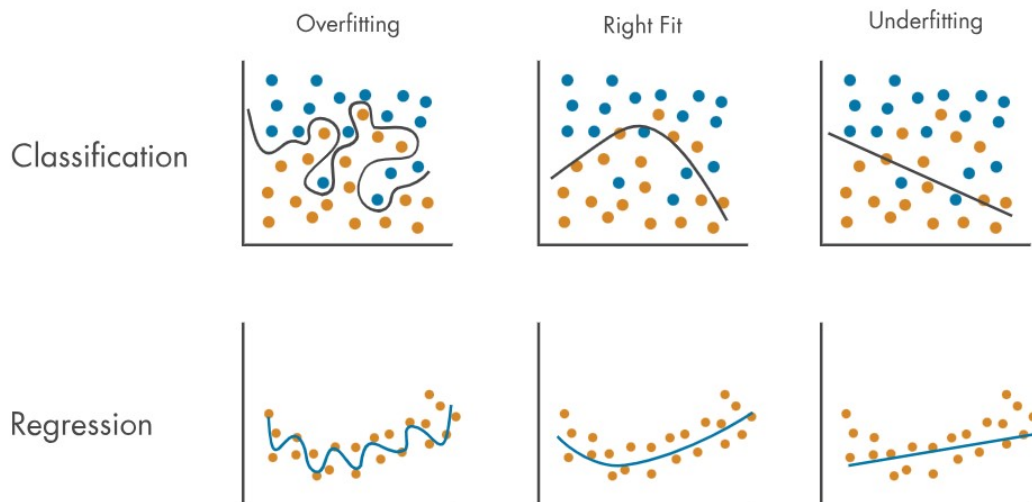


FIGURE 2 – Underfitting et Overfitting

1.3.2 P nalisation ou r gularisation

La p nalisation est une technique couramment utilis e dans les probl mes de r gression pour  viter l'overfitting. Il existe plusieurs techniques de r gularisation ou p nalisation :

- P nalisation Lasso (L1) qui permet de d tecter les param tres moins importants donc le mod le est plus simple ce qui implique qu'il est moins sujet   l'overfitting.
- P nalisation Ridge (L2) qui permet de construire un mod le plus homog ne et donc moins sujet   l'overfitting.
- P nalisation Elastic Net qui   la fois d tecte les param tres moins importants et qui les homog n ise.

2 Estimateur Lasso

L'estimateur LASSO (Least Absolute Selection and Shrinkage Operator) est défini pour $\lambda > 0$ par :

$$\hat{x} = \underset{x \in \mathbb{R}^n}{\operatorname{Argmin}} \left\{ \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1 \right\}.$$

Minimiser cette expression a un intérêt majeur si la matrice A est une matrice creuse (sparse) c'est à dire qu'elle contient un nombre important de valeurs nulles.

2.1 Convexité et non-différentiabilité

On pose l'application $f_\lambda : x \mapsto \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1$. L'objectif de cette partie sera donc d'étudier la fonction f_λ , qui va nous permettre de donner quelques propriétés de l'estimateur LASSO " $\hat{x}(\lambda)$ "

La fonction f_λ est convexe, non différentiable.

On montrera également que la solution du problème de minimisation ne peut pas être unique.

Démonstration : Soient $x_1, x_2 \in \mathbb{R}^n$ et $t \in [0, 1]$,

$$f_\lambda(tx_1 + (1-t)x_2) = \frac{1}{2} \|y - A(tx_1 + (1-t)x_2)\|_2^2 + \lambda \|tx_1 + (1-t)x_2\|_1.$$

Alors par application de l'inégalité triangulaire pour le deuxième terme, on a que :

$$\|tx_1 + (1-t)x_2\|_1 \leq t\|x_1\|_1 + (1-t)\|x_2\|_1$$

Et d'autre part,

$$\begin{aligned} & \|y - A(tx_1 + (1-t)x_2)\|_2^2 \\ &= \|t(y - Ax_1) + (1-t)(y - Ax_2)\|_2^2 \\ &= t^2 \|y - Ax_1\|_2^2 + (1-t)^2 \|y - Ax_2\|_2^2 + 2t(1-t) \langle y - Ax_1, y - Ax_2 \rangle \\ &\leq t^2 \|y - Ax_1\|_2^2 + (1-t)^2 \|y - Ax_2\|_2^2 + 2t(1-t) \|y - Ax_1\|_2 \|y - Ax_2\|_2 \\ &\leq t^2 \|y - Ax_1\|_2^2 + (1-t)^2 \|y - Ax_2\|_2^2 + t(1-t) (\|y - Ax_1\|_2^2 + \|y - Ax_2\|_2^2) \\ &= t \|y - Ax_1\|_2^2 + (1-t) \|y - Ax_2\|_2^2 \end{aligned}$$

Dans les inégalités, nous avons utilisé respectivement l'inégalité de Cauchy-Schwarz ainsi que l'inégalité suivante : $uv \leq \frac{1}{2}(u^2 + v^2)$ pour tout $u, v \in \mathbb{R}$. Ainsi :

$$\|y - A(tx_1 + (1-t)x_2)\|_2^2 \leq t\|y - Ax_1\|_2^2 + (1-t)\|y - Ax_2\|_2^2$$

Il en vient que :

$$f_\lambda(tx_1 + (1-t)x_2) \leq tf_\lambda(x_1) + (1-t)f_\lambda(x_2).$$

\implies la fonction f_λ est bien convexe. Puisque $n > m$ alors $rg(A) < n$ (et donc la matrice $A^T A$ est non inversible). En outre, la matrice Hessienne de f_λ qui est $A^T A$, n'est pas définie-positive, donc on en déduit que la fonction f_λ n'est pas strictement convexe. De plus, étant donné que pour tout $x \in \mathbb{R}^n$ la fonction norme 1, $x \mapsto \|x\|_1$, n'est pas différentiable sur \mathbb{R}^n alors on en déduit que f_λ ne l'est également pas.

2.2 Dualité de Fenchel

2.2.1 Problème dual

Posons : $g(z) = \frac{1}{2}\|y - z\|_2^2$ and $f(x) = \lambda\|x\|_1$. Cherchons donc :

$$\inf_{x \in \mathbb{R}^n} (f(x) + g(Ax))$$

Par Fenchel Rockafeller [3], comme g et f sont des fonctions convexes, propres, semi-continues inférieures (même continues car ils sont des normes) et $0 \in \text{int}(\mathbb{R}^n - A\mathbb{R}^n) = \text{int}(\mathbb{R}^n) = \mathbb{R}^n$

$$\inf_{x \in \mathbb{R}^n} (f(x) + g(Ax)) = - \inf_{\phi \in \mathbb{R}^n} (f^*(A^*\phi) + g^*(-\phi)) \quad (1)$$

Nous avons :

$$\begin{aligned} g^*(u) &= \sup_x \left(\langle x | u \rangle - \frac{1}{2}\|y - x\|_2^2 \right) \\ &= \sup_x \left(\langle x | u \rangle - \frac{1}{2}\|y\|_2^2 - \frac{1}{2}\|x\|_2^2 + \langle y, x \rangle \right) \\ &= \sup_x \left(\langle x | u + y \rangle - \frac{1}{2}\|x\|_2^2 \right) - \frac{1}{2}\|y\|_2^2 \\ &= \sup_x \left(-\frac{1}{2}\|x - (u + y)\|_2^2 + \frac{1}{2}\|u + y\|_2^2 \right) - \frac{1}{2}\|y\|_2^2 \\ &= \frac{1}{2}\|u + y\|_2^2 - \frac{1}{2}\|y\|_2^2 \text{ en prenant } x = u + y \end{aligned}$$

Pour le cas de la fonction f , nous supposons d'abord que

$$f(x) = \|x\|_1$$

Nous voulons montrer que : $f^*(y) = I_{\|\cdot\|_\infty \leq 1}(y) = \begin{cases} 0, & \|y\|_\infty \leq 1 \\ +\infty, & \|y\|_\infty > 1 \end{cases}$

- Si $\|y\|_\infty \leq 1$,

$$\langle x | y \rangle \leq \|y\|_\infty \|x\|_1 \leq \|x\|_1$$

Donc $f^*(y) = \sup_x (\langle x | y \rangle - \|x\|_1) = 0$

- Sinon, il existe i_0 tel que $\|y\|_\infty = |y_{i_0}| > 1$, On pose :

$$y_0 = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ y_{i_0} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad i_0\text{-eme position}$$

Soit $t > 0$, on a :

$$f^*(y) \geq \langle ty_0 \mid y \rangle - \|ty_0\|_1 = t(y_{i_0}^2 - |y_{i_0}|)$$

qui tend vers $+\infty$ lorsque t tend vers $+\infty$

$$\text{Donc : } f^*(y) = \begin{cases} 0, & \|y\|_\infty \leq 1 \\ +\infty, & \|y\|_\infty > 1 \end{cases} = I_{\|\cdot\|_\infty \leq 1}(y)$$

Pour le cas,

$$f(x) = \lambda \|x\|_1$$

Comme $\lambda > 0$, on a :

$$f^*(u) = \sup_x (\langle x \mid u \rangle - \lambda \|x\|_1) = \lambda \sup_x (\langle x \mid \lambda^{-1}u \rangle - \|x\|_1) = \lambda I_{\|\cdot\|_\infty \leq 1}(\lambda^{-1}u)$$

Donc le problème Lasso se transforme en :

$$\begin{aligned} \inf_{x \in \mathbb{R}^n} (f(x) + g(Ax)) &= - \inf_{\phi \in \mathbb{R}^n} (f^*(A^*\phi) + g^*(-\phi)) \\ &= - \inf_{\phi \in \mathbb{R}^n} (I_{\|\cdot\|_\infty \leq \lambda}(A^*\phi) + \frac{1}{2} \|\phi - y\|_2^2 - \frac{1}{2} \|y\|_2^2) \\ &\Leftrightarrow \inf_{\phi \in \mathbb{R}^n} \frac{1}{2} \|y - \phi\|_2^2 \text{ tel que } \|A^*\phi\| \leq \lambda \end{aligned}$$

2.2.2 Relation entre $\hat{\phi}$ et \hat{x}

Soit $\hat{\phi}$ et \hat{x} les minimiseurs. Comme la dualité est forte, nous avons alors l'égalité de Fenchel-Young pour f et pour g . On écrit le cas d'égalité Fenchel-Young pour g ,

$$\begin{aligned} g(A\hat{x}) + g^*(-\hat{\phi}) &= \langle A\hat{x}, -\hat{\phi} \rangle \\ \Leftrightarrow \frac{1}{2} \|y - A\hat{x}\|_2^2 + \frac{1}{2} \|y - \hat{\phi}\|_1^2 - \frac{1}{2} \|y\|_2^2 &= -\langle A\hat{x}, \hat{\phi} \rangle \\ \Leftrightarrow \langle y, -\hat{\phi} - A\hat{x} \rangle + \frac{1}{2} \|y\|_2^2 + \frac{1}{2} \|A\hat{x}\|_2^2 + \frac{1}{2} \|\hat{\phi}\|_2^2 &= -\langle A\hat{x}, \hat{\phi} \rangle \\ \Leftrightarrow \left\| -A\hat{x} + y - \hat{\phi} \right\|_2^2 &= 0 \\ \Leftrightarrow \hat{\phi} &= y - A\hat{x} \end{aligned}$$

2.3 R gles de screening

Le but est minimiser la fonction d finie par $f_\lambda : x \mapsto \frac{1}{2}\|y - Ax\|_2^2 + \lambda\|x\|_1$. Pour minimiser des fonctions usuelles, il faut passer par la d riv e de ces fonctions ou par leur diff rentielles, sauf que dans le cas  tudi , la fonction n'est pas diff rentiable, d'o  la n cessit  de d finir une notion qui est proche de la notion de la diff rentiabilit .

2.3.1 Sous-diff rentiel

Quand une fonction convexe $f : I \subset \mathbb{R}^n \mapsto \mathbb{R}$ est diff rentiable en un point $x_0 \in I$ on a l'in galit  suivante :

$$f(x) \geq f(x_0) + df(x_0)(x - x_0), \forall x \in I$$

o  $df(x_0) \in \mathbb{R}^n$ est la diff rentielle de f au point x_0 .

Notons que la fonction f_λ n'est pas diff rentiable   cause de la norme 1, introduisons la notion de sous-diff rentiel.

Soit f une fonction convexe d finie sur I , un ouvert de \mathbb{R}^n ,   valeur dans \mathbb{R} . Un sous-gradient de $f : I \mapsto \mathbb{R}$ en un point x_0 de I est un vecteur $S \in \mathbb{R}^n$ tel que, pour tout x appartenant   I :

$$f(x) \geq f(x_0) + S(x - x_0)$$

L'ensemble des tous les sous-gradients est appel  sous-diff rentiel de la fonction f en x_0 , not  $\partial f(x_0)$.

Voici une illustration du concept du sous-diff rentiel :

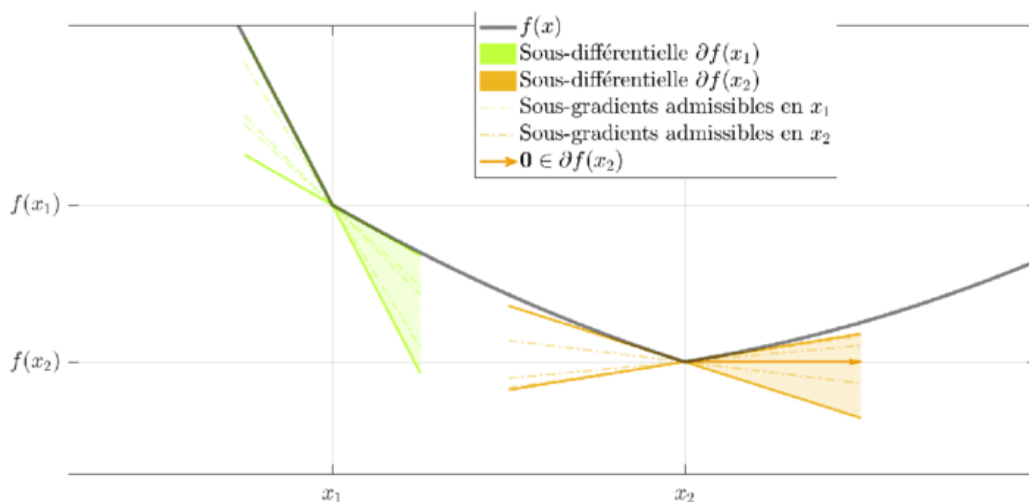


FIGURE 3 – Sous-diff rentiel

2.3.2 Sous-diff rentiel de la fonction f_λ

Il suffit de trouver le sous-diff rentiel de la norme 1 pour trouver le sous diff rentiel de la fonction f_λ vu que la norme 2 est diff rentiable et nous savons calculer sa diff rentielle.

Nous nous focalisons maintenant sur le calcul du sous-diff rentiel de la norme 1, pour cela nous admettons le r sultat affirmant que le sous-diff rentiel d'une somme finie de fonctions sous-diff rentiables est la somme des sous-diff rentiels de ces fonctions.

Nous pouvons donc nous restreindre   calculer le sous-diff rentiel des fonctions projections sur la base canonique : $f_i : \mathbb{R}^n \longrightarrow \mathbb{R}$, $\forall i \in \{1, 2, \dots, n\}$.

$$x \longmapsto |x_i|$$

Soit $\tilde{x} \in \mathbb{R}^n$, nous cherchons $S \in \mathbb{R}^n$ tel que :

$$|x_i| \geq |\tilde{x}_i| + S(x - \tilde{x}), \quad \forall x \in \mathbb{R}^n.$$

En particulier pour $x = \tilde{x} + e_j$ o  e_j est le vecteur canonique et $j \in \{1, 2, \dots, n\}$, $j \neq i$, on a donc :

$$\begin{aligned} |\tilde{x}_i| &\geq |\tilde{x}_i| + S(x - \tilde{x}) \\ \Rightarrow 0 &\geq S \times e_j \\ \Rightarrow 0 &\geq S_j. \end{aligned}$$

Avec le m me raisonnement pour $x := \tilde{x} - e_j$, nous avons aussi :

$$0 \leq S_j$$

Donc $S_j = 0, \forall j \in \{1, 2, \dots, n\}, j \neq i$.

Il suffit donc de calculer le sous diff rentiel de la fonction valeur absolue

Cette fonction est d rivable pour tout x non nul et sa d riv  vaut $\text{sign}(x)$. Calculons son sous-diff rentiel en z ro, nous cherchons $s \in \mathbb{R}$ tel que :

$$\begin{aligned} |x| &\geq |0| + s(x - 0), \quad \forall x \in \mathbb{R} \\ \Leftrightarrow |x| &\geq sx, \quad \forall x \in \mathbb{R} \\ \Rightarrow \begin{cases} x \geq sx, \forall x \in \mathbb{R}^+ \\ -x \geq sx, \forall x \in \mathbb{R}^- \end{cases} \\ \Rightarrow \begin{cases} s \leq 1, \quad \forall x \in \mathbb{R}^+ \\ s \geq -1, \quad \forall x \in \mathbb{R}^- \end{cases} \\ \Rightarrow s &\in [-1, 1]. \end{aligned}$$

Nous avons donc le sous-diff rentiel de $f_\lambda(x)$ est :

$$\partial f_\lambda(x) = -A^T(y - Ax) + \lambda \left\{ u \in \mathbb{R}^n : u_l = \begin{cases} [-1, 1], & \text{si } x_l = 0 \\ \text{sign}(x_l), & \text{sinon} \end{cases}, l \in \{1, 2, \dots, n\} \right\}.$$

2.3.3 Règles de screening

Soit \hat{x} solution admissible, en annulant le terme $\partial f_\lambda(\hat{x})$, nous avons :

$$\partial f_\lambda(\hat{x}) = 0 \Rightarrow \forall l \in [1, n], a_l^\top (y - A\hat{x}) = \lambda u_l.$$

$$\Rightarrow a_l^\top (y - A\hat{x}) \in \begin{cases} [-\lambda, \lambda] & \text{si } \hat{x}_l = 0 \\ \lambda \text{sign}(\hat{x}_l) & \text{si } \hat{x}_l \neq 0 \end{cases}$$

Nous concluons les règles de screening suivantes :

$$\text{si } \hat{x}_l = 0 : |a_l^\top (y - A\hat{x})| \leq \lambda \text{ et si } \hat{x}_l \neq 0 : |a_l^\top (y - A\hat{x})| = \lambda$$

2.4 Expression de λ_{max}

0 est un minimiseur de $f_\lambda \Leftrightarrow (\forall h > 0)(\forall v \neq 0) f_\lambda(hv) \geq f_\lambda(0)$.

$$(\forall h > 0)(\forall v \neq 0)$$

$$f_\lambda(hv) - f_\lambda(0) = \frac{1}{2} \|hAv - y\|_2^2 + \lambda \|hv\|_1 - \frac{1}{2} \|y\|_2^2$$

Donc,

$$\forall h > 0 \quad \forall v \neq 0$$

$$\frac{f_\lambda(hv) - f_\lambda(0)}{h} = \frac{1}{2} (h^2 \|Av\|_2^2 - 2hv^\top A^\top y) + \lambda \|v\|_1$$

pour $v \neq 0$ donné,

$$\begin{aligned} (\forall h > 0) \frac{f_\lambda(hv) - f_\lambda(0)}{h} &\Leftrightarrow (\forall h > 0) \frac{1}{2} (h^2 \|Av\|_2^2 - 2v^\top A^\top y) + \lambda \|v\|_1 \geq 0 \\ &\Leftrightarrow \frac{1}{2} (0 - 2v^\top A^\top y) + \lambda \|v\|_1 \geq 0 \end{aligned}$$

(L'implication directe se justifie en faisant tendre h vers 0, l'autre implication découle de la positivité du premier terme).

Donc pour $v \neq 0$ donné,

$$\begin{aligned} (\forall h > 0 \quad f_\lambda(hv) - f_\lambda(0) \geq 0) &\Leftrightarrow \lambda \|v\|_1 \geq v^\top A^\top y \\ &\Leftrightarrow \lambda \geq \frac{v^\top A^\top y}{\|v\|_1} \end{aligned}$$

Nous en déduisons donc que

$$0 \text{ minimise } f_\lambda \Leftrightarrow \forall v \neq 0 \quad \lambda \geq \frac{V^\top A^\top y}{\|v\|_1}$$

On pose $\psi : v \mapsto \frac{V^\top A^\top y}{\|v\|_1}$

$$\begin{aligned} \forall v \neq 0 \quad |\psi(v)| &\leq \frac{|\sum_i v_i (A^\top y)_i|}{\sum_i |v_i|} \\ &\leq \frac{\sum_i |v_i| \cdot \|A^\top y\|_\infty}{\sum_i |v_i|} \leq \|A^\top y\|_\infty \end{aligned}$$

Et en prenant le vecteur $v = \text{sign}(A^\top y)$, nous obtenons que $\max_v \psi(v) = \|A^\top y\|_\infty$, et donc $\lambda_{\max} = \|A^\top y\|_\infty$.

2.5 Choix du pas de descente

Dans cette sous-section, nous discuterons du choix optimal du pas de descente en choisissant un pas de descente approprié en fonction de la constante de Lipschitz de la fonction $f(x) = \frac{1}{2} \|Ax - y\|_2^2$.

Commençons par établir l'expression de cette constante de Lipschitz. Le gradient de la fonction $f(x)$ est donné par :

$$\nabla f(x) = A^\top (Ax - y) \quad (2)$$

Nous nous apercevons que la constante de Lipschitz du gradient est justement égale à $\sup_{\|x\|_2=1} \|A^\top Ax\|_2$, ou encore à la norme opérateur (avec des normes 2) de la hessienne de la fonction f qui est donné par :

$$\nabla^2 f(x) = A^\top A \quad (3)$$

La constante de Lipschitz du gradient de la fonction est égale à la plus grande valeur propre de la matrice hessienne :

$$\sigma_f = \rho_{\max}(A^\top A) \quad (4)$$

Considérons une étape de la descente de gradient à partir d'un point x vers un point z :

$$z = x - \alpha \nabla f(x) \quad (5)$$

Nous avons :

$$f(z) - f(x) \leq \langle \nabla f(x), z - x \rangle + \frac{\sigma_f}{2} \|z - x\|^2 \quad (6)$$

En substituant z par l'expression de l'étape de la descente de gradient, nous obtenons :

$$f(z) - f(x) \leq -\alpha \|\nabla f(x)\|^2 + \frac{\sigma_f}{2} \alpha^2 \|\nabla f(x)\|^2 \quad (7)$$

La différence $f(z) - f(x)$ est négative, à condition que :

$$\frac{\sigma_f}{2}\alpha^2 - \alpha \leq 0 \quad (8)$$

Ce qui donne la condition suivante pour le pas de descente α :

$$\alpha \leq \frac{2}{\sigma_f} \quad (9)$$

Cette condition assure que le pas de descente garantit une diminution suffisante de la fonction objectif à chaque étape de l'algorithme de descente de gradient.

$$\alpha \leq \frac{2}{\sigma_f} \quad (10)$$

où σ_f est la constante de Lipschitz du gradient de la fonction.

Pour fixer un choix de pas de descente optimal, nous pouvons considérer la valeur $\alpha = \frac{2t}{\sigma_f}$, où $t \in (0, 1)$.

$$f(z) - f(x) \leq -\alpha \|\nabla f(x)\|^2 + \frac{\sigma_f}{2}\alpha^2 \|\nabla f(x)\|^2 \quad (11)$$

En substituant α par $\frac{2t}{\sigma_f}$, nous avons :

$$f(z) - f(x) \leq -\frac{2t}{\sigma_f} \|\nabla f(x)\|^2 + \frac{\sigma_f}{2} \left(\frac{2t}{\sigma_f}\right)^2 \|\nabla f(x)\|^2 \quad (12)$$

Ou encore :

$$f(z) - f(x) \leq -\frac{2t}{\sigma_f} \|\nabla f(x)\|^2 + \frac{2t^2}{\sigma_f} \|\nabla f(x)\|^2 \quad (13)$$

i.e :

$$f(x) - f(z) \geq \frac{2t}{\sigma_f} \|\nabla f(x)\|^2 - \frac{2t^2}{\sigma_f} \|\nabla f(x)\|^2 \quad (14)$$

NOus avons donc :

$$f(x) - f(z) \geq \frac{2}{\sigma_f} t(1-t) \|\nabla f(x)\|^2 \quad (15)$$

La valeur de t peut être ajustée pour contrôler la vitesse de convergence de l'algorithme.

La parabole $t(1-t)$ atteint son maximum en $t = \frac{1}{2}$, et donc un choix optimal du pas de descente consiste à prendre $\alpha_{\text{optimal}} = \frac{1}{\sigma_f}$.

2.6 Descente du gradient proximal

Nous avons implémenté un solveur Lasso en utilisant l'algorithme de descente de gradient proximal. La descente de gradient proximal est une méthode d'optimisation pour résoudre des problèmes de minimisation convexes avec des termes de régularisation non différentiables.

2.6.1 Opérateur proximal

L'opérateur proximal joue un rôle clé dans l'algorithme de descente de gradient proximal. Pour une fonction convexe f , l'opérateur proximal est défini comme suit :

$$\text{prox}_f(v) = \arg \min_x \left\{ f(x) + \frac{1}{2} \|x - v\|_2^2 \right\} \quad (16)$$

L'opérateur proximal est utilisé pour effectuer des mises à jour des coefficients dans l'algorithme de descente de gradient proximal en minimisant une combinaison linéaire de la fonction objectif et de la distance euclidienne au point courant.

Dans le cas du Lasso, la fonction de régularisation est la norme ℓ_1 . L'opérateur proximal de la norme ℓ_1 est donné par :

$$\text{prox}_{\lambda \|\cdot\|_1}(v) = \text{sign}(v) \cdot \max(|v| - \lambda, 0) \quad (17)$$

2.6.2 Algorithme de descente de gradient proximal

L'algorithme de descente de gradient proximal pour le Lasso est implémenté comme suit :

1. Initialiser les coefficients x à zéro et définir les paramètres de l'algorithme (nombre d'itérations et pas).
2. Pour chaque itération :
 - (a) Calculer le gradient de la fonction objectif sans terme de régularisation par rapport à x
 - (b) Mettre à jour les coefficients x en utilisant la fonction proximale de la norme ℓ_1 . Pour ce faire, nous calculons d'abord $x - \alpha \cdot \text{grad}$, où α est le pas de gradient et grad est le gradient de la fonction objectif (en norme 2, sans terme de régularisation). Ensuite, nous appliquons la fonction proximale à ce résultat en utilisant $\alpha \cdot \lambda$ comme paramètre de régularisation :

$$x = \text{prox}_{\alpha \lambda \|\cdot\|_1}(x - \alpha \cdot \text{grad}) \quad (18)$$

- (c) Calculer la valeur de la fonction objectif pour les coefficients x mis à jour.

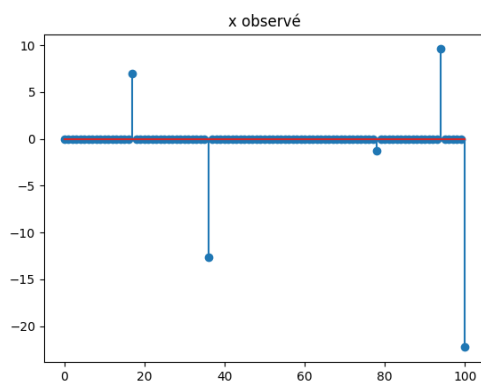
Le code Python correspondant est :

```
def proximal_Lasso(A, y, lambda):
    nb_iter = 200
    alpha = 0.005
    x = np.zeros(n + 1)
    fct_objective = np.zeros(nb_iter)
    for i in range(nb_iter):
        grad = A.T @ (A @ x - y)
        x = proximal_l1(x - alpha * grad, alpha * lambda)
        fct_objective[i] = 0.5 * np.linalg.norm(A @ x - y)**2 + lambda *
            np.sum(np.abs(x))
    return x, fct_objective
```

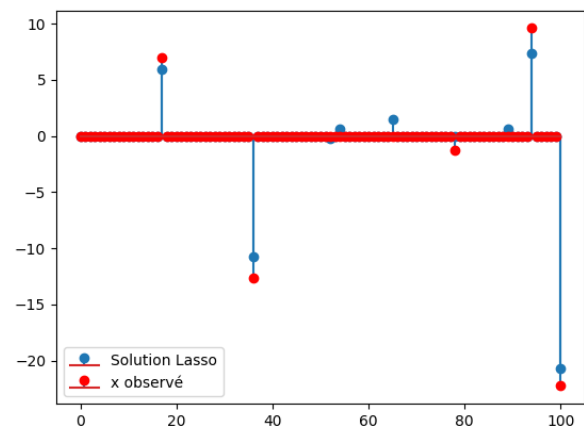
2.6.3 R sultats interm diaires

Nous avons appliqu  notre solveur Lasso   un ensemble de donn es synth tiques al atoires et un certain nombre de composantes non nulles dans le vrai mod le. Les figures ci-dessous montrent les r sultats obtenus avec cette impl mentation, avec $n = 100$ et $m = 25$, et 5 composantes non nulles.

Remarque : dans ce cas, nous n'avons pas accord  beaucoup d'importance au crit re portant sur le nombre de composantes non nulles (vu que la valeur de n est petite, et le coefficient de la p nalit  associ e aux composantes non nulles dans notre impl mentation est seulement de 0,01. Mais cela peut  tre ajust  si besoin.



(a) Vrai mod le



(b) Solution Lasso

FIGURE 4 – Convergence de la fonction objective

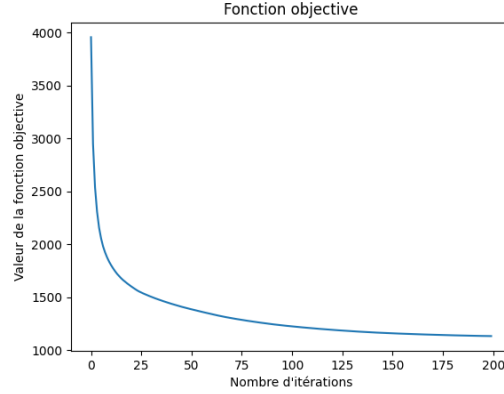


FIGURE 5 – Convergence de la fonction objective

2.7 Étude du Gap

Dans cette partie, nous allons introduire la notion de "duality gap" et en particulier la dualité forte définie dans le poly du cours théorème 15. Introduisons tout d'abord le changement de variable suivant : $\phi_t = \kappa_t(y - Ax_t)$ avec $\kappa_t \in \mathbb{R}^+$, tel que $\lim_{t \rightarrow +\infty} \phi_t = \hat{\phi}$, et $\lim_{t \rightarrow +\infty} x_t = \hat{x}$.

Cette variable aura une importance algorithmique car nous avons :

$$\forall (x, \phi) \in \mathbb{R}^2, f_\lambda(x) \geq -f_\lambda^*(\phi) \quad (19)$$

d'après l'inégalité de Fenchel-Young.

De plus, nous avons la condition d'admissibilité qui est la suivante :

$$0 \text{ solution} \Leftrightarrow \lambda \geq \|A^\top * y\|_\infty \quad (20)$$

On appellera par la suite : $\lambda_{\max} = \|A^\top * y\|_\infty$

Il s'agira par la suite de trouver l'expression de κ_t en utilisant la condition suivante :

$$\lambda \leq \kappa_t \|A^\top * \phi_t\|_\infty = \lambda_{\max} \quad (21)$$

Dès lors, nous pouvons poser : $\kappa_t = \frac{\lambda}{\|A^\top * \phi_t\|_\infty}$ avec $m \in \mathbb{N}$.

Il faut par ailleurs vérifier que cette expression convient. Nous avons tout d'abord,

$$\|A^\top * \hat{\phi}\|_\infty = \max_l |a_l^\top * \hat{\phi}| = \lambda \quad (22)$$

et $\lim_{t \rightarrow +\infty} \phi_t = \hat{\phi}$, qui doit être vérifiée lorsque $\lim_{t \rightarrow +\infty} x_t = \hat{x}$.

Le calcul de limite donne à la fin : $\phi_t(\hat{x}) = \frac{\hat{\phi}}{m}$ ce qui donne $m = 1$. La condition d'admissibilité est vérifiée par construction de κ_t .

Nous voyons donc que l'expression de ϕ_t peut être simplifiée et donne donc

$$\phi_t = \frac{\lambda}{\|A^\top * \phi_t\|_\infty} * (y - Ax_t) \quad (23)$$

L'algorithme a donc pour but de trouver la solution qui nous intéresse à savoir $f_\lambda(\hat{x})$ en minimisant l'écart entre cette valeur à savoir

$$\inf_{x \in \mathbb{R}^n} (\|f_\lambda(\hat{x}) - f_\lambda(x)\|_\infty) \quad (24)$$

Or, nous n'avons pas accès à \hat{x} donc nous utiliserons la relation 19 afin de contourner le problème. Nous obtenons donc

$$\inf_{x \in \mathbb{R}^n} (\|f_\lambda^*(\hat{\phi}) + f_\lambda(x)\|_\infty) \quad (25)$$

équivalent au problème suivant :

$$\inf_{x \in \mathbb{R}^n} (\left\| \lim_{t \rightarrow +\infty} f_\lambda^*(\phi_t) + f_\lambda(x) \right\|_\infty) \quad (26)$$

Sur python, nous pouvons simplifier le problème encore un peu en posons un critère ϵ aussi petit que l'on veut afin d'arrêter la boucle while ce qui se traduit donc par :

$$|f_\lambda^*(\phi_t) + f_\lambda(x)| \leq \epsilon \quad (27)$$

La condition d'admissibilité est par ailleurs vérifiée à chaque boucle grâce à la construction de κ_t .

3 Problème des moindres carrés non-négatifs

3.1 Utilité de ce problème

Le problème des moindres carrés non négatifs est une variante du problème des moindres carrés classique. Il est utilisé lorsque l'on souhaite trouver une solution qui est contrainte d'être non négative. L'objectif du problème des moindres carrés non négatifs est de minimiser la somme des carrés des écarts entre les valeurs observées et les valeurs prédites, tout en respectant la contrainte de non-négativité. Ce problème trouve des applications dans de nombreux domaines. Par exemple, en analyse de données, il peut être utilisé pour estimer les coefficients d'un modèle linéaire avec des contraintes de non-négativité. Dans le domaine de l'imagerie, il peut être utilisé pour la reconstruction d'images à partir de données incomplètes ou bruitées, en imposant que les valeurs des pixels soient non négatives.

3.2 Convexité et non-différentiabilité

Dans cette section, nous nous intéressons maintenant au problème des moindres carrés non-négatifs :

$$\hat{x} = \underset{x \in \mathbb{R}^n}{\text{Argmin}} \left\{ \frac{1}{2} \|y - Ax\|_2^2 + I_{[0, +\infty[^n}(x) \right\}.$$

Comme dans le cas du Lasso, la fonction objective est convexe ce qui permet d'appliquer les résultats de l'optimisation convexe, et elle est non-différentiable vu le deuxième terme de la fonction objective.

3.3 Dualité de Fenchel

Posons : $g(z) = \frac{1}{2} \|y - z\|_2^2$ et $f(x) = I_{[0, +\infty[^n}(x)$. Cherchons donc :

$$\inf_{x \in \mathbb{R}^n} (g(Ax) + f(x))$$

Par Fenchel Rockafeller [3], comme g et f sont des fonctions convexes, propres et semi-continues inférieures car \mathbb{R}_+^n est convexe fermé et $\text{dom}(Af) \cap \text{dom}(g) = \text{dom}(Af) \cap \mathbb{R}_+^n \neq \emptyset$

$$\begin{aligned} \inf_{x \in \mathbb{R}^n} (f(x) + g(Ax)) &= - \inf_{\phi \in \mathbb{R}^n} (f^*(A^*\phi) + g^*(-\phi)) \\ &= - \inf_{\phi \in \mathbb{R}^n} (I_{]-\infty, 0]^n}(A^*\phi) + \frac{1}{2} \|y - \phi\|_2^2 - \frac{1}{2} \|y\|_2^2) \end{aligned}$$

Comme la dualité est forte, nous avons le cas d'égalité dans l'inégalité de Fenchel-Young pour g (en utilisant la même preuve faite dans le problème Lasso) :

$$\begin{aligned} g(A\hat{x}) + g^*(-\hat{\phi}) &= \langle A\hat{x}, -\hat{\phi} \rangle \\ \Leftrightarrow \hat{\phi} &= y - A\hat{x} \end{aligned}$$

3.4 Sous-différentiel

Comme f n'est pas différentiable, on nous calculons son sous-gradient qui est bien défini car f est convexe. Le sous-gradient S de $f(x)$ en x est défini par l'inégalité suivante pour tout $y \in \mathbb{R}^n$:

$$f(y) \geq f(x) + S(y - x)$$

Si $x_i > 0$ pour tout i , on a $f(x) = 0$. Dans ce cas, nous voulons que

$$f(y) \geq S(y - x)$$

— Si pour tout $y_i \geq 0$ (avec $y \in \mathbb{R}_+^n$), $f(y) = 0$, et nous cherchons S tel que

$$0 \geq S(y - x)$$

Choisissons donc $S_i = 0$ pour tout i .

— S'il existe i tel que $y_i < 0$, $f(y) = +\infty$, et l'inégalité est vraie pour n'importe quelle valeur de S . Nous choisissons donc $S = 0$.

S'il existe i tel que $x_i = 0$,

— Si pour tout $y_i \geq 0$ (avec $y \in \mathbb{R}_+^n$), $f(y) = 0$, donc nous cherchons S tel que

$$0 \geq S(y - x)$$

Et nous choisissons donc n'importe quelle $S_i \leq 0$.

— S'il existe i tel que $y_i < 0$, $f(y) = +\infty$, et l'inégalité est vraie pour n'importe quelle valeur de S .

Au final,

$$\partial f(x) = \left\{ u \in \mathbb{R}^n : u_l = \begin{cases} \mathbb{R}_- & \text{si } x_l = 0 \\ 0 & \text{sinon} \end{cases}, l \in \{1, 2, \dots, n\} \right\}$$

3.5 Test de screening

Soit \hat{x} solution admissible, en annulant le sous-gradient de la fonction objective :

$$\forall l \in [1, n], a_l^\top (y - A\hat{x}) \in \begin{cases} \mathbb{R}_- & \text{si } \hat{x}_l = 0 \\ 0 & \text{si } \hat{x}_l > 0 \end{cases}$$

Nous déduisons le test de screening :

$$a_l^\top (y - A\hat{x}) < 0 \Rightarrow x_l = 0$$

3.6 Sphère Gap pour les moindres carrés non-négatifs

Pour cette partie, nous avons choisi la méthode du gradient projeté au lieu du gradient proximal, dont la formule est la suivante :

$$x_{t+1} = x_t + \eta * (P_{r_j}(x_t - \gamma * grad) - x_t). \quad (28)$$

Dans notre algorithme, nous avons choisi $\eta = 1$ pour le coefficient de relaxation et P_{r_j} correspond à la projection de $x_t - \gamma * grad$ sur \mathbb{R}_+^n (cf. Cours d'optimisation). De plus, l'étude mathématique de la sphère Gap est semblable à celle réalisée précédemment pour l'exemple du Lasso et nous aboutissons donc au même résultat.

Ce qui devra être modifié est l'expression de ϕ_t qui rappelons-le, dans le problème du Lasso, dépendait de λ afin de la rendre admissible, d'où l'introduction dans l'algorithme du paramètre k_t .

Démonstration :

Pour que ϕ_t soit admissible, c'est-à-dire $[A^\top * \phi_t]_i \leq 0 \forall i \in \{1, 2, \dots, n\}$, il faudra introduire un facteur tel que

$$\phi_t = y - k_t * A * x_t \quad (29)$$

Soit $y \in \mathbb{R}^m$ et $A \in \mathbb{R}_+^{m \times n}$,

Premier cas :

Supposons, $\forall i \in \{1, 2, \dots, n\}$, $[A * x_t]_i \neq 0$,

Dès lors, nous avons

$$[y]_i - k_t * [A * x_t]_i \leq 0 \Leftrightarrow \frac{y_i}{[A * x_t]_i} \leq k_t \quad (30)$$

Nous pouvons donc généraliser en prenant

$$k_t = \max\left(\frac{\|y\|_\infty}{\|A * x_t\|_\infty}, 1\right) \quad (31)$$

De plus, $\hat{\phi}$ est admissible, donc $[A^\top * \hat{\phi}]_i \leq 0 \Leftrightarrow \|y\|_\infty \leq \|A * \hat{x}\|_\infty$,

d'où le choix de 1 dans l'expression de k_t permettant la convergence de ϕ_t vers $\hat{\phi}$.

Deuxième cas :

$\exists i \in \{1, 2, \dots, n\}$ tel que $[A * x_t]_i = 0$.

L'algorithme devra donc s'arrêter car $\exists i$ tel que $[A^\top * \hat{\phi}]_i \geq 0$ ce qui va renvoyer $+\infty$ dans l'expression de la fonction indicatrice.

4 Implémentation

4.1 Pas de descente

4.1.1 Théorie autour du pas

Dans cette sous-section, nous discuterons du choix optimal du pas de descente en choisissant un pas de descente approprié en fonction de la constante de Lipschitz de la fonction $f(x) = \frac{1}{2} \|Ax - y\|_2^2$.

Commençons par établir l'expression de cette constante de Lipschitz. Le gradient de la fonction $f(x)$ est donné par :

$$\nabla f(x) = A^T(Ax - y) \quad (32)$$

Nous nous apercevons que la constante de Lipschitz du gradient est justement égale à $\sup_{\|x\|_2=1} \|A^T Ax\|_2$, ou encore à la norme opérateur (avec des normes 2) de la hessienne de la fonction $f(x)$ qui est donné par :

$$\nabla^2 f(x) = A^T A \quad (33)$$

La constante de Lipschitz du gradient de la fonction est égale à la plus grande valeur propre de la matrice hessienne :

$$\sigma_f = \lambda_{\max}(A^T A) \quad (34)$$

Considérons une étape de la descente de gradient à partir d'un point x vers un point z :

$$z = x - \alpha \nabla f(x) \quad (35)$$

Nous avons :

$$f(z) - f(x) \leq \langle \nabla f(x), z - x \rangle + \frac{\sigma_f}{2} \|z - x\|^2 \quad (36)$$

En substituant z par l'expression de l'étape de la descente de gradient, nous obtenons :

$$f(z) - f(x) \leq -\alpha \|\nabla f(x)\|^2 + \frac{\sigma_f}{2} \alpha^2 \|\nabla f(x)\|^2 \quad (37)$$

La différence $f(z) - f(x)$ est négative, à condition que :

$$\frac{\sigma_f}{2} \alpha^2 - \alpha \leq 0 \quad (38)$$

Ce qui donne la condition suivante pour le pas de descente α :

$$\alpha \leq \frac{2}{\sigma_f} \quad (39)$$

Cette condition assure que le pas de descente garantit une diminution suffisante de la fonction objectif à chaque étape de l'algorithme de descente de gradient.

$$\alpha \leq \frac{2}{\sigma_f} \quad (40)$$

où σ_f est la constante de Lipschitz du gradient de la fonction.

Pour fixer un choix de pas de descente optimal, nous pouvons considérer la valeur $\alpha = \frac{2t}{\sigma_f}$, où $t \in (0, 1)$.

$$f(z) - f(x) \leq -\alpha \|\nabla f(x)\|^2 + \frac{\sigma_f}{2} \alpha^2 \|\nabla f(x)\|^2 \quad (41)$$

En substituant α par $\frac{2t}{\sigma_f}$, nous avons :

$$f(z) - f(x) \leq -\frac{2t}{\sigma_f} \|\nabla f(x)\|^2 + \frac{\sigma_f}{2} \left(\frac{2t}{\sigma_f}\right)^2 \|\nabla f(x)\|^2 \quad (42)$$

Ou encore :

$$f(z) - f(x) \leq -\frac{2t}{\sigma_f} \|\nabla f(x)\|^2 + \frac{2t^2}{\sigma_f} \|\nabla f(x)\|^2 \quad (43)$$

i.e :

$$f(x) - f(z) \geq \frac{2t}{\sigma_f} \|\nabla f(x)\|^2 - \frac{2t^2}{\sigma_f} \|\nabla f(x)\|^2 \quad (44)$$

Nous avons finalement :

$$f(x) - f(z) \geq \frac{2}{\sigma_f} t(1-t) \|\nabla f(x)\|^2 \quad (45)$$

La valeur de t peut être ajustée pour contrôler la vitesse de convergence de l'algorithme.

La parabole $t(1-t)$ atteint son maximum en $t = \frac{1}{2}$, et donc un choix optimal du pas de descente consiste à prendre $\alpha_{\text{optimal}} = \frac{1}{\sigma_f}$.

4.1.2 Actualisation des pas

Nous avons testé plusieurs versions de l'algorithme en tenant compte du pas : Une première version consiste à considérer un pas constant tout au long de l'exécution du code, ce pas constant était $\alpha = \frac{1}{\sigma_f}$. Les résultats n'étaient pas à la hauteur de ce que nous espérions en terme de nombre d'itérations, donc nous avons opté pour un pas variable qui varie intrinsèquement avec la dimension de la matrice hessienne de la fonction objective. Ce pas décroît naturellement vu que la dimension de la matrice décroît après chaque itération de l'algorithme. Cette modification permet d'avoir des résultats plus performants en nombre d'itérations à précision fixe. Mais nous avons

pensé que nous pouvions faire une amélioration de notre implémentation. Cette dernière consiste à varier les pas par des blocs d'itérations, nous avons choisi par exemple de les varier par blocs de 20 itérations, avec cette méthode nous obtenons une convergence plus rapide en nombre d'itération avec une précision donnée.

4.2 Une région safe - Sphère Gap

Une fois cet algorithme de gradient proximal (sans screening) implémenté, nous nous intéressons à une implémentation avec screening, pour procéder ensuite à une comparaison entre les performances du Gradient proximal avec et sans screening. Les tests de screening faisant intervenir les solutions du primal ou du dual auxquelles nous n'avons pas accès, nous pouvons construire des régions safe contenant $\hat{\phi}$.

En passant par la formulation duale et en notant $-D(\phi) = \frac{1}{2}\|y\|^2 - \frac{1}{2}\|y - \phi\|^2$, $-D$ est 1-fortement concave, nous avons donc pour ϕ admissible :

$$-D(\phi) \leq -D(\hat{\phi}) + \langle \nabla(-D)(\hat{\phi}), \phi - \hat{\phi} \rangle - \frac{1}{2}\|\hat{\phi} - \phi\|^2.$$

Par dualité faible, nous avons :

$$-D(\hat{\phi}) \leq P(x).$$

De plus, par concavité et optimalité de $\hat{\phi}$:

$$\langle \nabla - D(\hat{\phi}), \phi - \hat{\phi} \rangle \leq 0.$$

En combinant ces inégalités, nous obtenons :

$$-D(\phi) \leq P(x) + \langle \nabla - D(\hat{\phi}), \phi - \hat{\phi} \rangle - \frac{1}{2}\|\hat{\phi} - \phi\|^2.$$

Comme $\langle \nabla - D(\hat{\phi}), \phi - \hat{\phi} \rangle \leq 0$, on a :

$$-D(\phi) \leq P(x) - \frac{1}{2}\|\hat{\phi} - \phi\|^2.$$

Nous obtenons finalement :

$$2 \cdot \text{gap}(x, \phi) \geq \|\hat{\phi} - \phi\|^2.$$

i.e :

$$\sqrt{2 \cdot \text{gap}(x, \phi)} \geq \|\hat{\phi} - \phi\|.$$

Ce qui prouve que la sphère de centre ϕ et de rayon $R = \sqrt{2 \cdot \text{gap}(x, \phi)}$ contient $\hat{\phi}$. En notant $B(\phi, R)$ cette boule, nous avons la propriété suivante :

$$\max\{|\langle a_i, v \rangle| : v \in B(\phi, R)\} = |\langle a_i, \phi \rangle| + R\|a_i\|_2$$

Il suffit pour s'en apercevoir de prendre $v = \phi + R\mathbf{e}$, où \mathbf{e} est un vecteur dans le disque unitaire $D(0, 1)$. Comme $v \in B(\phi, R)$, on a :

$$\|v - \phi\|_2 = R\|\mathbf{e}\|_2 \leq R.$$

Ainsi,

$$\begin{aligned} |\langle a_i, v \rangle| &= |\langle a_i, \phi + R\mathbf{e} \rangle| \\ &\leq |\langle a_i, \phi \rangle| + R|\langle a_i, \mathbf{e} \rangle|. \end{aligned}$$

le maximum de $|\langle a_i, v \rangle|$ est atteint pour $\mathbf{e} = \frac{a_i}{\|a_i\|_2}$, et donc :

$$\max\{|\langle a_i, v \rangle| : v \in B(\phi, R)\} = |\langle a_i, \phi \rangle| + R\|a_i\|_2.$$

Cette propriété va être utile pour effectuer le test de screening, sachant que nous ne disposons pas de la solution $\hat{\phi}$. Pour des régions safe contenant $\hat{\phi}$, nous pouvons utiliser cette propriété pour vérifier si $\max\{|\langle a_i, v \rangle| : v \in R\} < \lambda$. Si cette condition est vérifiée, alors nous pouvons éliminer la i -ème variable.

4.3 Gradient Proximal avec screening

Suivant la discussion sur les régions safe, nous implémentons la méthode Proximale Lasso avec screening. Nous suivons une démarche similaire pour l'algorithme sans screening ;

```
def proximal_Lasso_screening(A, y, lambda, max_iter=100, epsi=1e-4):
    m, n = A.shape
    original_n = n
    x = np.zeros(original_n)
    alpha = 1 / np.linalg.norm(A.T @ A, 2)
    final_x = np.zeros(n)
    objective_history = []
    A_screened = np.vstack([A, np.arange(original_n)])
    R = np.inf

    for iter in range(max_iter):
        res = y - A_screened[:-1, :] @ x
        grad = A_screened[:-1, :].T @ (-res)
        if iter % 20 == 0:
            alpha = 1 / np.linalg.norm(A_screened[:-1, :].T @ A_screened[:-1, :], 2)
        x = proximal_l1(x - alpha * grad, alpha * lambda)
        u = res
```

```

Atres = A_screened[:-1, :].T @ u
norm_Au_inf = np.linalg.norm(Atres, ord=np.inf)
kap = (lmbda / norm_Au_inf)
u = kap * u
P_x = 0.5 * np.linalg.norm(res)**2 + lmbda * np.sum(np.abs(x))
D_u = 0.5 * np.linalg.norm(y)**2 - 0.5 * np.linalg.norm(y - u)**2
gap = P_x - D_u
objective_history.append(P_x)
centre = u
R = np.sqrt(2 * gap)
print(R)
if np.abs(gap) < epsi:
    break

a_reduit = A_screened[:-1, :]
maxx_valeurs = np.abs(a_reduit.T @ centre) + R * np.linalg.norm(a_reduit, axis=0)
to_delete = maxx_valeurs < lmbda
A_screened = A_screened[:, ~to_delete]
x = x[~to_delete]
n -= np.sum(to_delete)

final_indices = A_screened[-1, :].astype(int)
for i, index in enumerate(final_indices):
    final_x[index] = x[i]

screened_indices = np.setdiff1d(np.arange(original_n), final_indices)
return final_x, centre, R, screened_indices, objective_history

```

Dans cette fonction, après chaque mise à jour de x , nous calculons u , la solution duale courante. Ensuite, nous calculons le gap entre la fonction objectif duale et la fonction objectif primale, et utilisons cela pour calculer R . Nous éliminons ensuite les variables dont la valeur maximale, calculée en utilisant la propriété présentée dans la section précédente, est inférieure à λ . Enfin, nous mettons à jour A_{screened} , x et n pour refléter les variables restantes.

En utilisant ensuite la propriété présentée dans la section précédente, nous éliminons les colonnes qui passent le test, en faisant attention aux indices correspondants dans la matrice originale (moyennant une dernière ligne ajoutée qui sauvegarde l'indice initial).

4.4 Résultats de l'implémentation

La version finale de l'implémentation montre les résultats présentés dans cette partie.

4.4.1 Nombre d'it rations

La figure suivante pr sente une comparaison du nombre des it rations n cessaires pour atteindre la fonction objective :

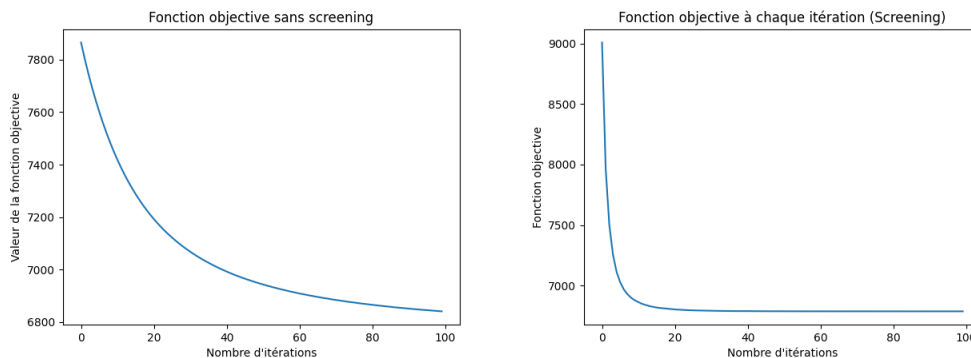


FIGURE 6 – Comparaison de la convergence de la fonction objective

Nous remarquons que la fonction objective est atteinte apr s plus ou moins 15 it rations dans la version avec screening, cependant, elle n cessite 80 it rations pour pouvoir atteindre la valeur minimale dans la version sans screening. Donc en raisonnant   une pr cision fixe, nous pouvons montrer que le screening permet de r duire drastiquement le nombre d'it ration. Ce qui laisse supposer que le screening permet de r duire le temps d'ex cution de l'algorithme.

4.4.2 Temps d'ex cution

Comme mentionn  pr c demment, la diminution des nombres d'it rations nous laisse penser   la diminution du temps de l'ex cution. Malheureusement, nous n'avons pas pu aboutir   cela, vu que notre code n cessite d' tre optimiser. Il faut aussi pr ciser que la comparaison des temps de l'ex cution n'est pas  quitable vu que le mod le sans screening ne contient pas l'impl mentation du gap et son calcul   chaque it ration, cependant le mod le avec screening contient cette  tape qui est tr s c teuse en complexit  temporelle et spatiale. Ajoutons   cela que l'impl mentation avec screening contient des actualisations des pas de la descente de gradient donc des calculs du rayon spectrale. C'est pour ces raisons que le temps de l'ex cution de l'algorithme avec et sans screening ne donne pas les r sultats qu'on voulait.

Voil  une comparaison des temps d'ex cution avec et sans screening :

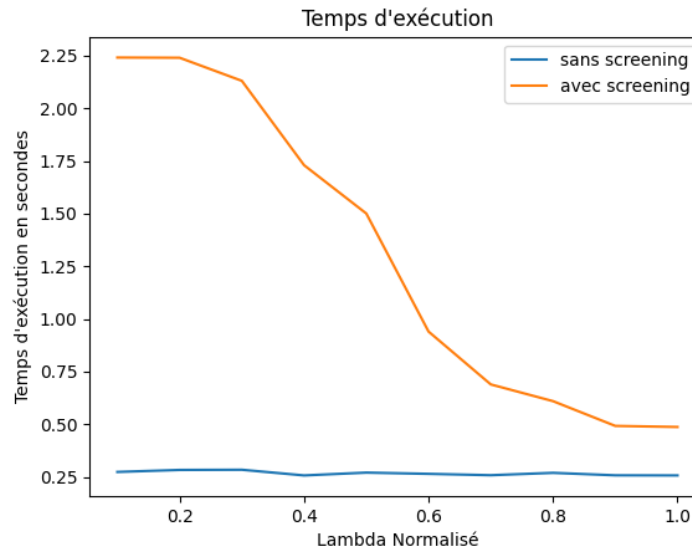


FIGURE 7 – Comparaison des temps d'ex cution

4.4.3 Empreinte carbone

L'empreinte carbone est obtenu par le calcul des quantit s de CO2  mises par le CPU et la RAM, qui sont calcul s par un carbone tracker du module psutil. Les valeurs obtenues pour l'utilisation du CPU et de la RAM sont pond r es et additionn es pour obtenir l'empreinte carbone totale de l'ex cution du code. Comme l'empreinte carbone est directement li  aux temps de l'ex cution, les r sultats sont toujours contradictoires   notre objectif qu'on a annonc  au d but, c'est   dire que le screening causera une diminution de l'empreinte carbone. Voil  une comparaison entre les quantit s de carbone en Kg  mises lors de l'ex cution du code :

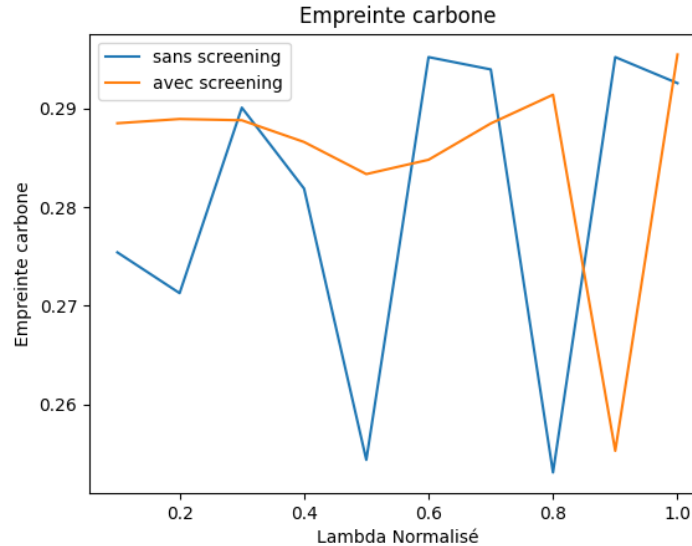


FIGURE 8 – Comparaison des empreintes carbone

4.5 Impl mentation Moindres Carr s non-n gatifs

Nous impl mentons l'algorithme pour les moindres carr s non-n gatifs en suivant une d marche analogue avec celle employ e pour le gradient proximal ;

```
def mcnnscreening(A, y, max_iter=1000, epsi=1e-4):
    m, n = A.shape
    original_n = n
    alpha = 0.1 / np.linalg.norm(A.T @ A, 2)
    x = np.ones(original_n)
    objective_history = []
    A_screened = np.vstack([A, np.arange(original_n)])
    for iter in range(max_iter):
        grad = A_screened[:-1, :].T @ (A_screened[:-1, :] @ x - y)
        if iter % 20 == 0:
            alpha = 1 / np.linalg.norm(A_screened[:-1, :].T @ A_screened[:-1, :], 2)
        x = np.maximum(x - alpha * grad, 0)
        Ax = A_screened[:-1, :] @ x
        if np.all(Ax > 0):
            kap = np.max([1, np.max(y / Ax)])
        else:
            break
        u = y - kap * A_screened[:-1, :] @ x
        P_x = 0.5 * np.linalg.norm(y - A_screened[:-1, :] @ x)**2
        D_u = 0.5 * np.linalg.norm(y)**2 - 0.5 * np.linalg.norm(y - u)**2
        gap = P_x - D_u
```



```

print(gap)
objective_history.append(P_x)
R = np.sqrt(2 * np.abs(gap))
if R < epsi:
    break
to_delete= []
for i in range(n):
    a_i = A_screened[:-1, i]
    maxx = np.dot(a_i, u) + R * np.linalg.norm(a_i)
    if maxx < 0:
        to_delete.append(i)
A_screened = np.delete(A_screened, to_delete, axis=1)
x = np.delete(x, to_delete)
n -= len(to_delete)
final_indices = A_screened[-1, :].astype(int)
final_x = np.zeros(original_n)
screened_indices = np.setdiff1d(np.arange(original_n), final_indices)
final_x[final_indices] = x
return final_x, screened_indices, objective_history

```

En se restreignant au cas où la matrice A est positive, et en adaptant les tests de screening (la valeur absolue disparaît et λ est remplacé par zéro, conformément aux nouveaux tests de screening obtenus, et le calcul du maximum de la quantité considérée sur les sphères gap découle d'un raisonnement similaire à celui pour le Lasso). Contrairement à ce que nous observons pour le cas du Lasso avec le Proximal avec screening, très peu de colonnes passent le test, ce qui se répercute sur la vitesse de convergence de cette implémentation avec screening par opposition à celle du Lasso ; le nombre d'itérations requises pour atteindre la convergence est très supérieur à celui qu'on avait précédemment.

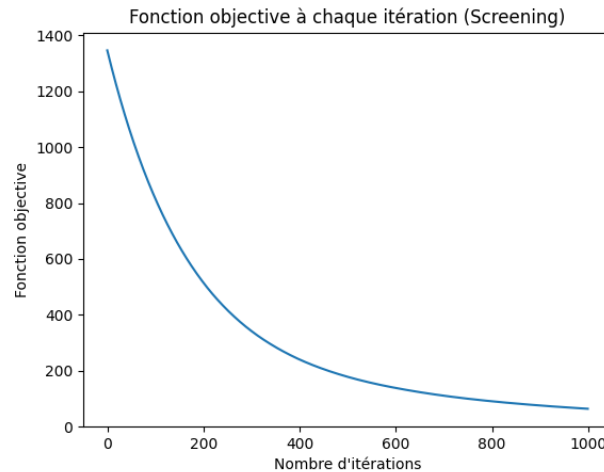


FIGURE 9 – Convergence pour l’algorithme avec screening dans le cas des Moindres Carres non négatifs

5 Conclusion

Nous avons pu obtenir des résultats concluants quant à la convergence du modèle avec screening qui s’avère être plus rapide en nombre d’itérations à précision donnée mais qui peut être plus coûteuse en temps et en énergie (et donc un impact plus grand sur l’environnement). En revanche, pour y remédier, nous pouvons dans un premier temps calculer le temps d’exécution sur un nombre d’itérations adéquats pour chaque méthode et comparer. Ensuite, régler des problèmes techniques dans les algorithmes par exemple optimiser le nombre de fois où nous calculons certaines variables mais aussi choisir peut-être une nouvelle méthode pour évaluer le pas de descente de gradient. Enfin, nous avons travaillé avec une matrice A dont les coefficients appartiennent à \mathbb{R}_+ ce qui n’est pas forcément une contrainte que nous pourrions trouver dans des problèmes plus concrets.

Références

- [1] Elvira Clément et Herzet Cédric, Safe Rules for the Identification of zeros in the Solutions of the Slope Problem, 2020
- [2] Ghaoui L. Viallon V. et Rabbani T., Safe feature elimination in sparse supervised learning, *Pacific Journal of Optimization*. 2010
- [3] Pesquet Jean-Christophe, Cours d’optimisation CentraleSupélec, 2022.