



HAL
open science

Étude de l'estimateur SLOPE par le prisme du schéma : Propriétés de parcimonie et d'appariement et calcul du chemin des solutions

Patrick J C Tardivel

► **To cite this version:**

Patrick J C Tardivel. Étude de l'estimateur SLOPE par le prisme du schéma : Propriétés de parcimonie et d'appariement et calcul du chemin des solutions. 2024. hal-04528428v1

HAL Id: hal-04528428

<https://hal.science/hal-04528428v1>

Preprint submitted on 1 Apr 2024 (v1), last revised 20 Sep 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Étude de l'estimateur SLOPE par le prisme du schéma : Propriétés
de parcimonie et d'appariement et calcul du chemin des solutions

Patrick Tardivel, Institut de Mathématiques de Bourgogne,
UMR 5584, CNRS & Université de Bourgogne, F-21000 Dijon, France

Résumé

L'estimateur SLOPE (acronyme signifiant « Sorted L One Penalized Estimation ») est défini comme une solution d'un problème d'optimisation convexe où le terme de pénalité est la norme ℓ_1 ordonnée. Cette norme, non-différentiable en un point ayant des composantes nulles ou égales en valeur absolue, induit des propriétés de parcimonie et d'appariement à l'estimateur SLOPE (certaines composantes de cet estimateur peuvent être nulles ou égales en valeur absolue). Ce mémoire d'habilitation illustre la pertinence de la notion de schéma du SLOPE pour l'étude de cet estimateur. La première partie de ce travail prouve que les schémas sont des classes d'équivalence pour la relation « avoir le même sous-différentiel pour la norme ℓ_1 ordonnée » et établit une bijection entre ces schémas et les faces de la boule unité de la norme ℓ_1 ordonnée duale. La seconde partie de ce mémoire donne des conditions théoriques garantissant la récupération du schéma des coefficients de régression via l'estimateur SLOPE. La dernière partie de ce manuscrit fournit un algorithme, basé sur les conditions du sous-différentiel et du schéma, pour calculer le chemin des solutions du SLOPE lorsque le paramètre de régularisation varie.

Mots-clés : SLOPE, Norme ℓ_1 ordonnée, Schéma du SLOPE, Parcimonie, Appariement, Chemin des solutions

Table des matières

Introduction	2
1 Estimateurs favorisant la parcimonie et l'appariement	8
1.1 Introduction	8
1.2 Schéma du SLOPE	10
1.3 Annexe : preuves	12
2 Récupération du schéma par l'estimateur SLOPE	18
2.1 Notions liées au schéma du SLOPE	18
2.2 Propriétés théoriques de récupération du schéma du SLOPE	20
2.2.1 Recupération du schéma dans le cas non-bruité : condition d'irreprésentabilité	21
2.3 Relâchement de la condition d'irreprésentabilité du SLOPE	24
2.4 Expériences numériques	25
2.5 Annexes : preuves	27
3 Chemin des solutions et chemin des valeurs ajustées de l'estimateur SLOPE	36
3.1 Introduction	36
3.2 Chemin multidimensionnel des solutions et valeurs ajustées du SLOPE	37
3.2.1 Chemin du gradient et groupes d'appariement	38
3.3 Calcul du chemin unidimensionnel de l'estimateur SLOPE	39
3.4 Expériences numériques	41
3.4.1 Calcul des chemins et minimisation de la somme des carrés résiduels sur l'échantillon de validation	43
3.4.2 Limite de cet algorithme de calcul du chemin des solutions SLOPE	44
3.5 Annexes : preuves	44
Conclusion et perspectives	50
A Généralisation de certaines propriétés de l'estimateur SLOPE	52

Introduction

Mon parcours scientifique depuis mes débuts en doctorat

Les estimateurs pénalisés ont été le fil conducteur de mes travaux depuis mes débuts en thèse jusqu'à aujourd'hui. En doctorat, j'ai exploré l'utilisation de l'estimateur LASSO (acronyme pour « Least Absolute Shrinkage and Selection Operator ») (Tibshirani, 1996), dont le terme de pénalité est la norme ℓ_1 , dans le contexte de la métabolomique. Bien que l'étude de cet estimateur ait été importante, elle n'était pas au centre de mes recherches doctorales. Mes principales contributions à l'époque comprenaient un article appliqué en métabolomique (Tardivel *et al.*, 2017) ainsi qu'un article théorique portant sur l'optimisation non convexe (Tardivel *et al.*, 2018). Mes recherches sur les estimateurs pénalisés ont pris de l'ampleur après mon doctorat, et je me suis spécialisé dans l'étude de l'estimateur SLOPE (acronyme pour « Sorted L One Penalized Estimation ») (Bogdan *et al.*, 2015; Zeng et Figueiredo, 2014), qui généralise l'estimateur LASSO en utilisant une norme ℓ_1 ordonnée comme terme de pénalité. Les paragraphes suivants présentent mes travaux de recherche de manière chronologique, mettant en lumière la façon dont mon sujet d'habilitation découle naturellement des thèmes que j'ai abordés tout au long de ma carrière.

Mes années de doctorat à l'INRA de Toulouse

Entre novembre 2014 et novembre 2017, j'ai effectué ma thèse intitulée « Représentation parcimonieuse et procédures de tests multiples : Application à la métabolomique » à l'Institut National de la Recherche Agronomique (INRA) au centre Toxalim à Toulouse sous la direction de Didier Concordet et Rémi Servien. L'un des objectifs de ma thèse était d'identifier, dans un mélange complexe de métabolites obtenu par Résonance Magnétique Nucléaire (RMN), les métabolites pertinents.

Sans rentrer dans les détails des nombreux pré-traitements, on peut modéliser un signal RMN comme la réponse d'un modèle de régression linéaire, où les variables explicatives, non-aléatoires, sont les spectres RMN des métabolites purs. L'identification des métabolites revient donc à déterminer les coefficients de régression non nuls de ce modèle. Après avoir tenté en vain de développer une méthode utilisant la programmation linéaire, j'ai opté pour une approche alternative basée sur l'estimateur LASSO. Ce choix était motivé par la parcimonie de l'estimateur LASSO (certaines composantes du LASSO sont nulles), sa popularité croissante à l'époque, et par ma formation antérieure sur le LASSO lors de mes études en master. Mon objectif était d'identifier les métabolites pertinents en utilisant le support du LASSO. Les experts recommandaient d'éviter les faux positifs, quitte à ne pas identifier certains métabolites présents en faible quantité dans le mélange. Ainsi, j'ai cherché à contrôler la probabilité d'avoir au moins un faux positif, c'est-à-dire la probabilité qu'une composante non nulle du LASSO soit associée à un métabolite absent du mélange.

Théoriquement, la récupération exacte du support des coefficients de régression avec le LASSO nécessite une hypothèse contraignante appelée condition d'irreprésentabilité (Zhao et Yu, 2006; Zou, 2006). Cette condition peut être assouplie en utilisant le LASSO adaptatif, où les poids sont inversement proportionnels aux

composantes, en valeur absolue, d'un estimateur convergent (Zou, 2006). Cependant, pour le LASSO ainsi que pour le LASSO adaptatif, le calibrage du paramètre de régularisation pour contrôler le taux de faux positif est techniquement très difficile.

Au final, bien que l'estimateur LASSO soit souvent mentionné dans ma thèse (Tardivel, 2017) et dans l'article en métabolomique (Tardivel *et al.*, 2017), cet estimateur s'est avéré peu utile pour mes travaux. Quelques années après ma soutenance, j'ai épuré mes travaux de doctorat de la technicité due à l'introduction de l'estimateur LASSO, et j'ai développé une procédure de tests multiples contrôlant la probabilité d'obtenir un ou plusieurs faux positifs, basée sur la construction d'une région de confiance rectangulaire de volume minimal (Tardivel *et al.*, 2021).

Mes années en tant qu'enseignant-chercheur à l'université de Wrocław

Après ma rencontre avec Małgorzata Bogdan en juillet 2017 lors d'un congrès, j'ai commencé à travailler à l'université de Wrocław en Pologne, d'abord en tant que chercheur invité de février à août 2018, puis en tant qu'adiunkt¹ de septembre 2018 à août 2020. Pendant cette période, mes travaux ont principalement été théoriques. L'une de mes contributions était en continuité directe avec ma thèse. En effet, j'ai fourni une preuve originale d'un résultat déjà connu : la condition d'irreprésentabilité implique, dans le cas idéal où le bruit est nul, que la minimisation de la norme ℓ_1 permet de récupérer les coefficients de régression (Fuchs, 2004, théorème 4). Sous cette dernière condition, l'estimateur LASSO, lorsque le bruit est faible, sépare les coefficients de régression nuls des coefficients non nuls. Ainsi, un seuillage appliqué à l'estimateur LASSO permet de récupérer le support des coefficients de régression (Tardivel et Bogdan, 2022). Une autre contribution a été l'obtention d'une formule pour l'estimateur SLOPE lorsque les colonnes de la matrice de régression sont orthogonales, c'est-à-dire une formule pour l'opérateur proximal de la norme ℓ_1 ordonnée. Ce fut mon premier article sur l'estimateur SLOPE. Le calcul de l'opérateur proximal pouvait se faire de manière algorithmique (Bogdan *et al.*, 2015), et ma principale motivation fut de proposer une formule concise pour cet opérateur (Tardivel *et al.*, 2020). En décembre 2018, lors d'un congrès, j'ai rencontré Ulrike Schneider, qui présentait des travaux sur une condition nécessaire et suffisante pour l'unicité de l'estimateur LASSO. D'après cette condition, il existe un vecteur réponse pour lequel la solution du problème LASSO n'est pas unique si et seulement si l'espace vectoriel engendré par les lignes de la matrice de régression coupe une face d'un hypercube dont la dimension est strictement inférieure à la dimension du noyau de cette matrice. Nous avons généralisé cette propriété à des estimateurs dont le terme de pénalité est une norme polyédrique (Schneider et Tardivel, 2022). En particulier, il existe un vecteur réponse pour lequel la solution du problème SLOPE n'est pas unique si et seulement si l'espace vectoriel engendré par les lignes de la matrice de régression coupe une face d'un permutoèdre signé dont la dimension est strictement inférieure à la dimension du noyau. Cette approche très géométrique de l'unicité nous a amené à dégager la notion de schéma du SLOPE ; en effet, il y a une bijection entre l'ensemble des schémas du SLOPE et les faces du permutoèdre signé obtenues via le sous-différentiel de la norme ℓ_1 ordonnée.

Mes dernières années en tant que maître de conférences à Dijon

J'ai été recruté à l'université de Bourgogne en tant que maître de conférences en septembre 2020. Depuis lors, j'ai collaboré avec un collègue, Xavier Dupuis, sur l'estimateur SLOPE (Dupuis et Tardivel, 2022, 2024). J'ai également co-encadré le doctorat de Tomasz Skalski, intitulé « Aspects géométriques et combinatoires des modèles statistiques » (Skalski, 2023), entre juin 2021 et octobre 2023. Cette thèse m'a donné l'occasion de revisiter la condition d'irreprésentabilité, car nous avons formulé une hypothèse similaire pour l'estimateur SLOPE. Initialement sceptique à l'idée de travailler sur une condition garantissant la récupération du schéma

1. Cette position correspond approximativement à un poste de maître de conférences avec une charge annuelle d'environ 180 heures d'enseignement.

des coefficients de régression par l'estimateur SLOPE en raison de mon expérience avec le LASSO lors de mon doctorat, j'ai finalement pris beaucoup de plaisir à redécouvrir ce sujet d'un point de vue mathématique. En effet, la condition d'irreprésentabilité peut être formulée de manière moins technique qu'un calcul opaque à vérifier, et Tomasz a trouvé une interprétation géométrique élégante de cette condition. Il est important de souligner que les résultats obtenus par Tomasz sont principalement théoriques ; la condition d'irreprésentabilité de l'estimateur SLOPE est très forte, et même lorsqu'elle est satisfaite, le calibrage des paramètres de régularisation de cet estimateur pour récupérer le schéma des coefficients de régression est techniquement très difficile. Cependant, les concepts théoriques développés par Tomasz au cours de sa thèse se sont révélés extrêmement utiles pour mon dernier article traitant du chemin des solutions de l'estimateur SLOPE (Dupuis et Tardivel, 2024).

Choix thématique et organisation du manuscrit

Ce mémoire d'habilitation se concentre sur l'établissement des propriétés d'appariement et de parcimonie de l'estimateur SLOPE. Ainsi, certains de mes travaux qui ne sont pas directement liés à cette thématique ne sont pas mentionnés. Mon objectif fut de rédiger un document de synthèse unifiant certains résultats connus pour l'estimateur SLOPE, accompagné d'une liste de mes articles et pré-publications énumérés ci-dessous :

- Patrick Tardivel, Rémi Servien, et Didier Concordet : Simple expressions of the lasso and slope estimators in low-dimension. *Statistics*, 2020 (Tardivel *et al.*, 2020).
- Ulrike Schneider et Patrick Tardivel : The geometry of uniqueness, sparsity, and clustering in penalized estimation. *Journal of Machine Learning Research*, 2022 (Schneider et Tardivel, 2022).
- Xavier Dupuis et Patrick Tardivel : Proximal operator for the sorted ℓ_1 norm : Application to testing procedures based on slope. *Journal of Statistical Planning and Inference*, 2022 (Dupuis et Tardivel, 2022).
- Xavier Dupuis et Patrick Tardivel : The solution path of slope. À paraître dans *Artificial Intelligence and Statistics*, 2024 (Dupuis et Tardivel, 2024).
- Małgorzata Bogdan, Xavier Dupuis, Piotr Graczyk, Bartosz Kolodziejek, Tomasz Skalski, Patrick Tardivel, et Maciej Wilczyński : Pattern recovery by slope. Pré-publication (Bogdan *et al.*, 2022).
- Piotr Graczyk, Ulrike Schneider, Tomasz Skalski, et Patrick Tardivel : Pattern recovery in penalized and thresholded estimation and its geometry. Pré-publication (Graczyk *et al.*, 2023).

Idéalement, je souhaiterais que ce manuscrit devienne un document de référence pour les doctorants et chercheurs s'initiant aux propriétés théoriques de l'estimateur SLOPE. Ainsi, tous les théorèmes et propositions énoncés dans ce mémoire sont démontrés ; de plus, j'ai intégralement revisité et simplifié les preuves de ces résultats. J'espère que ces efforts d'écriture rendront la lecture de ce mémoire plus accessible et agréable.

Choix linguistique

Souhaitant profiter de ma soutenance d'habilitation pour développer mon réseau en France j'ai rédigé ce texte en français afin que ce document soit agréable à lire pour les lectrices et lecteurs francophones ou francophiles. Par la suite, je prévois de traduire ce document en utilisant des logiciels de traduction automatique, afin de permettre à mes collègues non-francophones de bénéficier de son contenu.

Organisation des chapitres

Le premier chapitre présente les estimateurs OSCAR « Octagonal Shrinkage and Clustering Algorithm for Regression », PACS « Pairwise Absolute Clustering and Sparsity » et SLOPE, qui favorisent la parcimonie et l'appariement : ces estimateurs peuvent avoir des composantes nulles ou égales en valeur absolue. Plus spécifiquement, l'objet d'étude de cette habilitation est l'estimateur SLOPE, défini comme une solution du

problème d'optimisation convexe

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \sum_{i=1}^p \lambda_i |b|_{\downarrow i} \right\}, \quad (1)$$

où $X \in \mathbb{R}^{n \times p}$ (appelée matrice de régression en référence au modèle de régression linéaire), $y \in \mathbb{R}^p$, $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ avec $\lambda_1 > 0$ sont les paramètres de pénalité et $|b|_{\downarrow 1} \geq \dots \geq |b|_{\downarrow p}$ sont les composantes de b ordonnées par valeur absolue décroissante. Par ailleurs, la notion de schéma du SLOPE, notée $\text{schm}(b) \in \mathbb{Z}^p$ avec $b \in \mathbb{R}^p$, introduite dans ce chapitre est cruciale pour l'étude des propriétés de parcimonie et d'appariement de cet estimateur. Le principal résultat de ce chapitre est la bijection entre les faces du permutoèdre signé (c'est-à-dire la boule unité de la norme ℓ_1 ordonnée duale) et l'ensemble des schémas du SLOPE.

Le deuxième chapitre traite de la récupération du schéma des coefficients de régression, noté $\text{schm}(\beta)$, par l'estimateur SLOPE. On introduit la condition d'accessibilité, qui est nécessaire et suffisante pour l'existence d'un vecteur y pour lequel l'estimateur SLOPE et les coefficients de régression β ont le même schéma. Le principal résultat du chapitre est la caractérisation des vecteurs y pour lesquels le schéma de l'estimateur SLOPE coïncide avec un schéma donné, à l'aide des conditions du schéma et du sous-différentiel. Par ailleurs, on introduit la condition d'irreprésentabilité, qui est nécessaire pour la récupération du schéma $\text{schm}(\beta)$ par l'estimateur SLOPE avec une probabilité supérieure à un demi lorsque le vecteur y est aléatoire. Ce chapitre, qui correspond en grande partie aux travaux du doctorat de Tomasz (Skalski, 2023, chapitre 4), est un peu plus technique et théorique que les autres parties de cette habilitation. Néanmoins, les techniques développées dans cette partie seront très utiles pour le dernier chapitre.

Le troisième chapitre traite du chemin des solutions de l'estimateur SLOPE. Ce chemin est affine par morceaux, et les intervalles sont caractérisés par les conditions du schéma et du sous-différentiel introduites au chapitre précédent. Un algorithme pour calculer ce chemin est décrit dans cette partie et est illustré sur des données réelles. Par ailleurs, on montre que le calcul du chemin permet de sélectionner le paramètre de régularisation en minimisant la somme des carrés résiduels sur un échantillon de validation.

La quatrième partie est une annexe regroupant certains résultats démontrés pour des estimateurs pénalisés plus généraux que l'estimateur SLOPE. Pour ces estimateurs, la norme ℓ_1 ordonnée, qui est le terme de pénalité du SLOPE, est substituée par une fonction convexe s'exprimant comme le maximum d'une famille finie de formes linéaires.

Principales notions et notations

Les notions suivantes sont étudiées en détail dans les ouvrages de Hiriart-Urruty et Lemaréchal (2004) (pour la norme duale et le sous-différentiel) et de Ben-Israel et Greville (2003) (pour la pseudo inverse). Quelques rappels utiles pour ce manuscrit sont donnés ci-dessous :

Norme duale : Soit $\|\cdot\|$ une norme sur \mathbb{R}^p ; la norme duale notée $\|\cdot\|^*$ est définie par

$$\|x\|^* = \sup\{s^\top x \mid \|s\| \leq 1\}.$$

En particulier, la norme ℓ_1 ordonnée duale jouera un rôle crucial dans ce manuscrit.

Sous-différentiel : Soit $f : \mathbb{R}^p \rightarrow \mathbb{R}$ une fonction convexe. Le sous-différentiel de f au point $x \in \mathbb{R}^p$ est l'ensemble

$$\partial f(x) = \{s \in \mathbb{R}^p \mid f(y) \geq f(x) + s^\top (y - x) \quad \forall y \in \mathbb{R}^p\}.$$

Comme la norme ℓ_1 ordonnée est polyédrique son sous-différentiel pourra être déduit des expressions suivantes :

— Lorsque $f = \|\cdot\|$ est une norme alors

$$\partial\|\cdot\|(x) = \{s \in \mathbb{R}^p \mid \|s\|^* \leq 1 \text{ et } s^\top x \leq \|x\|\}.$$

— Lorsque f est le maximum d'une famille finie de formes linéaires : pour $x \in \mathbb{R}^p$, $f(x) = \max\{u_1^\top x, \dots, u_l^\top x\}$ avec $u_1, \dots, u_l \in \mathbb{R}^p$ alors

$$\partial f(x) = \text{conv}\{u_i \mid i \in \{1, \dots, l\}, u_i^\top x = f(x)\}.$$

Pseudo inverse : Soit $A \in \mathbb{R}^{n \times p}$ et $r = \text{rang}(A)$. La décomposition en valeurs singulières de A est $A = U\Sigma V^\top$ où

— $U \in \mathbb{R}^{n \times r}$ est une matrice ayant des colonnes orthogonales : $U^\top U = I_r$.

— $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ est une matrice dont les coefficients diagonaux $\sigma_1 \geq \dots \geq \sigma_r > 0$ sont les racines carrées des valeurs propres strictement positives de $A^\top A$ (ou de façon équivalente de AA^\top).

— $V \in \mathbb{R}^{p \times r}$ est une matrice ayant des colonnes orthogonales : $V^\top V = I_r$.

La pseudo inverse est définie par $A^+ = V\Sigma^{-1}U^\top$. On vérifie que la transposition et la pseudo inverse commute : $(A^+)^\top = (A^\top)^+$ et on notera $A^{\top+}$ cette matrice. Par ailleurs, $AA^\top = A^{\top+}A^\top$ est la projection sur l'espace vectoriel $\text{im}(A)$.

Principales notations

Précisons que l'ensemble \mathbb{R}^p est identifié à l'ensemble des matrices colonnes $\mathbb{R}^{p \times 1}$. Ainsi, $x = (x_1, \dots, x_p) \in$

\mathbb{R}^p est toujours représenté, matriciellement, par la matrice colonne $\begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$.

Notations liées aux paramètres de pénalité

— La notation \mathbb{R}^{p+} représente l'ensemble des vecteurs de \mathbb{R}^p dont les composantes sont décroissantes et positives : $\mathbb{R}^{p+} = \{\lambda \in \mathbb{R}^p \mid \lambda_1 \geq \dots \geq \lambda_p \geq 0\}$. Pour le problème (1), on supposera toujours que le paramètre de pénalité λ est un élément de \mathbb{R}^{p+} .

— La notation \mathbb{R}^{p++} est l'intérieur de \mathbb{R}^{p+} : $\mathbb{R}^{p++} = \{\lambda \in \mathbb{R}^p \mid \lambda_1 > \dots > \lambda_p > 0\}$. Supposer que $\lambda \in \mathbb{R}^{p++}$ permet d'écartier certains estimateurs pénalisés, comme le LASSO, qui est solution du problème (1) avec $\lambda_1 = \dots = \lambda_p > 0$ (donc $\lambda \in \mathbb{R}^{p+} \setminus \mathbb{R}^{p++}$).

Notations liées à la norme ℓ_1 ordonnée

Soit $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ et $\lambda_1 > 0$ (i.e. $\lambda \in \mathbb{R}^{p+} \setminus \{0\}$).

— La norme ℓ_1 ordonnée est définie par

$$J_\lambda(b) = \sum_{i=1}^p \lambda_i |b|_{\downarrow i} \quad \forall b \in \mathbb{R}^p,$$

où $|b|_{\downarrow 1} \geq \dots \geq |b|_{\downarrow p}$ sont les composantes de b ordonnées par valeur absolue décroissante.

— La norme ℓ_1 ordonnée duale vérifie

$$J_\lambda^*(b) = \max \left\{ \frac{|b|_{\downarrow 1}}{\lambda_1}, \frac{|b|_{\downarrow 1} + |b|_{\downarrow 2}}{\lambda_1 + \lambda_2}, \dots, \frac{|b|_{\downarrow 1} + \dots + |b|_{\downarrow p}}{\lambda_1 + \dots + \lambda_p} \right\} \quad \forall b \in \mathbb{R}^p.$$

- Le permutoèdre signé, noté P_λ^\pm , est un polytope connu qui correspond à la boule unité de la norme ℓ_1 ordonnée duale

$$P_\lambda^\pm = \text{conv}\{(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon_1, \dots, \epsilon_p \in \{-1, 1\}, \pi \in S_p\},$$

où conv et S_p représentent respectivement l'enveloppe convexe et l'ensemble des permutations sur $\{1, \dots, p\}$.

- Le permutoèdre P_λ est une face du permutoèdre signé. Ce polytope est définie par

$$P_\lambda = \text{conv}\{(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}) \mid \pi \in S_p\}.$$

Notations liées aux propriétés de parcimonie et d'appariement

- Le schéma du SLOPE de $b \in \mathbb{R}^p$, $\text{schem}(b) \in \mathbb{Z}^p$, est défini par

$$\text{schem}(b)_i = \text{signe}(b_i) \text{rang}(|b|)_i, \quad \forall i \in \{1, \dots, p\}$$

où $\text{rang}(|b|)_i \in \{0, 1, \dots, k\}$ avec k est le nombre de valeurs distinctes non nulles dans $\{|b_1|, \dots, |b_p|\}$, $\text{rang}(|b|)_i = 0$ si $b_i = 0$, $\text{rang}(|b|)_i > 0$ si $|b_i| > 0$ et $\text{rang}(|b|)_i < \text{rang}(|b|)_j$ si $|b_i| < |b_j|$. La notion de schéma du SLOPE permet de décrire les propriétés de parcimonie et d'appariement de cet estimateur.

- Les schémas du SLOPE, $\mathcal{P}_p^{\text{slope}} = \text{schem}(\mathbb{R}^p)$, est un ensemble fini de représentants canoniques pour la relation d'équivalence « avoir le même sous-différentiel pour la norme ℓ_1 ordonnée² ».
- Soit $m \in \mathcal{P}_p^{\text{slope}} \setminus \{0\}$ et $k = \|m\|_\infty \geq 1$. La matrice du schéma $U_m \in \mathbb{R}^{p \times k}$ est définie par

$$(U_m)_{ij} = \text{signe}(m_i) \mathbf{1}_{(|m_i|=k+1-j)} \quad \forall i \in \{1, \dots, p\} \forall j \in \{1, \dots, k\}.$$

L'ensemble $U_m \mathbb{R}^{p++}$ est la classe d'équivalence du schéma m pour la relation « avoir le même sous-différentiel pour la norme ℓ_1 ordonnée ».

- Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^p$, $m \in \mathcal{P}_p^{\text{slope}} \setminus \{0\}$ et $k = \|m\|_\infty \geq 1$. La matrice d'appariement de X est $\tilde{X}_m = XU_m \in \mathbb{R}^{n \times k}$; le paramètre d'appariement de λ est $\tilde{\lambda}_m = U_{|m|_\downarrow}^\top \lambda \in \mathbb{R}^k$. Par exemple, si $\hat{\beta}$ est un estimateur SLOPE tel que $\text{schem}(\hat{\beta}) = m$ alors la valeur ajustée $X\hat{\beta}$ vérifie $\tilde{X}_m s$, où $s \in \mathbb{R}^{k++}$ représente les composantes distinctes non-nulles de $\hat{\beta}$ ordonnées par valeur absolue décroissante, enfin $J_\lambda(\hat{\beta}) = \tilde{\lambda}_m^\top s$.

2. Sous réserve que le paramètre de pénalité λ de la norme ℓ_1 ordonnée J_λ soit un élément quelconque de \mathbb{R}^{p++} .

Chapitre 1

Estimateurs favorisant la parcimonie et l'appariement

1.1 Introduction

Soient $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^p$, $\lambda_1 \geq 0$ et $\lambda_2 \geq 0$ (avec $(\lambda_1, \lambda_2) \neq (0, 0)$). L'estimateur OSCAR (acronyme signifiant « Octagonal Shrinkage and Clustering Algorithm for Regression ») (Bondell et Reich, 2008) est défini comme une solution du problème d'optimisation suivant :

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \lambda_1 \sum_{i=1}^p |b_i| + \lambda_2 \sum_{1 \leq i < j \leq p} \frac{|b_i + b_j| + |b_i - b_j|}{2} \right\}. \quad (1.1)$$

L'estimateur LASSO (acronyme signifiant « Least Absolute Shrinkage and Selection Operator ») (Tibshirani, 1996), également connu sous le nom de poursuite de base débruitée dans la communauté étudiant l'acquisition comprimée (Chen et Donoho, 1994) est un cas particulier de l'estimateur OSCAR obtenu en prenant $\lambda_2 = 0$ dans la formule (1.1). De nombreux articles montrent que la norme ℓ_1 est un terme de pénalité favorisant la parcimonie (voir par exemple les livres (Giraud, 2021; Hastie *et al.*, 2015) et les références citées à l'intérieur). Ainsi, la spécificité de la méthode OSCAR réside dans le terme d'appariement $\frac{1}{2}(|b_i + b_j| + |b_i - b_j|) = \max\{|b_i|, |b_j|\}$, qui promeut l'égalité des composantes en valeur absolue pour une solution du problème (1.1). Les propriétés d'appariement (composantes égales en valeur absolue) et de parcimonie (composantes nulles) peuvent être illustrées de manière intuitive à l'aide de la Figure 1.1, qui est classique dans la littérature.

Une autre façon, plus analytique, de se convaincre des propriétés de parcimonie et d'appariement de l'estimateur OSCAR est d'observer que le terme de pénalité

$$\lambda_1 \sum_{i=1}^p |b_i| + \lambda_2 \sum_{1 \leq i < j \leq p} \max\{|b_i|, |b_j|\}$$

n'est pas différentiable en un point b admettant au moins une composante nulle ou au moins deux composantes égales en valeur absolue. Remarquons également que ce terme de pénalité peut se réécrire comme une somme pondérée et ordonnée de la façon suivante :

$$\lambda_1 \sum_{i=1}^p |b_i| + \lambda_2 \sum_{1 \leq i < j \leq p} \max\{|b_i|, |b_j|\} = \sum_{i=1}^p (\lambda_1 + \lambda_2(p-i)) |b|_{\downarrow i},$$

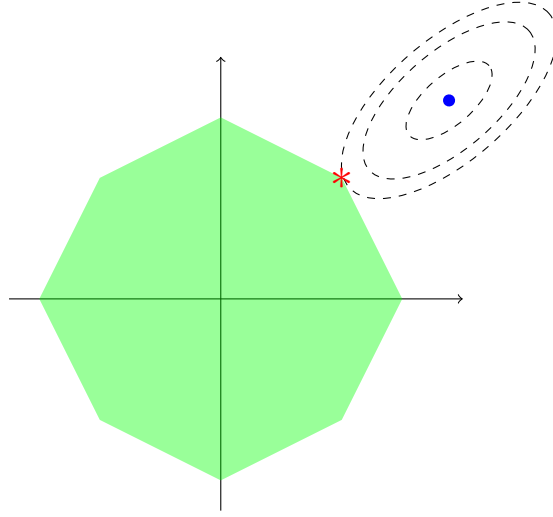


FIGURE 1.1 – Le polytope vert représente une boule pour la norme $b \in \mathbb{R}^2 \mapsto \lambda_1(|b_1| + |b_2|) + \lambda_2 \max\{|b_1|, |b_2|\}$, les ellipses représentent les lignes de niveau de la somme des carrés résiduels dont le centre, en bleu, représente l'estimateur des moindres carrés. Comme les sommets du polytope ont des composantes nulles ou égales en valeur absolue, cette figure illustre intuitivement que l'estimateur OSCAR est parcimonieux et possède des composantes appariées.

où $|b|_{\downarrow 1} \geq \dots \geq |b|_{\downarrow p} \geq 0$ sont les composantes de b ordonnées par valeur absolue décroissante. Ainsi, le terme de pénalité de l'estimateur OSCAR est un cas particulier de norme ℓ_1 ordonnée qui est une norme polyédrique définie à la Proposition 1.1.

Proposition 1.1 (Norme ℓ_1 ordonnée). *Soit $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ et $\lambda_1 > 0$. La fonction $J_\lambda(b) = \sum_{i=1}^p \lambda_i |b|_{\downarrow i}$, pour $b \in \mathbb{R}^p$ est une norme (dont la boule unité est un polyèdre) appelée norme ℓ_1 ordonnée, pouvant se reformuler de la façon suivante :*

$$J_\lambda(b) = \max \left\{ \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i \mid \pi \in S_p, \epsilon \in \{-1, 1\}^p \right\}.$$

L'estimateur SLOPE (Bogdan *et al.*, 2015; Zeng et Figueiredo, 2014) (acronyme signifiant « Sorted L One Penalized Estimation ») doit son nom à la norme ℓ_1 ordonnée qui apparaît comme le terme de pénalité du problème d'optimisation suivant

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + J_\lambda(b) \right\}. \quad (1.2)$$

Remarquons que par définition de la norme ℓ_1 ordonnée, l'estimateur OSCAR est un cas particulier de l'estimateur SLOPE dont le paramètre de pénalité λ est une suite arithmétique. La motivation première à l'usage d'une norme ℓ_1 ordonnée pour la méthode SLOPE est de pénaliser plus fortement les composantes larges de cet estimateur. La pénalité ℓ_1 ordonnée fait écho aux procédures de test multiples séquentielles seuillant plus fortement les grandes composantes de l'estimateur des moindres carrés. En particulier, un choix typique de paramètre de pénalité pour l'estimateur SLOPE est $\lambda_i^{\text{BH}} = \sigma \Phi^{-1}(1 - \frac{q_i}{2p})$, $i = 1, \dots, p$, où Φ est la fonction de répartition d'une loi $\mathcal{N}(0, 1)$ et $q \in]0, 1[$. Le vecteur λ^{BH} , appelé paramètre de pénalité de « Benjamini-Hochberg », fait référence à la très célèbre procédure de test multiples contrôlant le taux de faux positifs (Benjamini et Hochberg, 1995). Plus précisément, lorsque y est la réponse d'un modèle de régression linéaire gaussien où σ^2 est la variance des résidus, choisir λ^{BH} pour l'estimateur SLOPE permet de contrôler le taux de faux positifs au niveau q lorsque les colonnes de la matrice de régression X sont orthogonales (*i.e.* lorsque $X^\top X = I_p$) (Bogdan *et al.*, 2015). Ce dernier résultat est cité à titre indicatif car le contrôle du taux de faux positifs via l'estimateur

SLOPE est activement étudié dans la littérature et que cet estimateur est très largement connu grâce à cette propriété. Néanmoins, les notions de test multiples et de taux de faux positifs ne sont pas centrales dans ce mémoire.

Une autre généralisation de l'estimateur OSCAR qui se distingue de l'estimateur SLOPE est l'estimateur PACS (Sharma *et al.*, 2013) (acronyme signifiant « Pairwise Absolute Clustering and Sparsity »). Cet estimateur, qui est une version pondérée de l'estimateur OSCAR, est défini de la façon suivante :

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \sum_{i=1}^p w_i |b_i| + \sum_{1 \leq i < j \leq p} w_{ij}^+ |b_i + b_j| + \sum_{1 \leq i < j \leq p} w_{ij}^- |b_i - b_j| \right\}, \quad (1.3)$$

où $w_i > 0$ pour $i \in \{1, \dots, p\}$, $w_{ij}^+ > 0, w_{ij}^- > 0$ pour $1 \leq i < j \leq p$. Remarquons que lorsque $w_i = \lambda_1$ pour $i \in \{1, \dots, p\}$ et $w_{ij}^+ = w_{ij}^- = \lambda_2/2$ pour $1 \leq i < j \leq p$, on retrouve l'estimateur OSCAR défini à l'équation (1.1). Les poids peuvent être construit à partir d'un estimateur initial $\hat{\beta}$ (par exemple l'estimateur des moindres carrés) en posant $w_i = 1/|\hat{\beta}_i|$, $w_{ij}^+ = 1/|\hat{\beta}_i + \hat{\beta}_j|$ et $w_{ij}^- = 1/|\hat{\beta}_i - \hat{\beta}_j|$ de façon à renforcer les propriétés de parcimonie et d'appariement de l'estimateur PACS (Sharma *et al.*, 2013).

Dans la lignée des travaux de Figueiredo et Nowak (2016); Zeng et Figueiredo (2014), ce mémoire traite des propriétés de parcimonie et d'appariement de l'estimateur SLOPE. La notion de schéma du SLOPE, détaillée à la section suivante, sera cruciale pour cette étude.

1.2 Schéma du SLOPE

Une première étude sur la parcimonie et l'appariement de la méthode PACS a été réalisée dans l'article de Sharma *et al.* (2013). Lorsque β est un paramètre inconnu, ce travail établit des conditions théoriques permettant de récupérer, via la méthode PACS, les ensembles suivants :

$$\begin{cases} \{i \in \{1, \dots, p\} \mid \beta_i = 0\} & \text{composantes nulles de } \beta \\ \{1 \leq i < j \leq p \mid \beta_i - \beta_j = 0\} & \text{composantes de } \beta \text{ égales} \\ \{1 \leq i < j \leq p \mid \beta_i + \beta_j = 0\} & \text{composantes de } \beta \text{ égales en valeur absolue et de signes opposés} \end{cases} \quad (1.4)$$

La notion de schéma du SLOPE, introduite dans l'article de Schneider et Tardivel (2022), permet d'étudier les propriétés de parcimonie et d'appariement et fournit une description plus précise que les ensembles décrits à l'équation (1.4), spécifiant notamment les signes des composantes et la hiérarchie des groupes d'appariement (groupes de composantes égales en valeur absolue).

Définition 1.1. Soit $b \in \mathbb{R}^p$. Le schéma du SLOPE de b , noté $\text{schm}(b) \in \mathbb{Z}^p$, est défini par

$$\text{schm}(b)_i = \text{signe}(b_i) \text{rang}(|b|)_i, \quad \forall i \in \{1, \dots, p\}$$

où $\text{rang}(|b|)_i \in \{0, 1, \dots, k\}$ avec k le nombre de valeurs distinctes non nulles dans $\{|b_1|, \dots, |b_p|\}$, $\text{rang}(|b|)_i = 0$ si $b_i = 0$, $\text{rang}(|b|)_i > 0$ si $|b_i| > 0$ et $\text{rang}(|b|)_i < \text{rang}(|b|)_j$ si $|b_i| < |b_j|$.

Par exemple, pour $a = (4.2, -1.3, 0, 1.3, -4.2)$, on a $\text{schm}(a) = (2, -1, 0, 1, 2)$ et pour $b = (1.5, -2.8, 1.5, 2.8)$, on a $\text{schm}(b) = (1, -2, 1, 2)$. Analytiquement, nous montrerons que les schémas du SLOPE sont des représentants canoniques pour la relation d'équivalence « avoir le même sous-différentiel pour la norme ℓ_1 ordonnée ». Par ailleurs, géométriquement, nous prouverons que les schémas du SLOPE réalisent une bijection avec les faces de la boule unité de la norme ℓ_1 ordonnée duale. Dans un premier temps, nous introduisons la norme ℓ_1 ordonnée duale à la Proposition 1.2.

Proposition 1.2 (Norme ℓ_1 ordonnée duale). *Soit $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ et $\lambda_1 > 0$. La norme ℓ_1 ordonnée duale définie par $J_\lambda^*(b) = \max\{v^\top b \mid J_\lambda(v) \leq 1\}$, pour $b \in \mathbb{R}^p$, a une expression explicite :*

$$J_\lambda^*(b) = \max \left\{ \frac{|b|_{\downarrow 1}}{\lambda_1}, \frac{|b|_{\downarrow 1} + |b|_{\downarrow 2}}{\lambda_1 + \lambda_2}, \dots, \frac{|b|_{\downarrow 1} + \dots + |b|_{\downarrow p}}{\lambda_1 + \dots + \lambda_p} \right\}.$$

La boule unité de la norme ℓ_1 ordonnée duale est un polytope bien connu en géométrie appelé permutoèdre signé. Ce polytope peut se réécrire comme l'enveloppe convexe d'une famille de points, en particulier, on a l'équivalence suivante :

$$J_\lambda^*(b) \leq 1 \Leftrightarrow b \in \text{conv} \{(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon \in \{-1, 1\}^p, \pi \in S_p\}.$$

Dans la suite de ce manuscrit, pour $\lambda = (\lambda_1, \dots, \lambda_p)$ avec $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ et $\lambda_1 > 0$, on notera P_λ^\pm la boule unité de J_λ^* , c'est-à-dire le permutoèdre signé.

Dans le cas particulier où X est la matrice identité, $\widehat{\beta}$, l'unique solution du problème (1.2), est la différence entre y et son projeté sur la boule unité de J_λ^* (voir, par exemple, la Proposition 11 de l'article Schneider et Tardivel (2022) ou le Corollaire 2.3.1 de l'article Skalski *et al.* (2022)). En particulier, lorsque $X = I_2$, la Figure 1.2 illustre que la notion de schéma du SLOPE est adaptée pour décrire les propriétés de parcimonie et d'appariement de cet estimateur.

Sur la Figure 1.2 on remarque que chaque région colorée (en rouge, vert ou bleu) est associée à une face du permutoèdre signé (le polytope rouge, une arête ou un sommet) et à un schéma du SLOPE de \mathbb{R}^2 ($\{(0, 0), \pm(1, 0), \pm(0, 1), \pm(1, 1), \pm(1, -1), \pm(1, 2), \pm(1, -2), \pm(2, 1), \pm(2, -1)\}$). Le Théorème 1.1 confirme cette remarque et prouve que les schémas du SLOPE sont en bijection avec les faces du permutoèdre signé; cette bijection se réalise via le sous-différentiel de la norme ℓ_1 ordonnée. Notons que pour une norme $\|\cdot\|$ sur \mathbb{R}^p son sous-différentiel satisfait la formule suivante (voir par exemple l'ouvrage de Hiriart-Urruty et Lemaréchal (2004) page 180) :

$$\partial\|\cdot\|(b) = \{v \in \mathbb{R}^p \mid \|v\|^* \leq 1 \text{ et } v^\top b = \|b\|\}, \quad b \in \mathbb{R}^p,$$

où $\|\cdot\|^*$ représente la norme duale. Afin de formuler le Théorème 1.1 on introduit la notation $\mathcal{P}_p^{\text{slope}} = \text{schem}(\mathbb{R}^p)$ qui représente l'ensemble des schémas du SLOPE dans \mathbb{R}^p .

Théorème 1.1. *Soit $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 > \dots > \lambda_p > 0$.*

1. *Soit $a, b \in \mathbb{R}^p$ alors $\partial J_\lambda(a) = \partial J_\lambda(b)$ si et seulement si $\text{schem}(a) = \text{schem}(b)$.*
2. *L'application $m \in \mathcal{P}_p^{\text{slope}} \mapsto \partial J_\lambda(m)$ est une bijection entre les schémas du SLOPE et les faces du permutoèdre signé.*

Notons que l'hypothèse $\lambda_1 > \dots > \lambda_p > 0$ est importante pour établir le Théorème 1.1. Par exemple si $\lambda_1 = \dots = \lambda_p > 0$ alors le permutoèdre signé P_λ^\pm est un hypercube et clairement il n'y a pas de bijection entre $\mathcal{P}_p^{\text{slope}}$ et les faces de l'hypercube (les faces de l'hypercube sont en bijection avec $\{-1, 0, 1\}^p$). Notons qu'un résultat très récent de Godland et Kabluchko (2023) montrant que, lorsque $\lambda_1 > \dots > \lambda_p > 0$, les faces du permutoèdre signé sont en bijection avec la famille des sous-ensembles disjoints ordonnés et signés de $\{1, \dots, p\}$ confirme la deuxième assertion du Théorème 1.1. En effet, étant donné (\mathcal{B}, η) où $\mathcal{B} = (B_1, \dots, B_k)$ sont des ensembles disjoints ordonnés de $\{1, \dots, p\}$ et $\eta : B_1 \cup \dots \cup B_k \mapsto \{-1, 1\}$ on construit un schéma du SLOPE m en posant

$$m_i = \begin{cases} \eta(i)j & \text{si } i \in B_j \\ 0 & \text{si } i \notin B_k \cup \dots \cup B_1 \end{cases} \quad \forall i \in \{1, \dots, p\}.$$

Inversement étant donné un schéma du SLOPE on construit aisément des sous-ensembles disjoints ordonnés et signés de $\{1, \dots, p\}$. Une illustration de la première assertion du Théorème 1.1 est donnée à la Figure 1.3.

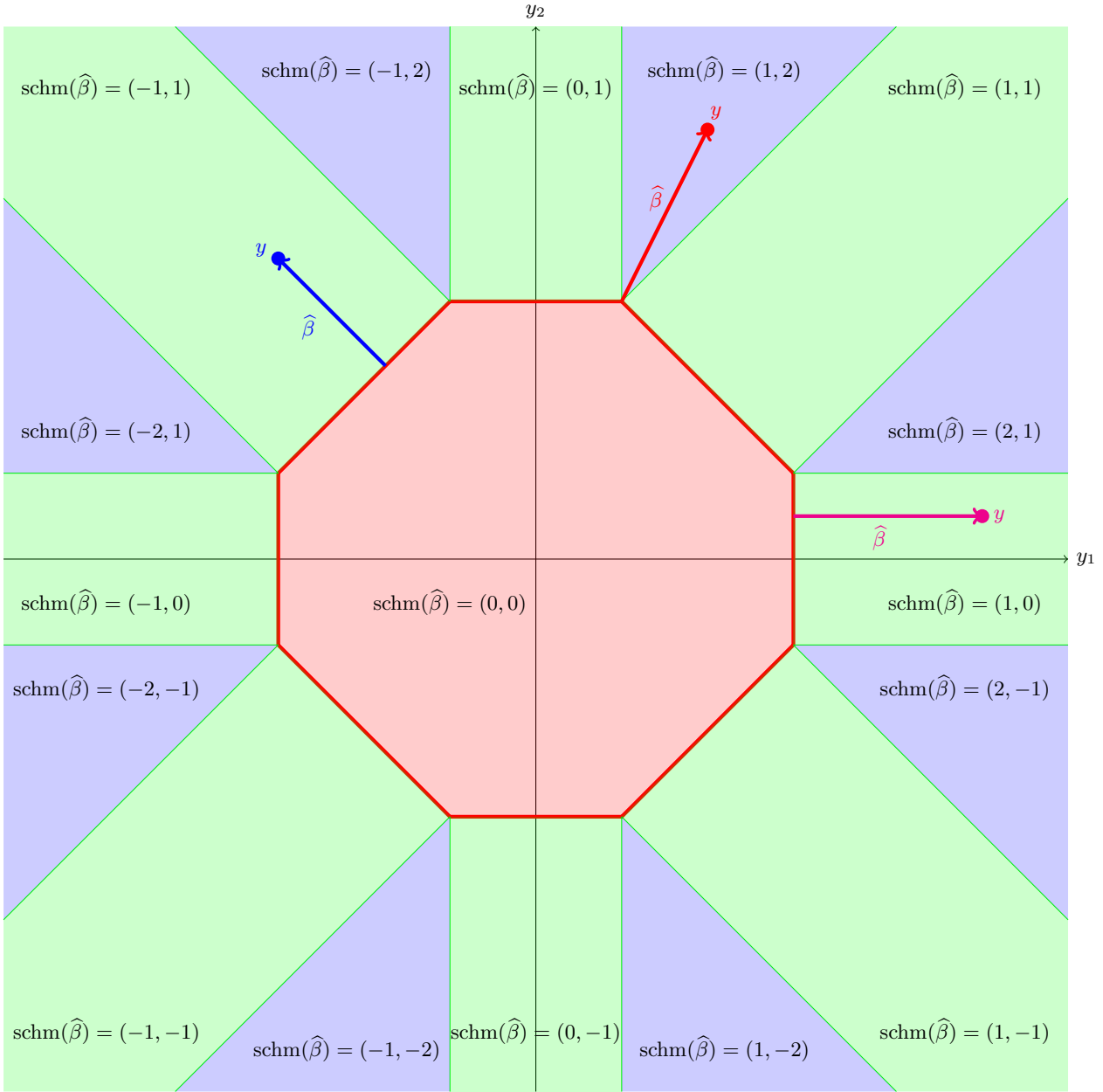


FIGURE 1.2 – Dans le cas particulier où $J_\lambda(b) = \lambda_1|b|_{\downarrow 1} + \lambda_2|b|_{\downarrow 2}$ avec $\lambda_1 > \lambda_2 > 0$ et $X = I_2$ cette image illustre que $\hat{\beta}$, l'unique solution du problème (1.2), est la différence entre y et son projeté sur le permutoèdre signé P_λ^\pm en rouge (la boule unité de J_λ^*). En particulier lorsque y est le point rose (resp. point rouge ou point bleu) $\hat{\beta}$ est représenté par le vecteur rose (resp. vecteur rouge ou vecteur bleu) et on peut observer que $\text{schm}(\hat{\beta}) = (1, 0)$ (resp. $\text{schm}(\hat{\beta}) = (1, 2)$ ou $\text{schm}(\hat{\beta}) = (-1, 1)$). De plus cette figure fournit $\text{schm}(\hat{\beta}) \in \{(0, 0), \pm(1, 0), \pm(0, 1), \pm(1, 1), \pm(1, -1), \pm(1, 2), \pm(1, -2), \pm(2, 1), \pm(2, -1)\}$ selon la localisation de $y \in \mathbb{R}^2$.

1.3 Annexe : preuves

L'inégalité de réarrangement

L'inégalité de réarrangement rappelée au Lemme 1.1 (voir par exemple l'ouvrage de Hardy *et al.* (1952)) sera très utile pour établir les preuves de ce chapitre.

Lemme 1.1. Soient $a \in \mathbb{R}^p$ et $b \in \mathbb{R}^p$ alors pour toute permutation $\pi \in S_p$ on a $\sum_{i=1}^p a_i b_{\pi(i)} \leq \sum_{i=1}^p a_{\downarrow i} b_{\downarrow i}$.

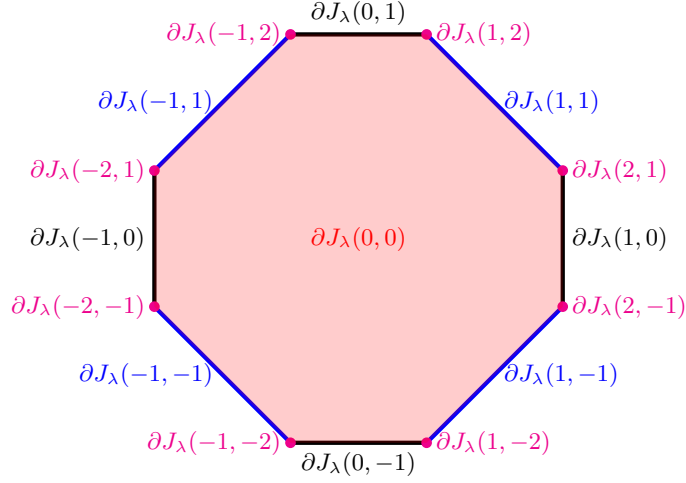


FIGURE 1.3 – Cette figure illustre la bijection $m \in \mathcal{P}_2^{\text{slope}} \mapsto \partial J_\lambda(m)$ entre les schéma du SLOPE de \mathbb{R}^2 ($\mathcal{P}_2^{\text{slope}} : \{(0,0), \pm(1,0), \pm(0,1), \pm(1,1), \pm(1,-1), \pm(1,2), \pm(1,-2), \pm(2,1), \pm(2,-1)\}$) et les faces du permutohédre signé P_λ^\pm avec $\lambda_1 > \lambda_2 > 0$.

Bien que classique, la preuve du Lemme 1.1 est donnée ci-dessous :

Démonstration. Dans un premier temps supposons que $a_1 \geq \dots \geq a_p$ et $b_1 \geq \dots \geq b_p$. On pose $A(\pi) = \sum_{i=1}^p a_i b_{\pi(i)}$, pour $\pi \in S_p$, et on note ψ une permutation de S_p pour laquelle la somme $A(\psi)$ est maximale et ayant un nombre maximal de points fixes. Nous allons montrer par l'absurde que ψ est l'identité. Si ψ n'est pas l'identité alors on pose $k = \min\{i \in \{1, \dots, p\} \mid \psi(i) \neq i\}$. Comme $\psi(1) = 1, \dots, \psi(k-1) = k-1$ on en déduit que $\psi(k) > k$ et qu'il existe $j > k$ tel que $\psi(j) = k$. On pose φ la permutation suivante obtenue en transposant les valeurs de ψ en k et j :

$$\forall i \in \{1, \dots, p\} \quad \varphi(i) = \begin{cases} \psi(i) & \text{si } i \notin \{k, j\} \\ k & \text{si } i = k \\ \psi(k) & \text{si } i = j \end{cases} .$$

Notons que φ a un point fixe de plus que ψ (le point k) et comme ces deux permutations sont partout égales sauf sur l'ensemble $\{k, j\}$ on en déduit l'inégalité suivante

$$A(\varphi) - A(\psi) = a_k b_k + a_j b_{\psi(k)} - a_k b_{\psi(k)} - a_j b_k = \underbrace{(a_k - a_j)}_{\geq 0 \text{ car } j > k} \times \underbrace{(b_k - b_{\psi(k)})}_{\geq 0 \text{ car } \psi(k) > k} \geq 0$$

Cette inégalité contredit que ψ est une permutation pour laquelle la somme $A(\psi)$ est maximale et ayant un nombre maximal de points fixes. Enfin, sans aucune hypothèse sur $a, b \in \mathbb{R}^p$, on pose ψ et φ des permutations de S_p pour lesquelles $a_{\psi(1)} \geq \dots \geq a_{\psi(p)}$ et $b_{\varphi(1)} \geq \dots \geq b_{\varphi(p)}$. La première partie de la preuve permet d'établir l'inégalité suivante pour toute permutation ϕ de S_p :

$$\sum_{i=1}^p a_{\downarrow i} b_{\downarrow i} = \sum_{i=1}^p a_{\psi(i)} b_{\varphi(i)} \geq \sum_{i=1}^p a_{\psi(i)} b_{\varphi \circ \phi(i)} \geq \sum_{i=1}^p a_i b_{\varphi \circ \phi \circ \psi^{-1}(i)} .$$

On obtient l'inégalité souhaité en prenant $\phi = \varphi^{-1} \circ \pi \circ \psi$. □

Preuves que J_λ et J_λ^* sont des normes

La preuve que J_λ est une norme polyédrique est inspirée de la preuve donnée dans (Bogdan *et al.*, 2015).

Démonstration de la Proposition 1.1. Clairement $J_\lambda(b) = 0$ si et seulement si $b = 0$ et pour tout $b \in \mathbb{R}^p$ et $t \in \mathbb{R}$ on a $J_\lambda(tb) = |t|J_\lambda(b)$. Pour achever la preuve il suffit de montrer que $J_\lambda(b) = \sum_{i=1}^p \lambda_i |b|_{\downarrow i}$ est une fonction convexe. D'après l'inégalité de réarrangement on a

$$J_\lambda(b) = \max \left\{ \sum_{i=1}^p \lambda_{\pi(i)} |b_i| \mid \pi \in S_p \right\} = \max \left\{ \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i \mid \pi \in S_p, \epsilon \in \{-1, 1\}^p \right\}.$$

Ainsi, comme J_λ est le maximum de fonctions linéaires, on en déduit que J_λ est convexe. \square

Démonstration de la Proposition 1.2. Clairement $J_\lambda^*(b) = 0$ si et seulement si $b = 0$ et pour tout $b \in \mathbb{R}^p$ et $t \in \mathbb{R}$ on a $J_\lambda^*(tb) = |t|J_\lambda^*(b)$. Montrons l'inégalité triangulaire. Soit π une permutation de $\{1, \dots, p\}$ telle que $|a_{\pi(1)} + b_{\pi(1)}| \geq \dots \geq |a_{\pi(p)} + b_{\pi(p)}|$ alors

$$\begin{aligned} J_\lambda^*(a+b) &= \max \left\{ \frac{\sum_{i=1}^j |a_{\pi(i)} + b_{\pi(i)}|}{\sum_{i=1}^j \lambda_i} \mid j \in \{1, \dots, p\} \right\}, \\ &\leq \max \left\{ \frac{\sum_{i=1}^j |a_{\pi(i)}|}{\sum_{i=1}^j \lambda_i} + \frac{\sum_{i=1}^j |b_{\pi(i)}|}{\sum_{i=1}^j \lambda_i} \mid j \in \{1, \dots, p\} \right\}, \\ &\leq \max \left\{ \frac{\sum_{i=1}^j |a|_{\downarrow i}}{\sum_{i=1}^j \lambda_i} + \frac{\sum_{i=1}^j |b|_{\downarrow i}}{\sum_{i=1}^j \lambda_i} \mid j \in \{1, \dots, p\} \right\}, \\ &\leq \underbrace{\max \left\{ \frac{\sum_{i=1}^j |a|_{\downarrow i}}{\sum_{i=1}^j \lambda_i} \mid j \in \{1, \dots, p\} \right\}}_{=J_\lambda^*(a)} + \underbrace{\max \left\{ \frac{\sum_{i=1}^j |b|_{\downarrow i}}{\sum_{i=1}^j \lambda_i} \mid j \in \{1, \dots, p\} \right\}}_{=J_\lambda^*(b)}. \end{aligned}$$

Ainsi J_λ^* est une norme. Montrons que J_λ^* est la norme ℓ_1 ordonnée duale. Soit $v \in \mathbb{R}^p$ tel que $J_\lambda^*(v) \leq 1$. D'après l'inégalité de réarrangement on a :

$$\sum_{i=1}^p v_i b_i \leq \sum_{i=1}^p |v_i| |b_i| \leq \sum_{i=1}^p |v|_{\downarrow i} |b|_{\downarrow i}$$

Par ailleurs, comme $|v|_{\downarrow 1} + \dots + |v|_{\downarrow j} \leq \lambda_1 + \dots + \lambda_j$ on en déduit l'inégalité suivante :

$$\begin{aligned} \sum_{i=1}^p |v|_{\downarrow i} |b|_{\downarrow i} &= \sum_{j=1}^{p-1} (|b|_{\downarrow j} - |b|_{\downarrow j+1}) (|v|_{\downarrow 1} + \dots + |v|_{\downarrow j}) + |b|_{\downarrow p} \sum_{i=1}^p |v|_{\downarrow i}, \\ &\leq \sum_{j=1}^{p-1} (|b|_{\downarrow j} - |b|_{\downarrow j+1}) (\lambda_1 + \dots + \lambda_j) + |b|_{\downarrow p} \sum_{i=1}^p \lambda_i = J_\lambda(b). \end{aligned}$$

Ainsi, $J_\lambda(b) \geq \max\{v^\top b \mid J_\lambda^*(v) \leq 1\}$. Enfin, en pour une permutation π de $\{1, \dots, p\}$ et $\epsilon \in \{-1, 1\}^p$ bien choisis on a $J_\lambda(b) = \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i$ donc $\bar{v}^\top b = J_\lambda(b)$ en prenant $\bar{v} = (\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)})$. Comme $J_\lambda^*(\bar{v}) \leq 1$ on en déduit que $J_\lambda(b) = \max\{v^\top b \mid J_\lambda^*(v) \leq 1\}$. Cette dernière égalité montre que J_λ est la norme duale de J_λ^* (ou, de façon équivalente, que J_λ^* est la norme duale de J_λ). \square

Preuve du Théorème 1.1

Preuve de l'assertion 1

Comme la norme ℓ_1 ordonnée est polyédrique, son sous-différentiel est donné par la formule suivante (voir par exemple le livre de Hiriart-Urruty et Lemaréchal (2004) page 180)

$$\partial J_\lambda(b) = \text{conv} \left\{ (\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon \in \{-1, 1\}^p, \pi \in S_p \text{ et } \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = J_\lambda(b) \right\}.$$

Le Lemme 1.2 donne une formulation plus concise du sous-différentiel de la norme ℓ_1 ordonnée permettant de vérifier que $\text{schr}(a) = \text{schr}(b)$ implique $\partial J_\lambda(a) = \partial J_\lambda(b)$.

Lemme 1.2. Soient $b \in \mathbb{R}^p$ et $\varphi \in S_p$ tel que $|b_{\varphi(1)}| \geq \dots \geq |b_{\varphi(p)}|$ et $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 > 0$ et $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. On définit le sous-groupe $H_{|b|}$ de S_p et le sous-ensemble E_b par

$$H_{|b|} = \{\pi \in S_p \mid (|b_{\pi(1)}|, \dots, |b_{\pi(p)}|) = (|b_1|, \dots, |b_p|)\} \text{ et } E_b = \{\epsilon \in \{-1, 1\}^p \mid \epsilon_i x_i = |x_i| \quad \forall i \in \{1, \dots, p\}\}.$$

Alors le sous-différentiel de J_λ au point b est donné par l'expression suivante :

$$\partial J_\lambda(b) = \text{conv}\{(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon \in E_b, \varphi \circ \pi \in H_{|b|}\}.$$

Cette formule est plus facile à analyser lorsque $b \in \mathbb{R}^p$ vérifie $b_1 \geq \dots \geq b_p \geq 0$ puisque φ est l'identité. Donnons quelques exemples :

— Lorsque $b_1 = \dots = b_p = 0$ alors $H_{|b|} = S_p$ et $E_b = \{-1, 1\}^p$ ainsi

$$\partial J_\lambda(b) = \text{conv}\{(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon \in \{-1, 1\}^p, \pi \in S_p\} = P_\lambda^\pm.$$

On retrouve le permutoèdre signé.

— Lorsque $b \in \mathbb{R}^p$ vérifie $b_1 = \dots = b_p > 0$ alors $H_{|b|} = S_p$ et E_b a un unique élément $(1, \dots, 1)$ ainsi

$$\partial J_\lambda(b) = \text{conv}\{(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}) \mid \pi \in S_p\}.$$

On obtient un polytope très connu également : le permutoèdre (noté P_λ).

— Lorsque $b \in \mathbb{R}^p$ vérifie $b_1 = \dots = b_k > b_{k+1} = \dots = b_p \geq 0$ alors $H_{|b|}$ est le sous-groupe des permutations : $\{\pi \in S_p \mid \pi(\{1, \dots, k\}) = \{1, \dots, k\} \text{ et } \pi(\{k+1, \dots, p\}) = \{k+1, \dots, p\}\}$. Par ailleurs, $E_b = \{(1, \dots, 1)\}$ si $b_p > 0$ ou $E_b = \{1\}^k \times \{-1, 1\}^{p-k}$ si $b_p = 0$ d'où

$$\partial J_\lambda(b) = \begin{cases} P_{\lambda_1, \dots, \lambda_k} \times P_{\lambda_{k+1}, \dots, \lambda_p} & \text{si } b_p > 0 \\ P_{\lambda_1, \dots, \lambda_k} \times P_{\lambda_{k+1}, \dots, \lambda_p}^\pm & \text{si } b_p = 0 \end{cases}.$$

Cette formule se généralise lorsque $b_1 \geq \dots \geq b_p \geq 0$; en effet $\partial J_\lambda(b)$ est un produit cartésien de permutoèdres avec éventuellement un permutoèdre signé lorsque $b_p = 0$ (Dupuis et Tardivel, 2022; Schneider et Tardivel, 2022).

Démonstration. Soient $\epsilon \in E_b$ et $\pi \in S_p$ tel que $\varphi \circ \pi \in H_{|b|}$. Montrons que $(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \in \partial J_\lambda(b)$. Comme $\epsilon \in E_b$ on en déduit les égalités suivantes :

$$\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = \sum_{i=1}^p \lambda_{\pi(i)} |b_i| = \sum_{i=1}^p \lambda_i |b_{\pi^{-1}(i)}| = \sum_{i=1}^p \lambda_i |b_{\pi^{-1} \circ \varphi^{-1} \circ \varphi(i)}|.$$

Comme $\pi^{-1} \circ \varphi^{-1} \in H_{|b|}$, car $H_{|b|}$ est un sous-groupe de S_p , on en déduit que pour tout i on a $|b_{\pi^{-1} \circ \varphi^{-1} \circ \varphi(i)}| = |b_{\varphi(i)}|$. Ainsi,

$$\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = \sum_{i=1}^p \lambda_i |b_{\varphi(i)}| = J_\lambda(b).$$

Inversement soit $\epsilon \in \{-1, 1\}^p$, $\pi \in S_p$ tels que $\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = J_\lambda(b)$. Pour tout $i \in \{1, \dots, p\}$ on a $\epsilon_i b_i = |b_i|$. En effet, s'il existe $i_0 \in \{1, \dots, p\}$ tel que $\epsilon_{i_0} b_{i_0} < -\epsilon_{i_0} b_{i_0} = |b_{i_0}|$ alors

$$\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i < -\epsilon_{i_0} \lambda_{\pi(i_0)} b_{i_0} + \sum_{i \neq i_0} \epsilon_i \lambda_{\pi(i)} b_i \leq J_\lambda(b),$$

ce qui contredit que $\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = J_\lambda(b)$. Ainsi, $J_\lambda(b) = \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = \sum_{i=1}^p \lambda_{\pi(i)} |b_i| = \sum_{i=1}^p \lambda_i |b_{\pi^{-1}(i)}|$. On en déduit, d'après l'inégalité de réarrangement, que $|b_{\pi^{-1}(1)}| \geq \dots \geq |b_{\pi^{-1}(p)}|$. Par conséquent

$$\forall i \in \{1, \dots, p\} \quad |b_{\pi^{-1}(i)}| = |b_{\varphi(i)}| \Leftrightarrow \forall i \in \{1, \dots, p\} \quad |b_i| = |b_{\varphi(\pi(i))}|.$$

Donc, $\varphi \circ \pi \in H_{|b|}$ ce qui termine la preuve. □

On remarque que la permutation $\varphi \in S_p$ qui ordonne les composantes de b par valeur absolue décroissante (*i.e.* telle que $|b_{\varphi(1)}| \geq \dots \geq |b_{\varphi(p)}|$) et le sous-groupe $H_{|b|}$ ne dépendent que de $\text{rang}(|b|)$. Par ailleurs, le sous-ensemble E_b ne dépend que de $\text{signe}(b)$. Ainsi, d'après le Lemme 1.2 si $\text{schm}(a) = \text{schm}(b)$ on a $\partial J_\lambda(a) = \partial J_\lambda(b)$. Au Lemme 1.4 on montre que si $\lambda_1 > \dots > \lambda_p > 0$ alors la réciproque est également vraie. Pour établir cette réciproque on utilisera le Lemme 1.3.

Lemme 1.3. *Soient $\lambda_1 > \dots > \lambda_p > 0$ et $b \in \mathbb{R}^p$. Si $\lambda \in \partial J_\lambda(b)$ alors $b_1 \geq \dots \geq b_p \geq 0$.*

Démonstration. Supposons que $b_{i_0} < 0$ pour certain $i_0 \in \{1, \dots, p\}$ alors

$$\sum_{i=1}^p \lambda_i b_i < -\lambda_{i_0} b_{i_0} + \sum_{\substack{i=1 \\ i \neq i_0}}^p \lambda_i b_i \leq J_\lambda(b).$$

Ce qui contredit que $\lambda \in \partial J_\lambda(b)$. Supposons que $b_{i_0} < b_{i_1}$ pour certains $1 \leq i_0 < i_1 \leq p$. Comme $\lambda_{i_0} > \lambda_{i_1}$ on a $\lambda_{i_0} b_{i_0} + \lambda_{i_1} b_{i_1} < \lambda_{i_0} b_{i_1} + \lambda_{i_1} b_{i_0}$ ainsi, en posant π la transposition (i_0, i_1) , on obtient l'inégalité :

$$\sum_{i=1}^p \lambda_i b_i < \sum_{i=1}^p \lambda_{\pi(i)} b_i \leq J_\lambda(b).$$

Ce qui contredit que $\lambda \in \partial J_\lambda(b)$. Par conséquent $b_1 \geq \dots \geq b_p \geq 0$. □

Lemme 1.4. *Soit $\lambda_1 > \dots > \lambda_p > 0$, $a \in \mathbb{R}^p$ et $b \in \mathbb{R}^p$. Si $\partial J_\lambda(a) = \partial J_\lambda(b)$ alors $\text{schm}(a) = \text{schm}(b)$.*

Démonstration. Dans un premier temps supposons que $a = |a|_\downarrow$, c'est-à-dire que $a_1 \geq a_2 \geq \dots \geq a_p \geq 0$. D'après le Lemme 1.2 on $\lambda \in \partial J_\lambda(a)$. Comme $\partial J_\lambda(a) = \partial J_\lambda(b)$ on en déduit, d'après le lemme 1.3, que $b_1 \geq \dots \geq b_p \geq 0$. Montrons que $\text{schm}(a)_p = \text{schm}(b)_p$, c'est-à-dire montrons que $a_p = b_p = 0$ ou $a_p > 0$ et $b_p > 0$. En effet si $a_p = 0$ et $b_p > 0$ alors

$$J_\lambda(a) = \sum_{i=1}^{p-1} \lambda_i a_i - \lambda_p a_p \text{ et } J_\lambda(b) > \sum_{i=1}^{p-1} \lambda_i b_i - \lambda_p b_p$$

ainsi $(\lambda_1, \dots, \lambda_{p-1}, -\lambda_p) \in \partial J_\lambda(a)$ et $(\lambda_1, \dots, \lambda_{p-1}, -\lambda_p) \notin \partial J_\lambda(b)$, ce qui contredit l'égalité des sous-différentiels (on montre de façon analogue que $a_p > 0$ et $b_p = 0$ mène à une contradiction). Pour prouver que $\text{schm}(a) =$

$\text{schm}(b)$, établissons que pour tout $i \in \{1, \dots, p-1\}$ on a $a_i = a_{i+1}$ et $b_i = b_{i+1}$ ou $a_i > a_{i+1}$ et $b_i > b_{i+1}$. En effet, s'il existe $i_0 \in \{1, \dots, p-1\}$ tel que $a_{i_0} = a_{i_0+1}$ et $b_{i_0} > b_{i_0+1}$ alors, en posant π la transposition $(i_0, i_0 + 1)$, on obtient

$$J_\lambda(a) = \sum_{i=1}^p \lambda_{\pi(i)} a_i \text{ et } J_\lambda(b) > \sum_{i=1}^p \lambda_{\pi(i)} b_i$$

ainsi $(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}) \in \partial J_\lambda(a)$ et $(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}) \notin \partial J_\lambda(b)$ ce qui contredit l'égalité des sous-différentiels (on montre de façon analogue que $a_{i_0} > a_{i_0+1}$ et $b_{i_0} = b_{i_0+1}$ mène à une contradiction).

Enfin, si $a \neq |a|_\downarrow$ alors choisissons une transformation orthogonale $\psi(x) = (\epsilon_1 x_{\pi(1)}, \dots, \epsilon_p x_{\pi(p)})$, pour $x \in \mathbb{R}^p$, où $\epsilon \in \{-1, 1\}^p$ et $\pi \in S_p$ sont choisis de telle sorte que $\psi(a) = |a|_\downarrow$. Comme la transformation ψ préserve les normes J_λ et J_λ^* on a l'équivalence suivante :

$$v \in \partial J_\lambda(x) \Leftrightarrow \underbrace{J_\lambda^*(v)}_{=J_\lambda^*(\psi(v))} \leq 1 \quad \text{et} \quad \underbrace{v^\top x}_{=\psi(v)^\top \psi(x)} = \underbrace{J_\lambda(x)}_{=J_\lambda(\psi(x))} \Leftrightarrow \psi(v) \in \partial J_\lambda(\psi(x)).$$

Ainsi, $\partial J_\lambda(a) = \partial J_\lambda(b)$ implique que $\partial J_\lambda(\psi(a)) = \partial J_\lambda(\psi(b))$. Comme $\psi(a) = |a|_\downarrow$, la première partie de la preuve montre que $\text{schm}(\psi(a)) = \text{schm}(\psi(b))$ où de façon équivalente $\psi(\text{schm}(a)) = \psi(\text{schm}(b))$ d'où $\text{schm}(a) = \text{schm}(b)$. \square

Preuve de l'assertion 2

La preuve de la seconde assertion du Théorème 1.1 est une conséquence du Lemme 1.5.

Lemme 1.5. *Soit $\lambda \in \mathbb{R}^p$ tel que $\lambda_1 > \dots > \lambda_p > 0$. Une face quelconque du permutoèdre signé P_λ^\pm peut s'exprimer comme le sous-différentiel $\partial J_\lambda(x)$ de la norme ℓ_1 ordonnée en un point $x \in \mathbb{R}^p$.*

D'après le Lemme 1.2, $\partial J_\lambda(x) = \partial J_\lambda(\text{schm}(x))$, ainsi on déduit du Lemme 1.5 que l'application $m \in \mathcal{P}_p^{\text{slope}} \mapsto \partial J_\lambda(m)$ est une surjection entre les schémas du SLOPE et les faces du permutoèdre signé. Enfin, cette application est injective puisque, d'après le Lemme 1.4, pour $m, \tilde{m} \in \mathcal{P}_p^{\text{slope}}$ avec $m \neq \tilde{m}$ on a $\partial J_\lambda(m) \neq \partial J_\lambda(\tilde{m})$.

Le Lemme 1.5 sera démontré dans un cadre plus générale, dans la dernière partie de ce manuscrit, en substituant la norme ℓ_1 ordonnée par le maximum d'une famille finie de formes linéaires.

Chapitre 2

Récupération du schéma par l'estimateur SLOPE

On considère le modèle de régression linéaire $Y = X\beta + \varepsilon$, où $X \in \mathbb{R}^{n \times p}$ est la matrice de régression, $\beta \in \mathbb{R}^p$ est le vecteur inconnu des coefficients de régression et $\varepsilon \in \mathbb{R}^n$ représente les résidus aléatoires. Le support de β ($\{i \in \{1, \dots, p\} \mid \beta_i \neq 0\}$) a une interprétation très simple : les variables explicatives (colonnes de X) pertinentes sont associées à des composantes non-nulles de β . L'estimateur LASSO étant parcimonieux, il est naturel d'utiliser cet estimateur pour identifier les coefficients de régression non-nuls ; en particulier, pléthore d'articles traitent de l'estimation du support de β via le LASSO (voir par exemple les ouvrages de Bühlmann et Van De Geer (2011); Hastie *et al.* (2015) et les références citées). Néanmoins, la littérature ne se résume pas qu'à la récupération du support de β , et d'autres structures sont également pertinentes à identifier (voir par exemple les articles de Vaïter *et al.* (2015, 2017)).

L'objectif de cette partie est d'étudier les propriétés théoriques garantissant la récupération du schéma $\text{schem}(\beta)$ par l'estimateur SLOPE. Notons qu'à l'instar du support, le schéma du SLOPE a une interprétation statistique. En effet, lorsque la matrice de régression est normalisée, les composantes de β égales en valeur absolue sont associées à des variables explicatives ayant le même impact sur la réponse Y du modèle de régression linéaire (Sharma *et al.*, 2013).

2.1 Notions liées au schéma du SLOPE

Supposons que la $i^{\text{ème}}$ composante de l'estimateur SLOPE soit nulle, alors la colonne X_i est non pertinente et peut être supprimée de la matrice de régression X . De même, supposons que les $i^{\text{ème}}$ et $j^{\text{ème}}$ composantes de l'estimateur SLOPE soient égales, alors les colonnes X_i et X_j peuvent être fusionnées. Plus généralement, la connaissance du schéma de l'estimateur SLOPE mène à construire une matrice d'appariement obtenue en modifiant la matrice de régression X via la prise en compte des composantes nulles et des groupes d'appariement non-nuls. La notion de matrice d'appariement associée à un schéma sera définie dans cette section. Au préalable, nous introduisons la notion de matrice du schéma.

Définition 2.1. Soit $m \in \mathcal{P}_p^{\text{slope}}$ un schéma du SLOPE non-nul et $k = \|m\|_\infty \geq 1$. La matrice du schéma $U_m \in \mathbb{R}^{p \times k}$ est définie par

$$(U_m)_{ij} = \text{signe}(m_i) \mathbf{1}_{(|m_i|=k+1-j)} \quad \forall i \in \{1, \dots, p\} \quad \forall j \in \{1, \dots, k\}.$$

Pour $k \geq 1$, on note $\mathbb{R}^{k+} = \{s \in \mathbb{R}^k \mid s_1 \geq \dots \geq s_k \geq 0\}$ et \mathbb{R}^{k++} l'intérieur de \mathbb{R}^{k+} , c'est-à-dire

$\mathbb{R}^{k^{++}} = \{s \in \mathbb{R}^k \mid s_1 > \dots > s_k > 0\}$. La matrice du schéma permet de caractériser les vecteurs partageant un même schéma du SLOPE : Soit $m \in \mathcal{P}_p^{\text{slope}} \setminus \{0\}$ et $k = \|m\|_\infty \geq 1$, alors

$$\{b \in \mathbb{R}^p \mid \text{schm}(b) = m\} = U_m \mathbb{R}^{k^{++}}.$$

Exemple 2.1. Soit $m = (2, -1, 0, 1, 2)$, on a

$$U_m = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 \\ 0 & -1 & 0 & 1 & 0 \end{pmatrix}^\top \quad \text{et} \quad U_{|m|_\downarrow} = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \end{pmatrix}^\top.$$

Définition 2.2. Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^p$, $m \in \mathcal{P}_p^{\text{slope}}$ un schéma du SLOPE non-nul et $k = \|m\|_\infty \geq 1$. La matrice d'appariement $\tilde{X}_m \in \mathbb{R}^{n \times k}$ de X est définie par $\tilde{X}_m = XU_m$; le paramètre d'appariement $\tilde{\lambda}_m \in \mathbb{R}^k$ de λ est défini par $\tilde{\lambda}_m = U_{|m|_\downarrow}^\top \lambda$.

La matrice d'appariement \tilde{X}_m du schéma m a moins de colonnes que la matrice de régression X . En effet, une composante nulle $m_i = 0$ mène à supprimer la colonne X_i de la matrice de régression X , et un groupe d'appariement $K \subseteq \{1, \dots, p\}$ de m (ensemble de composantes de m égales en valeur absolue) mène à substituer les colonnes $(X_i)_{i \in K}$ par une colonne égale à la somme signée : $\sum_{i \in K} \text{signe}(m_i) X_i$.

Exemple 2.2. Soient $X = (X_1|X_2|X_3|X_4|X_5)$, $m = (2, -1, 0, 1, 2)$ et $\lambda = (\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5) \in \mathbb{R}^5$. La matrice d'appariement et le paramètre d'appariement sont donnés par :

$$\tilde{X}_m = (X_1 + X_5 \mid -X_2 + X_4) \quad \text{et} \quad \tilde{\lambda}_m = \begin{pmatrix} \lambda_1 + \lambda_2 \\ \lambda_3 + \lambda_4 \end{pmatrix}.$$

Les notions de matrice du schéma et paramètre d'appariement permettent de donner une expression analytique du sous-différentiel de la norme ℓ_1 ordonnée. On a vu au chapitre 1 que $\partial J_\lambda(b)$ est l'enveloppe convexe d'une famille de sommets de la boule unité de la norme ℓ_1 ordonnée duale J_λ^* ; c'est donc un sous-ensemble du permutoèdre signé. La Proposition 2.1 précise ce résultat en montrant que $\partial J_\lambda(b)$ est une face du permutoèdre signé et exprime ce sous-différentiel comme l'intersection d'un espace affine avec la boule unité de la norme duale J_λ^* (Bogdan *et al.*, 2022).

Proposition 2.1. Soit $b \in \mathbb{R}^p \setminus \{0\}$. Le sous-différentiel de la norme ℓ_1 ordonnée vérifie les assertions suivantes.

1. Soit $\lambda \in \mathbb{R}^{p^+} \setminus \{0\}$ alors

$$\partial J_\lambda(b) = \left\{ v \in \mathbb{R}^p \mid J_\lambda^*(v) \leq 1 \text{ et } U_m^\top v = \tilde{\lambda}_m \right\} \quad \text{où } m = \text{schm}(b). \quad (2.1)$$

2. Si $\lambda \in \mathbb{R}^{p^{++}}$ alors l'espace affine donné dans l'expression (2.1) a une dimension minimale :

$$\text{aff}(\partial J_\lambda(b)) = \{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\}. \quad (2.2)$$

De la Proposition 2.1 découle une caractérisation des solutions du problème SLOPE : $\hat{\beta}$ est une solution non-nulle si et seulement si $J_\lambda^*(X^\top(y - X\hat{\beta})) \leq 1$ et $U_m^\top X^\top(y - X\hat{\beta}) = \tilde{\lambda}_m$, où $m = \text{schm}(\hat{\beta})$. Cette caractérisation est une réécriture très synthétique du Théorème 1 de la pré-publication de Nomura (2020).

Les articles de Vaïter *et al.* (2015, 2017) introduisent la notion de sous-espace modèle : le supplémentaire orthogonal de l'espace vectoriel parallèle au sous-différentiel. D'après l'expression (2.2), pour la norme ℓ_1 ordonnée lorsque $\lambda \in \mathbb{R}^{p^{++}}$, le sous-espace modèle est $\text{im}(U_m)$. Pour la norme ℓ_1 ordonnée, cette notion est difficile à interpréter. Par exemple, les schémas $\{\pm(2, 2, 1, 1), \pm(2, 2, -1, -1), \pm(1, 1, 2, 2), \pm(1, 1, -2, -2)\}$ ont le même sous-espace modèle alors que le schéma $(2, -2, 1, 1)$, qui possède pourtant les mêmes groupes d'appariement

($\{1, 2\}$ et $\{3, 4\}$), a un sous-espace modèle différent. Néanmoins, la notion de sous-espace modèle sera très utile pour décrire géométriquement, à la section suivante, la condition d'irreprésentabilité du SLOPE.

Pour conclure cette partie, la Proposition 2.2 montre que le problème d'optimisation SLOPE possède toujours au moins une solution dont le nombre de groupes d'appariement non-nuls est inférieur au rang de la matrice X .

Proposition 2.2. *Soient $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$ et $\lambda \in \mathbb{R}^{p++}$. Il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\|\text{schm}(\hat{\beta})\|_\infty \leq \text{rang}(X)$.*

Ce résultat issu de l'article de Dupuis et Tardivel (2024) généralise légèrement le Théorème 2.1 de l'article de Kremer *et al.* (2022) et le Corollaire 9 de l'article de Schneider et Tardivel (2022), puisque l'hypothèse d'unicité n'est pas requise. En particulier, au moins une solution du problème SLOPE a moins de n groupes d'appariement, ce qui est très faible lorsque le nombre de variables explicatives p dépasse très largement le nombre d'observation n . Cette remarque motive l'étude de la récupération du schéma d'un vecteur de régression ayant un très faible nombre de groupes d'appariement.

2.2 Propriétés théoriques de récupération du schéma du SLOPE

Soient $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, et $\lambda \in \mathbb{R}^{p++}$. On note $\mathcal{S}(y)$ l'ensemble des solutions du problème d'optimisation SLOPE :

$$\mathcal{S}(y) = \arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + J_\lambda(b) \right\}. \quad (2.3)$$

La notion de vecteur signe accessible pour le LASSO, sous l'hypothèse d'unicité, a été introduite par Sepehri et Harris (2017). Dans la Définition 2.3, nous introduisons une notion similaire pour le schéma du SLOPE qui ne requiert pas l'unicité de l'estimateur.

Définition 2.3 (Schéma accessible). *Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^{p++}$, et $m \in \mathcal{P}_p^{\text{slope}}$. On dit que le schéma du SLOPE m est accessible relativement à X et à la norme J_λ s'il existe $y \in \mathbb{R}^n$ et $\hat{\beta} \in \mathcal{S}(y)$ tels que $\text{schm}(\hat{\beta}) = m$.*

Des caractérisations d'un schéma accessible sont données par la Proposition 2.3.

Proposition 2.3. *Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^{p++}$, et $m \in \mathcal{P}_p^{\text{slope}}$. Le schéma du SLOPE m est accessible relativement à X et à la norme J_λ si et seulement si les assertions suivantes sont vérifiées :*

Caractérisation analytique : *Pour tout $b \in \mathbb{R}^p$ tel que $Xb = Xm$ on a $J_\lambda(b) \geq J_\lambda(m)$.*

Caractérisation géométrique : *L'espace vectoriel $\text{im}(X^\top)$ intersecte l'ensemble $\partial J_\lambda(m)$.*

D'après la Proposition 2.3, on vérifie facilement que : a) le schéma nul $m = 0$ est accessible ; et b) lorsque $\ker(X) = \{0\}$, tout schéma $m \in \mathcal{P}_p^{\text{slope}}$ est accessible (voir l'article de Skalski *et al.* (2022) pour le cas particulier où X est la matrice identité).

Il est facile de vérifier que $0 \in \mathcal{S}(y)$ si et seulement si $J_\lambda^*(X^\top y) \leq 1$. Par la suite, on appellera polyèdre nul du SLOPE l'ensemble $A_0 = \{y \in \mathbb{R}^n \mid J_\lambda^*(X^\top y) \leq 1\}$. Plus généralement, pour un schéma $m \in \mathcal{P}_p^{\text{slope}}$ non-nul, le Théorème 2.1 décrit l'ensemble des vecteurs $y \in \mathbb{R}^n$ pour lesquels au moins un élément de l'ensemble des solutions du problème SLOPE $\mathcal{S}(y)$ a un schéma m .

Théorème 2.1. *Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^{p++}$, et $m \in \mathcal{P}_p^{\text{slope}}$ un schéma du SLOPE non-nul accessible. On considère l'ensemble non vide $A_m = \left\{ y \in \mathbb{R}^n \mid \exists \hat{\beta} \in \mathcal{S}(y), \text{schm}(\hat{\beta}) = m \right\}$, alors les assertions suivantes sont satisfaites.*

1. *L'ensemble A_m est convexe et a une expression explicite :*

$$A_m = \left\{ y = z + \tilde{X}_m s \mid s \in \mathbb{R}^{k++} \text{ et } X^\top z \in \partial J_\lambda(m) \right\}.$$

1. Comme deux solutions du problème SLOPE ont la même norme ℓ_1 ordonnée, on a l'équivalence $0 \in \mathcal{S}(y) \Leftrightarrow \mathcal{S}(y) = \{0\}$.

2. L'ensemble A_m satisfait la caractérisation suivante :

$$y \in A_m \Leftrightarrow \begin{cases} \text{il existe } s \in \mathbb{R}^{k^{++}} \text{ tel que } \tilde{X}_m^\top Y - \tilde{\lambda}_m = \tilde{X}_m^\top \tilde{X}_m s & (\text{condition du schéma}) \\ X^\top \tilde{X}_m^\top + \tilde{\lambda}_m + X^\top (I_n - \tilde{X}_m \tilde{X}_m^\top) y \in \partial J_\lambda(m) & (\text{condition du sous-différentiel}) \end{cases}.$$

3. Si $\text{im}(X^\top) \cap \text{ri}(\partial J_\lambda(m)) \neq \emptyset$, alors A_m est un ensemble d'intérieur non-vide.

Lorsque A_m est d'intérieur non-vide et que y est un vecteur aléatoire ayant une densité strictement positive sur \mathbb{R}^n (par exemple lorsque y est la réponse d'un modèle de régression linéaire gaussien), alors la probabilité de l'événement $y \in A_m$, c'est-à-dire la probabilité de récupérer le schéma m avec l'estimateur SLOPE, est strictement positive. Par ailleurs, d'après la Proposition 5.3 de Gilbert (2017), la condition $Xb = Xm$ et $b \neq m$ implique $J_\lambda(b) > J_\lambda(m)$, qui est suffisante pour que A_m soit d'intérieur non-vide. Lorsque $X = I_p$, on observe que $A_m = \partial J_\lambda(m) + U_m \mathbb{R}^{k^{++}}$. Cette formule permet de retrouver la Figure 1.2 du chapitre 1. Notons que lorsque $\mathcal{S}(y)$ est un singleton pour tout $y \in \mathbb{R}^n$, les ensembles A_m pour $m \in \mathcal{P}_p^{\text{slope}}$ accessible forment une partition de \mathbb{R}^n . Une condition nécessaire et suffisante pour que $\mathcal{S}(y)$ soit un singleton pour tout $y \in \mathbb{R}^n$, appelée unicité uniforme, est donnée par la Proposition 2.4.

Proposition 2.4. Soient $X \in \mathbb{R}^{n \times p}$ et $\lambda \in \mathbb{R}^{p^{++}}$. Il existe $y \in \mathbb{R}^n$ tel que $\mathcal{S}(y)$ ne soit pas réduit à un singleton si et seulement si $\text{im}(X^\top)$ coupe une face du permutoèdre signé P_λ^\pm dont la dimension est strictement inférieure à $\dim(\ker(X))$.

Au dernier chapitre, nous énoncerons et prouverons une condition nécessaire et suffisante pour l'unicité uniforme dans un cadre plus général que le problème d'optimisation du SLOPE.

Exemple 2.3. Soient $\lambda = (4, 2, 1)$ et $X \in \mathbb{R}^{2 \times 3}$ la matrice suivante

$$X = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}.$$

D'après la Proposition 2.4, pour tout $y \in \mathbb{R}^2$, l'ensemble $\mathcal{S}(y)$ est un singleton si et seulement si l'espace $\text{im}(X^\top)$ ne coupe pas un sommet du permutoèdre signé P_λ^\pm , c'est-à-dire si et seulement si la condition suivante est satisfaite :

$$\forall \epsilon \in \{-1, 1\}^3 \quad \forall \pi \in S_3 \quad \underbrace{\begin{vmatrix} 1 & 0 & \epsilon_1 \lambda_{\pi(1)} \\ 0 & 1 & \epsilon_2 \lambda_{\pi(2)} \\ 1 & 1 & \epsilon_3 \lambda_{\pi(3)} \end{vmatrix}}_{= -\epsilon_1 \lambda_{\pi(1)} - \epsilon_2 \lambda_{\pi(2)} + \epsilon_3 \lambda_{\pi(3)}} \neq 0.$$

La dernière expression est vraie, puisque $|\epsilon_1 \lambda_{\pi(1)} + \epsilon_2 \lambda_{\pi(2)} - \epsilon_3 \lambda_{\pi(3)}| \geq 1$. Ainsi, $\mathcal{S}(y)$ est un singleton et on note $\hat{\beta}$ son unique élément. La Figure 2.1 illustre le polyèdre nul du SLOPE $A_0 = (XX^\top)^{-1} X(\text{im}(X^\top) \cap P_\lambda^\pm)$, ainsi que la partition de \mathbb{R}^2 en fonction du schéma de l'unique élément $\hat{\beta}$ de $\mathcal{S}(y)$.

2.2.1 Recupération du schéma dans le cas non-bruité : condition d'irreprésentabilité

Comme illustré par Fuchs (2004) (Théorème 2) et par Bühlmann et Van De Geer (2011) (Théorème 7.1), la condition d'irreprésentabilité est nécessaire pour la récupération des signes des coefficients de régression par l'estimateur LASSO dans cas non-bruité. De façon similaire, nous allons étudier le cas non-bruité pour dégager une condition théorique permettant de récupérer le schéma du SLOPE des coefficients de régression. Notons que dans le cas non-bruité la réponse du modèle de régression est $y = X\beta$ ainsi $y \in \tilde{X}_m \mathbb{R}^{k^{++}}$ où $m = \text{schm}(\beta)$ et $k = \|m\|_\infty$. La Proposition 2.5 donne une condition nécessaire et suffisante pour la récupération du schéma dans le cas non-bruité.

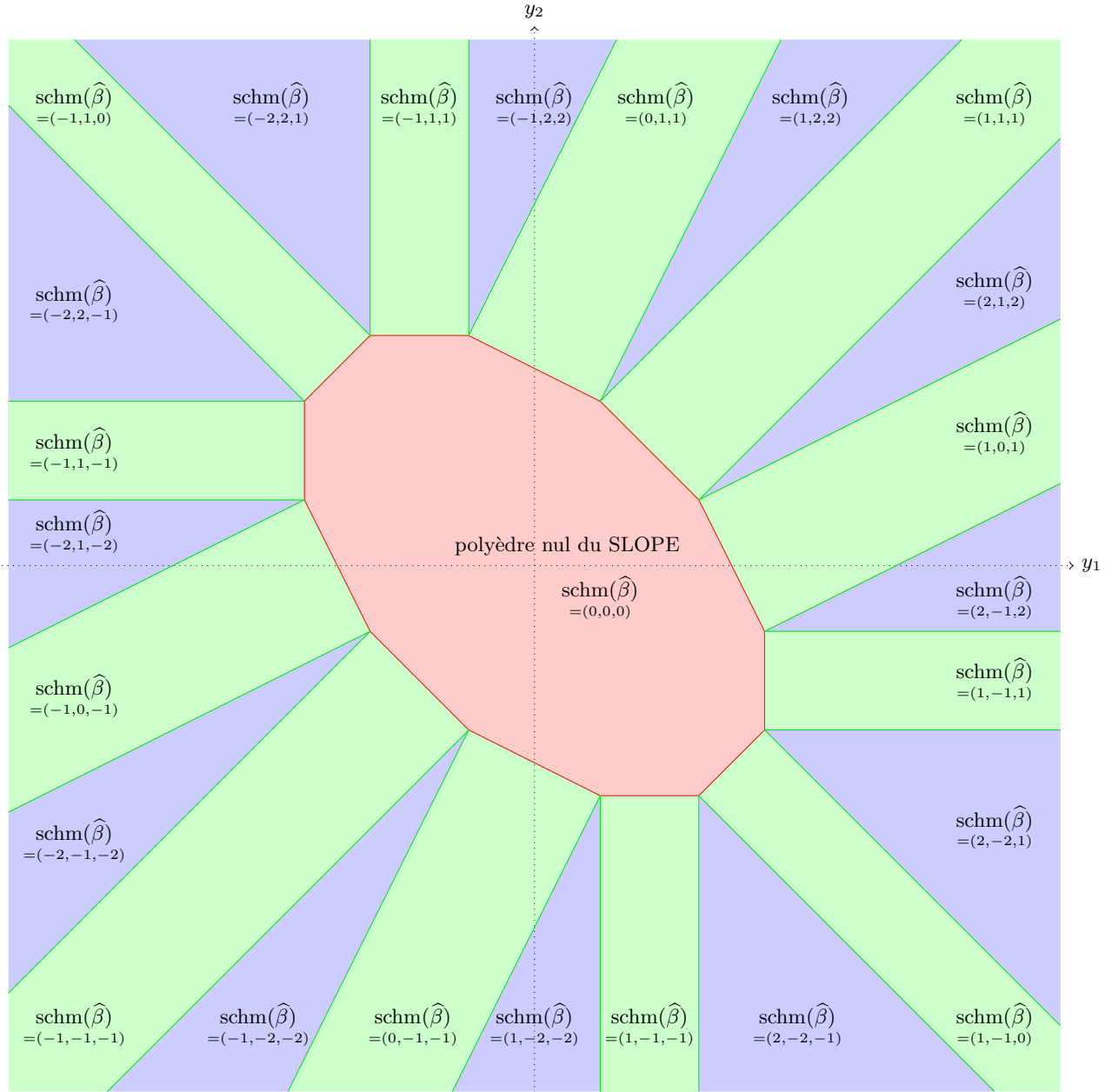


FIGURE 2.1 – Cette figure illustre le polyèdre nul du SLOPE, en rouge, et fournit $\text{schm}(\hat{\beta})$, où $\hat{\beta}$ est l'unique élément de $\mathcal{S}(y)$, selon la localisation de $y \in \mathbb{R}^2$.

Proposition 2.5. Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^{p^{++}}$, $m \in \mathcal{P}_p^{\text{slope}}$ avec $m \neq 0$ et $k = \|m\|_\infty \geq 1$. Il existe $y \in \tilde{X}_m \mathbb{R}^{k^{++}}$ et il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\text{schm}(\hat{\beta}) = m$ si et seulement si l'une des assertions suivantes est satisfaite.

Caractérisation analytique : On a $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \in \partial J_\lambda(m)$ (ou de façon équivalente $J_\lambda^*(X^\top \tilde{X}_m^\top + \tilde{\lambda}_m) \leq 1$ et $\tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top)$).

Caractérisation géométrique : L'espace vectoriel $X^\top X \text{im}(U_m) = \text{im}(X^\top \tilde{X}_m)$ coupe l'ensemble $\partial J_\lambda(m)$.

En référence aux premiers travaux traitant de la récupération du signe par l'estimateur LASSO (Zhao et Yu, 2006), on appellera condition d'irreprésentabilité² du SLOPE la condition $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \in \partial J_\lambda(m)$. Notons que

². Le nom « condition d'irreprésentabilité » est bien connu dans la littérature, cependant l'auteur du manuscrit ne comprend pas le sens de cette dénomination.

les conditions similaires à $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \in \partial J_\lambda(m)$ sont souvent exprimées comme des conditions dépendant de la norme duale de pénalité (voir par exemple (Zhao et Yu, 2006; Zou, 2006) pour le LASSO ou (Bach, 2008) pour le groupe LASSO). Ces simplifications résultent d'hypothèses mineures sur la matrice X , qui pour l'estimateur SLOPE, consisterait à supposer que $\tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top)$.

Pour se convaincre de la caractérisation géométrique, notons que si $\hat{\beta} \in \mathcal{S}(X\beta)$ a le même schéma que β , alors $\hat{\beta}, \beta \in U_m \mathbb{R}^{k++}$ et $X^\top X(\beta - \hat{\beta}) \in \partial J_\lambda(m)$. Comme $\beta - \hat{\beta}$ appartient au sous-espace modèle $\text{im}(U_m)$, on en déduit que $X^\top X \text{im}(U_m) \cap \partial J_\lambda(m) \neq \emptyset$.

Exemple 2.4. *On revisite l'Exemple 2.3 en considérant $m \in \{(1, 0, 1), (1, 1, 1)\}$. Pour $m = (1, 0, 1)$ la condition d'irreprésentabilité du SLOPE n'est pas satisfaite, en effet*

$$\tilde{X}_m = \begin{pmatrix} 2 & 1 \end{pmatrix}^\top, \tilde{\lambda}_m = 6 \text{ et } J_\lambda^*(X^\top \tilde{X}_m^\top + \tilde{\lambda}_m) = J_\lambda^*(X^\top \tilde{X}_m (\tilde{X}_m^\top \tilde{X}_m)^{-1} \tilde{\lambda}_m) = \frac{36}{35} > 1.$$

En revanche, pour $m = (1, 1, 1)$ la condition d'irreprésentabilité du SLOPE est satisfaite, en effet

$$\tilde{X}_m = \begin{pmatrix} 2 & 2 \end{pmatrix}^\top, \tilde{\lambda}_m = 7, J_\lambda^*(X^\top \tilde{X}_m^\top + \tilde{\lambda}_m) = J_\lambda^*(X^\top \tilde{X}_m (\tilde{X}_m^\top \tilde{X}_m)^{-1} \tilde{\lambda}_m) = 1 \text{ et } \tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top).$$

La Figure 2.2 illustre géométriquement la condition d'irreprésentabilité pour les schémas $m \in \{(1, 0, 1), (1, 1, 1)\}$.

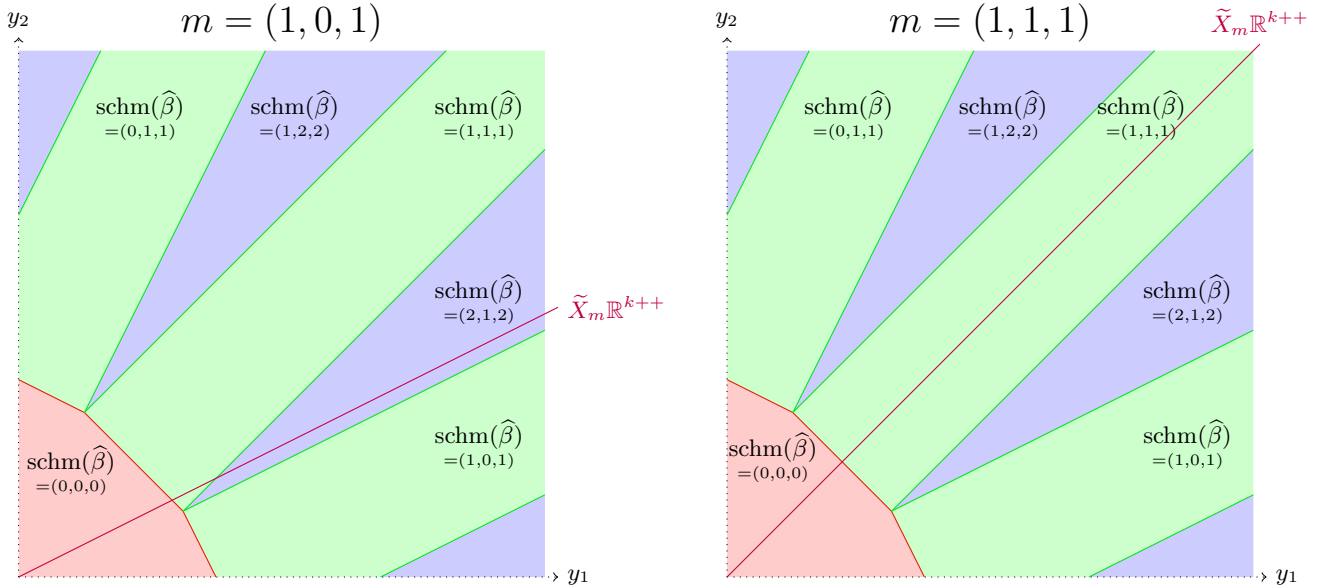


FIGURE 2.2 – La figure de gauche (resp. de droite) confirme que la condition d'irreprésentabilité du SLOPE n'est pas satisfaite (resp. est satisfaite) pour $m = (1, 0, 1)$ (resp. pour $m = (1, 1, 1)$). En effet pour la demi-droite pourpre $\tilde{X}_m \mathbb{R}^{k++}$ n'a aucun point commun avec A_m donc dans le cas sans bruit une solution de l'estimateur SLOPE ne peut pas avoir pour schéma m (resp. coupe l'ensemble A_m donc dans le cas sans bruit une solution de l'estimateur SLOPE peut avoir pour schéma m).

La Proposition 2.6 montre que le vecteur $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m$ a une interprétation géométrique : lorsque $\text{im}(X^\top \tilde{X}_m)$ coupe le plus petit espace affine contenant $\partial J_\lambda(m)$, alors l'intersection est réduite au singleton $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m$ (Bogdan *et al.*, 2022).

Proposition 2.6. *Soient $X \in \mathbb{R}^{n \times p}$, $m \in \mathcal{P}_p^{\text{slope}}$ avec $m \neq 0$ et $\lambda \in \mathbb{R}^{p++}$. Les assertions suivantes sont satisfaites.*

1. Si $\tilde{\lambda}_m \notin \text{im}(\tilde{X}_m^\top)$ alors $\text{aff}(\partial J_\lambda(m)) \cap \text{im}(X^\top \tilde{X}_m) = \emptyset$.

2. Si $\tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top)$ alors $\text{aff}(\partial J_\lambda(m)) \cap \text{im}(X^\top \tilde{X}_m) = \{X^\top \tilde{X}_m^\top + \tilde{\lambda}_m\}$.

À l'instar du LASSO (Wainwright, 2009), la Proposition 2.7 montre que sous une condition très réaliste sur les résidus ε , lorsque la condition d'irreprésentabilité du SLOPE n'est pas satisfaite, la probabilité de récupération du schéma est inférieure à 1/2.

Proposition 2.7. *Soit $Y = X\beta + \varepsilon$ où ε et $-\varepsilon$ ont la même loi et $m = \text{schm}(\beta)$. Si $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \notin \partial J_\lambda(m)$ alors*

$$\mathbb{P}(\exists \hat{\beta} \in \mathcal{S}(Y), \text{schm}(\hat{\beta}) = m) \leq 1/2$$

La Proposition 2.7 montre que la condition d'irreprésentabilité est nécessaire pour que le schéma de l'estimateur SLOPE soit égal au schéma de β avec une probabilité dépassant un demi. Par ailleurs, le Théorème 4.1 de Bogdan *et al.* (2022) montre que la condition $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \in \text{ri}(\partial J_\lambda(m))$, légèrement plus forte que la condition d'irreprésentabilité, est suffisante pour la récupération asymptotique de $\text{schm}(\beta)$ via l'estimateur SLOPE. Le Théorème 4.1 n'est pas détaillé dans ce manuscrit car nous allons nous focaliser, à la section suivante, sur le relâchement de la condition d'irreprésentabilité. Plus spécifiquement, nous montrons qu'il est possible de récupérer le schéma de β en affaiblissant la condition d'irreprésentabilité moyennant l'application de l'opérateur proximal de la norme ℓ_1 ordonnée à l'estimateur SLOPE.

2.3 Relâchement de la condition d'irreprésentabilité du SLOPE

La convergence d'un estimateur n'implique pas la convergence du schéma. En effet, même si $\hat{\beta}$ est un estimateur très proche de β , une composante nulle de β peut ne pas être exactement estimée à zéro et deux composantes égales pour β peuvent ne pas être estimées exactement à la même valeur. Nous verrons que l'utilisation de l'opérateur proximal de la norme ℓ_1 ordonnée J_λ défini comme l'unique solution du problème d'optimisation suivant :

$$\text{prox}_{\tau\lambda}(y) = \arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - b\|_2^2 + \tau J_\lambda(b) \right\} \text{ avec } \tau \geq 0.$$

permet de pallier le défaut de récupération du schéma par l'estimateur $\hat{\beta}$.

Concernant le calcul de l'opérateur proximal, on peut, sans perte de généralité, se ramener au cas où $y \in \mathbb{R}^{p+}$. Un algorithme résolvant ce problème d'optimisation est donné dans l'article de Bogdan *et al.* (2015), et une formule concise pour l'opérateur proximal est fournie à la Proposition 2.8 (Dupuis et Tardivel, 2022; Tardivel *et al.*, 2020).

Proposition 2.8. *Soit $y \in \mathbb{R}^{p+}$, $\lambda \in \mathbb{R}^{p+} \setminus \{0\}$ et $\tau \geq 0$. On considère la suite de Cesàro $(C_j)_{1 \leq j \leq p}$ où $C_j = \frac{1}{j} \sum_{i=1}^j (y_i - \tau \lambda_i)$, et on note $k \in \{1, \dots, p\}$ le plus grand entier pour lequel cette suite atteint son maximum alors :*

$$\text{prox}_{\tau\lambda}(y) = \begin{cases} (0, \dots, 0) & \text{si } C_k \leq 0 \\ (\underbrace{C_k, \dots, C_k}_{k \text{ éléments}}, \text{prox}_{(\tau\lambda_{k+1}, \dots, \tau\lambda_p)}(y_{k+1}, \dots, y_p)) & \text{sinon} \end{cases} \quad (2.4)$$

Une lecture intuitive de cette formule suggère que lorsque deux composantes de $\hat{\beta}$ sont proches, l'opérateur proximal $\text{prox}_{\tau\lambda}(\hat{\beta})$ tend à les apparier, et quand une composante de $\hat{\beta}$ est approximativement zéro, l'opérateur proximal $\text{prox}_{\tau\lambda}(\hat{\beta})$ tend à l'annuler. Le Théorème 2.2 donne une condition théorique sous laquelle l'application de l'opérateur proximal à l'estimateur SLOPE permet de récupérer asymptotiquement³ le schéma de β .

3. Dans le cadre du Théorème 2.2, on montre que l'estimateur SLOPE $\hat{\beta}^{(r)}$ est convergent.

Théorème 2.2. Soient $X \in \mathbb{R}^{n \times p}$, $\lambda \in \mathbb{R}^{p^{++}}$, $\beta \in \mathbb{R}^p$, $m = \text{schm}(\beta)$, et $y^{(r)} = X\beta + \varepsilon^{(r)}$ où $(\varepsilon^{(r)})_{r \in \mathbb{N}}$ est une suite de \mathbb{R}^n telle que $\lim_{r \rightarrow +\infty} \varepsilon^{(r)} = 0$. On pose $\widehat{\beta}^{(r)}$ et $\widehat{\beta}^{0(r)}$ les estimateurs définis ci-dessous :

$$\widehat{\beta}^{(r)} \text{ est une solution de } \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y^{(r)} - Xb\|_2^2 + \gamma_r J_\lambda(b) \right\} \text{ où } \gamma_r > 0, \quad (2.5)$$

$$\widehat{\beta}^{0(r)} = \begin{cases} \text{est une solution de } \min J_\lambda(b) \text{ sous la contrainte } Xb = XX^+y^{(r)} & \text{si } \dim(\ker(X)) \geq 1 \\ (X^\top X)^{-1}X^\top y^{(r)} & \text{si } \ker(X) = \{0\} \end{cases}.$$

Lorsque $\ker(X) = \{0\}$, les assertions suivantes sont satisfaites :

1. Si $\lim_{r \rightarrow +\infty} \gamma_r = 0$, alors il existe un seuil $\tau \geq 0$ tel que $\lim_{r \rightarrow +\infty} \text{schm}(\text{prox}_{\tau\lambda}(\widehat{\beta}^{(r)})) = m$.
2. Il existe un seuil $\tau \geq 0$ tel que $\lim_{r \rightarrow +\infty} \text{schm}(\text{prox}_{\tau\lambda}(\widehat{\beta}^{0(r)})) = m$.

Lorsque $\dim(\ker(X)) \geq 1$ et sous l'hypothèse que pour tout $b \neq \beta$, $Xb = X\beta$ implique $J_\lambda(b) > J_\lambda(\beta)$, les assertions suivantes sont satisfaites :

3. Si $\lim_{r \rightarrow +\infty} \gamma_r = 0$, alors il existe un seuil $\tau \geq 0$ tel que $\lim_{r \rightarrow +\infty} \text{schm}(\text{prox}_{\tau\lambda}(\widehat{\beta}^{(r)})) = m$.
4. Il existe un seuil $\tau \geq 0$ tel que $\lim_{r \rightarrow +\infty} \text{schm}(\text{prox}_{\tau\lambda}(\widehat{\beta}^{0(r)})) = m$.

Quelques remarques sur le Théorème 2.2 sont données ci-dessous :

- L'estimateur $\widehat{\beta}^{0(r)}$ est l'estimateur des moindres carrés lorsque $\ker(X) = \{0\}$ ou l'estimateur par récupération J_λ lorsque $\dim(\ker(X)) \geq 1$. Ces estimateurs peuvent être interprétés comme des cas limites de l'estimateur SLOPE $\widehat{\beta}^{(r)}$ lorsque $\gamma_r = 0$.
- La condition $Xb = X\beta$ et $b \neq \beta$ implique $J_\lambda(b) > J_\lambda(\beta)$ ne dépend que du schéma m de β (voir (Gilbert, 2017, Proposition 5.3)). Ainsi, cette condition est légèrement plus forte que la condition d'accessibilité du schéma m .
- Appliquer l'opérateur proximal à un estimateur $\widehat{\beta}$ est une construction particulière d'estimateur « seuillé » ; en effet, $\partial J_\lambda(\widehat{\beta}) \subseteq \partial J_\lambda(\text{prox}_{\tau\lambda}(\widehat{\beta}))$. Dans un cadre beaucoup plus général, le Théorème 5.3 de l'article de Graczyk *et al.* (2023) donne une condition théorique pour qu'un estimateur seuillé récupère le sous-différentiel de β (c'est-à-dire le schéma du SLOPE de β pour la cas particulier de la norme J_λ).
- On peut montrer que si m ne satisfait pas la condition d'accessibilité (c'est-à-dire s'il existe $v \in \mathbb{R}^p$ tel que $Xv = Xm$ et $J_\lambda(v) < J_\lambda(m)$), alors, indépendamment de $\tau \geq 0$, l'application de l'opérateur proximal à l'estimateur SLOPE ne permet pas de récupérer ce schéma (voir la Proposition 5.3 de l'article de Graczyk *et al.* (2023)). La condition $Xb = X\beta$ et $b \neq \beta$ implique $J_\lambda(b) > J_\lambda(\beta)$ est donc quasiment minimale pour la récupération du schéma. Nous ne savons pas si cet écart minime entre la condition nécessaire et la condition suffisante peut être comblé.

2.4 Expériences numériques

Pour nos simulations, nous fixons $\beta \in \{0, 1\}^{784}$. La Figure 2.3 donne une visualisation de β redimensionné sous la forme d'un graphique de taille 28×28 , où 0 représente un pixel blanc et 1 représente un pixel noir.

Pour les paramètres de la norme ℓ_1 ordonnée, nous choisissons $\lambda = (\sqrt{j} - \sqrt{j-1})_{1 \leq j \leq 784}$ comme suggéré dans l'article de Nomura (2020).⁴

4. Le code des expériences numériques de ce chapitre est disponible en ligne sur ma page internet

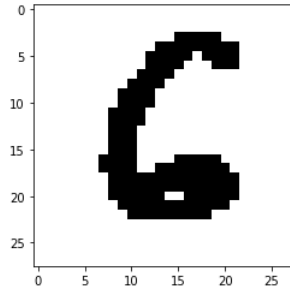


FIGURE 2.3 – Le vecteur $\beta \in \{0, 1\}^{784}$ redimensionné en une image de taille 28×28 , représentant le chiffre six.

Comparaison entre les conditions d'irreprésentabilité et d'accessibilité

Soit $X \in \mathbb{R}^{n \times 784}$ une matrice dont les coefficients sont indépendants et de même loi $\mathcal{N}(0, 1/n)$. D'après la Proposition 2.3, la probabilité que β soit accessible par rapport à X et J_λ vaut

$$\mathbb{P}_X(\min\{J_\lambda(b) \mid Xb = X\beta\} = J_\lambda(\beta)).$$

De plus, la probabilité que la condition d'irreprésentabilité pour β soit satisfaite vaut :

$$\mathbb{P}_X(X^\top \tilde{X}_\beta^\top + \tilde{\lambda}_\beta \in \partial J_\lambda(\beta)).$$

La Figure 2.4 fournit ces probabilités en fonction du nombre de lignes de la matrice X . Pour l'approximation de ces probabilités, nous avons utilisé 1000 réalisations de la matrice X .

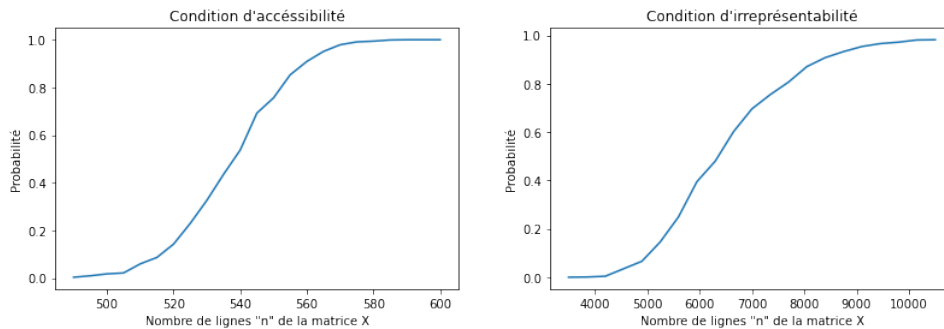


FIGURE 2.4 – Le graphique de gauche donne la probabilité que le schéma β soit accessible en fonction du nombre de lignes de la matrice X . On remarque que cette probabilité est quasiment nulle lorsque $n \leq 500$ et vaut quasiment 1 lorsque $n \geq 600$. Le graphique de droite donne la probabilité que la condition d'irreprésentabilité soit satisfaite pour le schéma m en fonction du nombre de lignes de la matrice X . On remarque que cette probabilité est quasiment nulle lorsque $n \leq 4000$ et vaut quasiment 1 lorsque $n \geq 10000$.

Récupération du schéma par seuillage du SLOPE

Pour l'expérience numérique suivante, nous considérons le modèle de régression linéaire gaussien $Y = X\beta + \varepsilon$, où $X \in \mathbb{R}^{600 \times 784}$ est une matrice dont les coefficients sont indépendants et de même loi $\mathcal{N}(0, 1/600)$, et les composantes de $\varepsilon \in \mathbb{R}^{600}$ sont indépendantes et de même loi $\mathcal{N}(0, 0.05^2)$. Pour des réalisations particulières Y et X , nous notons $\hat{\beta}_\gamma$ l'unique solution du problème suivant

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|Y - Xb\|_2^2 + \gamma J_\lambda(b) \right\}.$$

La notation γ_{sure} représente le paramètre sélectionné via la formule SURE du SLOPE (Minami, 2020) minimisant $\gamma > 0 \mapsto \|Y - X\hat{\beta}_\gamma\|_2^2 + 2 \times 0.05^2 \|\text{schm}(\hat{\beta}_\gamma)\|_\infty$. À la Figure 2.5, nous illustrons que, pour $n = 600$, SLOPE ne peut pas récupérer le schéma β , tandis que l'estimateur SLOPE seuillé peut récupérer ce schéma.

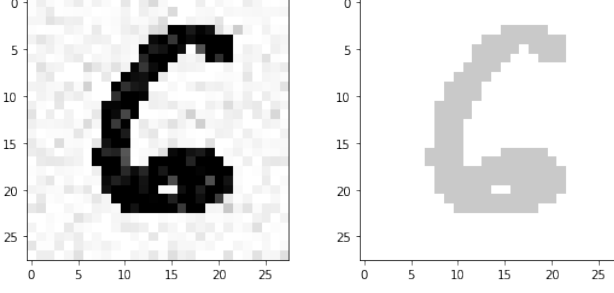


FIGURE 2.5 – Comme la condition d'irreprésentabilité n'est pas satisfaite pour la matrice particulière X et le schéma β , l'estimateur SLOPE $\hat{\beta}_{\gamma_{\text{sure}}}$ est peu susceptible de récupérer le schéma β . En effet, le graphique de gauche montre que $\hat{\beta}_{\gamma_{\text{sure}}}$, redimensionné comme une image de taille 28×28 , ne récupère pas le schéma β . En revanche, la condition d'accessibilité est satisfaite et donc l'estimateur SLOPE seuillé peut révéler le schéma β . En effet, le graphique de droite illustre que $\text{prox}_{\tau\lambda}(\hat{\beta}_{\gamma_{\text{sure}}})$, redimensionné comme une image de taille 28×28 , récupère le schéma β (ici, $\tau > 0$ est le plus petit nombre réel pour lequel $\text{prox}_{\tau\lambda}(\hat{\beta}_{\gamma_{\text{sure}}})$ a un unique groupe d'appariement non-nul).

2.5 Annexes : preuves

Preuve de la Proposition 2.1

Démonstration. 1) : Dans un premier temps montrons l'inclusion $\partial J_\lambda(b) \subseteq \{v \in \mathbb{R}^p \mid J_\lambda^*(v) \leq 1 \text{ et } U_m^\top v = \tilde{\lambda}_m\}$. Soit $v \in \partial J_\lambda(b)$ alors $J_\lambda^*(v) \leq 1$ ainsi, il reste à prouver que $U_m^\top v = \tilde{\lambda}_m$. On pose $q_l = \text{Card}(\{i \mid |m_i| \geq k+1-l\})$, pour $l \in \{1, \dots, k\}$; de plus on a l'inégalité suivante :

$$\sum_{i=1}^l [U_m^\top v]_i = \sum_{i \mid |m_i| \geq k+1-l} \text{signe}(m_i) v_i \leq \sum_{i \mid |m_i| \geq k+1-l} |v_i| \leq \sum_{i=1}^{q_l} |v|_{\downarrow i} \leq \sum_{i=1}^{q_l} \lambda_i = \sum_{i=1}^l [\tilde{\lambda}_m]_i. \quad (2.6)$$

On pose $b = U_m s$ avec $s \in \mathbb{R}^{k++}$ alors,

$$\begin{aligned} b^\top v &= s^\top U_m^\top v = \sum_{i=1}^k s_i [U_m^\top v]_i = \sum_{l=1}^{k-1} (s_l - s_{l+1}) \sum_{i=1}^l [U_m^\top v]_i + s_k \sum_{i=1}^k [U_m^\top v]_i \\ &\leq \sum_{l=1}^{k-1} (s_l - s_{l+1}) \sum_{i=1}^l [\tilde{\lambda}_m]_i + s_k \sum_{i=1}^k [\tilde{\lambda}_m]_i = \sum_{l=1}^k s_l [\tilde{\lambda}_m]_l = J_\lambda(b). \end{aligned}$$

De plus, pour $v \in \partial J_\lambda(b)$ on a $b^\top v = J_\lambda(b)$ ainsi

$$\sum_{i=1}^l [U_m^\top v]_i = \sum_{i=1}^l [\tilde{\lambda}_m]_i \quad \forall l \in \{1, \dots, k\}$$

donc les inégalités données dans (2.6) sont les égalités. Ainsi, pour tout $l \in \{1, \dots, k\}$ on a $[U_m^\top v]_l = [\tilde{\lambda}_m]_l$ d'où $U_m^\top v = \tilde{\lambda}_m$. Prouvons l'autre inclusion : $\partial J_\lambda(b) \supseteq \{v \in \mathbb{R}^p \mid J_\lambda^*(v) \leq 1 \text{ et } U_m^\top v = \tilde{\lambda}_m\}$. Supposons que $v \in \mathbb{R}^p$ satisfait $J_\lambda^*(v) \leq 1$ et $U_m^\top v = \tilde{\lambda}_m$. Pour prouver que $v \in \partial J_\lambda(b)$ il reste à établir que $b^\top v = J_\lambda(b)$. Comme

$b = U_m s$ avec $s \in \mathbb{R}^{k^{++}}$, on a

$$b^\top v = s^\top U_m^\top v = s^\top \tilde{\lambda}_m = J_\lambda(b).$$

2) : Dans un premier temps, pour simplifier, on suppose que $b \in \mathbb{R}^{p^+}$. Ainsi, il existe une subdivision $1 \leq q_1 < \dots < q_k \leq p$ telle que

$$\text{supp}(b) = \{1, \dots, q_k\} \text{ et } b_1 = \dots = b_{q_1} > b_{q_1+1} = \dots = b_{q_2} > \dots > b_{q_{k-1}+1} = \dots = b_{q_k} > 0.$$

On pose $\tilde{v} = U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m$, montrons que $\tilde{v} \in \partial J_\lambda(m)$. Clairement $U_m^\top \tilde{v} = \tilde{\lambda}_m$ ainsi, il reste à montrer que $J_\lambda(\tilde{v}) \leq 1$. Comme $U_m^\top U_m$ est une matrice diagonale dont les coefficients diagonaux sont $q_1, q_2 - q_1, \dots, q_k - q_{k-1}$ on peut réécrire \tilde{v} sous la forme explicite suivante :

$$\text{supp}(\tilde{v}) = \{1, \dots, q_k\} \text{ et } (\tilde{v}_1, \dots, \tilde{v}_{q_k}) = \underbrace{(\tilde{v}_1, \dots, \tilde{v}_{q_1})}_{=\frac{\lambda_1 + \dots + \lambda_{q_1}}{q_1}}, \underbrace{(\tilde{v}_{q_1+1}, \dots, \tilde{v}_{q_2})}_{=\frac{\lambda_{q_1+1} + \dots + \lambda_{q_2}}{q_2 - q_1}}, \dots, \underbrace{(\tilde{v}_{q_{k-1}+1}, \dots, \tilde{v}_{q_k})}_{=\frac{\lambda_{q_{k-1}+1} + \dots + \lambda_{q_k}}{q_k - q_{k-1}}} \in \mathbb{R}^{q_k^{++}}. \quad (2.7)$$

Comme $\lambda \in \mathbb{R}^{p^{++}}$, on a les inégalités suivantes :

$$\begin{cases} \|\tilde{v}\|_{(i)} = \tilde{v}_1 + \dots + \tilde{v}_i < \lambda_1 + \dots + \lambda_i & \text{si } i \notin \{q_1, q_2, \dots, q_k\} \\ \|\tilde{v}\|_{(i)} = \tilde{v}_1 + \dots + \tilde{v}_i = \lambda_1 + \dots + \lambda_i & \text{si } i \in \{q_1, q_2, \dots, q_k\} \end{cases}$$

Ainsi, $J_\lambda^*(\tilde{v}) \leq 1$ donc $\tilde{v} \in \partial J_\lambda(m)$. Soit $h \in \ker(U_m^\top)$. Montrons que $\tilde{v} + \eta h \in \partial J_\lambda(m)$ pour $\eta > 0$ bien choisi. Soit $\eta > 0$ vérifiant

$$\eta < \frac{\min\{\tilde{v}_{q_1} - \tilde{v}_{q_2}, \dots, \tilde{v}_{q_{k-1}} - \tilde{v}_{q_k}, \tilde{v}_{q_k}\}}{2\|h\|_\infty}. \quad (2.8)$$

Si $\eta > 0$, satisfaisant (2.8), est suffisamment petit alors pour tout $i \notin \{q_1, q_2, \dots, q_k\}$ on a l'inégalité suivante

$$\|\tilde{v} + \eta h\|_{(i)} \leq \|\tilde{v}\|_{(i)} + \eta \|h\|_{(i)} \leq \lambda_1 + \dots + \lambda_i.$$

Par ailleurs, comme $h \in \ker(U_m^\top)$, pour tout $i \in \{q_1, q_2, \dots, q_k\}$ on a $h_1 + \dots + h_i = 0$ de plus comme η satisfait (2.8) on a $\{|\tilde{v} + \eta h|_{\downarrow 1}, \dots, |\tilde{v} + \eta h|_{\downarrow i}\} = \{\tilde{v}_1 + \eta h_1, \dots, \tilde{v}_i + \eta h_i\}$ d'où

$$\|\tilde{v} + \eta h\|_{(i)} = \tilde{v}_1 + \eta h_1 + \dots + \tilde{v}_i + \eta h_i = \lambda_1 + \dots + \lambda_i.$$

Donc, $J_\lambda^*(\tilde{v} + \eta h) \leq 1$. Enfin, comme $U_m^\top(\tilde{v} + \eta h) = \tilde{\lambda}_m$ on en déduit que $\tilde{v} + \eta h \in \partial J_\lambda(m)$; ce qui établit que $\text{aff}(\partial J_\lambda(b)) = \tilde{v} + \ker(U_m^\top) = \{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\}$. Pour le cas général si $b \neq |b|_\downarrow$ on considère une transformation orthogonale $\psi(x) = (\epsilon_1 x_{\pi(1)}, \dots, \epsilon_p x_{\pi(p)})$, pour $x \in \mathbb{R}^p$, où $\epsilon \in \{-1, 1\}^p$ et $\pi \in S_p$ sont choisis de telle sorte que $\psi(b) = |b|_\downarrow$. L'inclusion $\text{aff}(\partial J_\lambda(b)) \subseteq \{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\}$ est immédiate; pour l'égalité montrons que ces deux espaces ont la même dimension. Comme $\partial J_\lambda(b) = \psi^{-1}(\partial J_\lambda(|b|_\downarrow))$ (voir la preuve du Lemme 1.4 au chapitre 1) alors

$$\dim(\text{aff}(\partial J_\lambda(b))) = \dim(\text{aff}(\partial J_\lambda(|b|_\downarrow))) = p - k = \dim(\{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\})$$

Ainsi $\text{aff}(\partial J_\lambda(b)) = \{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\}$. □

Dans la preuve du 2) le sous-différentiel de la norme ℓ_1 ordonnée en $b \in \mathbb{R}^{p^+}$ est le produit cartésien :

$$\partial J_\lambda(b) = \begin{cases} P_{\lambda_1, \dots, \lambda_{q_1}} \times P_{\lambda_{q_1+1}, \dots, \lambda_{q_2}} \times \dots \times P_{\lambda_{q_{k-1}+1}, \dots, \lambda_{q_k}} \times P_{\lambda_{q_k+1}, \dots, \lambda_p}^\pm & \text{si } q_k < p \\ P_{\lambda_1, \dots, \lambda_{q_1}} \times P_{\lambda_{q_1+1}, \dots, \lambda_{q_2}} \times \dots \times P_{\lambda_{q_{k-1}+1}, \dots, \lambda_{q_k}} & \text{si } q_k = p \end{cases}.$$

On peut remarquer que \tilde{v} , défini à l'équation (2.7), est l'isobarycentre de cet ensemble. Plus généralement, pour $m \in \mathcal{P}_p^{\text{slope}}$, $U_m(U_m^\top U_m)^{-1} \tilde{\lambda}_m$ est l'isobarycentre de $\partial J_\lambda(m)$. Cette remarque sera utile pour la preuve du Théorème 2.2.

Preuve de la Proposition 2.2

Démonstration. Si $0 \in \mathcal{S}(y)$ alors la proposition est clairement vraie. Supposons que $0 \notin \mathcal{S}(y)$. Soit $\hat{\beta} \in \mathcal{S}(y)$ tel que le nombre de groupes d'appariement non-nuls $k = \|\text{schm}(\hat{\beta})\|_\infty \geq 1$ soit minimal. On pose $m = \text{schm}(\hat{\beta})$. Montrons que $\ker(\tilde{X}_m) = \{0\}$. Si $\dim(\ker(\tilde{X}_m)) \geq 1$ alors on choisit $h \in \ker(\tilde{X}_m)$, $h \neq 0$, on pose $\hat{\beta} = U_m s$ où $s \in \mathbb{R}^{k^{++}}$ et $c(t) = \hat{\beta} + t U_m h = U_m(s + th)$. Comme $\tilde{X}_m h = X U_m h = 0$ alors $X^\top(y - Xc(t)) = X^\top(y - X\hat{\beta})$. Soit $t_{\min} = \inf\{|t| \mid s + th \notin \mathbb{R}^{k^{++}}\} > 0$; par construction, pour $t \in]-t_{\min}, t_{\min}[$, $s + th \in \mathbb{R}^{k^{++}}$ donc $\text{schm}(c(t)) = m$. Par conséquent,

$$\begin{aligned} & \forall t \in]-t_{\min}, t_{\min}[\quad X^\top(y - Xc(t)) \in \partial J_\lambda(m) = \partial J_\lambda(c(t)), \\ \Rightarrow & \quad \forall t \in]-t_{\min}, t_{\min}[\quad c(t) \in \mathcal{S}(y). \end{aligned}$$

Comme $\mathcal{S}(y)$ est un ensemble fermé, on peut en déduire que $c(\pm t_{\min}) \in \mathcal{S}(y)$. Enfin, par construction de t_{\min} , un des vecteurs $s + t_{\min}h$ ou $s - t_{\min}h$ a moins de k composantes distinctes, donc $\|\text{schm}(c(t_{\min}))\|_\infty < k$ ou $\|\text{schm}(c(-t_{\min}))\|_\infty < k$ qui contredit le fait que $\hat{\beta} \in \mathcal{S}(y)$ a un nombre minimal de groupes d'appariement non-nuls. Ainsi $\text{rang}(\tilde{X}_m) = k$ ce qui achève la preuve car $\text{rang}(\tilde{X}_m) \leq \text{rang}(X)$. \square

Cette preuve reste valide lorsque $\lambda \in \mathbb{R}^{p^+} \setminus \{0\}$. Néanmoins, la Proposition 2.2 est énoncé avec $\lambda \in \mathbb{R}^{p^{++}}$ car pour le LASSO, lorsque $\lambda_1 = \dots = \lambda_p > 0$, ce résultat peut prêter à confusion puisque la notion d'appariement n'est pas pertinente. De plus, pour cet estimateur un résultat plus précis est connu : il existe toujours une solution au problème LASSO ayant un nombre de composantes non-nulles inférieur à $\text{rang}(X)$ (Osborne *et al.*, 2000) (a fortiori une telle solution à moins de $\text{rang}(X)$ groupes d'appariement non-nuls).

Preuve de la Proposition 2.3

Démonstration de la Proposition 2.3.

Caractérisation géométrique : Supposons que le schéma $m \in \mathcal{P}_p^{\text{slope}}$ soit accessible alors, il existe $y \in \mathbb{R}^n$ et il existe $\hat{\beta} \in \mathcal{S}(y)$ tels que $\text{schm}(\hat{\beta}) = m$ ainsi

$$X^\top(y - X\hat{\beta}) \in \partial J_\lambda(\hat{\beta}) = \partial J_\lambda(m).$$

Cette inclusion montre que $\text{im}(X^\top) \cap \partial J_\lambda(m) \neq \emptyset$. Inversement, supposons que $\text{im}(X^\top) \cap \partial J_\lambda(m) \neq \emptyset$ alors $X^\top z \in \partial J_\lambda(m)$ pour un certain $z \in \mathbb{R}^n$. On pose $y = z + Xm$ alors

$$X^\top(y - Xm) = X^\top z \in \partial J_\lambda(m).$$

donc $m \in \mathcal{S}(y)$ ainsi m est accessible.

Caractérisation analytique : Cette preuve est inspirée de la démonstration de la Proposition 4.1 de l'article de Gilbert (2017). Considérons la fonction $\chi_m : \mathbb{R}^p \rightarrow \{0, +\infty\}$ définie par

$$\chi_m(b) = \begin{cases} 0 & \text{si } Xb = Xm \\ +\infty & \text{sinon.} \end{cases}$$

Pour tout $b \in \mathbb{R}^p$ tel que $Xb = Xm$ on a $J_\lambda(b) \geq J_\lambda(m)$ si et seulement si la fonction $f : b \mapsto J_\lambda(b) + \chi_m(b)$ atteint son minimum en m . Comme $\partial\chi_m(b) = \text{im}(X^\top)$ dès que $Xb = Xm$, m est un minimiseur de f si et seulement si

$$0 \in \text{im}(X^\top) + \partial J_\lambda(m) \Leftrightarrow \text{im}(X^\top) \cap \partial J_\lambda(m) \neq \emptyset.$$

Ainsi, la caractérisation analytique est équivalent à la caractérisation géométrique. \square

Preuve du Théorème 2.1

Démonstration. **1)** On pose $y = z + XU_m s$ avec $z \in \mathbb{R}^n$ tel que $X^\top z \in \partial J_\lambda(m)$ et $s \in \mathbb{R}^{k++}$. Alors, $X^\top(y - XU_m s) = X^\top z \in \partial J_\lambda(m) = \partial J_\lambda(U_m s)$ ainsi $U_m s \in \mathcal{S}(y)$ d'où $y \in A_m$. Inversement, soit $y \in A_m$ alors il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\text{schm}(\hat{\beta}) = m$. En posant $z = y - X\hat{\beta}$ et comme $X\hat{\beta} = \tilde{X}_m s$ pour un certain $s \in \mathbb{R}^{k++}$ on en déduit que $y = z + \tilde{X}_m s$ de plus $\partial J_\lambda(m) \ni X^\top(y - X\hat{\beta}) = X^\top z$. Montrons à présent que A_m est un ensemble convexe. Soit $y \in A_m, \bar{y} \in A_m$ et $\alpha \in [0, 1]$ alors $y = z + XU_m s$ pour un certain $z \in \mathbb{R}^n$ tel que $X^\top z \in \partial J_\lambda(m)$ et un certain $s \in \mathbb{R}^{k++}$. De même $\bar{y} = \bar{z} + XU_m \bar{s}$ avec $\bar{z} \in \mathbb{R}^n$ tel que $X^\top \bar{z} \in \partial J_\lambda(m)$ et $\bar{s} \in \mathbb{R}^{k++}$. Ainsi $\alpha y + (1 - \alpha)\bar{y} = \alpha z + (1 - \alpha)\bar{z} + XU_m(\alpha s + (1 - \alpha)\bar{s})$. Comme $\partial J_\lambda(m)$ et \mathbb{R}^{k++} sont des ensembles convexes on en déduit que $X^\top(\alpha z + (1 - \alpha)\bar{z}) \in \partial J_\lambda(m)$ et $\alpha s + (1 - \alpha)\bar{s} \in \mathbb{R}^{k++}$ donc $\alpha y + (1 - \alpha)\bar{y} \in A_m$.

2) Nécessité. Soit $y \in A_m$ alors il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\text{schm}(\hat{\beta}) = m$. Par conséquent, $\hat{\beta} = U_m s$ pour un certain $s \in \mathbb{R}^{k++}$. Comme $\hat{\beta}$ est un élément de $\mathcal{S}(y)$ dont le schéma est m alors $X^\top(y - X\hat{\beta}) \in \partial J_\lambda(\hat{\beta}) = \partial J_\lambda(m)$. En multipliant cette inclusion par U_m^\top , grâce à (2.1), on obtient $\tilde{X}_m^\top(y - X\hat{\beta}) = \tilde{\lambda}_m$ et ainsi

$$\tilde{X}_m^\top y - \tilde{\lambda}_m = \tilde{X}_m^\top X\hat{\beta} = \tilde{X}_m^\top \tilde{X}_m s. \quad (2.9)$$

ce qui prouve la condition du schéma. On applique $\tilde{X}_m^{\top+}$ dans l'expression (2.9). Comme $X\hat{\beta} \in \text{im}(\tilde{X}_m)$ et que $\tilde{X}_m^{\top+}\tilde{X}_m^\top$ est la projection orthogonale sur $\text{im}(\tilde{X}_m)$, on en déduit que $(\tilde{X}_m^{\top+}\tilde{X}_m^\top)X\hat{\beta} = X\hat{\beta}$. Ainsi,

$$\tilde{X}_m^{\top+}\tilde{X}_m^\top y - \tilde{X}_m^{\top+}\tilde{\lambda}_m = X\hat{\beta}.$$

L'égalité précédente donne la condition sous-différentiel :

$$\begin{aligned} \partial J_\lambda(m) \ni X^\top(y - X\hat{\beta}) &= X^\top y - X^\top(\tilde{X}_m^{\top+}\tilde{X}_m^\top y - \tilde{X}_m^{\top+}\tilde{\lambda}_m) \\ &= X^\top \tilde{X}_m^{\top+}\tilde{\lambda}_m + X^\top \underbrace{(I_n - \tilde{X}_m^{\top+}\tilde{X}_m^\top)}_{=\tilde{X}_m\tilde{X}_m^{\top+}} y. \end{aligned}$$

Suffisance. Supposons que les conditions du schéma et sous-différentiel soient vraies. Alors, d'après la condition du schéma, il existe $s \in \mathbb{R}^{k++}$ tel que

$$\tilde{\lambda}_m = \tilde{X}_m^\top y - \tilde{X}_m^\top \tilde{X}_m s. \quad (2.10)$$

Montrons que $U_m s \in \mathcal{S}(y)$. Par définition de U_m , on a $\text{schm}(U_m s) = m$ donc $\partial J_\lambda(U_m s) = \partial J_\lambda(m)$. D'après la condition du sous-différentiel et en utilisant l'équation (2.10) on en déduit que

$$\begin{aligned} \partial J_\lambda(U_m s) &\ni X^\top(y - \tilde{X}_m^{\top+}\tilde{X}_m^\top y + \tilde{X}_m^{\top+}\tilde{\lambda}_m) \\ &\ni X^\top(y - \tilde{X}_m^{\top+}\tilde{X}_m^\top y + \tilde{X}_m^{\top+}(\tilde{X}_m y - \tilde{X}_m^\top \tilde{X}_m s)) \\ &\ni X^\top(y - XU_m s). \end{aligned}$$

Par conséquent $U_m s \in \mathcal{S}(y)$.

3) On pose $y = z + \tilde{X}_m s$ avec $z \in \mathbb{R}^n$ tel que $X^\top z \in \text{ri}(\partial J_\lambda(m))$ et $s \in \mathbb{R}^{k++}$ alors, d'après 1), $y \in A_m$. Montrons que pour $\epsilon \in \mathbb{R}^n$ ayant une norme suffisamment petite on a $y + \epsilon \in A_m$. On pose $\mathbb{R}^n = \text{im}(\tilde{X}_m)^\perp \oplus \text{im}(\tilde{X}_m) = \ker(U_m^\top X^\top) \oplus \text{im}(\tilde{X}_m)$. Via cette décomposition on peut écrire $\epsilon = \eta + \delta$ avec $\eta \in \ker(U_m^\top X^\top)$ et $\delta \in \text{im}(\tilde{X}_m)$ d'où $\delta = \tilde{X}_m \tilde{X}_m^+ \delta$. Ainsi $y + \epsilon = z + \eta + \tilde{X}_m (s + \tilde{X}_m^+ \delta)$. Comme, par construction, $\|\delta\|_2 \leq \|\epsilon\|_2$ et que \mathbb{R}^{k++} est ouvert pour $\|\epsilon\|_2$ suffisamment petit on a $s + \tilde{X}_m^+ \delta \in \mathbb{R}^{k++}$. Par ailleurs, comme $U_m^\top X^\top \eta = 0$ et que $\lambda \in \mathbb{R}^{p++}$, d'après l'assertion 2) de la Proposition 2.1 on a $X^\top \eta \in \overrightarrow{\text{aff}}(\partial J_\lambda(m))$. Ainsi, pour ϵ suffisamment petit on a $X^\top (z + \eta) \in \partial J_\lambda(m)$. Donc, d'après 1), $y + \epsilon \in A_m$.

□

Preuve de la Proposition 2.5

Démonstration. Supposons qu'il existe $y \in \tilde{X}_m \mathbb{R}^{k++}$ et qu'il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\text{schm}(\hat{\beta}) = m$. D'après la condition du sous-différentiel donnée au Théorème 2.1 on a $X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m + X^\top (I_n - \tilde{X}_m \tilde{X}_m^+) y \in \partial J_\lambda(m)$. Par ailleurs, comme $\tilde{X}_m \tilde{X}_m^+$ est la matrice de projection sur l'espace vectoriel $\text{im}(\tilde{X}_m)$ et que $y \in \text{im}(\tilde{X}_m)$ on en déduit que

$$\partial J_\lambda(m) \ni X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m + X^\top (I_n - \tilde{X}_m \tilde{X}_m^+) y = X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m.$$

Inversement, si $X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m \in \partial J_\lambda(m)$ alors $U_m^\top X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m = \tilde{X}_m^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m = \tilde{\lambda}_m$. Comme $\tilde{X}_m^\top \tilde{X}_m^{\top+}$ est la matrice de projection sur $\text{im}(\tilde{X}_m^\top)$ alors $\tilde{\lambda} \in \text{im}(\tilde{X}_m^\top) = \text{im}(\tilde{X}_m^\top \tilde{X}_m)$. On pose $\tilde{\lambda}_m = \tilde{X}_m^\top \tilde{X}_m v$ pour un certain $v \in \mathbb{R}^k$ et $y = \tilde{X}_m s$ avec $s \in \mathbb{R}^{k++}$ vérifiant $s_k > \|v\|_\infty$ et $s_i - s_{i+1} > 2\|v\|_\infty$ pour $i \in \{1, \dots, k-1\}$. Comme

$$\tilde{X}_m^\top y - \tilde{\lambda}_m = \tilde{X}_m^\top \tilde{X}_m (s - v)$$

et que $s - v \in \mathbb{R}^{k++}$ on en déduit que la condition du schéma est satisfaite : $\tilde{X}_m^\top y - \tilde{\lambda}_m \in \tilde{X}_m^\top \tilde{X}_m \mathbb{R}^{k++}$. Par ailleurs, comme $y \in \text{im}(\tilde{X}_m)$ et que $\tilde{X}_m \tilde{X}_m^+$ la matrice de projection sur $\text{im}(\tilde{X}_m)$ on en déduit

$$X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m + X^\top (I_n - \tilde{X}_m \tilde{X}_m^+) y = X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m$$

donc la condition du sous-différentiel est satisfaite. Ainsi, il existe $\hat{\beta} \in \mathcal{S}(y)$ tel que $\text{schm}(\hat{\beta}) = m$.

□

Preuve de la Proposition 2.6

Démonstration. 1) Si $\text{aff}(\partial J_\lambda(m)) \cap \text{im}(X^\top \tilde{X}_m) \neq \emptyset$ alors il existe $z \in \mathbb{R}^k$, où $k = \|m\|_\infty \geq 1$, tel que $X^\top \tilde{X}_m z \in \text{aff}(\partial J_\lambda(m))$. Comme $\lambda \in \mathbb{R}^{p++}$, d'après l'assertion 2 de la Proposition 2.1, on a $\tilde{\lambda}_m = U_m^\top X^\top \tilde{X}_m z = \tilde{X}_m^\top \tilde{X}_m^\top z$ donc $\tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top)$ qui établit la première assertion.

2) Si $\tilde{\lambda}_m \in \text{im}(\tilde{X}_m^\top)$ alors $X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m \in \text{aff}(\partial J_\lambda(m))$. En effet, comme $\tilde{X}_m^\top (\tilde{X}_m^\top)^+$ est la matrice de projection sur $\text{im}(\tilde{X}_m^\top)$ on a

$$U_m^\top X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m = \tilde{X}_m^\top (\tilde{X}_m^\top)^+ \tilde{\lambda}_m = \tilde{\lambda}_m.$$

De plus, comme $\text{im}(\tilde{X}_m^{\top+}) = \text{im}(\tilde{X}_m)$ on en déduit que $X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m \in \text{im}(X^\top \tilde{X}_m)$. Pour prouver que $X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m$ est l'unique point de l'intersection entre $\{v \in \mathbb{R}^p \mid U_m^\top v = \tilde{\lambda}_m\}$ et $\text{im}(X^\top \tilde{X}_m)$, montrons que $\text{im}(X^\top \tilde{X}_m) \cap \text{im}(U_m)^\perp = \{0\}$. En effet, si $w \in \text{im}(X^\top \tilde{X}_m) \cap \text{im}(U_m)^\perp$ alors $w = X^\top \tilde{X}_m z$ pour un certain $z \in \mathbb{R}^k$ et $U_m^\top w = 0$. Donc, $\tilde{X}_m^\top \tilde{X}_m z = 0$, par conséquent $\tilde{X}_m z = 0$ d'où $w = 0$.

□

Preuve de la Proposition 2.7

Démonstration. D'après le Théorème 2.1, A_m est un ensemble est convexe. Par ailleurs, comme $X^\top \tilde{X}_m^\top + \tilde{\lambda}_m \notin \partial J_\lambda(m)$, d'après la Proposition 2.5 le schéma de l'estimateur SLOPE dans le cas non-bruité n'est pas égal à m ; en d'autres termes on a $X\beta \notin A_m$. Supposons que ε soit défini sur l'espace probabilisé $(\Omega, \mathcal{F}, \mathbb{P})$ alors pour tout $\omega \in \Omega$ on a $X\beta + \varepsilon(\omega) \notin A_m$ ou $X\beta - \varepsilon(\omega) \notin A_m$; en effet si $X\beta + \varepsilon(\omega) \in A_m$ et $X\beta - \varepsilon(\omega) \in A_m$ alors la convexité de A_m impliquerait que $X\beta \in A_m$ ce qui contredit que $X\beta \notin A_m$. Enfin, comme ε et $-\varepsilon$ ont la même loi on en déduit l'inégalité suivante

$$1 = \mathbb{P}(X\beta + \varepsilon \notin A_m \cup X\beta - \varepsilon \notin A_m) \leq \mathbb{P}(X\beta + \varepsilon \notin A_m) + \mathbb{P}(X\beta - \varepsilon \notin A_m) = 2\mathbb{P}(X\beta + \varepsilon \notin A_m).$$

Donc $\mathbb{P}(Y \notin A_m) \geq 1/2$ d'où, $\mathbb{P}(Y \in A_m) = \mathbb{P}(\exists \hat{\beta} \in \mathcal{S}(Y), \text{schm}(\hat{\beta}) = m) \leq 1/2$. \square

Preuve de la Proposition 2.8

La formule de l'opérateur proximal donnée à la Proposition 2.8 est une conséquence immédiate du Lemme 2.1.

Lemme 2.1. *Soit $y \in \mathbb{R}^{p+}$, $\lambda \in \mathbb{R}^{p+} \setminus \{0\}$ et $(C_j)_{1 \leq j \leq p}$ la suite de Cesàro définie par $C_j = \frac{1}{j} \sum_{i=1}^j (y_i - \lambda_i)$. Pour simplifier les notations on note b^* l'opérateur proximal de la norme ℓ_1 ordonnée évalué en y : $b^* = \text{prox}_\lambda(y)$.*

Les assertions suivantes sont satisfaites :

1. On a $b^* \in \mathbb{R}^{p+}$.
2. On a $b^* = (0, \dots, 0)$ si et seulement si la suite de Cesàro est négative ou nulle.
3. Si $b_1^* > 0$ et $k = \max\{i \in \{1, \dots, p\} \mid b_i^* = b_1^*\}$ alors le plus grand entier pour lequel la suite de Cesàro atteint son maximum est k et $C_k = b_1^* = \dots = b_k^*$. Inversement si le plus grand entier pour lequel la suite de Cesàro atteint son maximum est k et $C_k > 0$ alors $b_1^* = \dots = b_k^* = C_k$ et $k = \max\{i \in \{1, \dots, p\} \mid b_i^* = b_1^*\}$.
4. Si $b_1^* > 0$ et $k = \max\{i \in \{1, \dots, p\} \mid b_i^* = b_1^*\}$ avec $k < p$ alors $(b_{k+1}^*, \dots, b_p^*) = \text{prox}_{\lambda_{k+1}, \dots, \lambda_p}(y_{k+1}, \dots, y_p)$.

Démonstration. 1) La preuve de cette assertion est similaire à celle donnée dans l'article de Bogdan *et al.* (2015). Soit $b \in \mathbb{R}^p$, montrons l'inégalité suivante :

$$\frac{1}{2} \|y - b\|_2^2 - J_\lambda(b) \geq \frac{1}{2} \|y - |b|_\downarrow\|_2^2 - J_\lambda(|b|_\downarrow). \quad (2.11)$$

Comme $J_\lambda(b) = J_\lambda(|b|_\downarrow)$ et que $\|b\|_2^2 = \||b|_\downarrow\|_2$ nous avons

$$\begin{aligned} & \frac{1}{2} \|y - b\|_2^2 + J_\lambda(b) \geq \frac{1}{2} \|y - |b|_\downarrow\|_2^2 + J_\lambda(|b|_\downarrow), \\ \Leftrightarrow & \|y - b\|_2^2 \geq \|y - |b|_\downarrow\|_2^2, \\ \Leftrightarrow & b^\top y \leq |b|_\downarrow^\top y. \end{aligned}$$

Clairement $b^\top y \leq |b|^\top y$ où $|b| = (|b_1|, \dots, |b_p|)$ et d'après l'inégalité de réarrangement on a $|b|^\top y \leq |b|_\downarrow^\top y$ ce qui établit (2.11). En appliquant l'inégalité (2.11) en b^* , unique minimiseur de l'expression $b \in \mathbb{R}^p \mapsto \frac{1}{2} \|y - b\|_2^2 - J_\lambda(b)$, on en déduit que $b^* = |b^*|_\downarrow$ d'où $b_1^* \geq \dots \geq b_p^* \geq 0$.

2) On a $b^* = 0$ si et seulement si $y - 0 \in \partial J_\lambda(0)$ si et seulement si $J_\lambda^*(y) \leq 1$ si et seulement si les inégalités

suyvantes sont satisfaites :

$$\forall j \in \{1, \dots, p\} \quad \sum_{i=1}^j y_i \leq \sum_{i=1}^j \lambda_i \Leftrightarrow \forall j \in \{1, \dots, p\} \quad C_j \leq 0.$$

3) Supposons que $b_1^* > 0$ et posons $k = \max\{i \in \{1, \dots, p\} \mid b_i^* = b_1^*\}$. On a $y - b^* \in \partial J_\lambda(m)$ où $m = \text{schm}(b^*)$. D'après la première équation de $U_m^\top(y - b^*) = \tilde{\lambda}_m$ on a

$$\sum_{i=1}^k (y_i - b_i^*) = \sum_{i=1}^k (y_i - b_1^*) = \sum_{i=1}^k \lambda_i \Rightarrow b_1^* = \frac{1}{k} \sum_{i=1}^k (y_i - \lambda_i) = C_k.$$

Comme $J_\lambda^*(y - b^*) \leq 1$ alors, pour tout $j \in \{1, \dots, p\}$ on a l'inégalité suivante

$$\sum_{i=1}^j y_i - b_i^* \leq \sum_{i=1}^j |y - b^*|_{\downarrow i} \leq \sum_{i=1}^j \lambda_i \Rightarrow \frac{\sum_{i=1}^j y_i - \lambda_i}{j} = C_j \leq \underbrace{\frac{\sum_{i=1}^j b_i^*}{j}}_{\text{car } b^* \in \mathbb{R}^{p+}}$$

Par ailleurs, si $k < p$ alors pour $j \in \{k+1, \dots, p\}$ on a $C_j < b_1^*$ (car $b_{k+1}^* < b_1^*$) ainsi, k est le plus grand entier pour lequel la suite de Cesàro atteint son maximum.

Inversement, soit k le plus grand entier pour lequel la suite de Cesàro atteint son maximum et supposons que $C_k > 0$. D'après 1), b^* a des composantes positives et décroissantes et d'après 2), comme $C_k > 0$, b^* est non-nul ainsi on en déduit que $b_1^* > 0$ et qu'il existe un entier $l \in \{1, \dots, p\}$ tel que $l = \max\{i \in \{1, \dots, p\} \mid b_i^* = b_1^*\}$. La première partie de la preuve de cette assertion montre que le plus grand entier pour lequel la suite de Cesàro atteint son maximum est l donc $k = l$ et ainsi $b_1^* = \dots = b_k^* = C_k$.

4) On pose $\epsilon = (b_k^* - b_{k+1}^*)/2 > 0$ alors pour $b \in \mathbb{R}^p$ tel que $\|b - b^*\|_\infty < \epsilon$ on a $\min\{|b_i| \mid i \in \{1, \dots, k\}\} > \max\{|b_i| \mid i \in \{k+1, \dots, p\}\}$ d'où

$$\|y - b\|_2^2 + J_\lambda(b) = \sum_{i=1}^k (y_i - b_i)^2 + J_{\lambda_1, \dots, \lambda_k}(b_1, \dots, b_k) + \sum_{i=k+1}^p (y_i - b_i)^2 + J_{\lambda_{k+1}, \dots, \lambda_p}(b_{k+1}, \dots, b_p).$$

Comme b^* est l'opérateur proximal de la norme ℓ_1 ordonnée en y donc, $(b_{k+1}^*, \dots, b_p^*)$ minimise localement la fonction $(b_{k+1}, \dots, b_p) \mapsto \sum_{i=k+1}^p (y_i - b_i)^2 + J_{\lambda_{k+1}, \dots, \lambda_p}(b_{k+1}, \dots, b_p)$ sur l'ensemble $\{b \in \mathbb{R}^p \mid \|b - b^*\|_\infty < \epsilon\}$. Cette fonction étant strictement convexe on en déduit que $(b_{k+1}^*, \dots, b_p^*) = \text{prox}_{\lambda_{k+1}, \dots, \lambda_p}(y_{k+1}, \dots, y_p)$. \square

Preuve du Théorème 2.2

Lemme 2.2. Soit $X \in \mathbb{R}^{n \times p}$. L'application $N : b \in \text{im}(X) \mapsto \min\{J_\lambda(z) \mid Xz = b\}$ est une norme sur $\text{im}(X)$.

Démonstration. N est définie positive : clairement, $N(0) = 0$. Réciproquement, si $N(b) = 0$ alors 0 est une solution du système linéaire $Xz = b$ impliquant que $b = 0$.

N est homogène : clairement, si $t = 0$ alors $0 = N(tb) = |t|N(b)$. Supposons maintenant que $t \neq 0$. Soit z_b une solution du système linéaire d'équation $Xz = b$ ayant une norme J_λ minimale (i.e. $J_\lambda(z_b) = N(b)$) alors tz_b est une solution de $Xz = tb$ impliquant que $N(tb) \leq J_\lambda(tz_b) = |t|N(b)$. Réciproquement, si z_{tb} est une solution du système linéaire d'équations $Xz = tb$ (i.e. $J_\lambda(z_{tb}) = N(tb)$) alors z_{tb}/t est une solution du système $Xz = b$ ce qui conduit à l'implication $N(b) \leq J_\lambda(z_{tb}/t) \Rightarrow |t|N(b) \leq J_\lambda(z_{tb}) = N(tb)$. Par conséquent, pour tout $t \in \mathbb{R}$ et pour tout $b \in \text{im}(X)$, $N(tb) = |t|N(b)$.

N satisfait l'inégalité triangulaire : soit z_b une solution du système $Xz = b$ ayant une norme J_λ minimale (i.e. $J_\lambda(z_b) = N(b)$) et $z_{b'}$ une solution du système $Xz = b'$ ayant une norme J_λ minimale (i.e. $J_\lambda(z_{b'}) = N(b')$). Comme $z_b + z_{b'}$ est une solution du système $Xz = b + b'$, on a $N(b + b') \leq J_\lambda(z_b + z_{b'}) \leq J_\lambda(z_b) + J_\lambda(z_{b'}) = N(b) + N(b')$. \square

Démonstration du Théorème 2.2. Dans un premier temps montrons que les estimateurs $\widehat{\beta}^{(r)}$ et $\widehat{\beta}^{0(r)}$ convergent vers β .

Estimateur $\widehat{\beta}^{(r)}$: Comme $\widehat{\beta}^{(r)}$ est une solution du problème (2.5) et en évaluant la fonction objective en $\widehat{\beta}^{0(r)}$ on obtient l'inégalité suivante :

$$\frac{1}{2} \|y^{(r)} - X\widehat{\beta}^{(r)}\|_2^2 + \gamma_r J_\lambda(\widehat{\beta}^{(r)}) \leq \frac{1}{2} \|y^{(r)} - X\widehat{\beta}^{0(r)}\|_2^2 + \gamma_r J_\lambda(\widehat{\beta}^{0(r)}). \quad (2.12)$$

Les affirmations suivantes sont immédiates :

- Par définition de $\widehat{\beta}^{0(r)}$ et comme XX^+ est la projection orthogonale sur $\text{im}(X)$ on a $\|y - X\widehat{\beta}^{(r)}\|_2^2 \geq \|y^{(r)} - XX^+y^{(r)}\|_2^2 = \|y^{(r)} - X\widehat{\beta}^{0(r)}\|_2^2$ d'où $J_\lambda(\widehat{\beta}^{(r)}) \leq J_\lambda(\widehat{\beta}^{0(r)})$.
- Par construction de $\widehat{\beta}^{0(r)}$ on a $J_\lambda(\widehat{\beta}^{0(r)}) = N(XX^+y^{(r)})$.
- Comme une norme est continue et que $y^{(r)}$ tend vers $X\beta$ alors $\lim_{r \rightarrow +\infty} N(XX^+y^{(r)}) = N(X\beta)$.
- On a $N(X\beta) = J_\lambda(\beta)$ (cette identité est immédiate lorsque $\ker(X) = \{0\}$ sinon, lorsque $\dim(\ker(X)) \geq 1$, cette identité découle de l'hypothèse faite sur β).

Par conséquent la limite supérieure de $J_\lambda(\widehat{\beta}^{(r)})$ satisfait l'inégalité

$$\limsup_{r \rightarrow +\infty} J_\lambda(\widehat{\beta}^{(r)}) \leq J_\lambda(\beta). \quad (2.13)$$

Ainsi, la suite $(\widehat{\beta}^{(r)})_{r \in \mathbb{N}}$ est bornée ; on note $l \in \mathbb{R}^p$ une valeur d'adhérence de cette suite. D'après (2.13), on a $J_\lambda(l) \leq J_\lambda(\beta)$. Par ailleurs, comme $\lim_{r \rightarrow 0} \frac{1}{2} \|y^{(r)} - X\widehat{\beta}^{0(r)}\|_2^2 + \gamma_r J_\lambda(\widehat{\beta}^{0(r)}) = 0$, d'après (2.12), on a

$$0 = \lim_{r \rightarrow +\infty} \frac{1}{2} \|y^{(r)} - X\widehat{\beta}^{(r)}\|_2^2 + \gamma_r J_\lambda(\widehat{\beta}^{(r)}) = \frac{1}{2} \|X\beta - Xl\|_2^2.$$

Comme $Xl = X\beta$ et $J_\lambda(l) \leq J_\lambda(\beta)$ on en déduit que $\beta = l$ (cette implication est immédiate lorsque $\ker(X) = \{0\}$ sinon, lorsque $\dim(\ker(X)) \geq 1$, cette implication découle de l'hypothèse faite sur β). Finalement $(\widehat{\beta}^{(r)})_{r \in \mathbb{N}}$ est une suite bornée ayant une unique valeur d'adhérence β d'où $\lim_{r \rightarrow +\infty} \widehat{\beta}^{(r)} = \beta$.

Estimateur $\widehat{\beta}^{0(r)}$: D'après les affirmations b, c et d on déduit l'inégalité suivante

$$\lim_{r \rightarrow +\infty} J_\lambda(\widehat{\beta}^{0(r)}) = J_\lambda(\beta). \quad (2.14)$$

Ainsi, la suite $(\widehat{\beta}^{0(r)})_{r \in \mathbb{N}}$ est bornée ; on note $l \in \mathbb{R}^p$ une valeur d'adhérence de cette suite. D'après (2.14), on a $J_\lambda(l) = J_\lambda(\beta)$. Par ailleurs, comme $X\widehat{\beta}^{0(r)} = XX^+y^{(r)}$, en prenant la limite de cette identité on en déduit que $Xl = X\beta$. De façon similaire que pour l'estimateur $\widehat{\beta}^{(r)}$ on en déduit que $\lim_{r \rightarrow +\infty} \widehat{\beta}^{0(r)} = \beta$.

Ci-après, la notation $\widetilde{\beta}^{(r)}$ représente indistinctement $\widehat{\beta}^{(r)}$ ou $\widehat{\beta}^{0(r)}$. Le schéma de l'opérateur proximal est caractérisé par les conditions du schéma et du sous-différentiel du Théorème 2.1 dans le cas particulier où $X = I_p$. On pose $\beta = U_m s$ avec $s \in \mathbb{R}^{k++}$ et $v^{(r)} = \widetilde{\beta}^{(r)} - \beta$. Examinons la condition du schéma lorsque $X = I_p$:

$$\widetilde{X}_m^\top \widetilde{\beta}^{(r)} - \tau \widetilde{\lambda}_m = U_m^\top \widetilde{\beta}^{(r)} - \tau \widetilde{\lambda}_m = U_m^\top (U_m s + v^{(r)}) - \tau \widetilde{\lambda}_m = U_m^\top U_m (s - \tau (U_m^\top U_m)^{-1} \widetilde{\lambda}_m) + (U_m^\top U_m)^{-1} U_m^\top v^{(r)}.$$

On pose $\tau \geq 0$ suffisamment petit pour que $s - \tau (U_m^\top U_m)^{-1} \widetilde{\lambda}_m \in \mathbb{R}^{k++}$. Montrons que pour r assez grand on a

$\text{schm}(\text{prox}_{\tau\lambda}(\tilde{\beta}^{(r)})) = m$. Comme $v^{(r)}$ tend vers 0 alors

$$\lim_{r \rightarrow +\infty} s - \tau(U_m^\top U_m)^{-1} \tilde{\lambda} + (U_m^\top U_m)^{-1} U_m^\top v^{(r)} = s - \tau(U_m^\top U_m)^{-1} \tilde{\lambda} \in \mathbb{R}^{k++}.$$

Ainsi, il existe $r_0 \in \mathbb{N}$ tel que pour $r \geq r_0$ la condition du schéma est satisfaite : $U_m^\top \tilde{\beta}^{(r)} - \tau \tilde{\lambda}_m \in U_m^\top U_m \mathbb{R}^{k++}$. Examinons la condition du sous-différentiel lorsque $X = I_p$:

$$\begin{aligned} X^\top \tilde{X}_m^\top + \tilde{\lambda}_m + \frac{1}{\tau} X^\top (I_p - \tilde{X}_m \tilde{X}_m^\top) \tilde{\beta} &= U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m + \frac{1}{\tau} (I_p - U_m (U_m^\top U_m)^{-1} U_m^\top) (U_m s + v^{(r)}), \\ &= U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m + \frac{1}{\tau} (I_p - U_m (U_m^\top U_m)^{-1} U_m^\top) v^{(r)}. \end{aligned}$$

Remarquons que $U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m + \frac{1}{\tau} (I_p - U_m (U_m^\top U_m)^{-1} U_m^\top) v^{(r)} \in \text{aff}(\partial J_\lambda(m))$; en effet

$$U_m^\top (U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m + \frac{1}{\tau} (I_p - U_m (U_m^\top U_m)^{-1} U_m^\top) v^{(r)}) = \tilde{\lambda}_m.$$

Par ailleurs, comme $v^{(r)}$ tend vers 0 alors on a la limite suivante

$$\lim_{r \rightarrow +\infty} U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m + \frac{1}{\tau} (I_p - U_m (U_m^\top U_m)^{-1} U_m^\top) v^{(r)} = U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m.$$

Notons que $U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m \in \text{ri}(\partial J_\lambda(m))$ (voir la preuve de la Proposition 2.1 ; plus précisément $U_m (U_m^\top U_m)^{-1} \tilde{\lambda}_m$ est l'isobarycentre des points extrémaux de $\partial J_\lambda(m)$) ainsi il existe $r_1 \in \mathbb{N}$ tel que pour $r \geq r_1$ la condition du sous-différentiel soit satisfaite. Par conséquent, pour $r \geq \max\{r_0, r_1\}$ les conditions du schéma et du sous-différentiel sont satisfaites donc $\text{schm}(\text{prox}_{\tau\lambda}(\tilde{\beta}^{(r)})) = m$ ce qui achève la preuve. \square

Chapitre 3

Chemin des solutions et chemin des valeurs ajustées de l'estimateur SLOPE

3.1 Introduction

Les estimateurs pénalisés comme le LASSO, le LASSO généralisé ou SLOPE dépendent des paramètres de pénalité qui doivent être sélectionnés de façon appropriée pour que ces estimateurs aient de bonnes propriétés. Lorsque l'objectif est la récupération des coefficients de régression non nuls, des valeurs pour ces paramètres de pénalité peuvent être spécifiées a priori (voir, par exemple, (Bogdan *et al.*, 2015; Candès et Plan, 2009; Lounici, 2008; Tardivel et Bogdan, 2022)). Néanmoins, les conditions pour déterminer les coefficients de régression non nuls sont très contraignantes, par exemple, que la matrice de régression ait des colonnes orthogonales ou que la condition d'irreprésentabilité du LASSO soit satisfaite. Un objectif différent est de choisir les paramètres de pénalité de façon à construire un estimateur facile à interpréter avec peu de coefficients non nuls (ou peu de groupes d'appariement) et ayant une erreur quadratique moyenne ou une erreur de prédiction¹ faible. Pour un tel objectif, les paramètres de pénalité ne sont spécifiés a priori et sont calculés en fonction des données. Cette approche nécessite donc de calculer le chemin des solutions; c'est-à-dire la fonction qui à un paramètre de pénalité associe l'ensemble des solutions du problème pénalisé. Les premiers travaux sur les chemins des solutions traitent des estimateurs LASSO et LASSO généralisé (voir, par exemple (Rosset et Zhu, 2007; Mairal et Yu, 2012; Tibshirani et Taylor, 2011; Arnold et Tibshirani, 2016)). Soit $D \in \mathbb{R}^{m \times p}$, le chemin des solutions du LASSO généralisé (resp. du LASSO lorsque $D = I_p$) est l'application suivante :

$$\gamma > 0 \mapsto \arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \gamma \|Db\|_1 \right\}.$$

Le paramètre de pénalité peut être choisi a posteriori en minimisant un critère comme, par exemple, la formule SURE (acronyme signifiant « Stein Unbiased Risk Estimate ») ou la somme des carrés résiduels sur un échantillon de validation (Bertrand *et al.*, 2020, 2022). Récemment quelques pré-publications ou articles ont traité du chemin des solutions des estimateurs OSCAR ou SLOPE (Dupuis et Tardivel, 2024; Gu *et al.*, 2017; Nomura, 2020; Takahashi et Nomura, 2020). Soit $\bar{\lambda} \in \mathbb{R}^{p+} \setminus \{0\}$, le chemin unidimensionnel des solutions du SLOPE (resp. d'OSCAR lorsque $\bar{\lambda}$ est une suite arithmétique) est l'application suivante :

$$\gamma > 0 \mapsto \arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \gamma J_{\bar{\lambda}}(b) \right\}. \quad (3.1)$$

1. Pour un estimateur $\hat{\beta}$ d'un paramètre β , l'erreur quadratique moyenne et l'erreur de prédiction sont respectivement $\mathbb{E}(\|\hat{\beta} - \beta\|_2^2)$ et $\mathbb{E}(\|X\hat{\beta} - X\beta\|_2^2)$.

Cette approche du chemin des solutions de l'estimateur SLOPE où $\bar{\lambda}$ est un paramètre de pénalité fixé et γ est le paramètre scalaire qui varie dans l'intervalle sur $]0, +\infty[$ est calquée sur les chemins unidimensionnels des estimateurs LASSO et LASSO généralisé. Comme l'estimateur SLOPE dépend d'un paramètre de pénalité de dimension p , il est plus adapté de considérer le chemin des solutions du SLOPE suivant :

$$\lambda \in \mathbb{R}^{p^+} \setminus \{0\} \mapsto \underbrace{\arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + J_\lambda(b) \right\}}_{=\mathcal{S}_\lambda}.$$

L'approche multidimensionnelle du chemin des solutions du SLOPE a néanmoins un inconvénient technique : dès que $\dim(\ker(X)) \geq 1$ on ne peut pas espérer que l'ensemble des solutions \mathcal{S}_λ soit réduit à un singleton pour tout $\lambda \in \mathbb{R}^{p^+} \setminus \{0\}$. Un exemple de chemin des solutions du SLOPE où l'ensemble \mathcal{S}_λ n'est pas réduit à un singleton pour certains paramètres de pénalité $\lambda \in \mathbb{R}^{p^+} \setminus \{0\}$ est donné ci-dessous.

Exemple 3.1. Soit $X = \begin{pmatrix} 3 & 2 \end{pmatrix}$. D'après la Proposition 4 du chapitre 2, il existe $y \in \mathbb{R}$ pour lequel l'ensemble des solutions \mathcal{S}_λ n'est pas réduit à un singleton dès que le paramètre de pénalité $\lambda \in \mathbb{R}^{2^+} \setminus \{0\}$ est proportionnel à $(3, 2)$. Pour $y = 1$, l'expression suivante donne un élément de \mathcal{S}_λ pour $\lambda \in \mathbb{R}^{2^+} \setminus \{0\}$:

$$\begin{cases} (0, 0) \in \mathcal{S}_\lambda & \text{dès que } \lambda_1 \geq 3 \text{ et } \lambda_1 + \lambda_2 \geq 5, \\ \left(\frac{3-\lambda_1}{9}, 0\right) \in \mathcal{S}_\lambda & \text{dès que } \lambda_1 \leq 3 \text{ et } 2\lambda_1 \leq 3\lambda_2, \\ \left(\frac{5-\lambda_1-\lambda_2}{25}, \frac{5-\lambda_1-\lambda_2}{25}\right) \in \mathcal{S}_\lambda & \text{dès que } 2\lambda_1 \geq 3\lambda_2 \text{ et } \lambda_1 + \lambda_2 \leq 5. \end{cases}$$

On peut ainsi remarquer que \mathcal{S}_λ n'est pas un singleton dès que λ appartient au segment $] (0, 0), (3, 2)[$. Par exemple, lorsque $\lambda = (3/2, 1)$, l'ensemble \mathcal{S}_λ est le segment $] (1/6, 0), (1/10, 1/10)[$.

Une façon de s'affranchir de ce problème d'unicité est de considérer le chemin des valeurs ajustées du SLOPE.

Proposition 3.1. Soit $X \in \mathbb{R}^{n \times p}$ et $\lambda \in \mathbb{R}^{p^+}$. Si $\hat{\beta}$ et $\bar{\beta}$ sont deux éléments de \mathcal{S}_λ alors $X\hat{\beta} = X\bar{\beta}$

Cette Proposition montre que la notion de valeur ajustée du SLOPE $X\hat{\beta}$, où $\hat{\beta}$ est un élément arbitraire de \mathcal{S}_λ , est bien définie. Dans la suite de ce chapitre on notera $\hat{v}_\lambda(\lambda)$ la valeur ajustée du SLOPE. La norme ℓ_1 ordonnée ne joue pas de rôle particulier pour établir la Proposition 3.1 ; en particulier ce résultat sera prouvé à l'annexe A dans un cadre plus général.

3.2 Chemin multidimensionnel des solutions et valeurs ajustées du SLOPE

Dans cette section on s'intéresse à la fonction qui à $\lambda \in \mathbb{R}^{p^+}$ associe la solution ou la valeur ajustée du problème SLOPE.

Théorème 3.1. Soient $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, $\lambda \in \mathbb{R}^{p^+}$ et $m \in \mathcal{P}_p^{\text{slope}} \setminus \{0\}$ où $k = \|m\|_\infty \geq 1$.

1. L'ensemble $E_m = \{\lambda \in \mathbb{R}^{p^+} \mid \exists \hat{\beta} \in \mathcal{S}_\lambda \text{ tel que } \text{schem}(\hat{\beta}) = m\}$ est convexe (plus précisément un polyèdre) et admet la caractérisation suivante :

$$\lambda \in E_m \Leftrightarrow \begin{cases} \exists s \in \mathbb{R}^{k^{++}} \text{ tel que } \tilde{X}_m^\top y - \tilde{\lambda}_m = \tilde{X}_m^\top \tilde{X}_m s & \text{(condition du schéma),} \\ X^\top \tilde{X}_m^\top + \tilde{\lambda}_m + X^\top (I_n - \tilde{X}_m^\top + \tilde{X}_m^\top) y \in \partial J_\lambda(m) & \text{(condition du sous-différentiel).} \end{cases}$$

2. Pour $\lambda = 0$ on pose, par convention, $\hat{v}_\lambda(\lambda) = \|y - XX^+y\|_2^2$. Le chemin des valeurs ajustées $\lambda \in \mathbb{R}^{p^+} \mapsto$

$\widehat{\text{v}}\text{a}(\lambda)$ est continu sur \mathbb{R}^{p+} avec l'expression affine sur E_m :

$$\widehat{\text{v}}\text{a}(\lambda) = \widetilde{X}_m^\top + \widetilde{X}_m^\top y - \widetilde{X}_m^\top + \widetilde{\lambda}_m, \quad \lambda \in E_m.$$

3. Si $\ker(X) = \{0\}$ on pose, par convention, $\widehat{\beta}(\lambda) = (X^\top X)^{-1} X^\top y$ si $\lambda = 0$. Le chemin des solutions $\lambda \in \mathbb{R}^{p+} \mapsto \widehat{\beta}(\lambda)$ est continue sur \mathbb{R}^{p+} avec l'expression affine suivante sur E_m :

$$\widehat{\beta}(\lambda) = U_m (\widetilde{X}_m^\top \widetilde{X}_m)^{-1} (\widetilde{X}_m^\top y - \widetilde{\lambda}_m), \quad \lambda \in E_m.$$

Le Théorème 3.1 ne caractérise pas E_m lorsque $m = 0$. Néanmoins, cet ensemble est un polyèdre dont les équations sont faciles à décrire. En effet

$$\begin{aligned} E_0 &= \{ \lambda \in \mathbb{R}^{p+} \setminus \{0\} \mid \exists \widehat{\beta} \in \mathcal{S}_\lambda \text{ tel que } \widehat{\beta} = 0 \}, \\ &= \{ \lambda \in \mathbb{R}^{p+} \setminus \{0\} \mid J_\lambda(X^\top y) \leq 1 \}, \\ &= \left\{ \lambda \in \mathbb{R}^{p+} \setminus \{0\} \mid \sum_{i=1}^j \lambda_i \geq \sum_{i=1}^j |X^\top y|_{\downarrow i} \quad \forall j \in \{1, \dots, p\} \right\}. \end{aligned}$$

Exemple 3.2. On pose $X = \begin{pmatrix} 3 & 2 \end{pmatrix}$ et $y = 1$.

3.2.1 Chemin du gradient et groupes d'appariement

Une solution du problème d'optimisation SLOPE est caractérisée par les deux conditions suivantes

$$\widehat{\beta} \in \mathcal{S}_\lambda \Leftrightarrow \begin{cases} J_\lambda^*(X^\top(y - X\widehat{\beta})) \leq 1 \\ \widehat{\beta}^\top X^\top(y - X\widehat{\beta}) = J_\lambda(\widehat{\beta}) \end{cases}$$

On peut remarquer $X^\top(y - X\widehat{\beta}) = X^\top(y - \widehat{\text{v}}\text{a}(\lambda))$ est l'opposé du gradient en $\widehat{\beta}$ de la somme des carrés résiduels $b \mapsto \frac{1}{2} \|y - Xb\|_2^2$. Par la suite, nous appelons chemin du gradient l'expression $\lambda \in \mathbb{R}^{p+} \mapsto X^\top(y - \widehat{\text{v}}\text{a}(\lambda))$. Afin de construire l'ensemble $\mathcal{A}(\gamma)$, qui joue un rôle proéminent au Théorème 3.2, nous introduisons $\|\cdot\|_{(i)}$ la i -norme qui est la somme des i plus grandes composantes d'un vecteur en valeur absolue (voir par exemple l'article de Gaudio et al. (2020) illustrant l'intérêt de cette norme en apprentissage statistique). Cette norme est un cas particulier de la norme ℓ_1 ordonnée lorsque $\lambda_1 = \dots = \lambda_i = 1$ et $\lambda_{i+1} = \dots = \lambda_p = 0$. Comme $X^\top(y - \widehat{\text{v}}\text{a}(\lambda))$ est un élément du permutoèdre signé on a $\|X^\top(y - \widehat{\text{v}}\text{a}(\lambda))\|_{(i)} \leq \sum_{j=1}^i \lambda_j$ pour tout $i \in \{1, \dots, p\}$; de plus l'ensemble des inégalités saturées par le gradient est :

$$\mathcal{A}(\lambda) = \left\{ i \in \{1, \dots, p\} \mid \frac{\|X^\top(y - \widehat{\text{v}}\text{a}(\lambda))\|_{(i)}}{\sum_{j=1}^i \lambda_j} = 1 \right\}.$$

D'après le Théorème 3.2 ci-dessous, l'ensemble $\mathcal{A}(\lambda)$ fournit à la fois le nombre de groupes d'appariement non nuls, la taille de ces groupes ainsi que le nombre de composantes non nulles .

Théorème 3.2. Soit $\lambda \in \mathbb{R}^{p+} \setminus \{0\}$, $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, > 0 et $\widehat{\beta} \in \mathcal{S}_\lambda$.

1. Soit $1 \leq k_1 \leq \dots \leq k_l \leq p$ une subdivision telle que :

$$\text{Card}(\text{supp}(\widehat{\beta})) = k_l \text{ et } |\widehat{\beta}|_{\downarrow 1} = \dots = |\widehat{\beta}|_{\downarrow k_1} > \dots > |\widehat{\beta}|_{\downarrow k_l - 1 + 1} = \dots = |\widehat{\beta}|_{\downarrow k_l} > 0$$

alors $\{k_1, \dots, k_l\} \subseteq \mathcal{A}(\lambda)$.

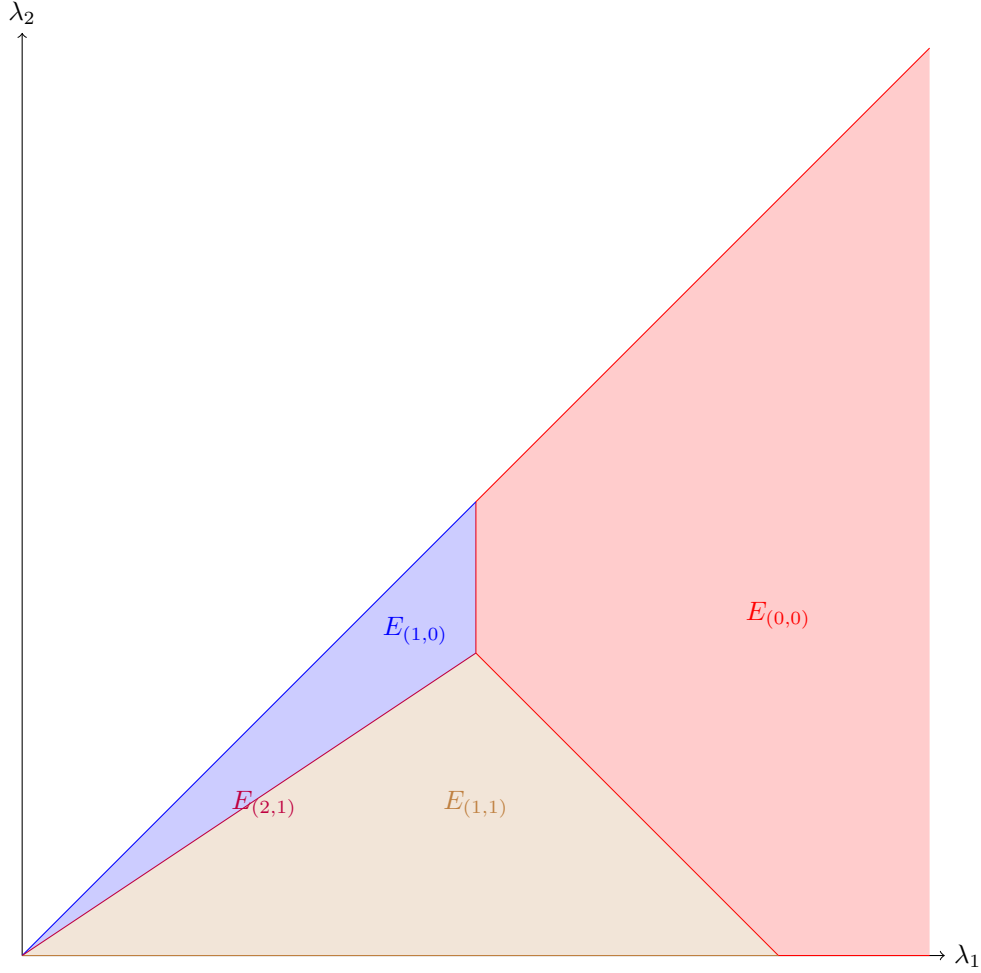


FIGURE 3.1 – Cette figure représente graphiquement le recouvrement de $\mathbb{R}^{2+} \setminus \{0\}$ par les sous-ensembles $E_{(0,0)}$ (en rouge), $E_{(1,0)}$ (en bleu), $E_{(1,1)}$ (en marron) et $E_{(2,1)}$ (en pourpre). On peut remarquer que les ensembles $E_{(1,0)}$, $E_{(1,1)}$ et $E_{(2,1)}$ ne sont pas disjoints. En effet, $E_{(1,0)} \cap E_{(1,1)} = E_{(2,1)}$.

2. Inversement, si $\{k_1, \dots, k_l\} = \mathcal{A}(\lambda)$ alors

$$|\widehat{\beta}|_{\downarrow 1} = \dots = |\widehat{\beta}|_{\downarrow k_1} \geq \dots \geq |\widehat{\beta}|_{\downarrow k_{l-1}+1} = \dots = |\widehat{\beta}|_{\downarrow k_l} \geq |\widehat{\beta}|_{\downarrow k_l+1} = \dots = |\widehat{\beta}|_{\downarrow p} = 0$$

Il existe des liens entre le Théorème 3.2 et les méthodes de dépistage de composantes nulles de l'estimateur SLOPE (Elvira et Herzet, 2023; Larsson *et al.*, 2020). Par exemple, l'exécution de l'algorithme 1 dans l'article de Larsson *et al.* (2020) avec $|X^\top(y - \widehat{v}_a(\lambda))|_{\downarrow}$ renvoie que la solution de l'estimateur SLOPE a au plus $\max\{\mathcal{A}(\lambda)\}$ composantes non nulles. Par ailleurs, le théorème 4.1 de l'article de Elvira et Herzet (2023) est étroitement lié à l'implication suivante : $|\widehat{\beta}(\lambda)|_{\downarrow i} \neq 0 \Rightarrow \exists k \geq i, k \in \mathcal{A}(\lambda)$.

3.3 Calcul du chemin unidimensionnel de l'estimateur SLOPE

À présent nous nous restreignons au chemin unidimensionnel des solutions du SLOPE qui est l'application multivaluée suivante :

$$\gamma > 0 \mapsto \arg \min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \gamma J_{\bar{\lambda}}(b) \right\}, \text{ où } \bar{\lambda} \in \mathbb{R}^{p+} \setminus \{0\}. \quad (3.2)$$

Une question ouverte pour l'auteur est l'existence d'une sélection continue du chemin des solutions du SLOPE. Une façon de contourner cette difficulté est de supposer que $\text{im}(X^\top)$ ne coupe pas une face du permutoèdre signé P_λ^\pm dont la dimension est strictement inférieure à $\dim(\ker(X))$. Cette hypothèse, satisfaite de manière générique d'après la proposition 3 de l'article de Schneider et Tardivel (2022), garanti que pour tout $y \in \mathbb{R}^n$, pour tout $\gamma > 0$ le problème d'optimisation donné dans l'expression (3.2) a une unique solution que l'on notera $\widehat{\beta}_\gamma$. Il convient de mentionner que le chemin unidimensionnel des solutions du SLOPE n'est pas un raffinement du chemin des solutions du LASSO généralisé. L'argument le plus immédiat étant que l'estimateur SLOPE n'est pas un estimateur LASSO généralisé dès que le paramètre de pénalité $\bar{\lambda} \in \mathbb{R}^{p+}$ n'est pas arithmétique; une preuve de cette affirmation est donnée en annexe de ce chapitre. Par ailleurs, même pour l'estimateur OSCAR où $\bar{\lambda} \in \mathbb{R}^{p+} \setminus \{0\}$ est arithmétique, qui est un estimateur SLOPE ou LASSO généralisé particulier, la méthode employée dans cette section pour le calcul du chemin des solutions du problème (3.1) est très différente de celle utilisée par Arnold et Tibshirani (2016). Par exemple, contrairement à l'article d'Arnold et Tibshirani (2016), l'hypothèse $\ker(X) = \{0\}$ n'est pas requise dans cette section pour le calcul de ce chemin; d'autres différences techniques sont également mentionnées dans l'article de Dupuis et Tardivel (2024).

Comme E_m est convexe on vérifie que l'ensemble $I_m = \{\gamma > 0 \mid \gamma\bar{\lambda} \in E_m\}$ est un intervalle. Ainsi, les intervalles non vides de la famille $(I_m)_{m \in \mathcal{P}_p^{\text{slope}}}$ forment une partition de $]0, \infty[$. Clairement $\widehat{\beta}_\gamma = 0$ si et seulement si $\gamma \geq J_{\bar{\lambda}}(X^\top y)$ ainsi, le calcul du chemin des solutions est pertinent pour $\gamma < J_{\bar{\lambda}}(X^\top y)$. Posons $J_{\bar{\lambda}}^*(X^\top y) = \gamma_0 > \gamma_1 > \dots > \gamma_r > \gamma_{r+1} = 0$ une subdivision telle que $\gamma \mapsto \widehat{\beta}_\gamma$ soit affine et de schéma $m^{(i)} \in \mathcal{P}_p^{\text{slope}}$ sur l'intervalle $]\gamma_{i+1}, \gamma_i[$ pour $i = 0, \dots, r$ (c'est-à-dire que l'intérieur de $I_{m^{(i)}}$ est $]\gamma_{i+1}, \gamma_i[$). Le but de l'algorithme suivant est de déterminer itérativement les noeuds du chemin des solutions du SLOPE ainsi que les expressions affines de $\widehat{\beta}_\gamma$ sur les intervalles $I_{m^{(0)}}, \dots, I_{m^{(r)}}$. A présent, avant de donner plus de détails sur l'algorithme nous supposons que $\bar{\lambda} \in \mathbb{R}^{p++}$ de tel sorte qu'il y ait une bijection entre les schémas du SLOPE et les faces du permutoèdre signé. Commençons par expliquer comment calculer le chemin des solutions de l'estimateur SLOPE sur $[\gamma_1, \gamma_0]$. Par construction de $m^{(0)}$, l'implication suivante est vérifiée

$$\forall \gamma \in]\gamma_1, \gamma_0[\quad \text{schm}(\widehat{\beta}_\gamma) = m^{(0)} \Rightarrow \frac{1}{\gamma} X^\top (y - X\widehat{\beta}_\gamma) \in \partial J_{\bar{\lambda}}(m^{(0)}).$$

De plus, comme $\gamma > 0 \mapsto X\widehat{\beta}_\gamma$ est une application continue, que $X\widehat{\beta}_{\gamma_0} = 0$ et que $\partial J_{\bar{\lambda}}(m^{(0)})$ est un ensemble fermé, nous obtenons

$$\frac{1}{\gamma_0} X^\top (y - X\widehat{\beta}_{\gamma_0}) = \frac{1}{\gamma_0} X^\top y \in \partial J_{\bar{\lambda}}(m^{(0)}).$$

L'algorithme 1 fournit le schéma $f(\frac{1}{\gamma_0} X^\top y)$ de la plus petite face du permutoèdre signé contenant $\frac{1}{\gamma_0} X^\top y$.

Algorithme 1 : schéma de la plus petite face contenant un vecteur

Entrées : $\lambda \in \mathbb{R}^{p++}$ et $z \in \mathbb{R}^p$ tel que $J_\lambda^*(z) \leq 1$

début

 Définir l'ensemble des inégalités saturées comme suit :

$$\mathcal{A}(z) = \left\{ i \in \{1, \dots, p\} \mid \frac{\|z\|_{(i)}}{\sum_{j=1}^i \lambda_j} = 1 \right\}.$$

si $\mathcal{A}(z) = \emptyset$ **alors**

$f(z) = (0, \dots, 0) \in \mathbb{R}^p$

sinon

$\forall j \in \{1, \dots, p\} \quad f_j(z) = \text{signe}(z_j) \sum_{i \in \mathcal{A}(z)} \mathbf{1}(|z_j| \geq \lambda_i).$

Sorties : $f(z)$

Exemple 3.3. Une illustration du calcul du chemin des solutions du SLOPE dans le voisinage du plus grand noeud est donné lorsque $y = (6, 2) \in \mathbb{R}^2$, $\bar{\lambda} = (4, 2) \in \mathbb{R}^{2++}$, et $X \in \mathbb{R}^{2 \times 2}$ est la matrice donnée ci-dessous

$$X = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}.$$

Plus grand noeud γ_0 : Comme $X^\top y = (7, 5)$ donc $\gamma_0 = J_{\bar{\lambda}}^*(X^\top y) = 2$.

Schéma $m^{(0)}$ dans le voisinage à gauche de γ_0 : Comme $\frac{1}{\gamma_0} X^\top y = (3.5, 2.5)$ est dans l'intérieur relatif du permutoèdre $\partial J_{\bar{\lambda}}(1, 1) = \text{conv}\{(4, 2), (2, 4)\}$, alors $m^{(0)} = f(\frac{1}{\gamma_0} X^\top y) = (1, 1)$.

Expression affine de $\hat{\beta}(\gamma)$ dans le voisinage à gauche de γ_0 : D'après l'assertion 3 du Théorème 3.1, lorsque $\gamma < \gamma_0 = 2$ est suffisamment proche de γ_0 , nous avons $\hat{\beta}(\gamma) = (\frac{8-4\gamma}{3}, \frac{8-4\gamma}{3})$.

L'algorithme 2 utilise la caractérisation de $I_{m^{(0)}}$, basée sur les conditions du schéma et du sous-différentiel, pour déterminer le noeud γ_1 ainsi que le schéma $m^{(1)}$. Partant d'un noeud γ_i le noeud γ_{i+1} est déduit de la condition du schéma, lorsque $\gamma_{i+1} = \gamma_{\text{schm}}$, ou de la condition du sous-différentiel, lorsque $\gamma_{i+1} > \gamma_{\text{schm}}$.

Algorithme 2 : Calcul itératif des noeuds et schémas

Données : $X \in \mathbb{R}^{n \times p}$, $\bar{\lambda} \in \mathbb{R}^{p++}$, $\gamma_i > 0$ et $m^{(i)} \in \mathcal{P}_p^{\text{slope}}$

début

└ on pose $k = \|m^{(i)}\|_\infty$
└ on calcule $s(\gamma) = (\tilde{X}_{m^{(i)}}^\top \tilde{X}_{m^{(i)}})^{-1} (\tilde{X}_{m^{(i)}}^\top y - \gamma \tilde{\lambda}_{m^{(i)}})$

si $s(\gamma) \in \mathbb{R}^{k+}$ **pour tout** $\gamma \in [0, \gamma_i]$ **alors**

└ on pose $\gamma_{\text{schm}} = 0$

sinon

└ on pose $\gamma_{\text{schm}} = \sup\{\gamma \in [0, \gamma_i] \mid s(\gamma) \notin \mathbb{R}^{k+}\}$.

si $X^\top (y - \tilde{X}_{m^{(i)}} s(\gamma)) \in \gamma \partial J_{\bar{\lambda}}(m^{(i)})$ **pour tout** $\gamma \in [\gamma_{\text{schm}}, \gamma_i]$ **alors**

└ on pose $\gamma_{i+1} = \gamma_{\text{schm}}$
└ on calcule $m^{(i+1)} = \text{schem}(U_{m^{(i)}} s(\gamma_{\text{schm}}))$

└ **Résultat :** $\gamma_{i+1}, m^{(i+1)}$

sinon

└ on pose $\gamma_{i+1} = \sup\{\gamma \in [\gamma_{\text{schm}}, \gamma_i] \mid X^\top (y - \tilde{X}_{m^{(i)}} s(\gamma)) \notin \gamma \partial J_{\bar{\lambda}}(m^{(i)})\}$
└ on calcule $m^{(i+1)} = f(\frac{1}{\gamma_{i+1}} X^\top (y - \tilde{X}_{m^{(i)}} s(\gamma_{i+1})))$ avec l'algorithme 1.

└ **Résultat :** $\gamma_{i+1}, m^{(i+1)}$

L'utilisation itérative de l'algorithme 2, jusqu'à ce que la valeur de γ_{i+1} soit nulle, permet de calculer entièrement le chemin des solutions de l'estimateur SLOPE. En particulier, la Figure 3.2 complète le chemin des solutions du SLOPE de l'Exemple 3.3.

3.4 Expériences numériques

Pour cette expérience numérique², nous utilisons le jeu de données réelles *Boston Housing* introduit par Harrison Jr et Rubinfeld (1978). Via un modèle de régression linéaire, les auteurs de cet article analysent le prix médian d'une maison dans un quartier de Boston en fonction de variables explicatives liées au quartier, telles que le taux de criminalité par habitant, le nombre moyen de pièces par logement, etc. Les notations $y \in \mathbb{R}^{506}$ et $X \in \mathbb{R}^{506 \times 13}$ représentent respectivement la réponse et la matrice de régression de ce modèle; les colonnes de la matrice X sont centrées et standardisées (pour tout $j \in \{1, \dots, 13\}$, $\sum_{i=1}^{506} X_{ij} = 0$ et $\sum_{i=1}^{506} X_{ij}^2 = 506$), et le vecteur y est centré ($\sum_{i=1}^{506} y_i = 0$). Pour l'analyse statistique, ces données sont divisées en un échantillon

2. Le code des expériences numériques de ce chapitre est disponible en ligne sur ma page internet

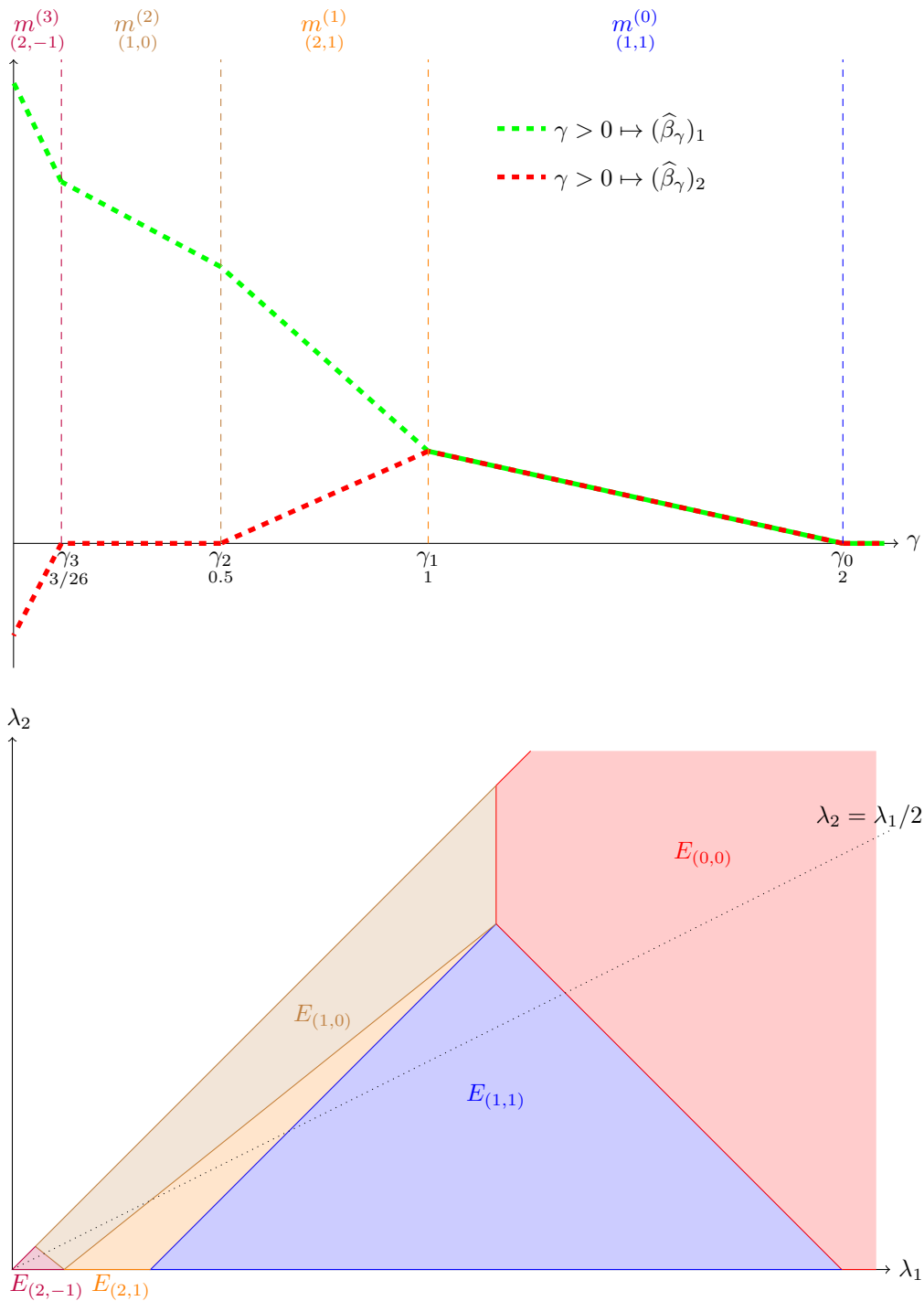


FIGURE 3.2 – La figure du haut représente le chemin des solutions de l’estimateur SLOPE en fonction du paramètre de régularisation $\gamma > 0$. La figure du bas représente la partition de $\mathbb{R}^{2+} \setminus \{0\}$ par les sous-ensembles $E_{(0,0)}$ (en rouge), $E_{(1,1)}$ (en bleu), $E_{(1,0)}$ (en marron), $E_{(2,1)}$ (en orange) et $E_{(2,-1)}$ (en pourpre). Par ailleurs, les régions coupées par la demi-droite d’équation $\lambda_2 = \lambda_1/2$ correspondent aux schémas observés sur le chemin des solutions du SLOPE.

d’apprentissage $X^{\text{app}}, y^{\text{app}}$ de taille 400 sur lequel le chemin des solutions est calculé et un échantillon de validation $X^{\text{val}}, y^{\text{val}}$ de taille 106 permettant de choisir le paramètre de régularisation γ .

3.4.1 Calcul des chemins et minimisation de la somme des carrés résiduels sur l'échantillon de validation

Ci-après, nous illustrons le chemin des solutions des estimateurs SLOPE et LASSO sur l'échantillon d'apprentissage du jeu de données *Boston Housing*. Pour l'estimateur SLOPE, nous prenons $\bar{\lambda} = (1, \sqrt{2} - 1, \sqrt{3} - \sqrt{2}, \dots, \sqrt{13} - \sqrt{12})$ comme paramètre de pénalité, de telle sorte que la boule unité de la norme ℓ_1 ordonnée soit la plus sphérique possible (Nomura, 2020). La Figure 3.3 illustre le chemin des solutions de l'estimateur SLOPE ainsi que le chemin des solutions du LASSO $\gamma > 0 \mapsto \arg \min_{b \in \mathbb{R}^{13}} \left\{ \frac{1}{2} \|y^{\text{app}} - X^{\text{app}} b\|_2^2 + \gamma \|b\|_1 \right\}$ calculé via l'algorithme d'homothopie de Mairal et Yu (2012).

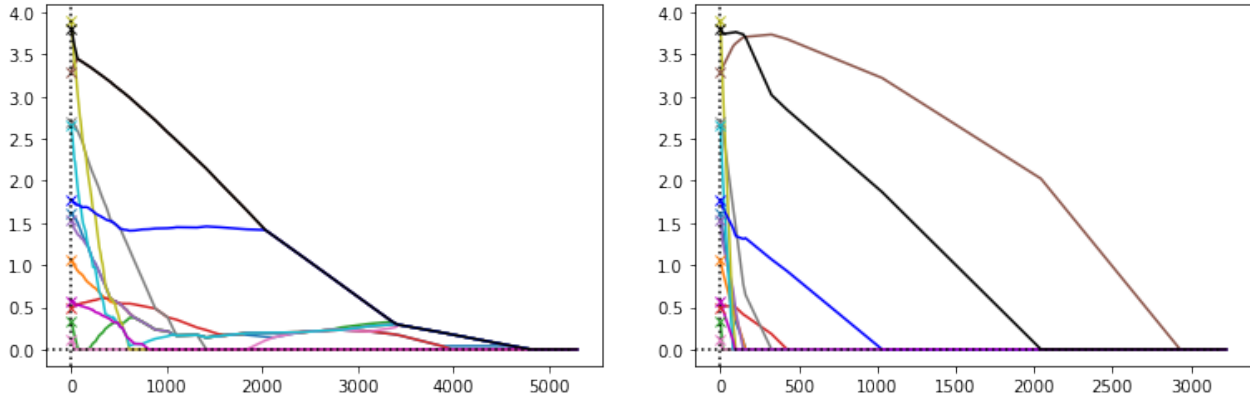


FIGURE 3.3 – Ces graphiques illustrent les chemins des solutions, en valeur absolue, du SLOPE (à gauche) et du LASSO (à droite) en fonction du paramètre de régularisation $\gamma > 0$. À gauche, on observe que certaines courbes se superposent ou coïncident partiellement avec l'axe des abscisses, illustrant les propriétés de appariement et de parcimonie de l'estimateur SLOPE. À droite, certaines courbes coïncident partiellement avec l'axe des abscisses, illustrant la propriété de parcimonie de l'estimateur LASSO. Ces courbes sont très similaires lorsque γ est petit car ces deux estimateurs convergent, lorsque γ tend vers 0, vers l'estimateur des moindres carrés, dont les composantes, en valeur absolue, sont représentées par des croix sur l'axe des ordonnées.

Le calcul du chemin des solutions des estimateurs SLOPE et LASSO sur l'échantillon d'apprentissage permet de choisir le paramètre γ en minimisant la somme des carrés résiduels sur l'échantillon de validation :

$$\text{scr}v : \gamma > 0 \mapsto \|y^{\text{val}} - X^{\text{val}} \widehat{\beta}_\gamma\|_2^2 \text{ où } \widehat{\beta}_\gamma \text{ représente indistinctement l'estimateur SLOPE ou LASSO.}$$

Indépendamment de l'estimateur SLOPE ou LASSO, on a $\text{scr}v(\gamma) = \|y^{\text{val}}\|_2^2 = 7895.05$ lorsque $\gamma > 0$ est suffisamment grand et $\lim_{\gamma \rightarrow 0} \text{scr}v(\gamma) = 2641.63$ (somme des carrés résiduels calculée pour l'échantillon de validation de l'estimateur des moindres carrés $(X^{\text{app}\top} X^{\text{app}})^{-1} X^{\text{app}\top} y^{\text{app}}$). Par ailleurs, la fonction $\text{scr}v$ étant quadratique entre deux nœuds adjacents, le minimum de cette expression est calculable de façon exacte. Le tableau 3.1 donne le minimum de la fonction $\text{scr}v$ pour les estimateurs SLOPE et LASSO ainsi que le paramètre $\gamma > 0$ pour lequel le minimum est atteint.

	$\min\{\text{scr}v(\gamma) \mid \gamma > 0\}$	$\arg \min_{\gamma > 0}\{\text{scr}v(\gamma)\}$
SLOPE	2159.05	357.39
LASSO	2498.52	42.61

TABLE 3.1 – Ce tableau illustre que lorsque le paramètre de régularisation est choisi de façon à minimiser la somme des carrés résiduels sur l'échantillon de validation, l'estimateur SLOPE est plus performant que l'estimateur LASSO.

L'estimateur SLOPE $\widehat{\beta}_\gamma$, pour $\gamma = 357.39$, a pour schéma $(-5, 4, -1, 4, -5, 8, 0, -7, 5, -3, -6, 2, -8)$. Par exemple, les variables explicatives associées au groupe d'appariement *huit* dont l'impact est le plus fort sur le

prix médian sont le nombre moyen de pièces par logement et le pourcentage de la population ayant un statut inférieur. La première variable a un impact positif sur le prix médian tandis que la seconde a un impact négatif; ce qui est très intuitif dans ce cas-là. Globalement, les composantes positives de l'estimateur SLOPE correspondent à des variables ayant un impact positif sur le prix médian (voir l'article Harrison Jr et Rubinfeld (1978) pour la liste de ces variables) et inversement pour les composantes négatives.

3.4.2 Limite de cet algorithme de calcul du chemin des solutions SLOPE

Ce chapitre propose un algorithme de calcul du chemin unidimensionnel des solutions de l'estimateur SLOPE et cette méthode est illustrée, avec succès, sur le jeu de données réelles *Boston Housing*. Néanmoins, il est facile de faire dysfonctionner cet algorithme sur des jeux de données plus volumineux; en particulier lorsque p est très large. De façon générale, au voisinage de l'origine, les noeuds tendent à être de plus en plus proches pouvant causer l'arrêt de l'algorithme. De même, les faces explorées par le gradient tendent à avoir une dimension de plus en plus faible menant à des problèmes de précisions numériques pour déterminer la face sur laquelle le gradient est localisé et faisant, in fine, dysfonctionner cette méthode. Il n'est pas rare que cet algorithme s'arrête lorsque γ est trop proche de 0 laissant une partie du chemin des solutions, souvent petite, non calculée. Enfin, donnons un exemple pathologique très simple pour lequel l'Algorithme 1 dysfonctionne. Reprenons l'Exemple 3.3 en substituant $\bar{\lambda} = (4, 2)$ par $\bar{\lambda} = X^\top y = (7, 5)$. Par construction, $\gamma_0 = J_{\bar{\lambda}}^*(X^\top y) = 1$ et $\frac{1}{\gamma_0} X^\top y = \partial J_{\bar{\lambda}}(2, 1)$ ainsi l'Algorithme 3.3 renvoie $m^{(0)} = (2, 1)$. Or, d'après le Figure 3.4 le schéma $m^{(0)}$ vaut $(1, 0)$; cette erreur sur le calcul du schéma $m^{(0)}$ empêche de calculer le chemin des solutions du SLOPE avec notre méthode.

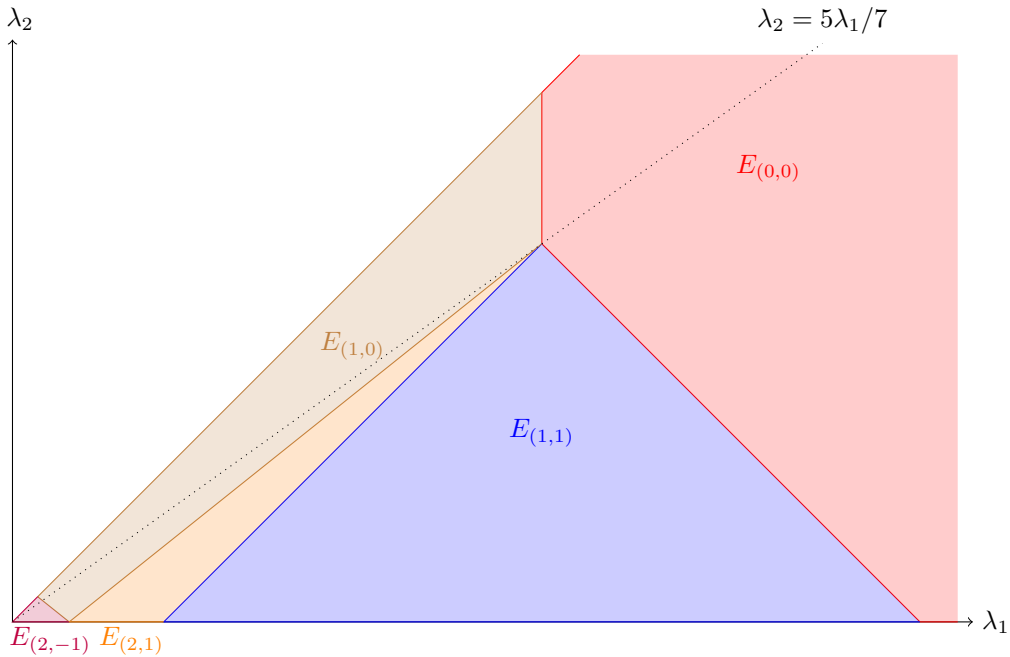


FIGURE 3.4 – Les régions coupées par la demi-droite d'équation $\lambda_2 = 5\lambda_1/7$ montre que l'estimateur SLOPE a pour schéma $m^{(0)} = (1, 0)$ lorsque $\gamma < 1$ est suffisamment proche de 1 et $m^{(1)} = (2, -1)$ lorsque γ est suffisamment proche de 0.

3.5 Annexes : preuves

Le Lemme 3.1 est utile pour montrer au Théorème 3.1 que l'ensemble $E_m = \{\lambda \in \mathbb{R}^{p^+} \mid \exists \hat{\beta} \in \mathcal{S}_\lambda \text{ tel que } \text{schm}(\hat{\beta}) = m\}$ est convexe

Lemme 3.1. Soient $\lambda \in \mathbb{R}^{p+}$, $\bar{\lambda} \in \mathbb{R}^{p+}$, $\alpha \geq 0$ et $\delta \geq 0$. La somme de Minkowski de sous-différentiels de norme ℓ_1 ordonnée satisfait l'identité suivante :

$$\alpha \partial J_\lambda(b) + \delta \partial J_{\bar{\lambda}}(b) = \partial J_{\alpha\lambda + \delta\bar{\lambda}}(b) \quad \forall b \in \mathbb{R}^p.$$

Démonstration. La preuve est une conséquence immédiate de l'identité :

$$\alpha J_\lambda(b) + \delta J_{\bar{\lambda}}(b) = J_{\alpha\lambda + \delta\bar{\lambda}}(b) \quad \forall b \in \mathbb{R}^p.$$

□

Ce lemme permet de retrouver un résultat connu en géométrie : la somme de Minkowski de deux permutoèdres est un permutoèdre (voir, par exemple, (Doker, 2011, Lemme 2.2.2)). En effet, soit $\hat{\lambda} \in \mathbb{R}^{p+}$ et $\lambda \in \mathbb{R}^{p+}$ alors :

$$\begin{aligned} & \underbrace{\text{conv}\{(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}) \mid \pi \in S_p\}}_{\partial J_\lambda(1, \dots, 1)} + \underbrace{\text{conv}\{(\bar{\lambda}_{\pi(1)}, \dots, \bar{\lambda}_{\pi(p)}) \mid \pi \in S_p\}}_{\partial J_{\bar{\lambda}}(1, \dots, 1)} \\ &= \underbrace{\text{conv}\{(\bar{\lambda}_{\pi(1)} + \lambda_{\pi(1)}, \dots, \bar{\lambda}_{\pi(p)} + \lambda_{\pi(p)}) \mid \pi \in S_p\}}_{\partial J_{\bar{\lambda} + \lambda}(1, \dots, 1)}. \end{aligned}$$

Preuve du Théorème 3.1

Démonstration.

1 : E_m est un ensemble convexe) Soit $\hat{\lambda} \in E_m$ et $\bar{\lambda} \in E_m$ alors, il existe $\hat{\beta} \in \mathcal{S}_{\hat{\lambda}}$ et $\bar{\beta} \in \mathcal{S}_{\bar{\lambda}}$ tels que $\text{schm}(\hat{\beta}) = \text{schm}(\bar{\beta}) = m$. On pose $\lambda = \alpha\hat{\lambda} + (1 - \alpha)\bar{\lambda}$ et $\beta = \alpha\hat{\beta} + (1 - \alpha)\bar{\beta}$ alors

$$X^\top(y - X\beta) = \alpha X^\top(y - X\hat{\beta}) + (1 - \alpha)X^\top(y - X\bar{\beta})$$

Comme $\text{schm}(\beta) = m$, $X^\top(y - X\hat{\beta}) \in \partial J_{\hat{\lambda}}(m)$ et $X^\top(y - X\bar{\beta}) \in \partial J_{\bar{\lambda}}(m)$ donc, d'après le Lemme 3.1, $X^\top(y - X\beta) \in \partial J_\lambda(m) = \partial J_\lambda(\beta)$. Ainsi, $\beta \in \mathcal{S}_\lambda$ d'où $\lambda \in E_m$, ce qui montre que E_m est un ensemble convexe.

1 : caractérisation de E_m) La preuve de cette caractérisation est similaire à a preuve de l'assertion 2 du Théorème 2.1 du chapitre 2.

Condition nécessaire. Soit $\lambda \in E_m$, alors il existe $\hat{\beta} \in \mathcal{S}_\lambda$ tel que $\text{schm}(\hat{\beta}) = m$. Par conséquent, $\hat{\beta} = U_m s$ pour un certain $s \in \mathbb{R}^{k+}$. Comme $\hat{\beta}$ est un élément de \mathcal{S}_λ dont le schéma est m alors $X^\top(y - \hat{\beta}) \in \partial J_\lambda(\hat{\beta}) = \partial J_\lambda(m)$. En multipliant cette inclusion par U_m^\top , on obtient $\tilde{X}_m^\top(y - \hat{\beta}) = \tilde{\lambda}_m$ d'où

$$\tilde{X}_m^\top y - \tilde{\lambda}_m = \tilde{X}_m^\top \hat{\beta} = \tilde{X}_m^\top U_m s, \quad (3.3)$$

ce qui prouve la condition du schéma. On applique $\tilde{X}_m^{\top+}$ à gauche de l'expression (3.3) et on utilise le fait que $\tilde{X}_m^{\top+} \tilde{X}_m^\top$ est la projection sur $\text{im}(\tilde{X}_m)$. Comme $\hat{\beta} \in \text{im}(\tilde{X}_m)$, on a $\tilde{X}_m^{\top+} \tilde{X}_m^\top \hat{\beta} = \hat{\beta}$. Ainsi,

$$\tilde{X}_m^{\top+} \tilde{X}_m^\top y - \tilde{X}_m^{\top+} \tilde{\lambda}_m = \hat{\beta}.$$

L'égalité ci-dessus donne la condition sous-différentielle :

$$\begin{aligned} \partial J_\lambda(m) \ni X^\top(y - \hat{\beta}) &= X^\top(y - (\tilde{X}_m^{\top+} \tilde{X}_m^\top y - \tilde{X}_m^{\top+} \tilde{\lambda}_m)) \\ &= X^\top \tilde{X}_m^{\top+} \tilde{\lambda}_m + X^\top(I\lambda_n - \tilde{X}_m^{\top+} \tilde{X}_m^\top)y. \end{aligned} \quad (3.4)$$

Condition Suffisante. Supposons que la condition de positivité et les conditions sous-différentielles soient vraies. Alors, d'après la condition de positivité, on peut choisir $s \in \mathbb{R}^{k^+}$ tel que

$$\tilde{\lambda}_m = \tilde{X}_m^\top y - \tilde{X}_m^\top \tilde{X}_m s. \quad (3.5)$$

Montrons que $U_m s \in \mathcal{S}_\lambda$. Par définition de U_m , on a $\text{schm}(U_m s) = m$ donc $\partial J_\lambda(U_m s) = \partial J_\lambda(m)$. De plus, en utilisant (3.4) et (3.5) on obtient

$$\begin{aligned} \partial J_\lambda(U_m s) &\ni X^\top (y - \tilde{X}_m^{\top+} \tilde{X}_m^\top y + \tilde{X}_m^{\top+} \tilde{\lambda}_m) \\ &\ni X^\top (y - \tilde{X}_m^{\top+} \tilde{X}_m^\top y + \tilde{X}_m^{\top+} (\tilde{X}_m y - \tilde{X}_m^\top \tilde{X}_m s)) \\ &\ni X^\top (y - XU_m s). \end{aligned}$$

Par conséquent $U_m s \in \mathcal{S}_\lambda$.

2 : continuité Montrons d'abord la continuité en $\lambda = 0$. Soit $\hat{\beta}(\lambda) \in \mathcal{S}_\lambda$. Comme XX^+ est la projection sur $\text{im}(X)$, que $\hat{\text{v}}\text{a}(\lambda) = X\hat{\beta}(\lambda) \in \text{im}(X)$ et que $\hat{\beta}(\lambda)$ est une solution du problème SLOPE, on a l'inégalité :

$$\frac{1}{2} \|y - XX^+ y\|_2^2 \leq \frac{1}{2} \|y - \hat{\text{v}}\text{a}(\lambda)\|_2^2 \leq \frac{1}{2} \|y - \hat{\text{v}}\text{a}(\lambda)\|_2^2 + J_\lambda(\hat{\beta}(\lambda)) \leq \frac{1}{2} \|y - XX^+ y\|_2^2 + J_\lambda(X^+ y) \quad \forall \lambda \in \mathbb{R}^{p^+}.$$

Puisque $\lim_{\lambda \rightarrow 0} J_\lambda(X^+ y) = 0$ on déduit de cette inégalité que $\lim_{\lambda \rightarrow 0} \hat{\text{v}}\text{a}(\lambda) = XX^+ y$. Soit $\lambda \in \mathbb{R}^{p^+} \setminus \{0\}$ et $(\lambda_n)_{n \in \mathbb{N}} \in \mathbb{R}^{p^+} \setminus \{0\}$ une suite convergeant vers λ et $\hat{\beta}(\lambda_n) \in \mathcal{S}_{\lambda_n}$. La suite $(\hat{\beta}(\lambda_n))_{n \in \mathbb{N}}$ est bornée. En effet, si $\eta \|b\|_\infty > 0.5 \|y\|_2^2$, où $\eta = \inf\{\|\lambda_n\|_\infty \mid n \in \mathbb{N}\} > 0$, alors $J_{\lambda_n}(b) \geq \eta \|b\|_\infty$ d'où $0.5 \|y - Xb\|_2^2 + J_{\lambda_n}(b) > 0.5 \|y - 0\|_2^2 + J_{\lambda_n}(0)$ donc b n'est pas un élément de \mathcal{S}_{λ_n} . Par conséquent pour tout $n \in \mathbb{N}$ on a $\|\hat{\beta}(\lambda_n)\|_\infty \leq 0.5 \|y\|_2^2 / \eta$. Quitte à extraire une sous-suite, on peut supposer que $(\hat{\beta}(\lambda_n))_{n \in \mathbb{N}}$ convergent vers une valeur d'adhérence $l \in \mathbb{R}^p$. Soit $\hat{\beta}(\lambda) \in \mathcal{S}_\lambda$. Comme $\hat{\beta}(\lambda_n) \in \mathcal{S}_{\lambda_n}$, l'inégalité suivante est satisfaite :

$$\frac{1}{2} \|y - \hat{\text{v}}\text{a}(\lambda_n)\|_2^2 + \lambda_n J_\lambda(\hat{\beta}(\lambda_n)) \leq \frac{1}{2} \|y - \hat{\text{v}}\text{a}(\lambda)\|_2^2 + \lambda_n J_\lambda(\hat{\beta}(\lambda)).$$

En prenant la limite de cette expression on obtient :

$$\frac{1}{2} \|y - Xl\|_2^2 + J_\lambda(l) \leq \frac{1}{2} \|y - \hat{\text{v}}\text{a}(\lambda)\|_2^2 + J_\lambda(\hat{\beta}(\lambda)).$$

Comme $\hat{\beta}(\lambda) \in \mathcal{S}_\lambda$, on en déduit que $l \in \mathcal{S}_\lambda$ d'où $Xl = \hat{\text{v}}\text{a}(\lambda)$. Par conséquent, $(\hat{\text{v}}\text{a}(\lambda_n))_{n \in \mathbb{N}}$ est une suite bornée (car $(\hat{\beta}(\lambda_n))_{n \in \mathbb{N}}$ est bornée) ayant une unique valeur d'adhérence : $\hat{\text{v}}\text{a}(\lambda)$ d'où $\lim_{n \rightarrow +\infty} \hat{\text{v}}\text{a}(\lambda_n) = \hat{\text{v}}\text{a}(\lambda)$. Donc, la fonction $\lambda \in \mathbb{R}^{p^+} \mapsto \hat{\text{v}}\text{a}(\lambda)$ est continue.

2) Lorsque $\lambda \in E_m$ multiplie alors les deux côtés de la condition du schéma par $\tilde{X}_m^{\top+}$. En utilisant le fait que $\tilde{X}_m^{\top+} \tilde{X}_m^\top$ est la projection sur $\text{im}(\tilde{X}_m)$ et que $U_m s \in \mathcal{S}_\lambda$ on obtient

$$\tilde{X}_m^{\top+} \tilde{X}_m^\top y - \tilde{X}_m^{\top+} \tilde{\lambda}_m = \tilde{X}_m^{\top+} \tilde{X}_m^\top \tilde{X}_m s = \tilde{X}_m s = \hat{\text{v}}\text{a}(\lambda).$$

3 : continuité) Montrons d'abord la continuité en $\lambda = 0$. Comme $\widehat{\beta}(\lambda)$ est une solution du problème SLOPE alors, pour tout $\lambda \in \mathbb{R}^{p+}$ on a l'inégalité :

$$\begin{aligned} \frac{1}{2} \|y - X(X^\top X)^{-1} X^\top y\|_2^2 &\leq \frac{1}{2} \|y - X\widehat{\beta}(\lambda)\|_2^2, \\ &\leq \frac{1}{2} \|y - X\widehat{\beta}(\lambda)\|_2^2 + J_\lambda(\widehat{\beta}(\lambda)), \\ &\leq \frac{1}{2} \|y - X(X^\top X)^{-1} X^\top y\|_2^2 + J_\lambda((X^\top X)^{-1} X^\top y). \end{aligned}$$

Comme $\lim_{\lambda \rightarrow 0} J_\lambda((X^\top X)^{-1} X^\top y) = 0$ on déduit de cette inégalité que $\lim_{\lambda \rightarrow 0} \widehat{\beta}(\lambda) = (X^\top X)^{-1} X^\top y$. Soit $\lambda \in \mathbb{R}^{p+} \setminus \{0\}$ et $(\lambda_n)_{n \in \mathbb{N}} \in \mathbb{R}^{p+} \setminus \{0\}$ une suite convergeant vers λ . La suite $(\widehat{\beta}(\lambda_n))_{n \in \mathbb{N}}$ est bornée. Quitte à extraire une sous-suite, on peut supposer que $(\widehat{\beta}(\lambda_n))_{n \in \mathbb{N}}$ convergent vers une valeur d'adhérence $l \in \mathbb{R}^p$. Comme $\widehat{\beta}(\lambda_n)$ est une solution du problème SLOPE, l'inégalité suivante est satisfaite :

$$\frac{1}{2} \|y - X\widehat{\beta}(\lambda_n)\|_2^2 + J_{\lambda_n}(\widehat{\beta}(\lambda_n)) \leq \frac{1}{2} \|y - X\widehat{\beta}(\lambda)\|_2^2 + J_{\lambda_n}(\widehat{\beta}(\lambda)).$$

En prenant la limite de cette expression on obtient :

$$\frac{1}{2} \|y - Xl\|_2^2 + J_\lambda(l) \leq \frac{1}{2} \|y - X\widehat{\beta}(\lambda)\|_2^2 + J_\lambda(\widehat{\beta}(\lambda)).$$

D'où $l = \widehat{\beta}(\lambda)$. Donc, la fonction $\lambda \in \mathbb{R}^{p+} \mapsto \widehat{\beta}(\lambda)$ est continue.

3) Soit $\lambda \in E_m$. Comme $\widetilde{X}_m^\top \widetilde{X}_m$ est inversible, la condition du schéma donne

$$\widehat{\beta}(\lambda) = U_m s = U_m (\widetilde{X}_m^\top \widetilde{X}_m)^{-1} (\widetilde{X}_m^\top y - \widetilde{\lambda}_m).$$

□

Notons que la convexité de E_m est une conséquence de la caractérisation de cet ensemble. La preuve donnée de l'assertion 1) à l'avantage d'établir la convexité de E_m sans utiliser cette caractérisation.

Preuve du Théorème 3.2

Dans cette preuve nous utiliserons le fait qu'un élément v du permutoèdre P_λ , avec $\lambda \in \mathbb{R}^{p+}$, vérifie $v_1 + \dots + v_p = \lambda_1 + \dots + \lambda_p$.

Démonstration. Soit $\pi \in S_p$ et $\epsilon \in \{-1, 1\}^p$ tels que

$$|\widehat{\beta}|_\downarrow = (\epsilon_1 \widehat{\beta}_{\pi(1)}, \dots, \epsilon_p \widehat{\beta}_{\pi(p)}),$$

et soit ϕ la transformation orthogonale définie par :

$$\phi(x) = (\epsilon_1 x_{\pi(1)}, \dots, \epsilon_p x_{\pi(p)}) \quad \forall x \in \mathbb{R}^p.$$

1) Comme $\widehat{\beta} \in S_\lambda$ est une solution du problème d'optimisation SLOPE, l'équivalence suivante est vérifiée :

$$X^\top (y - \widehat{v}_\lambda) \in \partial J_\lambda(\widehat{\beta}) \Leftrightarrow \phi(X^\top (y - \widehat{v}_\lambda)) \in \phi(\partial J_\lambda(\widehat{\beta})) = \partial J_\lambda(|\widehat{\beta}|_\downarrow).$$

Les composantes de $|\widehat{\beta}|_{\downarrow}$ sont décroissantes ainsi, $\partial J_{\lambda}(|\widehat{\beta}|_{\downarrow})$ est un produit cartésien de permutoèdres avec potentiellement un permutoèdre signé (si $\widehat{\beta}$ a au moins une composante nulle). Plus précisément

$$\partial J_{\lambda}(|\widehat{\beta}|_{\downarrow}) = \begin{cases} P_{\lambda_1, \dots, \lambda_{k_1}} \times \dots \times P_{\lambda_{k_{l-1}+1}, \dots, \lambda_{k_l}} & \text{si } k_l = p, \\ P_{\lambda_1, \dots, \lambda_{k_1}} \times \dots \times P_{\lambda_{k_{l-1}+1}, \dots, \lambda_{k_l}} \times P_{\lambda_{k_l+1}, \dots, \lambda_p}^{\pm} & \text{si } k_l < p. \end{cases}$$

Si $b \in P_{\lambda_1, \dots, \lambda_{k_1}} \times \dots \times P_{\lambda_{k_{l-1}+1}, \dots, \lambda_{k_l}}$, alors les égalités suivantes sont vérifiées :

$$\forall i \in \{k_1, \dots, k_l\} \quad \sum_{j=1}^i b_j = \|b\|_{(i)} = \sum_{j=1}^i \lambda_j.$$

Enfin, comme la i -norme $\|\cdot\|_{(i)}$ est invariante pour la transformation ϕ , on en déduit les égalités suivantes :

$$\forall i \in \{k_1, \dots, k_l\} \quad \frac{\|\phi(X^{\top}(y - \widehat{v}a(\lambda)))\|_{(i)}}{\sum_{j=1}^i \lambda_j} = \frac{\|X^{\top}(y - \widehat{v}a(\lambda))\|_{(i)}}{\sum_{j=1}^i \lambda_j} = 1.$$

2) Établissons d'abord l'inclusion suivante lorsque $b \in \mathbb{R}^p$ vérifie $b_1 \geq \dots \geq b_p > 0$:

$$\partial J_{\lambda}(b) \subseteq \text{conv} \{(\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)}), \pi \in S_p\} = P_{\lambda}. \quad (3.6)$$

Comme la norme ℓ_1 ordonnée est polyédrique, à savoir

$$J_{\lambda}(b) = \max \left\{ \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i \mid \epsilon_1, \dots, \epsilon_p \in \{-1, 1\}, \pi \in S_p \right\},$$

son sous-différentiel est donné par

$$\partial J_{\lambda}(b) = \text{conv} \left\{ (\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \mid \epsilon_1, \dots, \epsilon_p \in \{-1, 1\}, \pi \in S_p, \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i = J_{\lambda}(b) \right\}.$$

De plus, si $\epsilon_{i_0} \lambda_{\pi(i_0)} < 0$ pour un certain $i_0 \in \{1, \dots, p\}$, alors

$$\sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i < \lambda_{\pi(i_0)} b_{i_0} + \sum_{i \neq i_0} \epsilon_i \lambda_{\pi(i)} b_i \leq J_{\lambda}(b).$$

On déduit de cette inégalité que $(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \notin \partial J_{\lambda}(b)$. Par conséquent, $(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) \in \partial J_{\lambda}(b)$ implique que $(\epsilon_1 \lambda_{\pi(1)}, \dots, \epsilon_p \lambda_{\pi(p)}) = (\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)})$ ce qui prouve l'inclusion (3.6).

Supposons maintenant qu'il existe $i \notin \mathcal{A}(\lambda)$ tel que

$$\begin{cases} |\widehat{\beta}|_{\downarrow i} > |\widehat{\beta}|_{\downarrow i+1} & \text{si } i \leq p-1, \\ |\widehat{\beta}|_{\downarrow i} > 0 & \text{si } i = p. \end{cases}$$

Alors, on a $\partial J_{\lambda_1, \dots, \lambda_i}(|\widehat{\beta}|_{\downarrow 1}, \dots, |\widehat{\beta}|_{\downarrow i}) \subseteq P_{\lambda_1, \dots, \lambda_i}$. Par conséquent

$$\frac{\|\phi(X^{\top}(y - \widehat{v}a(\lambda)))\|_{(i)}}{\sum_{j=1}^i \lambda_j} = \frac{\|X^{\top}(y - \widehat{v}a(\lambda))\|_{(i)}}{\sum_{j=1}^i \lambda_j} = 1.$$

Donc $i \in \mathcal{A}(\lambda)$, ce qui conduit à une contradiction. □

Condition sous laquelle l'estimateur SLOPE est un estimateur LASSO généralisé

Soit $D \in \mathbb{R}^{m \times p}$. Le sous-différentiel en 0 de la fonction $b \in \mathbb{R}^p \mapsto \|Db\|_1$ est le zonotope³ $D^\top[-1, 1]^m$. D'autre part, le permutoèdre signé, c'est-à-dire le sous-différentiel en 0 de J_λ , est un zonotope si et seulement si $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ est arithmétique (Godland et Kabluchko, 2023, Théorème 4.13). Par conséquent, lorsque λ n'est pas arithmétique, on ne peut pas choisir une matrice $D \in \mathbb{R}^{m \times p}$ telle que $J_\lambda(\cdot) = \|D \cdot\|_1$, donc l'estimateur SLOPE n'est pas un LASSO généralisé. À l'inverse l'estimateur OSCAR, c'est-à-dire l'estimateur SLOPE lorsque $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ est arithmétique, est un exemple d'estimateur LASSO généralisé. Une preuve de ce commentaire est donnée au Lemme 3.2 dans le cas particulier où $p = 3$.

Lemme 3.2. *Soit $\lambda \in \mathbb{R}^{3+} \setminus \{0\}$. Il existe une matrice $D \in \mathbb{R}^{m \times 3}$ telle que pour tout $b \in \mathbb{R}^3$, $J_\lambda(b) = \|Db\|_1$ si et seulement si $\lambda_2 = (\lambda_1 + \lambda_3)/2$.*

Démonstration. (\Rightarrow) Supposons que, pour tout $b \in \mathbb{R}^3$, $J_\lambda(b) = \|Db\|_1$ pour une certaine matrice $D \in \mathbb{R}^{m \times 3}$. On va établir que $\lambda_2 = (\lambda_1 + \lambda_3)/2$ en analysant les sous-différentiels de $\|D \cdot\|_1$ et $J_\lambda(\cdot)$ au point $\bar{b} = (1, 1, 1)$. Comme $\partial\|D \cdot\|_1(\bar{b}) = D^\top \partial\|\cdot\|_1(D\bar{b})$ et que $\text{signe}(D\bar{b})$ est le centre de symétrie de $\partial\|\cdot\|_1(D\bar{b})$ alors $c = D^\top \text{signe}(D\bar{b})$ est le centre de symétrie de $\partial\|D \cdot\|_1(\bar{b})$. Par ailleurs, d'après le Lemme 2 du chapitre 1, on a $\partial J_\lambda(\bar{b}) = \text{conv}\{(\lambda_{\pi(1)}, \lambda_{\pi(2)}, \lambda_{\pi(3)}) \mid \pi \in S_3\}$. Le centre de symétrie de $\partial J_\lambda(\bar{b})$ est $c = (c_1, c_2, c_3)$ donc, pour toute permutation $\pi \in S_3$, le point $(2c_1 - \lambda_{\pi(1)}, 2c_2 - \lambda_{\pi(2)}, 2c_3 - \lambda_{\pi(3)})$ est un sommet du permutoèdre $\text{conv}\{(\lambda_{\pi(1)}, \lambda_{\pi(2)}, \lambda_{\pi(3)}) \mid \pi \in S_3\}$. Ainsi,

$$\frac{1}{6} \sum_{\pi \in S_3} (\lambda_{\pi(1)}, \lambda_{\pi(2)}, \lambda_{\pi(3)}) = \frac{1}{6} \sum_{\pi \in S_3} (2c_1 - \lambda_{\pi(1)}, 2c_2 - \lambda_{\pi(2)}, 2c_3 - \lambda_{\pi(3)}),$$

d'où $c_1 = c_2 = c_3 = \bar{\lambda}$ où $\bar{\lambda} = (\lambda_1 + \lambda_2 + \lambda_3)/3$. Comme $\lambda \in \partial J_\lambda(\bar{b})$ alors $2c - \lambda \in \partial J_\lambda(\bar{b})$ d'où

$$J_\lambda^*(2c - \lambda) = \max \left\{ \frac{2\bar{\lambda} - \lambda_3}{\lambda_1}, \frac{4\bar{\lambda} - \lambda_3 - \lambda_2}{\lambda_1 + \lambda_2}, \underbrace{\frac{6\bar{\lambda} - \lambda_1 - \lambda_2 - \lambda_3}{\lambda_1 + \lambda_2 + \lambda_3}}_{=1} \right\} \leq 1.$$

Par conséquent

$$\begin{cases} \frac{2\bar{\lambda} - \lambda_3}{\lambda_1} \leq 1 \\ \frac{4\bar{\lambda} - \lambda_3 - \lambda_2}{\lambda_1 + \lambda_2} \leq 1 \end{cases} \Rightarrow \begin{cases} \lambda_2 \leq \frac{\lambda_1 + \lambda_3}{2} \\ \lambda_2 \geq \frac{\lambda_1 + \lambda_3}{2} \end{cases} \Rightarrow \lambda_2 = \frac{\lambda_1 + \lambda_3}{2}.$$

(\Leftarrow) Si $\lambda_2 = \frac{\lambda_1 + \lambda_3}{2}$ alors

$$J_\lambda(b) = \lambda_3(|b_1| + |b_2| + |b_3|) + \frac{\lambda_1 - \lambda_3}{2} \left(\frac{|b_1 + b_2| + |b_1 - b_2|}{2} + \frac{|b_1 + b_3| + |b_1 - b_3|}{2} + \frac{|b_2 + b_3| + |b_2 - b_3|}{2} \right),$$

donc $J_\lambda(b) = \|Db\|_1$ où $D \in \mathbb{R}^{9 \times 3}$ vérifie

$$D^\top = \begin{pmatrix} \lambda_3 & 0 & 0 & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_1 - \lambda_3}{4} & 0 & 0 \\ 0 & \lambda_3 & 0 & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_3 - \lambda_1}{4} & 0 & 0 & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_1 - \lambda_3}{4} \\ 0 & 0 & \lambda_3 & 0 & 0 & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_3 - \lambda_1}{4} & \frac{\lambda_1 - \lambda_3}{4} & \frac{\lambda_3 - \lambda_1}{4} \end{pmatrix}.$$

□

Cette preuve, inspirée de la démonstration du Théorème 4.13 de Godland et Kabluchko (2023), peut-être généralisée lorsque $p \geq 3$ moyennant l'introduction de notations un peu plus techniques.

3. Un zonotope est l'image d'un cube par une transformation affine.

Conclusion et perspectives

Ce manuscrit illustre la notion de schéma du SLOPE et son intérêt pour l'étude de cet estimateur. Le chapitre 1 introduit cette notion et démontre que les schémas sont des représentants canoniques pour la relation d'équivalence « avoir le même sous-différentiel pour la norme ℓ_1 ordonnée ». Par ailleurs, géométriquement, il y a une bijection entre les schémas et les faces du permutoèdre signé (la boule unité de la norme ℓ_1 ordonnée duale). Le chapitre 2 introduit la condition d'irreprésentabilité du SLOPE qui est nécessaire pour la récupération du schéma des coefficients de régression avec une probabilité supérieure à un demi. De plus, cette partie illustre que cette condition peut être relâchée en appliquant l'opérateur proximal de la norme ℓ_1 ordonnée à l'estimateur SLOPE. Enfin, le chapitre 3 propose un algorithme pour le calcul du chemin des solutions de l'estimateur SLOPE basé sur les conditions du schéma et du sous-différentiel. L'achèvement de ce manuscrit ouvre quelques perspectives pour des travaux futurs.

Chemin multidimensionnel des solutions du SLOPE

Le chapitre 3 de ce manuscrit évoque le chemin multidimensionnel des solutions du SLOPE lorsque le paramètre de pénalité λ varie dans \mathbb{R}^{p++} . Cette ébauche de travail, inédite, n'est pas basée sur une pré-publication. Dans le futur, il serait pertinent d'étudier ce chemin multidimensionnel afin de développer une méthode numérique permettant de sélectionner le paramètre de pénalité $\lambda \in \mathbb{R}^{p++}$ minimisant, par exemple, la formule SURE ou la somme des carrés résiduels sur un échantillon de validation.

Généralisation de la notion de schéma et applications

La notion de schéma, cruciale pour ce manuscrit, est facilement généralisable en considérant la relation d'équivalence « avoir le même sous-différentiel » relativement à un terme de pénalité. D'après l'article de Graczyk *et al.* (2023), lorsque le terme de pénalité est une famille finie de formes linéaires, les classes d'équivalence de cette relation sont les intérieurs relatifs des cônes normaux des faces du sous-différentiel en zéro du terme de pénalité. Dans la continuité du chapitre 3, une notion plus générale de schéma pourrait être très utile pour développer des algorithmes de calcul de chemin des solutions d'estimateurs pénalisés (par exemple, pour calculer le chemin des solutions de l'estimateur PACS). Une perspective future serait de développer de tels algorithmes calqués sur celui calculant le chemin des solutions de l'estimateur SLOPE. Une autre perspective serait d'étendre ces travaux obtenus dans le cadre de la régression linéaire à d'autres modèles. Par exemple, certains collaborateurs travaillent actuellement sur la condition d'irreprésentabilité du SLOPE dans le cadre du modèle graphique. Il pourrait également être pertinent de considérer le chemin des solutions du SLOPE dans le cadre des modèles de régression logistique, graphique etc.

Collaborations interdisciplinaires

Récemment, j'ai collaboré avec des médecins du Centre Hospitalier Universitaire d'Angers, exerçant dans le domaine de la médecine du travail. L'étude que nous avons menée portait sur une modélisation de la probabi-

lité de développer une maladie musculo-squelettique en fonction variables liées à la trajectoire professionnelle (Deltreil *et al.*, 2022). D'un point de vue mathématique, la modélisation mettait en jeu un modèle de régression logistique. Bien que nous n'ayons pas utilisé l'estimateur SLOPE logistique (une autre méthode était envisagée au moment où j'ai rejoint le projet), cet estimateur aurait été pertinent dans le contexte de cette étude pour réduire le nombre de variables explicatives et construire des groupes de variables ayant un impact similaire sur cette probabilité. Dans le futur, je souhaite renforcer ces collaborations, ce qui me permettrait d'illustrer mes travaux théoriques sur des jeux de données réelles ou encore de valoriser mes compétences en apprentissage statistique par des travaux interdisciplinaires.

Annexe A

Généralisation de certaines propriétés de l'estimateur SLOPE

Les propriétés suivantes, liées à l'estimateur SLOPE, n'ont pas encore été démontrées :

- Une face quelconque du permuttoèdre signé s'exprime comme le sous-différentiel de la norme ℓ_1 ordonnée en un point (Lemme 1.5 du chapitre 1).
- La valeur ajustées ne dépend pas de la solution du problème SLOPE (Proposition 3.1 du chapitre 3).
- La condition nécessaire et suffisante pour l'unicité uniforme de l'estimateur SLOPE (Proposition 2.4 du chapitre 2).

Les propriétés mentionnées ont peu d'intérêt à être prouvées en utilisant la norme ℓ_1 ordonnée et gagnent en clarté à être prouvées dans un cadre plus général. Nous avons vu que la norme ℓ_1 ordonnée est polyédrique, sa boule unité est un polyèdre, et nous rappelons son expression

$$J_\lambda(b) = \max \left\{ \sum_{i=1}^p \epsilon_i \lambda_{\pi(i)} b_i \mid \pi \in S_p, \epsilon \in \{-1, 1\}^p \right\}.$$

On peut noter que J_λ est le maximum de formes linéaires; c'est l'expression utilisée pour le terme de pénalité du Lemme A.1 ou du Théorème A.1 par ailleurs, J_λ est convexe; c'est la propriété retenue pour le terme de pénalité à la Proposition A.1.

Notations : Soit $u_1, \dots, u_k \in \mathbb{R}^p$ on note

- $\text{aff}\{u_1, \dots, u_k\}$ le plus petit espace affine contenant u_1, \dots, u_k , c'est à dire

$$\text{aff}\{u_1, \dots, u_k\} = \{\alpha_1 u_1 + \dots + \alpha_k u_k \mid \alpha_1, \dots, \alpha_k \in \mathbb{R} \text{ et } \alpha_1 + \dots + \alpha_k = 1\}.$$

- $\text{vect}\{u_1, \dots, u_k\}$ le plus petit espace vectoriel contenant u_1, \dots, u_k , c'est à dire

$$\text{vect}\{u_1, \dots, u_k\} = \{\alpha_1 u_1 + \dots + \alpha_k u_k \mid \alpha_1, \dots, \alpha_k \in \mathbb{R}\}.$$

- $\overrightarrow{\text{aff}}\{u_1, \dots, u_k\}$ l'espace vectoriel obtenu en translatant $\text{aff}\{u_1, \dots, u_k\}$ vers l'origine, c'est-à-dire

$$\overrightarrow{\text{aff}}\{u_1, \dots, u_k\} = \{z = \alpha_1 u_1 + \dots + \alpha_k u_k \mid \alpha_1, \dots, \alpha_k \in \mathbb{R} \text{ et } \alpha_1 + \dots + \alpha_k = 0\}.$$

- Soit $h \in \mathbb{R}^p$, h^\perp représente l'hyperplan vectoriel $\{z \in \mathbb{R}^p \mid h^\top z = 0\}$.

Faces d'un polytope et sous-différentiel

Le Lemme A.1 montre que les faces d'un polytope (polyèdre borné) s'expriment comme le sous-différentiel d'un maximum d'une famille finie de formes linéaires.

Lemme A.1. Soient $u_1, \dots, u_k \in \mathbb{R}^p$, P le polytope $P = \text{conv}\{u_1, \dots, u_k\}$ et pen la fonction convexe définie comme le maximum d'une famille finie de formes linéaires : $\text{pen}(x) = \max\{u_1^\top x, \dots, u_k^\top x\}$ pour $x \in \mathbb{R}^p$. Le sous-différentiel de pen en un point $x \in \mathbb{R}^p$ est une face du polytope P ayant l'expression suivante :

$$\partial \text{pen}(x) = \{s \in P \mid x^\top s = \text{pen}(x)\}, \text{ où } x^\top s \leq \text{pen}(x) \text{ est une inégalité valide pour tout } s \in P.$$

Inversement, une face non vide F de P est le sous-différentiel de pen en un point $x \in \mathbb{R}^p$: $F = \partial \text{pen}(x)$.

Démonstration.

Le sous-différentiel est une face de P : On pose $s = \sum_{i=1}^k \alpha_i u_i$ un élément arbitraire de P où $\alpha_1 \geq 0, \dots, \alpha_k \geq 0$ et $\sum_{i=1}^k \alpha_i = 1$. Comme $u_i^\top x \leq \text{pen}(x)$ pour tout $i \in \{1, \dots, k\}$ on en déduit que $x^\top s \leq \text{pen}(x)$ est une inégalité valide pour tout $s \in P$.

Le sous-différentiel de pen satisfait la formule suivante $\partial \text{pen}(x) = \text{conv}\{u_i \mid i \in I_{\text{pen}}(x)\}$ où $I_{\text{pen}}(x) = \{i \in \{1, \dots, k\} \mid u_i^\top x = \text{pen}(x)\}$ (Hiriart-Urruty et Lemaréchal, 2004). Ainsi, pour $i \in I_{\text{pen}}(x)$, u_i est un élément de l'ensemble $\{s \in P \mid x^\top s = \text{pen}(x)\}$. Comme $\{s \in P \mid x^\top s = \text{pen}(x)\}$ est un ensemble convexe, on en déduit l'inclusion suivante :

$$\partial \text{pen}(x) = \text{conv}\{u_i \mid i \in I_{\text{pen}}(x)\} \subseteq \{s \in P \mid x^\top s = \text{pen}(x)\}.$$

Inversement, soit $s \in P$ tel que $s \notin \text{conv}\{u_i \mid i \in I_{\text{pen}}(x)\}$ alors $s = \sum_{i=1}^k \alpha_i u_i$ où $\alpha_1 \geq 0, \dots, \alpha_k \geq 0$, $\sum_{i=1}^k \alpha_i = 1$ et $\alpha_{i_0} > 0$ pour un certain $i_0 \notin I_{\text{pen}}(x)$. Comme $u_i^\top x \leq \text{pen}(x)$ pour tout $i \in \{1, \dots, k\}$ et $u_{i_0}^\top x < \text{pen}(x)$, d'où

$$x^\top s = \sum_{i=1}^k \alpha_i u_i^\top x < \text{pen}(x).$$

Par conséquent, $\partial \text{pen}(x) = \text{conv}\{u_i \mid i \in I_{\text{pen}}(x)\} = \{s \in P \mid x^\top s = \text{pen}(x)\}$ donc $\partial \text{pen}(x)$ est une face de P .

Une face de P s'exprime comme un sous-différentiel : soit $F = \{s \in P \mid a^\top s = c\}$ une face non vide de P où $a \in \mathbb{R}^p$, $c \in \mathbb{R}$ et $a^\top s \leq c$ est une inégalité valide pour tout $s \in P$. Montrons que $F = \partial \text{pen}(a)$. Soit $s \in F$ on a $a^\top s = c$ et $a^\top s \leq \text{pen}(a)$ donc $c \leq \text{pen}(a)$. Par ailleurs, pour tout $s \in \partial \text{pen}(a)$, on a $a^\top s = \text{pen}(a)$ ainsi que $a^\top s \leq c$ car $\partial \text{pen}(a) \subseteq P$, d'où $\text{pen}(a) \leq c$. Par conséquent $\text{pen}(a) = c$ donc $F = \partial \text{pen}(a)$. \square

Valeur ajustée

La Proposition A.1 montre que la notion de valeur ajustée est bien définie même dans le cas où l'ensemble des solutions n'est pas un singleton.

Proposition A.1. Soient $X \in \mathbb{R}^{n \times p}$, $y \in \mathbb{R}^n$, $\gamma > 0$ et pen la fonction convexe sur \mathbb{R}^p . Si $\widehat{\beta}$ et $\bar{\beta}$ sont des solutions du problème d'optimisation

$$\min_{b \in \mathbb{R}^p} \left\{ \frac{1}{2} \|y - Xb\|_2^2 + \gamma \text{pen}(b) \right\} \quad (\text{A.1})$$

alors $X\widehat{\beta} = X\bar{\beta}$ et $\text{pen}(\widehat{\beta}) = \text{pen}(\bar{\beta})$.

Cette proposition montre que la valeur ajustée $X\widehat{\beta}$ ne dépend pas de la solution $\widehat{\beta}$ du problème (A.1). Pour la suite on notera $\mathcal{S}_{\gamma \text{pen}}(y)$ l'ensemble des solution du problème (A.1). Notons que $\mathcal{S}_{\gamma \text{pen}}(y)$ est fermé mais que

sans faire d'hypothèses supplémentaires autres que la convexité du le terme de pénalité pen cet ensemble peut être non borné, borné (donc compact), un singleton ou encore vide. Ce dernier cas, n'est pas rencontré pour les termes de pénalité utilisés en pratique; par exemple lorsque le terme de pénalité est une norme (comme pour l'estimateur SLOPE) la fonction objective du problème (A.1) est coercive donc $\mathcal{S}_{\gamma\text{pen}}(y)$ admet au moins un élément. Même pour le LASSO généralisé, dont le terme de pénalité n'est pas coercif, $\mathcal{S}_{\gamma\text{pen}}(y)$ n'est pas vide (voir le Théorème 3.1 de Dupuis et Vaïter (2023)). La preuve de la Proposition A.1, donnée ci-dessous, est similaire à la preuve du Lemme 1 de l'article Tibshirani (2013) qui traite de l'estimateur LASSO.

Démonstration. Supposons qu'il existe $\widehat{\beta} \in \mathcal{S}_{\gamma\text{pen}}(y)$ et $\bar{\beta} \in \mathcal{S}_{\gamma\text{pen}}(y)$ tel que $X\widehat{\beta} \neq X\bar{\beta}$ et posons $\beta = (\widehat{\beta} + \bar{\beta})/2$. Comme la fonction $t \in \mathbb{R}^n \mapsto \|y - t\|_2^2$ est strictement convexe, on en déduit que

$$\|y - X\beta\|_2^2 < \frac{1}{2}\|y - X\widehat{\beta}\|_2^2 + \frac{1}{2}\|y - X\bar{\beta}\|_2^2.$$

Ainsi, comme pen est convexe, on en déduit l'inégalité suivante :

$$\frac{1}{2}\|y - X\beta\|_2^2 + \gamma\text{pen}(\beta) < \frac{1}{2}\left(\frac{1}{2}\|y - X\widehat{\beta}\|_2^2 + \gamma\text{pen}(\widehat{\beta}) + \frac{1}{2}\|y - X\bar{\beta}\|_2^2 + \gamma\text{pen}(\bar{\beta})\right),$$

par conséquent l'un des deux éléments $\widehat{\beta}$ ou $\bar{\beta}$ n'appartient pas à $\mathcal{S}_{\gamma\text{pen}}(y)$ ce qui fournit une contradiction. Ainsi $X\widehat{\beta} = X\bar{\beta}$ d'où $\text{pen}(\widehat{\beta}) = \text{pen}(\bar{\beta})$. \square

Condition nécessaire et suffisante d'unicité uniforme

La condition de matrice en position générale (Dossal, 2012) est suffisante pour l'unité uniforme, pour tout $y \in \mathbb{R}^n$ et pour tout paramètre de régularisation $\gamma > 0$, de la solution du problème LASSO (Tibshirani, 2013). Cette condition fut d'abord affaiblie en une condition nécessaire et suffisante d'unicité pour le LASSO (Ewald et Schneider, 2020) puis étendue aux estimateurs des moindres carrés pénalisés par une norme polyédrique (Schneider et Tardivel, 2022) ou une jauge polyédrique (Graczyk *et al.*, 2023). Le Théorème A.1 englobe ces résultats et donne une condition nécessaire et suffisante pour que, pour tout $y \in \mathbb{R}^n$ et pour tout $\gamma > 0$, l'ensemble $\mathcal{S}_{\gamma\text{pen}}(y)$ ait au plus un élément où pen est le maximum d'une famille finie de formes linéaires.

Théorème A.1. *Soient $u_1, \dots, u_k \in \mathbb{R}^p$ et pen la fonction convexe définie comme le maximum d'une famille finie de formes linéaires : $\text{pen}(x) = \max\{u_1^\top x, \dots, u_k^\top x\}$ pour $x \in \mathbb{R}^p$. Il existe $y \in \mathbb{R}^n$ et il existe $\gamma > 0$ pour lesquels $\mathcal{S}_{\gamma\text{pen}}(y)$ a au moins deux éléments si et seulement si $\text{im}(X^\top)$ coupe une face du polytope $B^* = \text{conv}\{u_1, \dots, u_k\}$ dont la dimension est strictement inférieure à $\dim(\ker(X))$.*

Démonstration. (\Leftarrow) Supposons que $\text{im}(X^\top)$ coupe une face F du polytope $B^* = \text{conv}\{u_1, \dots, u_k\}$ telle que $\dim(F) < \dim(\ker(X))$. D'après le Lemme A.1, $F = \partial\text{pen}(\widehat{\beta})$ pour un certain $\widehat{\beta} \in \mathbb{R}^p$. Par hypothèse il existe $z \in \mathbb{R}^n$ tel que $X^\top z \in F$. Posons $y = X\widehat{\beta} + \gamma z$ alors $\widehat{\beta} \in \mathcal{S}_{\gamma\text{pen}}(y)$ en effet

$$\frac{1}{\gamma}X^\top(y - X\widehat{\beta}) = X^\top z \in \partial\text{pen}(\widehat{\beta}).$$

Construisons maintenant $\bar{\beta} \in \mathcal{S}_{\gamma\text{pen}}(y)$ avec $\bar{\beta} \neq \widehat{\beta}$. On a $\partial\text{pen}(\widehat{\beta}) = \text{conv}\{u_i \mid i \in I\}$ où $I = \{i \in \{1, \dots, k\} \mid u_i^\top \widehat{\beta} = \text{pen}(\widehat{\beta})\}$. Montrons que $\dim(\ker(X) \cap \text{vect}\{u_i \mid i \in I\}^\perp) \geq 1$:

1) Si $0 \in \text{aff}\{u_i \mid i \in I\}$ alors $\text{aff}\{u_i \mid i \in I\} = \text{vect}\{u_i \mid i \in I\}$ et $\dim(F) = \dim(\text{vect}\{u_i \mid i \in I\}) < \dim(\ker(X))$. Par conséquent $\dim(\ker(X)) + \dim(\text{vect}\{u_i \mid i \in I\}^\perp) > p$, ce qui prouve l'affirmation.

2) Si $0 \notin \text{aff}\{u_i \mid i \in I\}$ alors on pose $v = X^\top z \in \text{im}(X^\top) \cap \text{conv}\{u_i \mid i \in I\}$ satisfaisant $X^\top z \neq 0$. Comme $\dim(\text{vect}\{u_i \mid i \in I\}) = \dim(\text{aff}\{u_i \mid i \in I\}) + 1 = \dim(F) + 1 \leq \dim(\ker(X))$ on a $\dim(\ker(X)) + \dim(\text{vect}\{u_i \mid i \in I\}^\perp) > p$, ce qui prouve l'affirmation.

$i \in I^\perp \geq p$. Si $\ker(X) \cap \text{vect}\{u_i \mid i \in I\}^\perp = \{0\}$, alors $\mathbb{R}^p = \ker(X) \oplus \text{vect}\{u_i \mid i \in I\}^\perp$. Or $\ker(X) \subseteq v^\perp$ et $\text{vect}\{u_i \mid i \in I\}^\perp \subseteq v^\perp$, ce qui contredit que $\mathbb{R}^p = \ker(X) \oplus \text{vect}\{u_i \mid i \in I\}^\perp$ et prouve l'affirmation.

Soit $h \in \ker(X) \cap \text{vect}\{u_i \mid i \in I\}^\perp$ avec $h \neq 0$ alors pour tout $i \in I$ on a $u_i^\top(\hat{\beta}+h) = u_i^\top\hat{\beta} = \text{pen}(\hat{\beta})$. De plus, pour tout $i \notin I$, on a $u_i^\top\hat{\beta} < \text{pen}(\hat{\beta})$ donc, si la norme de h est suffisamment petite, alors $u_i^\top(\hat{\beta}+h) \leq \text{pen}(\hat{\beta})$. Par conséquent $\text{pen}(\hat{\beta}+h) = \max\{u_i^\top(\hat{\beta}+h) \mid i \in \{1, \dots, k\}\} = \text{pen}(\hat{\beta})$. Cette égalité combinée avec $X\hat{\beta} = X(\hat{\beta}+h)$ montre que $\bar{\beta} = \hat{\beta} + h \in \mathcal{S}_{\gamma\text{pen}}(y)$.

(\implies) Supposons qu'il existe $y \in \mathbb{R}^n$ tel que $\hat{\beta}, \bar{\beta} \in \mathcal{S}_{\gamma\text{pen}}(y)$ avec $\hat{\beta} \neq \bar{\beta}$ alors

$$\frac{1}{\gamma}X^\top(y - X\hat{\beta}) \in \partial\text{pen}(\hat{\beta}) \quad \text{et} \quad \frac{1}{\gamma}X^\top(y - X\bar{\beta}) \in \partial\text{pen}(\bar{\beta}).$$

D'après la Proposition A.1, $X\hat{\beta} = X\bar{\beta}$, donc $\frac{1}{\gamma}X^\top(y - X\hat{\beta}) = \frac{1}{\gamma}X^\top(y - X\bar{\beta})$ par conséquent $\text{im}(X^\top)$ coupe la face $\partial\text{pen}(\hat{\beta}) \cap \partial\text{pen}(\bar{\beta})$. Soit $F^* = \text{conv}\{u_i \mid i \in I^*\}$ une face de $\partial\text{pen}(\hat{\beta}) \cap \partial\text{pen}(\bar{\beta})$ de dimension minimale parmi les faces de $\partial\text{pen}(\hat{\beta}) \cap \partial\text{pen}(\bar{\beta})$ coupées par $\text{im}(X^\top)$. Par minimalité de $\dim(F^*)$, $\text{im}(X^\top)$ coupe l'intérieur relatif de F^* , c'est-à-dire, il existe $z \in \mathbb{R}^n$ tel que $v = X^\top z$ appartienne à F^* , mais pas à une face propre de F^* . Montrons que si $\dim(F^*) \geq \dim(\ker(X))$, alors $\text{im}(X^\top)$ coupe une face propre de F^* , conduisant à une contradiction. On pose $h = \hat{\beta} - \bar{\beta} \neq 0$. Clairement, $h \in \ker(X)$. De plus, comme $\text{pen}(\hat{\beta}) = \text{pen}(\bar{\beta})$ d'après la Proposition A.1, et que $u_i \in \partial\text{pen}(\hat{\beta}) \cap \partial\text{pen}(\bar{\beta})$ pour tout $i \in I^*$, on a

$$u_i^\top h = u_i^\top\hat{\beta} - u_i^\top\bar{\beta} = \text{pen}(\hat{\beta}) - \text{pen}(\bar{\beta}) = 0 \quad \forall i \in I^*.$$

Par conséquent, $h \in \ker(X) \cap \text{vect}\{u_i \mid i \in I^*\}^\perp$. Supposons que $\dim(F^*) = \dim(\overrightarrow{\text{aff}}\{u_i \mid i \in I^*\}) \geq \dim(\ker(X))$ alors $\dim(\text{im}(X^\top)) + \dim(\overrightarrow{\text{aff}}\{u_i \mid i \in I^*\}) \geq \text{rang}(X) + \dim(\ker(X)) = p$. Si $\text{im}(X^\top) \cap \overrightarrow{\text{aff}}\{u_i \mid i \in I^*\} = \{0\}$, alors $\mathbb{R}^p = \text{im}(X^\top) \oplus \overrightarrow{\text{aff}}\{u_i \mid i \in I^*\}$. Cependant, la dernière relation ne peut pas être vraie car $\text{im}(X^\top) = \ker(X)^\perp \subseteq h^\perp$ ainsi que $\overrightarrow{\text{aff}}\{u_i \mid i \in I^*\} \subseteq \text{vect}\{u_i \mid i \in I^*\} \subseteq h^\perp$, où $h \neq 0$. Par conséquent, il existe $\bar{v} \in \text{im}(X^\top) \cap \overrightarrow{\text{aff}}\{u_i \mid i \in I^*\}$ avec $\bar{v} \neq 0$. La droite affine $D = \{X^\top z + t\bar{v} \mid t \in \mathbb{R}\} \subseteq \text{im}(X^\top)$ coupe l'intérieur relatif de F^* en $t = 0$ et est incluse dans $\text{aff}(F^*) = \text{aff}\{u_i \mid i \in I^*\}$, car $X^\top z \in F^*$ et $\bar{v} \in \overrightarrow{\text{aff}}\{u_i \mid i \in I^*\}$. Par conséquent, D coupe une face propre de F^* donc $\text{im}(X^\top)$ coupe une face propre de F^* , ce qui contredit que F^* est une face de $\partial\text{pen}(\hat{\beta}) \cap \partial\text{pen}(\bar{\beta})$ coupées par $\text{im}(X^\top)$ dont la dimension minimale. \square

Cette notion d'unicité uniforme est pertinente en statistique, où $y \in \mathbb{R}^n$ représente la réponse aléatoire d'un modèle de régression. En effet, certaines quantités très classiques, comme l'espérance $\mathbb{E}(\hat{\beta})$, où $\hat{\beta}$ est un élément de $\mathcal{S}_{\gamma\text{pen}}(y)$, n'ont pas de sens si la probabilité, relativement à y , que l'ensemble $\mathcal{S}_{\gamma\text{pen}}(y)$ soit un singleton n'est pas égale à un. La condition de matrice en position générale aussi bien que la condition donnée au Théorème A.1 n'ont pas vocation à être testées numériquement car une telle vérification est combinatoire. En revanche, l'ensemble des matrices $X \in \mathbb{R}^{n \times p}$ qui ne satisfont pas la condition de position générale, ou encore l'ensemble

$$\{X \in \mathbb{R}^{n \times p} \mid \exists y \in \mathbb{R}^n \exists \gamma > 0, \mathcal{S}_{\gamma\text{pen}}(y) \text{ n'est pas un singleton}\},$$

avec pen une norme polyédrique, est négligeable pour la mesure de Lebesgue sur $\mathbb{R}^{n \times p}$ (voir le Lemme 1.4 de Tibshirani (2013) et la Proposition 3 de Schneider et Tardivel (2022)). Ainsi, par exemple, faire l'hypothèse que le problème d'optimisation LASSO ou SLOPE a une solution unique, ou encore supposer l'unicité du chemin des solutions, relativement à $\gamma > 0$, de ces estimateurs est raisonnable. D'autres approches différentes du Théorème A.1 enrichissent l'étude sur l'unicité des solutions au problème des moindres carrés pénalisé. Par exemple, les articles de Mousavi et Shen (2019); Fadili *et al.* (2023) proposent des méthodes numériques pour vérifier l'unicité d'une solution du problème $\mathcal{S}_{\gamma\text{pen}}(y)$. D'autres travaux s'intéressent à la description géométrique de

cet ensemble de solutions (Barbara *et al.*, 2019; Boyer *et al.*, 2019); par exemple l'ensemble des solutions du problème LASSO généralisé est un polyèdre. Ce dernier résultat est pertinent lorsque l'ensemble des solutions n'est pas un singleton. Enfin, pour conclure sur une note originale, le jeu de données réelles *giset*te fournit une matrice $X \in \mathbb{R}^{5000 \times 6000}$ et un vecteur $y \in \mathbb{R}^{5000}$ pour lesquels le problème LASSO n'a pas une unique solution (Dupuis et Vaiter, 2023). Cela illustre que le cas particulier où $\mathcal{S}_{\gamma\text{pen}}(y)$ n'est pas un singleton ne peut pas être complètement mis sous le tapis.

Bibliographie

- Taylor B ARNOLD et Ryan J TIBSHIRANI : Efficient implementations of the generalized lasso dual path algorithm. *Journal of Computational and Graphical Statistics*, 25(1):1–27, 2016.
- Francis R BACH : Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research*, 9(6), 2008.
- Abdessamad BARBARA, Abderrahim JOURANI et Samuel VAITER : Maximal solutions of sparse analysis regularization. *Journal of Optimization Theory and Applications*, 180:374–396, 2019.
- Adi BEN-ISRAEL et Thomas NE GREVILLE : *Generalized inverses : theory and applications*, volume 15. Springer Science & Business Media, 2003.
- Yoav BENJAMINI et Yosef HOCHBERG : Controlling the false discovery rate : a practical and powerful approach to multiple testing. *Journal of the Royal statistical society : series B (Methodological)*, 57(1):289–300, 1995.
- Quentin BERTRAND, Quentin KLOPFENSTEIN, Mathieu BLONDEL, Samuel VAITER, Alexandre GRAMFORT et Joseph SALMON : Implicit differentiation of lasso-type models for hyperparameter optimization. *In International Conference on Machine Learning*, pages 810–821. PMLR, 2020.
- Quentin BERTRAND, Quentin KLOPFENSTEIN, Mathurin MASSIAS, Mathieu BLONDEL, Samuel VAITER, Alexandre GRAMFORT et Joseph SALMON : Implicit differentiation for fast hyperparameter selection in non-smooth convex learning. *The Journal of Machine Learning Research*, 23(1):6680–6722, 2022.
- Małgorzata BOGDAN, Xavier DUPUIS, Piotr GRACZYK, Bartosz KOŁODZIEJEK, Tomasz SKALSKI, Patrick TARDIVEL et Maciej WILCZYŃSKI : Pattern recovery by slope. *arXiv preprint arXiv :2203.12086*, 2022.
- Małgorzata BOGDAN, Ewout VAN DEN BERG, Chiara SABATTI, Weijie SU et Emmanuel J CANDÈS : Slope—adaptive variable selection via convex optimization. *The annals of applied statistics*, 9(3):1103, 2015.
- Howard D BONDELL et Brian J REICH : Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with oscar. *Biometrics*, 64(1):115–123, 2008.
- Claire BOYER, Antonin CHAMBOLLE, Yohann De CASTRO, Vincent DUVAL, Frédéric DE GOURNAY et Pierre WEISS : On representer theorems and convex regularization. *SIAM Journal on Optimization*, 29(2):1260–1281, 2019.
- Peter BÜHLMANN et Sara VAN DE GEER : *Statistics for high-dimensional data : methods, theory and applications*. Springer Science & Business Media, 2011.
- Emmanuel J CANDÈS et Yaniv PLAN : Near-ideal model selection by ℓ_1 minimization. *The Annals of Statistics*, 37:2145 – 2177, 2009.

- Shaobing CHEN et David DONOHO : Basis pursuit. *In Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers*, volume 1, pages 41–44. IEEE, 1994.
- Guillaume DELTREIL, Patrick TARDIVEL, Piotr GRACZYK, Mikael ESCOBAR-BACH et Alexis DESCATHA : How to use biomechanical job exposure matrices with job history to access work exposure for musculoskeletal disorders? application of mathematical modeling in severe knee pain in the constances cohort. *International Journal of Environmental Research and Public Health*, 19(23):16217, 2022.
- Jeffrey Samuel DOKER : *Geometry of generalized permutohedra*. Thèse de doctorat, University of California, Berkeley, 2011.
- Charles DOSSAL : A necessary and sufficient condition for exact sparse recovery by ℓ_1 minimization. *Comptes Rendus Mathématique*, 350(1-2):117–120, 2012.
- Xavier DUPUIS et Patrick TARDIVEL : The solution path of slope. *In à paraître dans Artificial Intelligence and Statistics*, 2024.
- Xavier DUPUIS et Patrick JC TARDIVEL : Proximal operator for the sorted ℓ_1 norm : Application to testing procedures based on slope. *Journal of Statistical Planning and Inference*, 221:1–8, 2022.
- Xavier DUPUIS et Samuel VAITER : The geometry of sparse analysis regularization. *SIAM Journal on Optimization*, 33(2):842–867, 2023.
- Clément ELVIRA et Cédric HERZET : Safe rules for the identification of zeros in the solutions of the slope problem. *SIAM Journal on Mathematics of Data Science*, 5(1):147–173, 2023.
- Karl EWALD et Ulrike SCHNEIDER : On the distribution, model selection properties and uniqueness of the lasso estimator in low and high dimensions. *Electronic Journal of Statistics*, 14:944 – 969, 2020.
- Jalal FADILI, Tran TA NGHIA et Trinh TT TRAN : Sharp, strong and unique minimizers for low complexity robust recovery. *Information and Inference : A Journal of the IMA*, 12(3):iaad005, 2023.
- Mario FIGUEIREDO et Robert NOWAK : Ordered weighted l1 regularized regression with strongly correlated covariates : Theoretical aspects. *In Artificial Intelligence and Statistics*, pages 930–938. PMLR, 2016.
- Jean-Jacques FUCHS : On sparse representations in arbitrary redundant bases. *IEEE transactions on Information theory*, 50(6):1341–1344, 2004.
- Manlio GAUDIOSO, Enrico GORGONE et J-B HIRIART-URRUTY : Feature selection in svm via polyhedral k-norm. *Optimization letters*, 14(1):19–36, 2020.
- Jean Charles GILBERT : On the solution uniqueness characterization in the l1 norm and polyhedral gauge recovery. *Journal of Optimization Theory and Applications*, 172:70–101, 2017.
- Christophe GIRAUD : *Introduction to high-dimensional statistics*. CRC Press, 2021.
- Thomas GODLAND et Zakhar KABLUCHKO : Projections and angle sums of belt polytopes and permutohedra. *Results in Mathematics*, 78(4):140, 2023.
- Piotr GRACZYK, Ulrike SCHNEIDER, Tomasz SKALSKI et Patrick TARDIVEL : Pattern recovery in penalized and thresholded estimation and its geometry. *arXiv preprint arXiv :2307.10158*, 2023.
- Bin GU, Guodong LIU et Heng HUANG : Groups-keeping solution path algorithm for sparse regression with automatic feature grouping. *In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 185–193, 2017.

- Godfrey Harold HARDY, John Edensor LITTLEWOOD et George PÓLYA : *Inequalities*. Cambridge university press, 1952.
- David HARRISON JR et Daniel L RUBINFELD : Hedonic housing prices and the demand for clean air. *Journal of environmental economics and management*, 5(1):81–102, 1978.
- Trevor HASTIE, Robert TIBSHIRANI et Martin WAINWRIGHT : *Statistical learning with sparsity : the lasso and generalizations*. CRC press, 2015.
- Jean-Baptiste HIRIART-URRUTY et Claude LEMARÉCHAL : *Fundamentals of convex analysis*. Springer Science & Business Media, 2004.
- Philipp J KREMER, Damian BRZYSKI, Małgorzata BOGDAN et Sandra PATERLINI : Sparse index clones via the sorted ℓ_1 -norm. *Quantitative finance*, 22(2):349–366, 2022.
- Johan LARSSON, Malgorzata BOGDAN et Jonas WALLIN : The strong screening rule for slope. *Advances in neural information processing systems*, 33:14592–14603, 2020.
- Karim LOUNICI : Sup-norm convergence rate and sign concentration property of lasso and dantzig estimators. *Electronic Journal of Statistics*, 2:90 – 102, 2008.
- Julien MAIRAL et Bin YU : Complexity analysis of the lasso regularization path. *In Proceedings of the 29th International Conference on Machine Learning*, pages 353–360, 2012.
- Kentaro MINAMI : Degrees of freedom in submodular regularization : A computational perspective of stein’s unbiased risk estimate. *Journal of Multivariate Analysis*, 175:104546, 2020.
- Seyedahmad MOUSAVI et Jinglai SHEN : Solution uniqueness of convex piecewise affine functions based optimization with applications to constrained ℓ_1 minimization. *ESAIM : Control, Optimisation and Calculus of Variations*, 25:56, 2019.
- Shunichi NOMURA : An exact solution path algorithm for slope and quasi-spherical oscar. *arXiv preprint arXiv :2010.15511*, 2020.
- Michael R OSBORNE, Brett PRESNELL et Berwin A TURLACH : On the lasso and its dual. *Journal of Computational and Graphical statistics*, 9(2):319–337, 2000.
- Saharon ROSSET et Ji ZHU : Piecewise linear regularized solution paths. *The Annals of Statistics*, pages 1012–1030, 2007.
- Ulrike SCHNEIDER et Patrick TARDIVEL : The geometry of uniqueness, sparsity and clustering in penalized estimation. *Journal of Machine Learning Research*, 23(331):1–36, 2022.
- Amir SEPEHRI et Naftali HARRIS : The accessible lasso models. *Statistics*, 51(4):711–721, 2017.
- Dhruv B SHARMA, Howard D BONDELL et Hao Helen ZHANG : Consistent group identification and variable selection in regression with correlated predictors. *Journal of Computational and Graphical Statistics*, 22(2):319–340, 2013.
- Tomasz SKALSKI : *Aspects Géométriques et Combinatoires de Modeles Statistiques*. Thèse de doctorat, Angers, 2023.
- Tomasz SKALSKI, Piotr GRACZYK, Bartosz KOŁODZIEJEK et Maciej WILCZYŃSKI : Pattern recovery and signal denoising by slope when the design matrix is orthogonal. *arXiv preprint arXiv :2202.08573*, 2022.

- Atsumori TAKAHASHI et Shunichi NOMURA : Efficient path algorithms for clustered lasso and oscar. *arXiv preprint arXiv :2006.08965*, 2020.
- Patrick TARDIVEL : *Représentation parcimonieuse et procédures de tests multiples : application à la métabolomique*. Thèse de doctorat, Université Paul Sabatier-Toulouse III, 2017.
- Patrick TARDIVEL, Rémi SERVIEN et Didier CONCORDET : Powerful multiple testing procedures derived from hyperrectangular confidence regions having a minimal volume. *Journal de la Société Française de Statistique*, 162(1):2–21, 2021.
- Patrick JC TARDIVEL et Małgorzata BOGDAN : On the sign recovery by least absolute shrinkage and selection operator, thresholded least absolute shrinkage and selection operator, and thresholded basis pursuit denoising. *Scandinavian Journal of Statistics*, 49(4):1636–1668, 2022.
- Patrick JC TARDIVEL, Cécile CANLET, Gaëlle LEFORT, Marie TREMBLAY-FRANCO, Laurent DEBRAUWER, Didier CONCORDET et Rémi SERVIEN : Asics : an automatic method for identification and quantification of metabolites in complex 1d 1 h nmr spectra. *Metabolomics*, 13:1–9, 2017.
- Patrick JC TARDIVEL, Rémi SERVIEN et Didier CONCORDET : Sparsest representations and approximations of an underdetermined linear system. *Inverse Problems*, 34(5):055002, 2018.
- Patrick JC TARDIVEL, Rémi SERVIEN et Didier CONCORDET : Simple expressions of the lasso and slope estimators in low-dimension. *Statistics*, 54(2):340–352, 2020.
- Robert TIBSHIRANI : Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society : Series B (Methodological)*, 58(1):267–288, 1996.
- Ryan J TIBSHIRANI : The lasso problem and uniqueness. *Electronic Journal of Statistics*, 7:1456 – 1490, 2013.
- Ryan J. TIBSHIRANI et Jonathan TAYLOR : The solution path of the generalized lasso. *The Annals of Statistics*, 39(3):1335 – 1371, 2011.
- Samuel VAITER, Mohammad GOLBABAEE, Jalal FADILI et Gabriel PEYRÉ : Model selection with low complexity priors. *Information and Inference : A Journal of the IMA*, 4(3):230–287, 2015.
- Samuel VAITER, Gabriel PEYRÉ et Jalal FADILI : Model consistency of partly smooth regularizers. *IEEE Transactions on Information Theory*, 64(3):1725–1737, 2017.
- Martin J WAINWRIGHT : Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.
- Xiangrong ZENG et Mário AT FIGUEIREDO : Decreasing weighted sorted ℓ_1 regularization. *IEEE Signal Processing Letters*, 21(10):1240–1244, 2014.
- Peng ZHAO et Bin YU : On model selection consistency of lasso. *The Journal of Machine Learning Research*, 7:2541–2563, 2006.
- Hui ZOU : The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.