

WROCLAW UNIVERSITY OF SCIENCE AND TECHNOLOGY

Faculty of Pure and Applied Mathematics

UNIVERSITÉ D'ANGERS

Laboratoire Angevin de REcherche en MATHématiques, UMR CNRS 6093

DOCTORAL DISSERTATION

Tomasz SKALSKI

Geometric and Combinatorial Aspects of Statistical Models

Supervisor: dr hab. Maciej WILCZYŃSKI, prof. WUST

Supervisor: prof. dr hab. Piotr GRACZYK

Co-advisor in France: dr Patrick TARDIVEL

Wrocław, March 2023

POLITECHNIKA WROCŁAWSKA

Wydział Matematyki

UNIVERSITÉ D'ANGERS

Laboratoire Angevin de REcherche en MATHématiques, UMR CNRS 6093

ROZPRAWA DOKTORSKA

Tomasz SKALSKI

Geometryczne i kombinatoryczne zagadnienia modeli statystycznych

Promotor rozprawy: dr hab. Maciej WILCZYŃSKI, prof. uczelni

Promotor rozprawy: prof. dr hab. Piotr GRACZYK

Promotor pomocniczy we Francji: dr Patrick TARDIVEL

Wrocław, marzec 2023

UNIVERSITÉ D'ANGERS

Laboratoire Angevin de REcherche en MATHématiques, UMR CNRS 6093

POLITECHNIKA WROCLAWSKA

Wydział Matematyki

THÈSE DE DOCTORAT

Tomasz SKALSKI

Aspects Géométriques et Combinatoires de Modeles Statistiques

Directeur de thèse: prof. Piotr GRACZYK

Co-directeur de thèse: prof. Maciej WILCZYŃSKI

Co-encadrant en France: Patrick TARDIVEL, MCF

Wrocław, Mars 2023

Acknowledgements

First of all, I would like to express the gratitude to my supervisors Professors Maciej Wilczyński and Piotr Graczyk for their guidance and support during the time of my PhD studies. Thank you for the opportunity of sharing the studies between Wrocław and Angers and for your encouragement to submit the results of our long and fruitful discussions for publication. I would like to thank Dr. Patrick Tardivel for agreeing to be the co-encadrant in France and to extend our research to discussions in Dijon. Many thanks to Professor Krzysztof Bogdan for his guidance from the beginning of the PhD studies, for many inspiring discussions and for initiating the co-tutelle agreement with Université d'Angers. I would like to express my sincere gratitude to Prof. Małgorzata Bogdan for introducing me to the SLOPE estimator. I owe special appreciation to the late Prof. H el ene Massam, who agreed to take care of my journey through the theory of graphical models.

I owe many thanks to other co-authors of my research articles: Dr. Hab. Bartosz Kołodziejek, Dr. Michał Bosy, Dr. Hab. Ulrike Schneider and Dr. Xavier Dupuis for broadening the scope of my research and for letting me gather more knowledge. Our co-operation was a pleasure. I would like to thank the organizers and participants of seminars related to the topic of the dissertation, which helped us to exchange the current knowledge, as well in Wrocław (Random Graphs and Discrete Structures, Mathematical Statistics, Theory of Markov Semigroups and Schr odinger Operators), as in Angers (Probability and Statistics seminar, Doctorants' seminar) and on-line (Statistical Learning Seminars).

Special thanks are due to Dr. Cl ement Elvira, Dr. Hab. C edric Herzet, Dr. Hab. Wojciech Rejchel, Dr. Tomasz Stroiński, Dr. Samuel Vaiter and Prof. Piotr Zwiernik for mind-opening discussions, which helped me seeing the other branches of the currently presented research. I would also want to thank the anonymous referees of our articles for their insightful comments and for improving them.

This dissertation would be not possible to complete without a good atmosphere in the academic community. I would like to thank my colleagues from 6.06 and KNM for our long discussions on mathematics, life et al. Special thanks are due to prof. Janusz G orniak for his enormous help during the beginning of my academic adventure and to Monika Kaczmarz for her invaluable support and care of PhD students and their initiatives.

I would also like to thank everyone, who helped my to adapt in Angers, especially to Małgorzata Graczyk and Alexandra Le Petitcorps.

I would like to thank from all my heart to my family, especially to my Parents, who taught me how to count and think and supported my through the PhD, even during the harder times. Thanks are also due to all the people whom I've met while pursuing my hobbies, especially to the Linguistics Olympiad community and the Vytautas' Light Cavalry.

The major part of the research was supported by a French Government Scholarship. The research in Wrocław University of Science and Technology was supported in part by grant 049M/0010/19 from WUST. The research in Université d'Angers was supported by Centre Henri Lebesgue, program ANR-11-LABX-0020-0. Thank you for helping me to focus on research.

Contents

Summary	vii
Streszczenie (Summary in Polish)	ix
Résumé (Summary in French)	xi
Wstęp (Introduction in Polish)	xiii
PL.1 Penalizowana regresja liniowa	xiii
PL.2 SLOPE	xv
PL.2.1 Zgodność	xvi
PL.2.2 Warunek niereprezentowalności	xvi
PL.2.3 Penalizowana regresja liniowa i jej geometria	xvii
PL.3 Dyskretne rodziny wykładnicze	xviii
PL.4 Modele graficzne	xix
PL.5 Plan rozprawy	xx
Introduction (Introduction in French)	xxiii
FR.1 Régression Linéaire Pénalisée	xxiii
FR.2 SLOPE	xxv
FR.2.1 Consistance	xxvi
FR.2.2 Condition d'irreprésentabilité	xxvii
FR.2.3 Géométrie de la régression linéaire pénalisée	xxvii
FR.3 Familles exponentielles discrètes	xxviii
FR.4 Modèles graphiques	xxix
FR.5 Plan de la thèse	xxx
1 Introduction	1
1.1 Penalized Linear Regression	1
1.2 SLOPE	3
1.2.1 Consistency	4
1.2.2 Irrepresentability condition	4
1.2.3 Geometry of penalized linear regression	4
1.3 Discrete exponential families	6
1.4 Graphical models	7
1.5 Plan of the dissertation	8

2	Preliminaries and basic notions on penalized linear regression	11
2.1	Notation	11
2.2	Penalized linear regression	12
2.2.1	Linear regression	12
2.3	Convex Polytopes and cones	14
2.3.1	Convex analysis	14
2.3.2	Polytopes	16
2.3.3	Pattern equivalence class	16
2.3.4	Normal cones	17
2.3.5	Polyhedral gauges	18
2.3.6	Thresholded penalized least squares estimation	19
2.3.7	Permutahedron	19
2.3.8	Basics on Moore-Penrose inverse	19
2.3.9	Functional analysis	20
2.3.10	Tools from optimization	20
2.3.11	Tools from probability	21
3	Pattern recovery and signal denoising by SLOPE when the design matrix is orthogonal	25
3.1	Introduction	25
3.1.1	Introduction and motivations	25
3.1.2	Outline	26
3.2	Approach by minimax theorem	27
3.2.1	Technical results	27
3.2.2	Saddle point	28
3.2.3	SLOPE solution when X has full column rank	28
3.3	Properties of SLOPE in the orthogonal design	30
3.3.1	SLOPE vs. OLS	30
3.4	Asymptotic properties of SLOPE	32
3.4.1	Strong consistency of the SLOPE estimator	33
3.4.2	Asymptotic pattern recovery in the orthogonal design	35
3.5	Numerical experiment	37
3.6	Appendix	39
4	Pattern recovery by SLOPE	41
4.1	Introduction	41
4.1.1	History of SLOPE	41
4.1.2	Our contribution	43
4.1.3	Motivation	43
4.2	Preliminaries on clustering properties by SLOPE	45
4.2.1	Clustered design matrix and clustered parameter	45
4.2.2	Sorted ℓ_1 norm, dual sorted ℓ_1 norm and subdifferential	46
4.2.3	Characterization of SLOPE solutions	47
4.3	Characterization of pattern recovery by SLOPE	47
4.3.1	SLOPE irrepresentability condition	48
4.4	Geometrical interpretation of Irrepresentability Condition	50
4.5	Asymptotics of pattern recovery and pattern consistency	55
4.5.1	X is a fixed matrix	55

4.5.2	X is random, p is fixed, n tends to infinity	57
4.6	Strong consistency of SLOPE and its pattern	59
4.6.1	Refined results on strong consistency of the SLOPE pattern	59
4.6.2	Strong consistency of the SLOPE estimator	61
4.7	Simulation study	63
4.7.1	Sharp upper bound when X is orthogonal	63
4.7.2	Limiting probability when X is asymptotically orthogonal	64
4.8	Discussion	64
4.9	Appendix — Proofs	66
4.9.1	Proof of Proposition 4.2.1	66
4.9.2	Proof of Proposition 4.2.2	67
4.9.3	Proof of Theorem 4.3.1	68
4.9.4	Proof of Corollary 4.3.1	69
4.9.5	Proof of Theorem 4.5.1	69
4.9.6	Proofs from Section 4.5.2	70
5	Geometry of Pattern Recovery by Penalized and Thresholded Estimators	75
5.1	Introduction	75
5.1.1	Pattern recovery by penalized least squares estimators	75
5.1.2	Pattern recovery by a thresholded estimator	77
5.2	Geometry of pattern equivalence	78
5.2.1	Pattern equivalence classes and normal cones	79
5.2.2	Model subspace recovery	80
5.3	Examples of polyhedral gauges and their patterns	81
5.4	Pattern recovery in penalized estimation	84
5.4.1	Accessibility condition	84
5.4.2	Noiseless recovery condition	85
5.4.3	Irrepresentability Condition for polyhedral gauges	85
5.5	Pattern recovery by thresholded estimators	87
5.6	Full characterization of the uniform uniqueness	89
5.7	Numerical experiments	90
5.7.1	Numerical experiments for LASSO	90
5.7.2	Numerical experiments when the penalty term is the supremum norm	91
5.8	Appendix	93
5.8.1	Facts about real-valued polyhedral gauges	93
5.8.2	Proofs	95
6	Maximum likelihood estimation for discrete exponential families and random graphs	105
6.1	Introduction and preliminaries	105
6.1.1	Discrete exponential family	106
6.1.2	Alternative setting	108
6.2	Main results	108
6.2.1	Nonexistence of MLE	110
6.2.2	Linear programming	111
6.3	Applications	112
6.3.1	Rademacher functions	113
6.4	Random graphs	117

6.5	Applications to Walsh functions	120
6.6	Appendix	123
6.6.1	Proof of Lemma 6.1.1	123
6.6.2	Control by oscillations	124
6.6.3	Proof of Lemma 6.4.1	124
6.6.4	Proof of Lemma 6.4.2	125
6.6.5	Proof of Lemma 6.5.1	126
6.6.6	Propagation of extrema, relative interior and the criterion of Barndorff-Nielsen	126
7	On Laplacian of Graphical Models in Various Graphs	129
7.1	Introduction	129
7.2	Trees	129
7.3	Discussion and non-tree graphs	131
7.3.1	Non-tree graphs	131
7.3.2	Eigenvalues of augmented Laplacian	134
7.3.3	Discussion	134
	Bibliography	137

Summary

This dissertation concerns new applications of discrete geometry and combinatorics in modern statistics. First of them focuses on one of widely used remedies to the inevitable growth of data, that is the use of penalized linear regression methods. With an aim to recover the needed properties possessed by the vector of regression coefficients, we start our discussion with the Sorted ℓ_1 Penalized Estimator (SLOPE), which was proposed almost a decade ago. Especially, we examine the notion of the SLOPE pattern, which maintains the information about the support, sign and ranking between the regression coefficients. In particular, it preserves the clusters of coefficients with the same absolute value. In Chapter 3 we provide the conditions, under which SLOPE recovers the set of relevant covariables and the clusters when the design matrix is orthogonal. We also derive new results on the strong consistency of the SLOPE estimator and its pattern. Chapter 4 extends the discussion on SLOPE to a general class of fixed design matrices. We provide the SLOPE irrepresentability condition, which is necessary and sufficient for the pattern recovery in the noiseless case and illustrate it geometrically. Later on, we consider the case of asymptotic growth of the number of explanatory variables and of the incremental error. In Chapter 5 we study the wider class of penalized estimators, called the polyhedral gauges. It allows one to use the notions from the geometry of polyhedra to generalize the notion of the pattern and the results on its recovery. Chapter 6 is articulated around the existence of the Maximum Likelihood Estimator (MLE) for discrete exponential families. We give its new characterization based on the notion of the set of uniqueness. Later on, we inspect the size of independent identically distributed samples which is needed to ensure its existence with high probability. For that reason we use the notions from the analysis of discrete hypercubes and apply our results in the environment of random graphs. Last of the chapters connects the theory of graphical models in statistics with the notion of graph Laplacian matrices and discretized Wiener processes. The thesis is based on three already published articles and two preprints, which are available on-line.

Streszczenie

Niniejsza rozprawa poświęcona jest nowym zastosowaniom geometrii dyskretnej i kombinatoryki w nowoczesnej statystyce. Pierwsze z nich skupione jest na jednym z popularniejszych rozwiązań na radzenie z ciągłym przyrostem danych, jest nim penalizowana regresja liniowa. Mając na celu odtworzenie potrzebnych nam własności wektora współczynników regresji, rozpoczynamy dyskusję od estymatora SLOPE (Sorted ℓ_1 Penalized Estimator), który został wprowadzony w poprzedniej dekadzie. Szczególną uwagę poświęcamy pojęciu wzorca SLOPE, który zachowuje informację o nośniku, znaku i rankingu między współczynnikami regresji. Informuje on również o klastrach współczynników o tej samej wartości bezwzględnej. W rozdziale trzecim podajemy warunki, dla których SLOPE poprawnie odtwarza nośnik oraz klastry wektora współczynników regresji przy ortogonalnej macierzy eksperymentu. Przy tym założeniu wyprowadzamy też nowe wyniki dotyczące mocnej zgodności estymatora SLOPE i jego wzorca. Rozdział czwarty rozszerza dyskusję na temat SLOPE, pomijając założenie o ortogonalności macierzy eksperymentu. Wprowadzamy warunek niereprezentowalności dla SLOPE, który jest konieczny i dostateczny do odtworzenia wzorca w przypadku braku szumu, po czym ilustrujemy ten warunek geometrycznie. Następnie rozważamy przypadek asymptotycznego przyrostu liczby zmiennych objaśniających i szumu rosnącego inkrementalnie. W rozdziale piątym omawiamy szerszą klasę penalizowanych estymatorów zwaną polyhedral gauges. Pozwala ona na wykorzystanie twierdzeń z geometrii wielościanów do uogólnienia pojęcia wzorca i wyników dotyczących jego odtwarzania. Rozdział szósty dotyczy istnienia estymatora największej wiarygodności (MLE) w dyskretnej rodzinach wykładniczych. Podajemy jego pełną charakteryzację za pomocą pojęcia zbioru jednoznaczności. Następnie badamy rozmiar próby niezależnych zmiennych losowych o tym samym rozkładzie, która zapewnia istnienie MLE z wysokim prawdopodobieństwem. W tym celu wykorzystujemy narzędzia z analizy hipersześcianów dyskretnej i stosujemy otrzymane wyniki w modelach wykładniczych grafów losowych. Ostatni z rozdziałów skupiony jest na połączeniu między teorią modeli graficznych w statystyce, a pojęciami laplasjanu grafu oraz dyskretyzacji procesów Wienera. Rozprawa jest oparta na trzech opublikowanych artykułach oraz dwóch preprintach dostępnych on-line.

Résumé

Cette thèse traite des applications de la combinatoire et de la géométrie discrète aux statistiques modernes. La première application porte sur l'une des approches fréquemment utilisées pour faire face au volume croissant des données, à savoir l'utilisation de méthodes de régression linéaire pénalisée. Dans le but de retrouver la structure que possède le vecteur des coefficients de régression, nous commençons notre discussion par l'estimateur SLOPE (Sorted ℓ_1 Penalized Estimator), qui a été proposé il y a presque dix ans. En particulier, nous examinons la notion de schéma SLOPE, qui conserve l'information sur le support, le signe et le classement des coefficients de régression. En particulier, il préserve les groupes de coefficients ayant la même valeur absolue. Dans le chapitre 3, nous fournissons les conditions, sous lesquelles SLOPE récupère l'ensemble des variables pertinentes et des groupes lorsque la matrice de planification est orthogonale. Nous déduisons également de nouveaux résultats sur la forte consistance de l'estimateur SLOPE et de son schéma. Le chapitre 4 étend la discussion sur l'estimateur SLOPE à une classe générale de matrices de planification fixes. Nous fournissons la condition d'irreprésentabilité du SLOPE, qui est nécessaire et suffisante pour la recouvrement du schéma dans le cadre non-bruité et nous illustrons cette propriété géométriquement. Dans un deuxième temps, nous considérons le cas asymptotique lorsque le nombre de variables explicatives tend vers l'infini et que l'erreur est incrémentale. Dans le chapitre 5, nous étudions la classe plus large des estimateurs pénalisés, appelés les jauges polyédriques. Elle permet d'utiliser les notions issues de la géométrie des polyèdres pour généraliser la notion du schéma et les résultats sur sa récupération. Le chapitre 6 s'articule autour de l'existence de l'estimateur du maximum de vraisemblance (MLE) pour les familles exponentielles discrètes. Nous donnons sa nouvelle caractérisation basée sur la notion d'ensemble d'unicité. Par la suite, nous inspectons la taille des échantillons indépendants identiquement distribués qui est nécessaire pour assurer son existence avec une grande probabilité. Pour cela, nous utilisons les notions issues de l'analyse des hypercubes discrets et appliquons nos résultats dans l'environnement des graphes aléatoires. Le dernier des chapitres relie la théorie des schémas graphiques en statistique avec la notion de matrices laplaciennes de graphes et de processus de Wiener discrétisés. La thèse est basée sur trois articles déjà publiés et deux prépublications, qui sont disponibles en ligne.

Wstęp

Obecnie można zauważyć szybki i nieunikniony przyrost danych, zarówno pod względem liczby obserwacji, jak i pod względem liczby sposobów, za pomocą których da się je zmierzyć. Duże zbiory danych i ich analiza rosną na znaczeniu w życiu codziennym, dzięki czemu statystyka matematyczna oraz nauki związane z analizą danych zyskują coraz większe zainteresowanie w innych działach matematyki i jej zastosowaniach. Jednak statystyka i analiza danych nie powinny być traktowane w oderwaniu od innych działów matematyki. Kluczowym punktem badań niniejszej rozprawy będzie znajdowanie nowych połączeń między współczesną statystyką, geometrią i kombinatoryką.

Dyskusja podjęta w tej rozprawie rozpoczyna się od rozważania przestrzeni euklidesowych o skończonym wymiarze. Zwracamy uwagę na przypadki, w których zbiór wartości danego estymatora można podzielić na skończoną liczbę podzbiorów w taki sposób, żeby istniała bijekcja między nimi, a interesującymi nas własnościami tego estymatora. Taki sposób predykcji obserwacji jest znany jako problem klasyfikacji [32], nad którym badania zostały rozpoczęte przez Fishera [77] i który jest stosowany w prawie każdym dziale nauk związanych z przetwarzaniem danych.

PL.1 Penalizowana regresja liniowa

Za jedno z istotnych zastosowań problemu klasyfikacji możemy uznać wersję regresji liniowej, w której zamiast dokładnego oszacowania wartości nieznanego wektora interesują nas jego wybrane właściwości. W modelu regresji liniowej wielorakiej mającym n obserwacji i p zmiennych objaśniających zakładamy, że zmienna objaśniana $\mathbf{Y} = (Y_1, \dots, Y_n)' \in \mathbb{R}^n$ jest postaci $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, gdzie $\mathbf{X} \in \mathbb{R}^{n \times p}$ jest macierzą eksperymentu, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$ jest nieznanym wektorem współczynników regresji, a $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)' \in \mathbb{R}^n$ jest losowym wektorem błędu (szumu). Głównym zagadnieniem regresji liniowej są oszacowanie $\boldsymbol{\beta}$ oraz wydobycie jego istotnych własności. Szeroka klasa zagadnień opartych na wyborze podzbioru współrzędnych $\boldsymbol{\beta}$, które będą najlepiej spełniać oczekiwania dotyczące estymatora i jego złożoności, znany jest jako problem wyboru modelu. Więcej informacji o tym zagadnieniu można znaleźć w przeglądowym artykule [49]. Przez większość dyskusji zakładamy, że wektor błędu $\boldsymbol{\varepsilon}$ jest losowy, a jego rozkład jest ciągły i symetryczny. Niektóre z przedstawionych wyników dotyczą również przypadku niezasumionego, w którym $\boldsymbol{\varepsilon} \equiv \mathbf{0}$, co pomoże w wyprowadzeniu wyników asymptotycznych w przypadku zasumionym. Dzięki istniejącym wynikom dotyczącym rozkładu normalnego możemy pokazać ulepszone wyniki dla szumu gaussowskiego $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Najbardziej klasyczna metoda regresji liniowej, metoda najmniejszych kwadratów (Ordinary Least Squares, OLS), została wprowadzona przez Legendre'a na początku XIX wieku [10, 97, 124]. Estymator OLS jest zdefiniowany jako wektor \mathbf{b} minimalizujący odległość euklidesową między \mathbf{Y} , a $\mathbf{X}\mathbf{b}$:

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2.$$

Popularność estymatora najmniejszych kwadratów ma miejsce ze względu na swoją prostotę, dokładność przy względnie małej liczbie zmiennych [33, str. 4] i wielu innych przydatnych statystycznych własności. Dla przykładu, jeżeli macierz $\mathbf{X}'\mathbf{X}$ jest odwracalna, a $\mathbf{Cov}(\mathbf{Y}) = \sigma^2\mathbf{I}_n$ dla $\sigma > 0$, to OLS jest najlepszym estymatorem w klasie nieobciążonych liniowych estymatorów (BLUE) [1, Theorem 2.7.1.] wektora β . W tym przypadku istnieje wzór jawny na $\hat{\beta}^{\text{OLS}}$, który można w dosyć łatwy sposób wyprowadzić [1, str. 28]:

$$\hat{\beta}^{\text{OLS}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}.$$

Po dodaniu założenia o gaussowskości wektora \mathbf{Y} estymator OLS staje się estymatorem największej wiarygodności (MLE) [1, str. 28]. W przypadku, gdy macierz eksperymentu \mathbf{X} jest ortogonalna, tzn. $\mathbf{X}'\mathbf{X} = c\mathbf{I}_p$, $c > 0$, powyższy wzór sprowadza się do $\hat{\beta}^{\text{OLS}} = \frac{1}{c}\mathbf{X}'\mathbf{Y}$. Przy powyższych założeniach OLS jest również mocno zgodnym estymatorem dla β [6]. Z drugiej strony, estymator ten nie jest jednoznacznie zdefiniowany, kiedy macierz $\mathbf{X}'\mathbf{X}$ jest nieodwracalna, co ma miejsce w wysokowymiarowym przypadku, gdy $p > n$. Ponadto, przy ogólnie przyjmowanych założeniach o błędzie ϵ , OLS nie jest estymatorem rzadkim, z prawdopodobieństwem 1 składa się on z p współrzędnych o parami różnych wartościach. Mniej klasyczne założenia o wektorze \mathbf{Y} i macierzy \mathbf{X} , przy których estymator nie jest jednoznaczny, można znaleźć między innymi w niedawno opublikowanym artykule Dupuisa i Vaitera [62]. W praktyce często wektor β składa się ze względnie małej liczby niezerowych współrzędnych, co w naturalny sposób sugeruje metody, które promują rzadkość wektora β rozumianą jako małą liczbę niezerowych współrzędnych lub opisywalność β za pomocą małej liczby parametrów. Istnieją różne propozycje rozwiązania powyższego problemu, m.in. porównanie pasujących modeli przy pomocy kryterium informacyjnego, np. BIC [158] lub AIC [2]. Innym podejściem jest penalizowana regresja liniowa postaci:

$$\hat{\beta} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}_\lambda(\mathbf{b}) \right],$$

gdzie $\text{pen}_\lambda(\mathbf{b})$ jest ustaloną nieujemną karą, którą można modyfikować za pomocą parametru $\lambda > 0$. To podejście jest wykorzystywane m.in. w neuroobrazowaniu [39], prognozowaniu cen energii elektrycznej [106, 180] oraz w matematyce finansowej w celu grupowaniu aktywów na podstawie ich korelacji częściowej z szeregiem czasowym stóp zwrotu z funduszy inwestycyjnych [116]. Pierwszym z zaproponowanych sposobów penalizowanej regresji była metoda wyboru najlepszego podzbioru [12, 102], gdzie kara $\text{pen}_\lambda(\mathbf{b}) = \lambda \|\mathbf{b}\|_0$ jest równa liczbie współrzędnych niezerowych wektora \mathbf{b} pomnożonej przez λ . Kłopot z jej zastosowaniem polega na tym, że dla dużych wartości p znalezienie najlepszego podzbioru w ten sposób jest problemem NP-trudnym [78]. Inną metodą na znalezienie rzadkiego rozwiązania jest metoda LASSO (skrót od Least Absolute Shrinkage and Selection Operator [47, 176]), w której karą nałożoną na estymator \mathbf{b} jest jego norma ℓ_1 pomnożona przez parametr $\lambda > 0$:

$$\hat{\beta}^{\text{LASSO}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_1 \right].$$

Estymator LASSO jest obciążony. Jako estymator ściągający (shrinkage estimator) sprowadza on współrzędne $\hat{\beta}^{\text{LASSO}}$ w stronę zera. Niektóre z nich zostają ściągnięte całkowicie do zera, co przynosi rezultat w postaci rzadszego estymatora. W przypadku, gdy macierz eksperymentu \mathbf{X} jest ortonormalna, tzn. $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$, jawny wzór na $\hat{\beta}^{\text{LASSO}}$, wyprowadzony w pracy Tibshiraniego [176], jest oparty na $\hat{\beta}^{\text{OLS}}$:

$$\hat{\beta}_i^{\text{LASSO}} = \text{sign}(\hat{\beta}_i^{\text{OLS}}) \max \{ |\hat{\beta}_i^{\text{OLS}}| - \lambda, 0 \}.$$

Aby zapewnić istnienie rozwiązania zagadnienia penalizowanej regresji liniowej, często używana jest kara wypukła, co umożliwia wykorzystanie narzędzi z analizy wypukłej. Zainteresowanym różnymi przykładami pomysłów na modyfikacje estymatora LASSO można polecić artykuł [78]. Warto zauważyć, że kary niewypukłe również są wykorzystywane m.in. przy wyborze najlepszego podzbioru [12, 102] lub estymatora SCAD (smoothly clipped absolute deviation [72]).

PL.2 SLOPE

Innym sposobem na zmniejszenie wymiaru jest estymator SLOPE (Sorted ℓ_1 Penalized Estimator [27, 26, 189]), który poza uogólnieniem metody LASSO, skleja takie same współczynniki regresji β oraz skorelowane kolumny macierzy \mathbf{X} . Zgodnie z nazwą, w estymatorze SLOPE norma ℓ_1 w funkcji kary została zastąpiona przez zdefiniowaną poniżej posortowaną normę ℓ_1 :

$$J_{\mathbf{\Lambda}}(\mathbf{b}) := \sum_{i=1}^p |b|_{(i)} \lambda_i,$$

gdzie $\lambda_1 > 0, \lambda_1 \geq \dots, \lambda_p \geq 0$, natomiast $|b|_{(1)} \geq \dots \geq |b|_{(p)}$ to wartości bezwzględne współrzędnych wektora \mathbf{b} posortowane malejąco. Szczególny przypadek SLOPE z $\mathbf{\Lambda}$ będącym ciągiem arytmetycznym jest znany pod nazwą OSCAR [29]. Zauważmy, że w parametr λ został zastąpiony w SLOPE przez nierosnący wektor $\mathbf{\Lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p)'$ parametrów (tuning vector). Możemy więc nie tylko modyfikować wielkość wektora $\mathbf{\Lambda}$, ale również jego kształt. Główną motywacją autorów metody SLOPE było testowanie p hipotez zerowych $H_i^0 : \beta_i = 0$ i kontrola współczynnika fałszywych odkryć (FDR control) zdefiniowanego jako oczekiwana proporcja między liczbą fałszywych odrzuceń hipotezy zerowej, a łączną liczbą jej odrzuceń (w przypadku braku odrzuceń przyjmuje się $\text{FDR} = 0$). Ponadto, SLOPE jest uogólnieniem poprzednich metod:

- $\lambda_1 = \dots = \lambda_p = 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OLS}}$,
- $\lambda_1 = \dots = \lambda_p > 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{LASSO}}$,
- $\mathbf{\Lambda}$ is an arithmetic sequence $\Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OSCAR}}$.

Dokładniejszy opis historii badań nad estymatorem SLOPE znajduje się w podrozdziale 4.1.1. W naszych badaniach skupiamy się na odtworzeniu istotnych właściwości wektora β za pomocą estymatora SLOPE, które są zakodowane w wektorze zwanym wzorcem SLOPE. Tak jak znak wektora β w przypadku LASSO, wzorec SLOPE można całkowicie opisać za pomocą sub-różniczkowej funkcji kary. Dokładniej mówiąc, niech k będzie liczbą klastrów wektora β , tzn. liczbą niezerowych różnych wartości współrzędnych $|\beta|$.

Definicja PL.2.1 (Wzorec SLOPE). *Wzorec SLOPE $\mathbf{patt} : \mathbb{R}^p \rightarrow \mathbb{Z}^p$ jest zdefiniowany następująco:*

$$\mathbf{patt}(\mathbf{b})_i = \text{sign}(b_i) \text{rank}(|b_i|),$$

gdzie $\text{rank}(|b_i|) \in \{1, 2, \dots, k\}$ jest zdefiniowany jako liczba takich $|c_j|$, dla których $|b_i| \geq |c_j|$, gdzie $|c_1|, |c_2|, \dots, |c_k|$, $k \leq p$ są niezerowymi różnymi wartościami spośród $|b_1|, \dots, |b_p|$. Przyjmujemy $\text{rank}(0) = 0$.

Wzorec SLOPE zachowuje informację nie tylko o znaku wektora, ale również o jego klastrach, tzn. zbiorach współrzędnych mających tę samą wartość bezwzględną oraz o hierarchii między tymi wartościami.

Fakt PL.2.1 (Własności wzorca SLOPE [156]).

- (a) Dla każdego $1 \leq l \leq \|\mathbf{patt}(\mathbf{b})\|_\infty$ istnieje takie j , że $|\mathbf{patt}(\mathbf{b})_j| = l$,
- (b) $\text{sign}(\mathbf{patt}(\mathbf{b})) = \text{sign}(\mathbf{b})$,
- (c) $|b_i| = |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| = |\mathbf{patt}(\mathbf{b})_j|$,
- (d) $|b_i| > |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| > |\mathbf{patt}(\mathbf{b})_j|$.

Example PL.2.2. $\mathbf{patt}((4, 0, -1.5, 1.5, -4)') = (2, 0, -1, 1, -2)'$.

Mówimy, że $\hat{\beta}^{SLOPE}$ odtwarza wzorzec β , kiedy

$$\mathbf{patt}\left(\hat{\beta}^{SLOPE}\right) = \mathbf{patt}(\beta).$$

W rozdziale trzecim i czwartym omawiamy nowe warunki konieczne i dostateczne na odtwarzanie wzorca SLOPE, jak również nowe wyniki o zgodności i mocnej zgodności SLOPE i jego wzorca w przypadku, gdy $n \geq p$.

PL.2.1 Zgodność

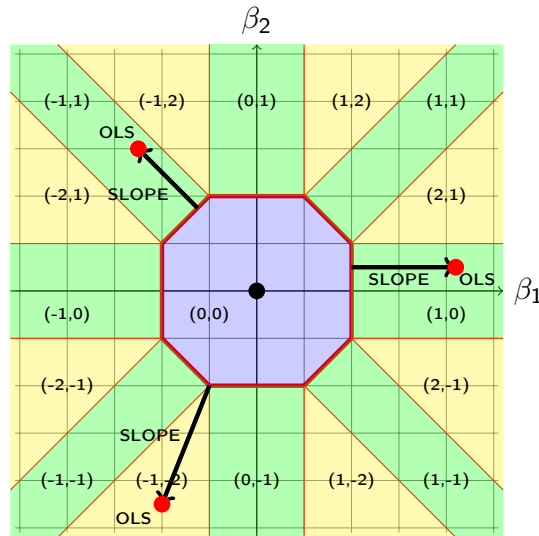
Głównym narzędziem do udowodnienia mocnej zgodności SLOPE w przypadku $n \geq p$ jest mocna zgodność estymatora największych kwadratów udowodniona m.in. w pracy Andersona i Taylora [6]. Główne wyniki dotyczące zgodności LASSO można znaleźć w następujących publikacjach:

- zgodność LASSO: Knight, Fu, 2000 [112],
- mocna zgodność LASSO: Chatterjee, Lahiri, 2011 [44],
Główne założenia: λ_n jest rzędu mniejszego niż n oraz $\mathbb{E}|\varepsilon| < \infty$.
- zgodność znaku LASSO: Zhao, Yu, 2006 [192].

Podrozdziały 3.4 oraz 4.6 niniejszej rozprawy dotyczą mocnej zgodności SLOPE i jego wzorca. Warto zauważyć, że w przypadku, gdy wektor \mathbf{A} jest stały (LASSO), wzorzec SLOPE nie jest zgodny, nawet jeśli założenia o zgodności wektora znaku są spełnione. Estymator LASSO nie odtwarza klastrów, które są elementem wzorca SLOPE. Zauważmy też, że mocna zgodność estymatora LASSO (SLOPE) nie implikuje mocnej zgodności jego znaku (ani jego wzorca). Jako kontrprzykład można rozważyć ciąg $((1/n, 1/n, -1/n, -1/n)')_{n \geq 1}$, który zbiega do $(0, 0, 0, 0)'$, podczas gdy jego znak i jego wzorzec SLOPE są równe $(1, 1, -1, -1)'$ dla wszystkich n .

PL.2.2 Warunek niereprezentowalności

Dla danych $n, p > 0$ zgodność zbioru niezerowych współrzędnych wektora $\hat{\beta}^{LASSO}$ jest prawie równoważna do warunku, w którym współrzędne spoza nośnika wektora β nie są reprezentowane przez współrzędne doń należące [192]. Z tego powodu warunek został nazwany warunkiem niereprezentowalności (irrepresentability condition). W rozdziale 4. rozprawy wprowadzony jest analogon powyższego warunku dla estymatora SLOPE.



$\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$ w przypadku ortogonalnym: $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ for $\Lambda = (2, 1)'$.

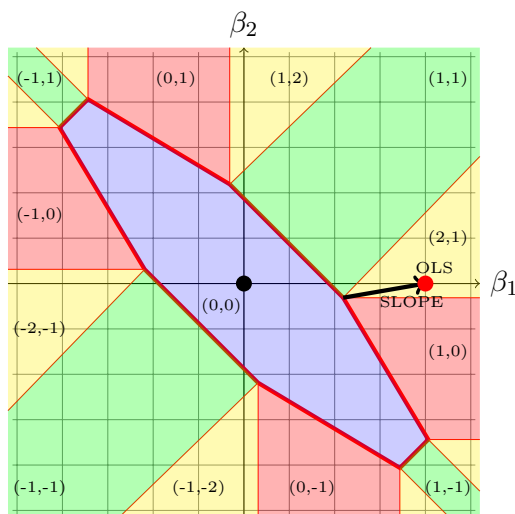
PL.2.3 Penalizowana regresja liniowa i jej geometria

Dla lepszego zrozumienia powiazań między $\hat{\beta}^{\text{SLOPE}}$, a $\hat{\beta}^{\text{OLS}}$, ilustrujemy je w przypadku niskowymiarowym $p \leq n$ dla $p = 2$. Tutaj wartości estymatora SLOPE jest równa różnicy między estymatorem najmniejszych kwadratów, a jednym z jego rzutów na $(\mathbf{X}'\mathbf{X})^{-1}C_\Lambda$, gdzie C_Λ jest kulą jednostkową w normie dualnej do J_Λ . Powyższy wynik ma łatwą interpretację w przypadku ortogonalnym, w którym zarówno \mathbf{X} , jak i rzut na C_Λ , są ortogonalne. Ta zależność ma rezultat w postaci jawnego wzoru na $\hat{\beta}^{\text{SLOPE}}$ w przypadku ortogonalnym, który został niedawno wprowadzony przez Tardivela, Serviena i Concordeta [175]. Więcej ilustracji związków między estymatorem SLOPE, a jego geometrią, można znaleźć w podrozdziale 4.4.

Zachodzi połączenie między estymatorami LASSO i SLOPE, a systemami pierwiastkowymi (root systems) wykorzystywanymi w analizie harmonicznej. Dokładniej mówiąc, można zauważyć, że kula jednostkowa w normie ℓ_∞ dualnej do kary stosowanej w LASSO oraz kula jednostkowa w normie J_Λ^* dualnej do normy SLOPE są proporcjonalne odpowiednio do powłok wypukłych orbit grup Weyla odpowiadających p -tej potędze systemu pierwiastkowego A_1 oraz systemowi B_p [90, 144]. Więcej na temat systemów pierwiastkowych można znaleźć m.in. w książce Helgasona [100] oraz w publikacjach [92, 63].

Geometrię penalizowanej regresji liniowej i odtwarzania wzorców jej estymatorów można rozważać w ogólniejszym przypadku, co robimy w rozdziale piątym. Mianowicie, estymator SLOPE można zaklasyfikować do jednej ze skończonej liczby klas, gdy kara jest postaci polyhedral gauge, tzn. nieujemną funkcją wypukłą równą maksimum skończonej liczby funkcji liniowych. Istotne właściwości estymatora, również nazwane jego wzorcem (pattern), można całkowicie opisać za pomocą subrózniczki kary. Ponieważ nierówności liniowe definiujące polyhedral gauge są spełnione przez przekrój skończonej liczby związanych z nimi półprzestrzeni, zachodzi tutaj silne powiązanie z teorią wielościanów, z którą można zapoznać się. m.in. w książkach Grubera [94], Grünbauma [95], Hiriart-Urrutiego i Lemarechala [101] oraz Zieglera [193].

W rozdziale piątym uogólniamy niektóre z nowo otrzymanych wyników SLOPE na powyższą klasę estymatorów. W tym celu rozważamy pojęcie klasy równoważności wzorca (pattern equivalence class), która pozwala nam opisać odtworzenie wzorca jako odtworzenie wartości sub-



$$\hat{\beta}^{\text{SLOPE}} \text{ oraz } \hat{\beta}^{\text{OLS}} \text{ dla } \mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 1 \end{bmatrix} \text{ and } \mathbf{\Lambda} = (2, 1)'.$$

różniczki kary pen.

Wprowadzamy i charakteryzujemy pojęcie osiągalności (accessibility) wzorca i wprowadzamy warunek niereprezentowalności dla polyhedral gauges. Z pomocą narzędzi z geometrii wielościanów, w szczególności pojęcia stożków normalnych ścian rozważanych wielościanów, rozważamy własności równoważności wzorca.

Następnie skupiamy się na odtwarzaniu modelu przez estymatory progowane. Są one uogólnieniem progowanego LASSO, jednak w odróżnieniu od swojego poprzednika ich celem jest zmniejszenie złożoności estymatora nie tylko poprzez sprowadzenie części współrzędnych do zera, ale również do innych jego uproszczeń związanych z danym wzorcem. Rozszerzamy też pojęcie jednorodnej jednoznaczności estymatora znanej dla norm wielościennych [156] do polyhedral gauges.

PL.3 Dyskretne rodziny wykładnicze

W rozdziale szóstym przechodzimy z przestrzeni euklidesowych do przestrzeni ze skończoną liczbą elementów (dyskretnych). Punktem naszych zainteresowań są rodziny wykładnicze zdefiniowane na skończonej przestrzeni stanów \mathcal{X} , $|\mathcal{X}| = K < \infty$. Z powiązaną dyskusją na temat przeliczalnych zbiorów można zapoznać się w pracy Jacobsena [105]. My rozważamy podprzestrzeń liniową \mathcal{B} przestrzeni funkcji $\mathbb{R}^{\mathcal{X}}$, do której należy funkcja stała dodatnia. Na zbiorze \mathcal{X} wprowadzamy też ściśle dodatnią wagę $\mu : \mathcal{X} \rightarrow (0, \infty)$. Rodzinę wykładniczą konstruujemy następująco:

Dla rzeczywistej funkcji ϕ definiujemy odpowiednio jej funkcję partycji oraz log-partycji:

$$Z(\phi) = \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x), \quad \psi(\phi) = \log Z(\phi),$$

a także gęstość wykładniczą

$$p = e(\phi) = e^{\phi - \psi(\phi)} = e^{\phi} / Z(\phi).$$

Rodzina wykładnicza rozpięta przez \mathcal{B} jest ukazana poniżej.

$$e(\mathcal{B}) := \{p = e(\phi) : \phi \in \mathcal{B}\}.$$

Dzięki temu możemy wprowadzić funkcję wiarygodności i log-wiarygodności. Ta ostatnia jest ściśle wklęsła, co zapewnia nam jednoznaczność MLE w przypadku jego istnienia. Istnienie estymatora nie jest zapewnione pomimo faktu, że funkcja wiarygodności jest ograniczona. W naszych badaniach wyprowadzamy nowe kryterium na istnienie MLE i wykorzystujemy je w rodzinach wykładniczych rozpiętych przez funkcje Rademachera i Walsha oraz w wykładniczych modelach grafów losowych.

Naszym głównym narzędziem jest pojęcie zbioru jednoznaczności, czyli takiego zbioru $U \subset \mathcal{X}$, że $\phi = 0$ jest jedyną funkcją z zadanej klasy, która jest równa zero na całym zbiorze U . Dokładniej, pokazujemy, że MLE istnieje dla $e(\mathcal{B})$ oraz próby i.i.d. $x_1, x_2, \dots, x_n \in \mathcal{X}$ wtedy i tylko wtedy, gdy zbiór $\{x_1, \dots, x_n\}$ jest jednoznaczności dla nieujemnego stożka funkcji $\mathcal{B}_+ := \{\phi \in \mathcal{B} : \phi \geq 0\}$. Powyższe kryterium podajemy również jako zagadnienie programowania liniowego.

Dodatkowo, w rozważanych zastosowaniach wyznaczamy asymptotyczne wielkości zbiorów jednoznaczności. W tym celu używamy klasycznych wyników z teorii grafów losowych oraz z problemu zbieracza kuponów (Coupon Collector's Problem) [118, str. 194-195], [68], and [143].

W rodzinie rozpiętej przez funkcje Walsha wykorzystujemy własności hipersześcianu $\{-1, 1\}^k$ ($k = \log_2 K$), jego podkości oraz związanego z nim grafu hipersześcianu.

PL.4 Modele graficzne

Jednym z działów nowoczesnej statystyki wykorzystujących zarówno MLE, jak i estymatorów penalizowanych jest teoria modeli graficznych (graphical models). Model graficzny jest rodziną rozkładów prawdopodobieństwa skończonego zbioru zmiennych losowych X_1, X_2, \dots, X_N , które są przedstawione za pomocą N wierzchołków grafu (skierowanego lub nieskierowanego, zależnie od zastosowania). Obecność krawędzi pomiędzy dwoma wierzchołkami świadczy w modelu o zależności warunkowej między powiązаныmi z nimi zmiennymi. Z gaussowskim modelem graficznym mamy do czynienia, kiedy zmienne X_1, X_2, \dots, X_N są z rozkładu normalnego. W tym przypadku pełna informacja o strukturze niezależności warunkowej między zmiennymi jest zawarta w macierzy odwrotnej do macierzy kowariancji (w macierzy precyzji) wektora losowego $\mathbf{X} = (X_1, \dots, X_N)'$. Czytelnika zainteresowanego głębszym zapoznaniem się z modelami graficznymi zachęcamy do książek wprowadzających do tej teorii [122, 187, 20]. Dyskusję nad istnieniem MLE w dyskretnych modelach hierarchicznych, w tym w modelach graficznych, można znaleźć w pracy Wanga, Rauha i Massam [185]. Odpowiednik estymatora LASSO w modelach graficznych został wprowadzony około piętnastu lat temu [130, 188, 79] i jest w dalszym ciągu używany [37]. Graficzny odpowiednik SLOPE został wprowadzony niezależnie w pracy doktorskiej Sobczyka [166] oraz w pracy Mazzy-Anthony'ego, Mazourego i Coatesa [129].

Wróćmy do struktury niezależności warunkowej w modelu graficznych. Można zauważyć, że położenie zer w macierzy precyzji wektora \mathbf{X} jest identyczne z położeniem zer w laplasjanie grafu opisującego tę strukturę. Nasza dyskusja skupiona jest na tych z procesów Markowa, których dyskretyzacje mają macierz kowariancji równą odwrotności macierzy zbliżonej do laplasjanu zadanego grafu. Jako przykład można podać sytuację, w której zadany graf jest n -wierzchołkową ścieżką. Wtedy dodanie 1 do lewego górnego wyrazu laplasjanu skutkuje macierzą odwrotną do $(\Sigma)_{i,j} = \min\{i, j\}$, która to jest dyskretyzacją macierzy kowariancji

klasycznego procesu Wienera.

PL.5 Plan rozprawy

Rozprawa jest oparta na pięciu artykułach. Trzy z nich są opublikowane, jeden został wysłany do czasopisma, a kolejny jest w przygotowaniu i niebawem zostanie wysłany. Preprinty nieopublikowanych artykułów są dostępne w repozytoriach arXiv oraz HAL.

1. Pattern recovery and signal denoising by SLOPE when the design matrix is orthogonal [165] (z P. Graczykiem, B. Kołodziejkiem i M. Wilczyńskim), *Probability and Mathematical Statistics*, 42(2):283–302, 2022
2. Pattern recovery by SLOPE [23] (z M. Bogdan, X. Dupuisem, P. Graczykiem, B. Kołodziejkiem, P. Tardivelem i M. Wilczyńskim), wysłany do recenzji <https://arxiv.org/pdf/2203.12086.pdf>
3. Pattern Recovery in Penalized and Thresholded Estimation and its Geometry [93] (z P. Graczykiem, U. Schneider i P. Tardivelem), w przygotowaniu, <https://hal.science/hal-03262087v2/document>
4. Maximum likelihood estimation for discrete exponential families and random graphs [22] (z K. Bogdanem i M. Bosym), *ALEA*, 19, 1045–1070 (2022)
5. Remarks on Laplacian of Graphical Models in Various Graphs [164], *Proceedings, GSI 2021, Paris, France, July 21–23, 2021*

Plan niniejszej rozprawy jest następujący:

Rozdział drugi zawiera podstawowe pojęcia dotyczące wyników używanych w dalszych rozdziałach.

Rozdział trzeci skupia się na metodzie SLOPE w przypadku ortogonalnej macierzy eksperymentu, tzn. $\mathbf{X}'\mathbf{X} = n\mathbf{I}_p$. W tym przypadku wprowadzamy również nowe wyniki dotyczące mocnej zgodności estymatora SLOPE i jego wzorca. Następnie ilustrujemy numerycznie zastosowanie SLOPE w klasteryzacji do odtwarzania sygnałów o wysokiej częstotliwości.

W rozdziale czwartym omawiamy odtwarzanie wzorca SLOPE bez ograniczeń z poprzedniego rozdziału. Pokazujemy, że odtwarzanie wzorca można scharakteryzować za pomocą dwóch warunków, które nazywamy 'positivity condition' oraz 'subdifferential condition'. Następnie wprowadzamy 'irrepresentability condition' dla SLOPE, który jest uogólnieniem 'irrepresentability condition' dla LASSO [82, 41] wraz z jego geometryczną interpretacją. W dalszej kolejności podajemy bardziej wyrafinowane warunki asymptotyczne na zgodność oraz mocną zgodność estymatora SLOPE i jego wzorca.

Rozdział piąty dotyczy odtwarzania wzorca w ogólniejszym przypadku kar postaci polyhedral gauges. Wyprowadzamy warunki konieczne i dostateczne na jednorodną jednoznaczność estymatora, tzn. jego jednoznaczność dla dowolnej wartości $\mathbf{Y} \in \mathbb{R}^n$. Również pokazujemy kryteria na odtwarzanie wzorca dla podanych estymatorów i ich progowanych (thresholded) odmian. Ponadto, uogólniamy irrepresentability condition na penalizację poprzez polyhedral gauges oraz opisujemy relacje pomiędzy wzorcem, a modelem rozważanym w publikacji Vaitera i in. [182].

W rozdziale szóstym podajemy konieczne i dostateczne warunki na istnienie estymatora największej wiarygodności (MLE) w dyskretnej rodzinie wykładniczych. Głównym narzędziem

używany w tym celu jest pojęcie zbiorów jednoznaczności, czyli takich podzbiorów przestrzeni stanów, na którym funkcja z danej klasy jest równa zero wtedy i tylko wtedy, gdy jest ona stale równa zero na całej przestrzeni. Następnie wykorzystujemy nasze kryterium do klasy funkcji Rademachera i Walsha oraz do modelu wykładniczych grafów losowych (ERGM).

W rozdziale siódmym rozważamy powiązania między gaussowskimi modelami graficznymi, a macierzami kowariancji zdyskretyzowanych procesów Wienera. W tym celu używamy laplasjanu grafu, który opisuje strukturę niezależności warunkowej zadanego modelu.

Résumé

Récemment nous pouvons observer l'émergence rapide de données massives, tant au niveau du nombre d'observations que du nombre de variables explicatives à mesurer. En raison de l'augmentation de l'influence des big data sur la vie quotidienne, la statistique mathématique et la science des données gagnent un intérêt croissant dans le domaine des mathématiques et des sciences appliquées. Cependant, elles ne doivent pas être considérées séparément des autres branches des mathématiques. Dans cette thèse, nous nous attachons à proposer de nouveaux liens entre la statistique moderne, la géométrie et la combinatoire. Nous commençons notre discussion par les espaces euclidiens de dimension finie. Nous prêtons attention aux situations où l'espace des valeurs possibles de l'estimateur peut être partitionné en un nombre fini de sous-ensembles tels qu'il existe une correspondance bijective entre les sous-ensembles et les propriétés spécifiques des estimations. Cette prédiction des propriétés d'une observation est mieux connue sous le nom de problème de classification [32], dont la recherche a été lancée par Fisher [77] et qui est appliquée dans presque tous les domaines liés aux données.

FR.1 Régression Linéaire Pénalisée

Parmi les applications importantes du problème de la classification, on peut citer les modifications de la régression linéaire, qui ne se concentrent pas sur la valeur exacte de l'estimation mais sur ses propriétés importantes.

Dans le modèle de régression linéaire multiple ayant n observations et p variables explicatives, nous supposons que le vecteur de réponse $\mathbf{Y} = (y_1, \dots, y_n)' \in \mathbb{R}^n$ est de la forme $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, où $\mathbf{X} \in \mathbb{R}^{n \times p}$ est une matrice de planification, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$ est un vecteur inconnu de coefficients de régression et $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)' \in \mathbb{R}^n$ est un bruit aléatoire. L'objectif principal de la régression est d'estimer $\boldsymbol{\beta}$ et de retrouver ses propriétés essentielles. La abondante classe de problèmes de choix du sous-ensemble de coordonnées de $\boldsymbol{\beta}$, qui est le plus approprié pour nos demandes sur l'estimation et sa parcimonie, est mieux connue sous le nom de problème de sélection de modèle. Pour une étude sur le problème du sélection de modèle, cf. [49]. Dans la plupart de nos discussions, nous supposons que l'erreur $\boldsymbol{\varepsilon}$ a une distribution symétrique et continue. Certains de nos résultats couvrent le cas sans erreur de $\boldsymbol{\varepsilon} \equiv \mathbf{0}$, qui est utile pour établir des résultats asymptotiques. En raison de l'abondance des connaissances existantes sur la distribution gaussienne, nous sommes également en mesure de présenter des résultats plus efficaces sur le comportement des estimateurs avec gaussienne $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. La méthode de régression linéaire la plus classique, la méthode des moindres carrés ordinaires (OLS), a été proposée par Legendre au début du XIXe siècle [10, 97, 124]. L'estimateur OLS est défini comme le minimiseur de la somme des carrés des résidus, à savoir

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2.$$

OLS a gagné sa popularité dans le cas $n \geq p$ en raison de sa simplicité, de son efficacité pour un nombre relativement bas de variables [33, str. 4] et de nombreuses propriétés statistiques utiles. Par exemple, si la matrice $\mathbf{X}'\mathbf{X}$ est inversible et $\mathbf{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n$ pour $\sigma > 0$, alors les OLS constituent le meilleur estimateur linéaire sans biais (BLUE) [1, Theorem 2.7.1.] de β . Dans ce cas, la formule exacte de $\hat{\beta}^{OLS}$ peut être calculée à partir de l'équation suivante pour $\hat{\beta}^{OLS}$ peut être facilement déduite [1, str. 28]:

$$\hat{\beta}^{OLS} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}.$$

Si nous supposons en plus la gaussianité de \mathbf{Y} , alors les OLS sont aussi l'estimateur du maximum de vraisemblance (MLE) de β [1, str. 28]. Lorsque la matrice de planification \mathbf{X} est orthogonale, c'est-à-dire que $\mathbf{X}'\mathbf{X} = c\mathbf{I}_p$, $c > 0$, la formule ci-dessus se résume à $\hat{\beta}^{OLS} = c\mathbf{X}'\mathbf{Y}$. Sous les hypothèses ci-dessus, les OLS est également un estimateur fortement consistant de β [6]. Cependant, cet estimateur n'est pas défini de manière unique lorsque la matrice $\mathbf{X}'\mathbf{X}$ est unité. défini lorsque la matrice $\mathbf{X}'\mathbf{X}$ n'est pas inversible, ce qui est le cas dans un cadre de haute dimension $p > n$. De plus, sous des hypothèses communes sur le terme d'erreur ε , l'estimation OLS n'est pas parcimonieux, avec la probabilité 1 elle contient p coordonnées de valeurs mutuellement différentes. Pour une exemple d'hypothèses inhabituelles sur \mathbf{Y} et \mathbf{X} et d'absence d'estimation unique, nous pouvons référer à l'article récent par Dupuis et Vaïter [62]. Comme le vrai vecteur $\beta \in \mathbb{R}^p$ des coefficients de régression pourrait contenir beaucoup moins de coordonnées non nulles, il existe une voie naturelle pour proposer des méthodes qui favoriseraient la parcimonie de β , c'est-à-dire un petit nombre de coordonnées non nulles, ou du moins favoriseraient les estimations de β qui sont descriptibles par un petit nombre de paramètres. Plusieurs solutions ont été proposées pour traiter ce problème. L'une d'entre elles consiste à comparer les modèles appropriés par un critère d'information, par exemple BIC [158] ou AIC [2]. Une autre approche consiste à utiliser la régression pénalisée de la forme

$$\hat{\beta} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}_\lambda(\mathbf{b}) \right],$$

où $\text{pen}_\lambda(\mathbf{b})$ est un pénalisateur arbitraire non négatif, modifiable en fonction de la valeur de $\lambda > 0$. Elle a été appliquée, par exemple, à l'imagerie cérébrale [39], à la prévision des prix de l'électricité [106, 180] ou encore dans les mathématiques financières pour regrouper les actifs en fonction de leur corrélation partielle avec la série chronologique des rendements des fonds spéculatifs [116]. La première présence de cette idée est la sélection du meilleur sous-ensemble [12, 102] avec $\text{pen}_\lambda(\mathbf{b}) = \lambda \|\mathbf{b}\|_0$ étant un nombre de coordonnées non nulles de \mathbf{b} multiplié par λ . Cependant, pour de grandes valeurs de p , trouver le meilleur sous-ensemble par cette méthode est un problème NP-difficile [78]. Pour une solution plus parcimonieuse, on peut utiliser l'estimateur LASSO (Least Absolute Shrinkage and Selection Operator [47, 176]), dans lequel la pénalité ajoutée à la somme des carrés des résidus $\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2$ est une norme ℓ_1 de $\hat{\beta}$ et le paramètre de régularisation est $\lambda > 0$:

$$\hat{\beta}^{\text{LASSO}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_1 \right].$$

L'estimateur LASSO n'est pas sans biais, en tant qu'estimateur de rétrécissement il apporte certains coefficients $\hat{\beta}_j^{\text{LASSO}}$ vers zéro. Certaines des coordonnées sont complètement réduites à zéro, ce qui donne une estimation plus parcimonie. Lorsque la matrice de planification \mathbf{X} est orthonormée, c'est-à-dire $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$, la formule exacte de $\hat{\beta}^{\text{LASSO}}$ trouvée par Tibshirani [176] est basée sur $\hat{\beta}^{OLS}$:

$$\hat{\beta}_i^{\text{LASSO}} = \text{sign}(\hat{\beta}_i^{\text{OLS}}) \max \{ |\hat{\beta}_i^{\text{OLS}}| - \lambda, 0 \}.$$

Pour garantir l'existence d'une solution du problème de la régression linéaire pénalisée, il y a nombreux propositions que la pénalité soit convexe. Cela permet également d'appliquer les outils de l'analyse convexe. Pour une comparaison plus large d'autres modifications d'un estimateur LASSO, nous renvoyons à l'article [78]. Cependant, il existe également des applications de pénalisateurs non convexes, par exemple la sélection du meilleur sous-ensemble [12, 102] ou l'estimateur SCAD (Smoothly Clipped Absolute Deviation [72]).

FR.2 SLOPE

Une autre approche pour réduire la dimensionnalité est l'estimateur pénalisé trié ℓ_1 (SLOPE [27, 26, 189]), qui, en plus de généraliser la méthode LASSO, regroupe les coefficients égaux de β et les colonnes corrélées de la matrice de planification \mathbf{X} . Comme son nom l'indique, dans SLOPE, la norme ℓ_1 comme pénalisateur est remplacée par la norme ℓ_1 triée :

$$J_{\mathbf{\Lambda}}(\mathbf{b}) := \sum_{i=1}^p |b|_{(i)} \lambda_i,$$

où $\lambda_1 > 0, \lambda_1 \geq \dots \geq \lambda_p \geq 0$ et $|\mathbf{b}|_{(1)} \geq \dots \geq |\mathbf{b}|_{(p)}$ sont les valeurs absolues des coordonnées de \mathbf{b} triées par ordre décroissant. Le sous-cas particulier de SLOPE avec $\mathbf{\Lambda}$ étant une séquence arithmétique est également connu sous le nom OSCAR [29]. Il est important de noter que dans SLOPE, un seul paramètre de régularisation λ est remplacé par un vecteur non croissant $\mathbf{\Lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p)'$ de p paramètres de régularisation (le vecteur de régularisation). Ainsi, il nous permet de modifier non seulement l'échelle de $\mathbf{\Lambda}$, mais aussi sa forme. La principale motivation des auteurs de SLOPE était de tester les p hypothèses nulles $H_0^i : \beta_i = 0$ et le contrôle du taux de fausses découvertes (le contrôle FDR), qui est défini par la proportion attendue entre le nombre de faux rejets d'hypothèses nulles et le nombre total de rejets d'hypothèses nulles (en cas d'absence de rejet, le FDR est défini comme égal à zéro). De plus, SLOPE généralise certaines des approches précédentes en régression linéaire :

- $\lambda_1 = \dots = \lambda_p = 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OLS}},$
- $\lambda_1 = \dots = \lambda_p > 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{LASSO}},$
- $\mathbf{\Lambda}$ est une suite arithmétique $\Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OSCAR}}.$

Pour une description plus détaillée d'histoire des recherches sur l'estimateur SLOPE, voir la Section 4.1.1.

Dans notre recherche, nous nous concentrons sur une autre propriété importante de SLOPE, à savoir la recouvrement des caractéristiques importantes d'un vecteur de coefficients de régression β , appelé son SLOPE schéma. Il s'agit d'un analogue du vecteur de signe dans LASSO et, comme son prédécesseur, il peut être entièrement décrit par le sous-différentiel de la norme pénalisante. Pour être plus spécifique, laissez k être le nombre de clusters de $\mathbf{patt}(\beta) = (m_1, \dots, m_p)'$ c'est-à-dire le nombre de composantes distinctes non nulles de $|\beta|$.

Définition FR.2.1 (SLOPE schéma). *Le SLOPE schéma est une fonction $\mathbf{patt} : \mathbb{R}^p \rightarrow \mathbb{Z}^p$ telle que*

$$\mathbf{patt}(\mathbf{b})_i = \text{sign}(b_i) \text{rank}(|b_i|),$$

où $\text{rank}(|b_i|) \in \{1, 2, \dots, k\}$ est défini comme le nombre de $|c_j|$'s satisfaisant $|b_i| \geq |c_j|$, où $|c_1|, |c_2|, \dots, |c_k|$, $k \leq p$, sont des valeurs distinctes non nulles parmi $|b_1|, \dots, |b_p|$. Nous adoptons la convention que $\text{rank}(0) = 0$.

Le SLOPE schéma ne préserve pas seulement le signe d'un vecteur, mais détecte également ses clusters, c'est-à-dire les ensembles de coordonnées partageant la même valeur absolue et l'ordre entre ces valeurs absolues (ordre entre les clusters).

Fait FR.2.1 (Propriétés de base du SLOPE schéma [156]).

- (a) pour chaque $1 \leq l \leq \|\mathbf{patt}(\mathbf{b})\|_\infty$ il existe j tel que $|\mathbf{patt}(\mathbf{b})_j| = l$,
- (b) $\text{sign}(\mathbf{patt}(\mathbf{b})) = \text{sign}(\mathbf{b})$ (préservation du signe),
- (c) $|b_i| = |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| = |\mathbf{patt}(\mathbf{b})_j|$ (préservation du cluster),
- (d) $|b_i| > |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| > |\mathbf{patt}(\mathbf{b})_j|$ (préservation de l'ordre).

Exemple FR.2.2. $\mathbf{patt}((4, 0, -1.5, 1.5, -4)') = (2, 0, -1, 1, -2)'$.

On dit que l'estimateur SLOPE $\hat{\beta}^{\text{SLOPE}}$ récupère le schéma de β lorsque

$$\mathbf{patt}\left(\hat{\beta}^{\text{SLOPE}}\right) = \mathbf{patt}(\beta).$$

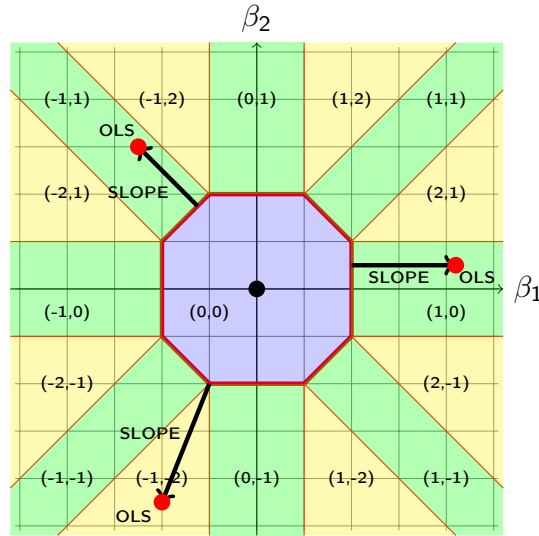
Dans les chapitres 3 et 4, nous discutons de nouvelles conditions nécessaires et suffisantes pour le recouvrement du SLOPE schéma, ainsi que de nouveaux résultats pour $n \geq p$ sur la consistance et la forte consistance de l'estimation SLOPE et de son schéma.

FR.2.1 Consistance

L'outil principal pour montrer la consistance forte de SLOPE est la forte consistance de l'estimateur des moindres carrés, dont la preuve a été présentée par exemple dans l'article de Anderson et Taylor [6]. Les principaux résultats sur la consistance de LASSO peuvent être trouvés dans les articles suivants:

- consistance de LASSO: Knight, Fu, 2000 [112],
- forte consistance de LASSO : Chatterjee, Lahiri, 2011 [44],
Les principales hypothèses sont que le paramètre de régularisation λ_n est d'un ordre inférieur à n et qu'il existe une espérance infinie de la valeur absolue d'un terme d'erreur.
- consistance du signe de LASSO: Zhao, Yu, 2006 [192].

Les sections 3.4 et 4.6 se concentrent sur la forte consistance de SLOPE et de son schéma. Il est important de rappeler que si le vecteur d'accord $\mathbf{\Lambda}$ est constant (LASSO comme cas particulier de SLOPE), alors le SLOPE schéma n'est pas consistant, même si l'hypothèse sur la consistance du vecteur de signe est satisfaite. A titre d'exemple, les clusters d'un vrai vecteur de paramètres, qui sont contenus dans un SLOPE schéma, ne sont pas préservés par LASSO. Comme autre remarque intéressante, on peut mentionner que la forte consistance de LASSO (ou SLOPE) n'implique pas nécessairement la forte consistance de son signe (de son schéma SLOPE). Comme contre-exemple, on peut remarquer que la séquence $((1/n, 1/n, -1/n, -1/n)')_{n \geq 1}$ converge vers $(0, 0, 0, 0)'$, alors que son signe et son schéma SLOPE sont égaux $(1, 1, -1, -1)'$ pour tout n positif.



$\hat{\beta}^{\text{SLOPE}}$ et $\hat{\beta}^{\text{OLS}}$ dans un plan orthogonal: $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ pour $\Lambda = (2, 1)'$.

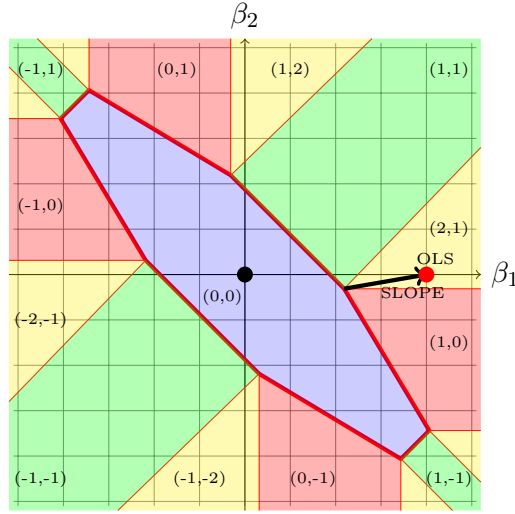
FR.2.2 Condition d'irreprésentabilité

Étant donné tout $n, p > 0$, dans LASSO, la consistance de la sélection du vrai sous-ensemble de β est presque équivalente à ce que les coordonnées en dehors du vrai support de β ne soient pas représentables par les coordonnées à l'intérieur du support [192]. Nous introduisons un analogue de cette condition d'irreprésentabilité de LASSO à SLOPE. Puis nous dérivons la caractérisation géométrique d'une estimation SLOPE.

FR.2.3 Géométrie de la régression linéaire pénalisée

Pour une meilleure compréhension des connexions entre $\hat{\beta}^{\text{SLOPE}}$ et $\hat{\beta}^{\text{OLS}}$, nous donnons une illustration d'un cas de bas dimension $p \leq n$ pour $p = 2$. Ici, l'estimation SLOPE est une différence entre l'estimation des moindres carrés et sa projection sur $(\mathbf{X}'\mathbf{X})^{-1}C_\Lambda$, où C_Λ est une boule unité dans une norme double de J_Λ . Ce résultat est particulièrement facile à interpréter dans le cas où \mathbf{X} est une matrice orthogonale. Il en résulte une formule plus simple, qui a été récemment proposée par Tardivel, Servien et Concordet [175]. Pour une présentation plus large de l'application de l'approche géométrique à l'estimateur SLOPE et à sa recouvrement de schémas, nous invitons à la Section 4.4.

LASSO et SLOPE présentent une connexion avec la théorie des systèmes racines. En particulier, on peut observer qu'une boule unité en norme ℓ_∞ , qui est duale de la norme ℓ_1 utilisée dans LASSO et une boule unité en norme J_Λ^* sont proportionnelles, respectivement, aux coques convexes des orbites des groupes de Weyl correspondant au produit cartésien $p^{\text{ième}}$ du système racine A_1 et au système racine B_p [90, 144]. Pour plus d'informations sur la théorie des systèmes racine et ses recherches actuelles, nous renvoyons à un cahier de cours de Helgason [100] et à [92, 63]. La géométrie des estimateurs de régression linéaire pénalisés et leur recouvrement peuvent être inspectés dans un cas plus général, ce qui constitue l'essence du chapitre 5. En effet, l'estimateur SLOPE peut être classé dans l'une des classes d'un nombre fini de classes lorsque la pénalité est une gauge polyédrique, c'est-à-dire une fonction convexe non négative, qui est le maximum d'un nombre fini de fonctions linéaires. Les caractéristiques importantes de cette estimation, aussi appelée son schéma, peuvent être entièrement décrites avec le sous-



$$\hat{\beta}^{\text{SLOPE}} \text{ et } \hat{\beta}^{\text{OLS}} \text{ for } \mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 1 \end{bmatrix} \text{ et } \boldsymbol{\Lambda} = (2, 1)'.$$

différentiel du pénaliseur. Puisque les inégalités pour les fonctions linéaires déterminant le gauge polyédrique sont satisfaites par une intersection du nombre fini de demi-espaces correspondants, ces estimateurs sont fortement liés à la théorie des polytopes, cf. les livres de Gruber [94], Grünbaum [95], Hiriart-Urruty et Lemarechal [101] ou Ziegler [193]. Dans le chapitre 5, nous généralisons certains des nouveaux résultats pour SLOPE à la classe d'estimateurs ci-dessus. Pour généraliser la notion du schéma à la classe des gauges polyédriques, nous considérons la classe d'équivalence des patrons, ce qui nous permet de désigner son recouvrement, comme le recouvrement de la sous-différentielle de pen.

Nous introduisons et caractérisons l'accessibilité du schéma et nous donnons une condition d'irreprésentabilité pour les gabarits polyédriques. Avec des outils issus de la géométrie polyédrique, en particulier avec les cônes normaux des faces des polyèdres considérés, nous discutons les propriétés de l'égalité des schémas entre vecteurs. Nous étudions également les relations entre les ensembles de vecteurs ayant le même schéma avec la notion de sous-espace modèle, cf. [182]. Plus tard, nous discutons de la recouvrement du schéma par les estimateurs pénalisés seuillés, qui sont une généralisation du LASSO seuillé, qui ne considère pas seulement le signe de l'estimateur, mais son schéma entier. Notre discussion étend également la caractérisation de l'unicité uniforme d'un estimateur de la classe des normes polyédriques [156] aux gauges polyédriques.

FR.3 Familles exponentielles discrètes

Dans le chapitre 6, nous passons des espaces euclidiens aux espaces discrets, c'est-à-dire les espaces avec un nombre fini d'éléments. Ici, le point majeur de notre intérêt sont les familles exponentielles discrètes, que nous comprenons comme des familles exponentielles définies sur un espace fini d'états \mathcal{X} , $|\mathcal{X}| = K < \infty$. Pour les familles sur les ensembles infinis dénombrables, voir Jacobsen [105]. Nous considérons ensuite un sous-espace linéaire \mathcal{B} de l'espace des fonctions linéaires $\mathbb{R}^{\mathcal{X}}$, qui contient une fonction constante strictement positive. Nous introduisons également une fonction de poids strictement positive $\mu : \mathcal{X} \rightarrow (0, \infty)$. Nous construisons la famille exponentielle de la façon suivante:

Pour une fonction à valeur réelle ϕ , nous définissons les fonctions de partition et de log-partition,

$$Z(\phi) = \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x), \quad \psi(\phi) = \log Z(\phi),$$

respectivement, et densité exponentielle

$$p = e(\phi) = e^{\phi - \psi(\phi)} = e^{\phi} / Z(\phi).$$

La famille exponentielle couverte par l'ensemble \mathcal{B} est

$$e(\mathcal{B}) := \{p = e(\phi) : \phi \in \mathcal{B}\}.$$

Cela nous permet de dériver la fonction de vraisemblance et de log-vraisemblance. La fonction de log-vraisemblance est strictement concave. Par conséquent, si l'estimateur du maximum de vraisemblance (MLE) existe, il est unique. Malgré le caractère borné de la fonction de vraisemblance, l'MLE peut ne pas exister. Notre objectif principal est d'établir une nouvelle caractérisation de l'existence de l'MLE et de l'appliquer à des familles spécifiques, les familles exponentielles couvertes par les fonctions de Rademacher et de Walsh, et les familles exponentielles de graphes aléatoires.

L'outil clé dans notre discussion est la notion nouvellement introduite d'ensemble d'unicité, c'est-à-dire un tel ensemble $U \subset X$, que $\phi = 0$ est la seule fonction d'une classe donnée de fonctions qui est égale à zéro sur U . Pour être plus précis, nous montrons que le MLE pour $e(\mathcal{B})$ et l'échantillon i.i.d. $x_1, x_2, \dots, x_n \in \mathcal{X}$ existe si et seulement si $\{x_1, \dots, x_n\}$ est l'ensemble d'unicité pour le cône non négatif $\mathcal{B}_+ := \{\phi \in \mathcal{B} : \phi \geq 0\}$. Nous reformulons également ce critère sous la forme d'un problème de programmation linéaire.

De plus, pour les applications proposées, nous utilisons les limites de probabilité établies par les outils classiques de la théorie des graphes aléatoires et du problème du collecteur de coupons, voir [118, p. 194-195], [68], et [143]. Dans la famille des fonctions de Walsh, nous utilisons les propriétés de l'hypercube $\{-1, 1\}^k$ ($k = \log_2 K$), de ses sous-cubes et du graphe hypercube correspondant.

FR.4 Modèles graphiques

L'une des branches de la statistique qui est connue à la fois pour l'utilisation de l'MLE et pour son efficacité et pour appliquer l'estimation pénalisée est la théorie des modèles graphiques. Un modèle graphique est une famille de distributions de probabilités d'une collection finie de variables aléatoires X_1, \dots, X_N , qui sont codées par N sommets d'un graphe (dirigé ou non dirigé). La présence (ou l'absence) d'une arête entre deux sommets renseigne sur la dépendance (ou l'indépendance) conditionnelle entre les sommets correspondants. Lorsque toutes les variables aléatoires sont gaussiennes, on parle de modèles graphiques gaussiens. Dans ce cas, la structure d'indépendance conditionnelle est entièrement codée par l'emplacement des zéros dans l'inverse de la matrice de covariance (la matrice de précision) d'un vecteur $\mathbf{X} = (X_1, \dots, X_N)'$. Par savoir plus sur les modèles graphiques nous invitons à consulter les introductions existantes sur ce sujet [122, 187, 20]. Pour la discussion sur l'existence de MLE dans les modèles hiérarchiques discrets, qui incluent les modèles graphiques, nous nous référons à l'article de Wang, Rauh et Massam [185]. La version modèle graphique d'un estimateur LASSO a été proposée il y a environ quinze ans [130, 188, 79] et gagne encore en popularité, cf. par exemple [37]. La SLOPE graphique a été proposée indépendamment par Sobczyk dans sa thèse de doctorat [166] et par

Mazza-Anthony, Mazoure et Coates [129]. Revenons-en à la structure d'indépendance conditionnelle. On peut observer que l'emplacement des zéros dans la matrice de précision de \mathbf{X} est le même que dans la matrice laplacienne d'un graphe sous-jacent. Dans notre discussion, nous examinons de plus près de tels graphes et recherchons ces processus de Markov, dont les discrétisations ont une matrice de covariance qui est l'inverse d'une matrice laplacienne un peu modifiée. À titre d'exemple, lorsqu'un graphique sous-jacent est un chemin de n sommets, l'ajout de 1 à l'entrée supérieure gauche de sa matrice laplacienne donne une inverse égale à $(\boldsymbol{\Sigma})_{i,j} = \min\{i, j\}$, qui est une version discrétisée de la matrice de covariance d'un processus de Wiener.

FR.5 Plan de la thèse

La thèse est basée sur cinq articles, dont trois sont publiés, un est en cours de révision, et un est une version étendue d'un preprint, qui peut être trouvé sur le site de HAL.

1. Pattern recovery and signal denoising by SLOPE when the design matrix is orthogonal [165] (avec P. Graczyk, B. Kołodziejek et M. Wilczyński), *Probability and Mathematical Statistics*, 42(2):283–302, 2022
2. Pattern recovery by SLOPE [23] (avec M. Bogdan, X. Dupuis, P. Graczyk, B. Kołodziejek, P. Tardivel et M. Wilczyński), soumis, <https://arxiv.org/pdf/2203.12086.pdf>
3. Pattern Recovery in Penalized and Thresholded Estimation and its Geometry [93] (avec P. Graczyk, U. Schneider et P. Tardivel), en préparation, <https://hal.science/hal-03262087v2/document>
4. Maximum likelihood estimation for discrete exponential families and random graphs [22] (avec K. Bogdan et M. Bosy), *ALEA*, 19, 1045–1070 (2022)
5. Remarks on Laplacian of Graphical Models in Various Graphs [164], *Proceedings, GSI 2021, Paris, France, July 21–23, 2021*

Le plan de la thèse est le suivant:

Le chapitre 2 contient les préliminaires de la recherche présentée et décrit les notions de base et nouvelles utilisées dans les chapitres suivants. Dans le chapitre 3, nous nous concentrons sur l'estimation du SLOPE dans le cas où la matrice de planification est orthogonale, c'est-à-dire, $\mathbf{X}'\mathbf{X} = n\mathbf{I}_p$. En ce cas, nous présentons également de nouveaux résultats sur la forte consistance des estimateurs SLOPE et sur la forte consistance de la recouvrement de schémas par SLOPE et nous illustrons les avantages du regroupement SLOPE dans le contexte du dénoisement de signaux à haute fréquence. Le chapitre 4 se concentre sur la recouvrement du SLOPE schéma sans restriction sur la matrice \mathbf{X} . Nous montrons que la recouvrement du schéma peut être déterminée par deux critères, appelés condition de positivité et condition subdifférentielle. Plus tard, nous introduisons une condition d'irreprésentabilité SLOPE, qui généralise la condition d'irreprésentabilité LASSO [82, 41] et nous l'illustrons ensuite géométriquement. Plus tard, nous dérivons les conditions asymptotiques raffinées sur la consistance et la consistance forte de l'estimation SLOPE et de son schéma. Le chapitre 5 concerne le problème de la recouvrement du schéma dans un cas général où la pénalité est une gauge polyédrique. Nous établissons la condition nécessaire et suffisante pour l'unicité uniforme de l'estimateur. Ensuite, nous dérivons des conditions pour la recouvrement du schéma des estimateurs pénalisés par une gauge polyédrique et pour leurs versions seuillées. Enfin, nous généralisons la condition d'irreprésentabilité pour les pénalisateurs étant des gauges

polyédriques et nous discutons des connexions entre un schéma et un modèle de l'article de Vaiter et al. [182].

Le chapitre 6 porte sur les conditions nécessaires et suffisantes pour l'existence de l'estimateur du maximum de vraisemblance (MLE) dans les familles exponentielles discrètes. L'outil principal de cet article est la notion d'ensembles d'unicité, c'est-à-dire les sous-ensembles d'un espace d'état sur lesquels une fonction d'une classe donnée est égale à zéro si et seulement si elle est égale à zéro sur un espace entier. Nous appliquons nos résultats à la classe des fonctions de Rademacher et de Walsh et aux modèles de graphes aléatoires exponentiels (ERGM). Dans le chapitre 7, nous abordons le lien entre les modèles graphiques gaussiens et les matrices de covariance des processus de Wiener discrétisés. Pour cela, nous utilisons la matrice laplacienne discrète d'une structure graphique sous-jacente.

Chapter 1

Introduction

In recent times we may observe the rapid and inevitable emergence of data collection, both in the number of observations and in the number of features to measure them. Due to the increase of the influence of big data on everyday life, mathematical statistics and data science gain growth of interest in mathematics and applied sciences. However, they should not be considered separately from other branches of mathematics. In this dissertation we focus on proposing new connections of modern statistics with geometry and combinatorics.

We start our discussion with finite-dimensional Euclidean spaces. We pay attention to those situations where the space of possible values of the estimator can be partitioned into a finite number of subsets such that there is a one-to-one correspondence between the subsets and specific properties of estimates. This prediction of the properties of an observation is better known as a classification problem [32], research on which was started by Fisher [77] and is being applied in almost all fields related to data processing.

1.1 Penalized Linear Regression

As one of important applications of the classification problem, we can point out the modifications of linear regression, which do not focus on the exact estimate value but on its important properties. In the multiple linear regression model having n observations and p explanatory variables, we assume that the response vector $\mathbf{Y} = (Y_1, \dots, Y_n)' \in \mathbb{R}^n$ is of the form $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$ is an unknown vector of regression coefficients and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)' \in \mathbb{R}^n$ is a random noise. The main objective of linear regression is to estimate $\boldsymbol{\beta}$ and recover its essential properties. The wide class of problems of choosing the subset of coordinates of $\boldsymbol{\beta}$, which is most suitable for our demands on the estimate and its sparsity, is better known as the model selection problem. For the survey on the model selection problem, cf. [49]. In most of our discussion, we assume that the error $\boldsymbol{\varepsilon}$ has a symmetric and continuous distribution. Some of our results cover the noiseless case of $\boldsymbol{\varepsilon} \equiv \mathbf{0}$, which is a strong tool for establishing asymptotic results. Due to the abundance of the existing knowledge on the Gaussian distribution, we are also able to present more efficient results on the behavior of considered estimators with Gaussian $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. The most classical linear regression method, the Ordinary Least Squares method (OLS), was proposed by Legendre at the beginning of the nineteenth century [10, 97, 124]. The OLS estimator is defined as the minimizer of the residual sum of squares, namely

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2.$$

The OLS estimator gained its popularity in the case $n \geq p$ due to its simplicity, efficiency for a relatively small number of variables [33, p. 4] and many useful statistical properties. For example, if the matrix $\mathbf{X}'\mathbf{X}$ is invertible and $\text{Cov}(\mathbf{Y}) = \sigma^2\mathbf{I}_n$ for $\sigma > 0$, then the OLS is the best linear unbiased estimator (BLUE) [1, Theorem 2.7.1.] of $\boldsymbol{\beta}$. In this case, the exact formula for $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ can be easily deduced [1, p. 28]:

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}.$$

If we additionally assume the gaussianity of \mathbf{Y} , then OLS is also the maximum likelihood estimator (MLE) of $\boldsymbol{\beta}$ [1, p. 28]. When the design matrix \mathbf{X} is orthogonal, i.e. $\mathbf{X}'\mathbf{X} = c\mathbf{I}_p$, $c > 0$, the above formula boils down to $\hat{\boldsymbol{\beta}}^{\text{OLS}} = \frac{1}{c}\mathbf{X}'\mathbf{Y}$. Under the above assumptions, the OLS is also a strongly consistent estimator of $\boldsymbol{\beta}$ [6]. However, this estimator is not uniquely defined when the matrix $\mathbf{X}'\mathbf{X}$ is not invertible, which is the case in a high-dimensional setting $p > n$. Additionally, under common assumptions on the error term $\boldsymbol{\varepsilon}$, OLS estimate is not sparse, with probability 1 it contains p coordinates of mutually different values. For an example of unusual assumptions on \mathbf{Y} and \mathbf{X} and no unique estimate, one may get acquainted with to a recent article by Dupuis and Vaïter [62]. As the true vector $\boldsymbol{\beta} \in \mathbb{R}^p$ of regression coefficients could contain much fewer nonzero coordinates, there is a natural pathway to propose methods that would promote the sparsity of $\boldsymbol{\beta}$, that is, a small number of nonzero coordinates, or at least to promote the estimates of $\boldsymbol{\beta}$ that are describable by a small number of parameters. Several solutions were proposed to deal with such problem. One of them is to compare the suitable models by an information criterion, for example BIC [158] or AIC [2]. Another approach is to use the penalized regression of the form

$$\hat{\boldsymbol{\beta}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}_\lambda(\mathbf{b}) \right],$$

where $\text{pen}_\lambda(\mathbf{b})$ is an arbitrary nonnegative penalizer, modifiable according to the value of $\lambda > 0$. It has been applied, for example, in brain imaging [39], forecasting electricity prices [106, 180] or in financial mathematics to group assets with respect to their partial correlation with the hedge fund return times series [116]. The first presence of this idea is the best subset selection [12, 102] with $\text{pen}_\lambda(\mathbf{b}) = \lambda\|\mathbf{b}\|_0$ being a number of nonzero coordinates of \mathbf{b} multiplied by λ . However, for large values of p , finding the best subset by this method is an NP-hard problem [78]. For a sparser solution, one may use the Least Absolute Shrinkage and Selection Operator (LASSO [47, 176]), in which the penalty added to the residual sum of squares $\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2$ is an ℓ_1 norm of $\hat{\boldsymbol{\beta}}$ and the tuning parameter is $\lambda > 0$:

$$\hat{\boldsymbol{\beta}}^{\text{LASSO}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda\|\mathbf{b}\|_1 \right].$$

The LASSO estimator is not unbiased, as a shrinkage estimator it brings some coefficients $\hat{\beta}_j^{\text{LASSO}}$ toward zero. Some of the coordinates are being shrunk completely to zero, which results in a sparser estimate. When the design matrix \mathbf{X} is orthonormal, i.e. $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$, the exact formula for $\hat{\boldsymbol{\beta}}^{\text{LASSO}}$ found by Tibshirani [176] is based on $\hat{\boldsymbol{\beta}}^{\text{OLS}}$:

$$\hat{\beta}_i^{\text{LASSO}} = \text{sign}(\hat{\beta}_i^{\text{OLS}}) \max \left\{ |\hat{\beta}_i^{\text{OLS}}| - \lambda, 0 \right\}.$$

To guarantee the existence of a solution of the penalized linear regression problem, many proposals offer the penalty to be convex. This also allows one to apply the tools from convex analysis. For a broader comparison of other modifications of a LASSO estimator, we refer to the article [78]. However, there are also applications of nonconvex penalizers, e.g., the best subset selection [12, 102] or the smoothly clipped absolute deviation (SCAD) [72].

1.2 SLOPE

Another approach to reduce dimensionality is the Sorted ℓ_1 Penalized Estimator (SLOPE [27, 26, 189]), which apart of generalizing the LASSO method, clusterizes the equal coefficients of β and correlated columns of the design matrix \mathbf{X} . As the name suggests, in SLOPE the ℓ_1 norm as penalizer is replaced by the sorted ℓ_1 norm:

$$J_{\mathbf{\Lambda}}(\mathbf{b}) := \sum_{i=1}^p |b|_{(i)} \lambda_i,$$

where $\lambda_1 > 0, \lambda_1 \geq \dots, \lambda_p \geq 0$ and $|b|_{(1)} \geq \dots \geq |b|_{(p)}$ are the absolute values of coordinates of \mathbf{b} sorted in descending order. The special subcase of SLOPE with $\mathbf{\Lambda}$ being an arithmetic sequence is also known as the OSCAR estimator [29]. It is important to note that in SLOPE a single tuning parameter λ is replaced by a non-increasing vector $\mathbf{\Lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p)'$ of p tuning parameters (the tuning vector). Thus, it allows us to modify not only the scaling of $\mathbf{\Lambda}$, but also its shape. The main motivation of the authors of SLOPE was the testing of p null hypotheses $H_i^0 : \beta_i = 0$ and the control of the rate of false discoveries (the FDR control), which is defined as the expected proportion between the amount of false rejections of null hypotheses and the total amount of rejections of null hypotheses (in case of no rejections, the FDR is defined to equal zero). Moreover, SLOPE generalizes some of the previous approaches in linear regression:

- $\lambda_1 = \dots = \lambda_p = 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OLS}},$
- $\lambda_1 = \dots = \lambda_p > 0 \Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{LASSO}},$
- $\mathbf{\Lambda}$ is an arithmetic sequence $\Rightarrow \hat{\beta}^{\text{SLOPE}} = \hat{\beta}^{\text{OSCAR}}.$

A more detailed description of the history of research on SLOPE can be found in Section 4.1.1. In our research, we focus on another important property of SLOPE, that is the recovery of important characteristics of a vector of regression coefficients β , called its SLOPE pattern. It is an analogue of the sign vector in LASSO and, like its predecessor, it is fully describable by the subdifferential of the penalizing norm. To be more specific, let k be the number of clusters of $\mathbf{patt}(\beta) = (m_1, \dots, m_p)'$ i.e., the number of distinct nonzero components of $|\beta|$.

Definition 1.2.1 (SLOPE pattern). *The SLOPE pattern is a function $\mathbf{patt} : \mathbb{R}^p \rightarrow \mathbb{Z}^p$ such that*

$$\mathbf{patt}(\mathbf{b})_i = \text{sign}(b_i) \text{rank}(|b_i|),$$

where $\text{rank}(|b_i|) \in \{1, 2, \dots, k\}$ is defined to be the number of $|c_j|$'s satisfying $|b_i| \geq |c_j|$, where $|c_1|, |c_2|, \dots, |c_k|$, $k \leq p$, are distinct nonzero values among $|b_1|, \dots, |b_p|$. We adopt the convention that $\text{rank}(0) = 0$.

The SLOPE pattern does not only preserve the sign of a vector, but also detects its clusters, i.e. the sets of coordinates sharing the same absolute value and the order between those absolute values (hierarchy between clusters).

Fact 1.2.1 (Basic properties of SLOPE pattern [156]).

- (a) for every $1 \leq l \leq \|\mathbf{patt}(\mathbf{b})\|_{\infty}$ there exists j such that $|\mathbf{patt}(\mathbf{b})_j| = l$,
- (b) $\text{sign}(\mathbf{patt}(\mathbf{b})) = \text{sign}(\mathbf{b})$ (sign preservation),
- (c) $|b_i| = |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| = |\mathbf{patt}(\mathbf{b})_j|$ (cluster preservation),

(d) $|b_i| > |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| > |\mathbf{patt}(\mathbf{b})_j|$ (hierarchy preservation).

Example 1.2.2. $\mathbf{patt}((4, 0, -1.5, 1.5, -4)') = (2, 0, -1, 1, -2)'$.

We say that the SLOPE estimator $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$ recovers the pattern of $\boldsymbol{\beta}$ when

$$\mathbf{patt}\left(\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\right) = \mathbf{patt}(\boldsymbol{\beta}).$$

In Chapter 3 and Chapter 4 we discuss novel necessary and sufficient conditions for the recovery of the SLOPE pattern, as well as the novel results for $n \geq p$ on the consistency and strong consistency of both the SLOPE estimate and its SLOPE pattern.

1.2.1 Consistency

The main tool to show the strong consistency of SLOPE is the strong consistency of the Ordinary Least Squares estimator, the proof of which was presented e.g. in the article of Anderson and Taylor [6]. Main results on the consistency of LASSO might be found in the following articles:

- consistency of LASSO: Knight, Fu, 2000 [112],
- strong consistency of LASSO: Chatterjee, Lahiri, 2011 [44],
Main assumptions are that the tuning parameter λ_n is of a smaller order than n and that there is a finite expectation of the absolute value of an error term.
- consistency of the sign of LASSO: Zhao, Yu, 2006 [192].

Sections 3.4 and 4.6 focus on the strong consistency of SLOPE and its pattern. It is important to recall that if the tuning vector $\boldsymbol{\Lambda}$ is constant (LASSO as the special case of SLOPE), then the SLOPE pattern is not consistent, even if the assumption on the consistency of the sign vector is satisfied. As an example, the clusters of a true vector of parameters, which are contained in a SLOPE pattern, are not preserved by LASSO. As another interesting remark, it may be mentioned that the strong consistency of LASSO (or SLOPE) does not necessarily imply the strong consistency of its sign (of its SLOPE pattern). As a counterexample, one may notice that the sequence $((1/n, 1/n, -1/n, -1/n)')_{n \geq 1}$ converges to $(0, 0, 0, 0)'$, while its sign and its SLOPE pattern are equal $(1, 1, -1, -1)'$ for every positive n .

1.2.2 Irrepresentability condition

Given any $n, p > 0$, in LASSO, the consistency of selection of the true subset of $\boldsymbol{\beta}$ is almost equivalent to the coordinates outside the true support of $\boldsymbol{\beta}$ not being representable by the coordinates inside the support [192]. We introduce an analogue of this Irrepresentability Condition from LASSO to SLOPE. Then we derive the geometric characterization of a SLOPE estimate.

1.2.3 Geometry of penalized linear regression

For a better understanding of the connections between $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$ and $\hat{\boldsymbol{\beta}}^{\text{OLS}}$, below we give an illustration of a low-dimensional case $p \leq n$ for $p = 2$. Here, the SLOPE estimate is a difference between the least squares estimate and its projection onto $(\mathbf{X}'\mathbf{X})^{-1}C_{\boldsymbol{\Lambda}}$, where $C_{\boldsymbol{\Lambda}}$ is a unit ball in a norm dual to $J_{\boldsymbol{\Lambda}}$. This result is especially easy to interpret in the case of \mathbf{X} being an orthogonal matrix. That resulted in an easier formula, which was recently proposed by Tardivel,

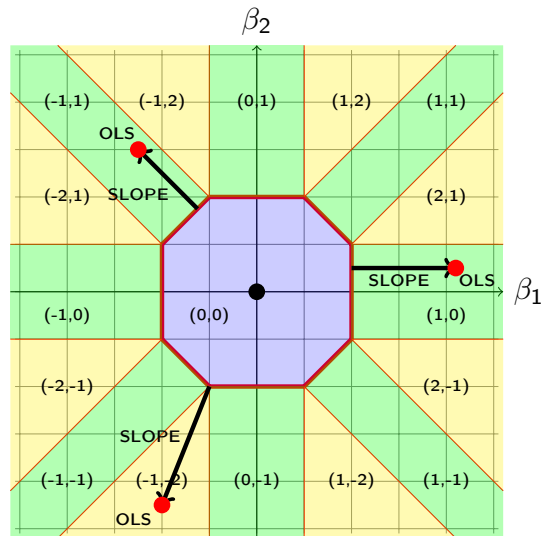


Figure 1.1: $\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$ in orthogonal design: $\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ for $\Lambda = (2, 1)'$.

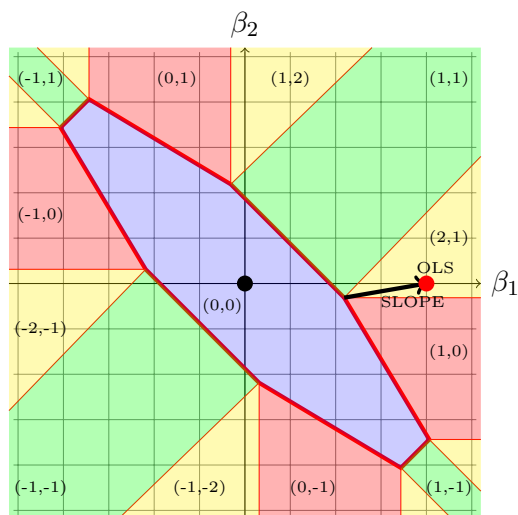


Figure 1.2: $\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$ for $\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 1 \end{bmatrix}$ and $\Lambda = (2, 1)'$.

Servien and Concordet [175]. For a wider presentation of applying the geometrical approach to the SLOPE estimator and its pattern recovery we invite the reader to Section 4.4.

LASSO and SLOPE estimators exhibit a connection to the theory of root systems. Namely, it may be observed that a unit ball in ℓ_∞ norm, which is dual to the ℓ_1 norm used in LASSO and a unit ball in J_Λ^* norm are proportional, respectively, to the convex hulls of orbits of Weyl groups corresponding to the p^{th} Cartesian product of the root system A_1 and of the root system B_p [90, 144]. For more on the theory of root systems and its current research, we refer to a coursebook of Helgason [100] and to [92, 63].

The geometry of penalized linear regression estimators and their pattern recovery may be inspected in more general case, which is the essence of Chapter 5. Indeed, the SLOPE estimator can be classified into one of finite number of classes when the penalty is a polyhedral gauge, i.e. a nonnegative convex function, which is a maximum of a finite number of linear functions. The important features of this estimate, also named its pattern, may be fully described with the subdifferential of the penalizer. Since the inequalities for linear functions determining the polyhedral gauge are satisfied by an intersection of the finite number of corresponding half-spaces, those estimators are strongly connected with the theory of polytopes, cf. the books of Gruber [94], Grünbaum [95], Hiriart-Urruty and Lemarechal [101] or Ziegler [193].

In Chapter 5 we generalize some of the new results for SLOPE to the above class of estimators. To generalize the notion of pattern to the class of polyhedral gauges, we consider the pattern equivalence class, which allows us to denote its recovery, as the recovery of the subdifferential of pen.

We introduce and characterize the accessibility of the pattern and we give an irrepresentability condition for polyhedral gauges. With tools from polyhedral geometry, especially with the normal cones of faces of considered polyhedra, we discuss the properties of the pattern equality between vectors. We also investigate the relations between sets of vectors having the same pattern with the notion of the model subspace, cf. [182].

Later on, we discuss the pattern recovery by thresholded penalized estimators, which are a generalization of the thresholded LASSO, which considers not only the sign of the estimator, but its entire pattern. Our discussion also extends the characterization of the uniform uniqueness of an estimator from the class of polyhedral norms [156] to polyhedral gauges.

1.3 Discrete exponential families

In Chapter 6 we move our focus from Euclidean to discrete spaces, i.e. spaces with a finite number of elements. Here, the major point of our interest are the discrete exponential families, which we understand as exponential families defined on a finite space of states \mathcal{X} , $|\mathcal{X}| = K < \infty$. For families on infinite countable sets, see Jacobsen [105]. Then we consider a linear subspace \mathcal{B} of the space of linear functions $\mathbb{R}^{\mathcal{X}}$, which contains a strictly positive constant function. We also introduce a strictly positive weight function $\mu : \mathcal{X} \rightarrow (0, \infty)$. We construct the exponential family in a following way:

For a real-valued function ϕ we define the partition and log-partition functions,

$$Z(\phi) = \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x), \quad \psi(\phi) = \log Z(\phi),$$

respectively, and *exponential density*

$$p = e(\phi) = e^{\phi - \psi(\phi)} = e^\phi / Z(\phi).$$

The exponential family spanned by \mathcal{B} is

$$e(\mathcal{B}) := \{p = e(\phi) : \phi \in \mathcal{B}\}.$$

This allows us to derive the likelihood and log-likelihood function. The log-likelihood function is strictly concave. Therefore, if the Maximum Likelihood Estimator (MLE) exists, then it is unique. Despite the boundedness of the likelihood function, MLE may not exist. Our main goal is to establish a new characterization of the existence of MLE and to apply it to specific families, exponential families spanned by Rademacher and Walsh functions, and exponential families of random graphs.

The key tool used in our discussion is the newly introduced notion of set of uniqueness, that is such set $U \subset \mathcal{X}$, that $\phi = 0$ is the only function of a given class of functions that equal zero on U . To be more specific, we show that the MLE for $e(\mathcal{B})$ and the i.i.d. sample $x_1, x_2, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ is the set of uniqueness for the nonnegative cone $\mathcal{B}_+ := \{\phi \in \mathcal{B} : \phi \geq 0\}$. We also restate this criterion as a linear programming problem.

Additionally, for proposed applications, we use the probability bounds established by classical tools from random graph theory and from the Coupon Collector's Problem, see [118, pp. 194-195], [68], and [143].

In the family spanned by Walsh functions, we use the properties of the hypercube $\{-1, 1\}^k$ ($k = \log_2 K$), of its subcubes and of the corresponding hypercube graph.

1.4 Graphical models

One of those branches of statistics that is known both for using the MLE and for applying penalized estimation is the theory of graphical models. A graphical model is a family of probability distributions of a finite collection of random variables X_1, X_2, \dots, X_N , which are encoded by N vertices of a graph (directed or undirected). The presence (or absence) of an edge between two vertices informs about the conditional dependence (or independence) between the corresponding vertices. When all random variables are Gaussian, we refer to Gaussian graphical models. In this case, the conditional independence structure is fully encoded by the location of zeros in the inverse of the covariance matrix (the precision matrix) of a vector $\mathbf{X} = (X_1, \dots, X_N)'$. We invite the reader interested in graphical models to the existing introductions for this topic [122, 187, 20]. For the discussion on the existence of MLE in discrete hierarchical models, which include graphical models, we refer to the article of Wang, Rauh and Massam [185]. The graphical model version of a LASSO estimator was proposed around fifteen years ago [130, 188, 79] and still gains popularity, cf. e.g. [37]. The graphical SLOPE was proposed independently by Sobczyk in his Ph.D. dissertation [166] and by Mazza-Anthony, Mazoure and Coates [129].

Let us come back to the conditional independence structure. It may be observed that the location of zeros in the precision matrix of \mathbf{X} is the same as in the Laplacian matrix of an underlying graph. In our discussion we take a closer look at such graphs and look for those Markov processes, the discretizations of which have a covariance matrix being an inverse of a slightly modified Laplacian matrix. As an example, when an underlying graph is a path of n vertices, then adding 1 to the top-left entry to its Laplacian matrix results in an inverse equal to $(\Sigma)_{i,j} = \min\{i, j\}$, which is a discretized version of a covariance matrix of a Wiener process.

1.5 Plan of the dissertation

The dissertation is based on five articles, three of which are published, one is currently under the review process, and one is an extended version of a preprint, which can be found on HAL and will be submitted soon.

1. Pattern recovery and signal denoising by SLOPE when the design matrix is orthogonal [165] (with P. Graczyk, B. Kołodziejek and M. Wilczyński), *Probability and Mathematical Statistics*, 42(2):283–302, 2022
2. Pattern recovery by SLOPE [23] (with M. Bogdan, X. Dupuis, P. Graczyk, B. Kołodziejek, P. Tardivel and M. Wilczyński), submitted, <https://arxiv.org/pdf/2203.12086.pdf>
3. Pattern Recovery in Penalized and Thresholded Estimation and its Geometry [93] (with P. Graczyk, U. Schneider and P. Tardivel), in preparation, <https://hal.science/hal-03262087v2/document>
4. Maximum likelihood estimation for discrete exponential families and random graphs [22] (with K. Bogdan and M. Bosy), *ALEA*, 19, 1045–1070 (2022)
5. Remarks on Laplacian of Graphical Models in Various Graphs [164], *Proceedings, GSI 2021, Paris, France, July 21–23, 2021*

The outline of the dissertation is the following:

Chapter 2 contains preliminaries of the presented research and describes the basic and novel notions used in latter chapters.

In Chapter 3 we focus on the SLOPE estimation in case of design matrix being orthogonal, that is, $\mathbf{X}'\mathbf{X} = n\mathbf{I}_p$. We also present new results on the strong consistency of the SLOPE estimators and on the strong consistency of pattern recovery by SLOPE when the design matrix is orthogonal and illustrate the advantages of the SLOPE clustering in the context of high frequency signal denoising.

Chapter 4 focuses on the SLOPE pattern recovery with no restrictions on the design matrix. Here, we show that recovery of the pattern can be determined by two criteria, called the positivity condition and the subdifferential condition. Later, we introduce a SLOPE irrepresentability condition, which generalizes the well known LASSO irrepresentability condition [82, 41] and then we illustrate it geometrically. Later on, we derive the refined asymptotic conditions on both consistency and strong consistency of the SLOPE estimate and of its pattern.

Chapter 5 concerns the problem of pattern recovery in a general case of a penalty being a polyhedral gauge. We establish the necessary and sufficient condition for the uniform uniqueness of the estimator. Then we derive conditions for the pattern recovery of polyhedral gauge penalized estimators and for their thresholded versions. Finally, we generalize the irrepresentability condition for penalizers being polyhedral gauges and discuss the connections between a pattern and a model from the article by Vaiter et al. [182].

Chapter 6 focuses on the necessary and sufficient conditions for the existence of the Maximum Likelihood Estimator (MLE) in discrete exponential families. The main tool in this article is the notion of sets of uniqueness, i.e. such subsets of a state space on which a function from a given class equals zero if and only if it equals zero on a whole space. We apply our results to the class of Rademacher and Walsh functions and to exponential random graph models (ERGM).

In Chapter 7 we discuss the connection between Gaussian graphical models and the covariance

matrices of discretized Wiener processes. For that reason, we use the discrete Laplacian matrix of an underlying graphical structure.

Chapter 2

Preliminaries and basic notions on penalized linear regression

2.1 Notation

The content of this section will be completed after unifying the notation in the whole dissertation.

- \mathbf{A}' — transpose of a matrix \mathbf{A}
- $|\mathbf{b}|_{\downarrow}$ — nonincreasing permutation of absolute values of coordinates of \mathbf{b}
- $\partial f(\mathbf{x})$ — subdifferential of f at \mathbf{x}
- $B_{\|\cdot\|}(\mathbf{b}, r)$ — ball of center \mathbf{b} and radius r in norm $\|\cdot\|$
- $cl(C)$ — closure of C
- $int(C)$ — interior of C
- $aff(C)$ — affine hull of (C)
- \vec{C} — unique linear space parallel to C
- $cl(C)$ — closure of C
- $bd(C)$ — boundary of C
- $ri(C)$ — relative interior of C
- $rb(C)$ — relative boundary of C
- $conv(C)$ — convex hull of C
- $cone(C)$ — conic hull of C
- $rec(C)$ — recession cone of C
- $P_C(\mathbf{x})$ — orthogonal projection of \mathbf{x} onto C
- C^* — polar set of C
- $\mathbf{1}_k = (1, \dots, 1)' \in \mathbb{R}^k$

- $\mathbb{1}_A(x) = \begin{cases} 1, & x \in A, \\ 0, & x \notin A. \end{cases}$ — characteristic function of a set A

For the convenience of the reader, when a set A is defined as a set of points satisfying the relation ϕ we reduce the notation from $\mathbb{1}_{\{x|\phi(x)\}}(x)$ to $\mathbb{1}_{\{\phi(x)\}}$, e.g. $\mathbb{1}_{\{x>0\}}$.

- $\mathbf{J}_n \in \mathbb{R}^{n \times n}$, $(\mathbf{J}_n)_{ij} = 1$ for every $1 \leq i, j \leq n$.
- $\text{supp}(\mathbf{v}) = \{i \in \{1, \dots, p\} \mid s_i \neq 0\}$ — support of a vector $\mathbf{v} \in \mathbb{R}^p$
- $\text{diag}(\mathbf{S}) \in \mathbb{R}^{p \times p}$ — diagonal matrix with $\mathbf{S} \in \mathbb{R}^p$ on the diagonal
($\text{diag}(\mathbf{S}))_{i,j} = s_i \mathbb{1}_{\{i=j\}}$)
- \mathcal{M}_p — set of all possible SLOPE patterns of $\mathbf{b} \in \mathbb{R}^p$
- $\text{sign}(x) := \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}$ — sign function
- $\text{sign}(\mathbf{x}) = (\text{sign}(x_1), \dots, \text{sign}(x_p))'$ — sign vector
- $\mathbb{R}^{k+} := \{\boldsymbol{\kappa} \in \mathbb{R}^k : \kappa_1 > \kappa_2 > \dots > \kappa_k > 0\}$
- \mathcal{S}_p — symmetric permutation group on the set $\{1, \dots, p\}$
- $Q^- := \max\{Q, 0\}$.

2.2 Penalized linear regression

2.2.1 Linear regression

We consider the following linear regression model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where $\mathbf{Y} \in \mathbb{R}^n$ is a response vector with n being a sample size, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix with p being the number of unknown parameters, $\boldsymbol{\beta} \in \mathbb{R}^p$ is an unknown parameter vector and $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ is a noise term.

The case of $\boldsymbol{\varepsilon} = \mathbf{0}$ ($\boldsymbol{\varepsilon} \neq \mathbf{0}$) is referred later as a noiseless (noisy) case.

In the penalized linear regression we want to find an estimator $\hat{\boldsymbol{\beta}}$ of a vector $\boldsymbol{\beta}$, which is of the form

$$\hat{\boldsymbol{\beta}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}_{\boldsymbol{\Lambda}}(\mathbf{b}) \right],$$

where, depending of the estimation method, $\boldsymbol{\Lambda}$ is a tuning parameter or a vector of tuning parameters (tuning vector). When dealing with a single tuning parameter, we denote it with a small letter λ , restricting a capital letter $\boldsymbol{\Lambda}$ for a vector $(\lambda_1, \dots, \lambda_p)'$.

As we are mostly interested with recovering a pattern of a vector $\boldsymbol{\beta}$, we are interested to define it in a general case. For that reason, we use a notion of a subdifferential and define two vectors to have the same pattern if the penalizer has the same subdifferential at those vectors

Definition 2.2.1 (Subdifferential). [41, p. 76] Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a convex function. A vector $\mathbf{d} \in \mathbb{R}^p$ is called a subgradient of f at point $\mathbf{x} \in \mathbb{R}^p$ if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})' \mathbf{d}, \quad \forall \mathbf{y} \in \mathbb{R}^p.$$

The set of all subgradients of the convex function f at $\mathbf{x} \in \mathbb{R}^p$ is called the subdifferential of f at \mathbf{x} , and it is denoted by $\partial f(\mathbf{x})$.

Remark 2.2.1. [17, p. 716] $\mathbf{x} \in \mathbb{R}^p$ is a minimum of a convex function f if and only if $\mathbf{0} \in \partial f(\mathbf{x})$.

Example 2.2.2. [89, Lemma D.5] The subdifferential of the ℓ_1 norm at $\mathbf{x} \in \mathbb{R}^p$ is given by

$$\partial \|\cdot\|_1(\mathbf{x}) = \partial |\cdot|(x_1) \times \cdots \times \partial |\cdot|(x_p) \text{ where } \partial |\cdot|(t) = \begin{cases} \{1\} & \text{if } t > 0 \\ [-1, 1] & \text{if } t = 0 \\ \{-1\} & \text{if } t < 0 \end{cases}$$

The subdifferential of the ℓ_∞ norm at $\mathbf{0}$ is the unit ball of the ℓ_1 norm and for $\mathbf{x} \in \mathbb{R}^p$ where $\mathbf{x} \neq \mathbf{0}$ this subdifferential is equal to

$$\partial \|\cdot\|_\infty(\mathbf{x}) = \left\{ \mathbf{s} \in \mathbb{R}^p : \|\mathbf{s}\|_1 = 1 \text{ and } \begin{cases} s_i x_i \geq 0 & \text{if } |x_i| = \|\mathbf{x}\|_\infty \\ s_i x_i = 0 & \text{otherwise} \end{cases} \right\}.$$

Definition 2.2.2 (Equality of patterns). Let $\text{pen} : \mathbb{R}^p \mapsto \mathbb{R}$ be a convex penalizer. We say that $\mathbf{x} \in \mathbb{R}^p$ and $\mathbf{z} \in \mathbb{R}^p$ have the same pattern with respect to pen when $\partial \text{pen}(\mathbf{x}) = \partial \text{pen}(\mathbf{z})$, where ∂pen represents the subdifferential of pen .

We say that the estimator $\hat{\boldsymbol{\beta}}$ recovers the pattern of $\boldsymbol{\beta}$ when

$$\partial \text{pen}(\hat{\boldsymbol{\beta}}) = \partial \text{pen}(\boldsymbol{\beta}).$$

Example 2.2.3 (LASSO pattern). In the LASSO regression, i.e., when $\text{pen}(\mathbf{x}) = \|\mathbf{x}\|_1$, the pattern of $\boldsymbol{\beta}$ identifies with its sign vector:

$$\mathbf{patt}(\boldsymbol{\beta}) = \mathbf{sign}(\boldsymbol{\beta}).$$

The proof of the above example goes straightforwardly from Example 2.2.2.

Theorem 2.2.4 (Subdifferential description of the SLOPE pattern [156]).

Let $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_p)$ satisfy $\lambda_1 > \dots > \lambda_p > 0$. Then

$$\mathbf{patt}(\mathbf{b}_1) = \mathbf{patt}(\mathbf{b}_2) \iff \partial J_{\boldsymbol{\Lambda}}(\mathbf{b}_1) = \partial J_{\boldsymbol{\Lambda}}(\mathbf{b}_2).$$

The pattern of SLOPE defined in Definition 2.2.2 coincides with the following definition.

Definition 2.2.3 (SLOPE pattern [23]). The SLOPE pattern is a function

$\mathbf{patt} : \mathbb{R}^p \rightarrow \mathbb{Z}^p$ such that

$$\mathbf{patt}(\mathbf{b})_i = \text{sign}(b_i) \text{rank}(|b_i|),$$

where $\text{rank}(|b_i|) \in \{1, 2, \dots, k\}$ is defined to be the number of $|c_j|$'s satisfying $|b_i| \geq |c_j|$, where $|c_1|, |c_2|, \dots, |c_k|$, $k \leq p$, are distinct non-zero values among $|b_1|, \dots, |b_p|$. We adopt the convention that $\text{rank}(0) = 0$.

Fact 2.2.5 (Basic properties of SLOPE pattern [156]).

- (a) for every $1 \leq l \leq \|\mathbf{patt}(\mathbf{b})\|_\infty$ there exists j such that $|\mathbf{patt}(\mathbf{b})_j| = l$,
- (b) $\text{sign}(\mathbf{patt}(\mathbf{b})) = \text{sign}(\mathbf{b})$ (sign preservation),
- (c) $|b_i| = |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| = |\mathbf{patt}(\mathbf{b})_j|$ (cluster preservation),
- (d) $|b_i| > |b_j| \Rightarrow |\mathbf{patt}(\mathbf{b})_i| > |\mathbf{patt}(\mathbf{b})_j|$ (hierarchy preservation).

Example 2.2.6. $\mathbf{patt}(4, 0, -1.5, 1.5, -4)' = (2, 0, -1, 1, -2)'$.

As a (nonzero) cluster in a vector $\mathbf{M} \in \mathbb{R}^p$ we denote a maximal (in terms of inclusion) set of indices i in $\text{supp}(\mathbf{M})$ such that the absolute value \mathbf{M}_i is the same. For example, the vector $(3, -3, 1, 3, 1, 4)'$ has three clusters: $\{1, 2, 4\}$, $\{3, 5\}$ and $\{6\}$.

Definition 2.2.4. Let $\mathbf{M} \neq \mathbf{0}$ be a pattern in \mathcal{M}_p with $k = \|\mathbf{M}\|_\infty$ nonzero clusters. The pattern matrix $\mathbf{U}_\mathbf{M} \in \mathbb{R}^{p \times k}$ is defined as follows

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

Example 2.2.7. Let $\boldsymbol{\beta} = (3, -4, -3, 0, 4, 10)'$. Then the pattern of $\boldsymbol{\beta}$ equals $\mathbf{M} = (1, -2, -1, 0, 2, 3)'$ and

$$\mathbf{U}_\mathbf{M} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}.$$

Note that $\boldsymbol{\beta}$ can be represented by a product of $\mathbf{U}_\mathbf{M}$ and the vector $(10, 4, 3)' \in \mathbb{R}^{3+}$. For every $\boldsymbol{\beta}$ with pattern \mathbf{M} such a vector in \mathbb{R}^{3+} exists.

Definition 2.2.5. Let $\mathbf{M} \neq \mathbf{0}$ be a pattern in \mathbb{R}^p and $k = \max\{\|\mathbf{M}\|_\infty, 1\}$.

For $\mathbf{X} \in \mathbb{R}^{n \times p}$ we define the clustered design matrix by $\widetilde{\mathbf{X}}_\mathbf{M} = \mathbf{X}\mathbf{U}_\mathbf{M} \in \mathbb{R}^{n \times k}$ and the clustered parameter by $\widetilde{\boldsymbol{\Lambda}}_\mathbf{M} = (\mathbf{U}_{|\mathbf{M}|_\downarrow})' \boldsymbol{\Lambda}$.

Example 2.2.8. Let $\mathbf{X} = (\mathbf{X}_1 | \mathbf{X}_2 | \mathbf{X}_3 | \mathbf{X}_4)$, $\mathbf{M} = (1, 2, -1, 0)'$

and $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_4)' \in \mathbb{R}^{4+}$. Then $|\mathbf{M}|_\downarrow = (2, 1, 1, 0)'$. The clustered matrix and the clustered parameter equal:

$$\widetilde{\mathbf{X}}_\mathbf{M} = (\mathbf{X}_2 | \mathbf{X}_1 - \mathbf{X}_3) \quad \text{and} \quad \widetilde{\boldsymbol{\Lambda}}_\mathbf{M} = \begin{pmatrix} \lambda_1 \\ \lambda_2 + \lambda_3 \end{pmatrix}.$$

2.3 Convex Polytopes and cones

2.3.1 Convex analysis

We consider \mathbb{R}^n with the norm $\|\cdot\|$.

Definition 2.3.1 (Affine set). [33, Sec. 2.1.2] A set $C \in \mathbb{R}^n$ is affine if for any $\mathbf{x}_1, \mathbf{x}_2 \in C$ and $\theta \in \mathbb{R}$ we have $\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in C$.

Definition 2.3.2 (Affine hull). [33, Sec. 2.1.2] The set of all affine combinations of points in a set $C \in \mathbb{R}^p$ is called the affine hull of C ($\text{aff}(C)$):

$$\text{aff}(C) := \{\theta_1 x_1 + \dots + \theta_k x_k : x_1, \dots, x_k \in C, \theta_1, \dots, \theta_k \in \mathbb{R}, \theta_1 + \dots + \theta_k = 1\}.$$

Remark 2.3.1. [33, Sec. 2.1.4] $\text{aff}(C)$ is the smallest affine set containing C .

Definition 2.3.3 (Relative interior). [33, Sec. 2.1.3] Let $C \subset \mathbb{R}^p$. The relative interior of a set C is defined as the interior of C with respect to its affine hull $\text{aff}(C)$:

$$\text{ri}(C) := \{\mathbf{x} \in C : B_{\|\cdot\|}(\mathbf{x}, r) \cap \text{aff}(C) \subset C \text{ for some } r > 0\}.$$

Definition 2.3.4 (Relative boundary). [33, Sec. 2.1.3]

$$\text{rb}(C) := \overline{C} \setminus \text{ri}C.$$

Definition 2.3.5 (Convex set). [33, Sec. 2.1.4] A set $C \in \mathbb{R}^n$ is convex if for any $\mathbf{x}_1, \mathbf{x}_2 \in C$ and $\theta \in [0, 1]$ we have $\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2 \in C$.

Definition 2.3.6 (Convex hull). [33, Sec. 2.1.4]

$$\text{conv}(C) := \{\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k : \mathbf{x}_i \in C, \theta_i \geq 0, \quad i = 1, \dots, k, \quad \theta_1 + \dots + \theta_k = 1\}.$$

Remark 2.3.2. [33, Sec. 2.1.4] $\text{conv}(C)$ is the smallest convex set containing C .

Definition 2.3.7 (Convex cone). [33, Sec. 2.1.5] A set $C \in \mathbb{R}^p$ is a convex cone if for any $\mathbf{x}_1, \mathbf{x}_2 \in C$ and $\theta_1, \theta_2 \geq 0$ we have $\theta_1 \mathbf{x}_1 + \theta_2 \mathbf{x}_2 \in C$.

Definition 2.3.8 (Conic hull). [33, Sec. 2.1.5]

$$\text{cone}(C) := \{\theta_1 \mathbf{x}_1 + \dots + \theta_k \mathbf{x}_k : \mathbf{x}_i \in C, \theta_i \geq 0, \quad i = 1, \dots, k\}.$$

Remark 2.3.3. [33, Sec. 2.1.5] $\text{cone}(C)$ is the smallest convex cone containing C .

Definition 2.3.9 (Recession cone). [193, Definition 1.11. 2.1.5] Let $C \in \mathbb{R}^p$ be a convex set. Then the recession cone of C is defined as

$$\text{rec}(C) := \{\mathbf{y} \in \mathbb{R}^p : \mathbf{x} + t\mathbf{y} \in C \text{ for all } \mathbf{x} \in C, t \geq 0\}.$$

Remark 2.3.4. [101, Prop. III. 1.2.1] A non-empty closed and convex set C is compact if and only if $\text{rec}(C) = \{\mathbf{0}\}$.

Definition 2.3.10 (Orthogonal projection). [33, Sec. 8.1] Let $C \subset \mathbb{R}^p$ be a closed set and let $\mathbf{x}_0 \in \mathbb{R}^p$. The orthogonal projection of \mathbf{x}_0 onto C is defined as

$$P_C(\mathbf{x}_0) := \arg \min_{\mathbf{x} \in C} \|\mathbf{x} - \mathbf{x}_0\|_2$$

Definition 2.3.11 (Half-space). [34, p. 9] Let $\mathbf{y} \in \mathbb{R}^p$ and $\alpha \in \mathbb{R}$. Then

$$K(\mathbf{y}, \alpha) := \{\mathbf{x} \in \mathbb{R}^p : \langle \mathbf{x}, \mathbf{y} \rangle \leq \alpha\}.$$

Definition 2.3.12 (Hyperplane). [34, p. 9] Let $\mathbf{y} \in \mathbb{R}^p$ and $\alpha \in \mathbb{R}$. Then

$$H(\mathbf{y}, \alpha) := \{\mathbf{x} \in \mathbb{R}^p : \langle \mathbf{x}, \mathbf{y} \rangle = \alpha\}.$$

Definition 2.3.13. Let $A \subset \mathbb{R}^p$ be an affine space, i.e.

$$\exists \mathbf{y}_1, \dots, \mathbf{y}_l \in \mathbb{R}^p, \alpha_1, \dots, \alpha_l \in \mathbb{R} \quad A = \bigcap_{i=1}^l H(\mathbf{y}_i, \alpha_i).$$

Then by \vec{A} we define the unique linear space parallel to A

$$\vec{A} := \bigcap_{i=1}^l H(\mathbf{y}_i, 0).$$

Definition 2.3.14 (Polar set). [34, p. 37] For any $M \subset \mathbb{R}^p$, the polar set M^* is defined by

$$M^* := \{\mathbf{y} \in \mathbb{R}^p : \forall \mathbf{x} \in M : \langle \mathbf{x}, \mathbf{y} \rangle \leq 1\}.$$

Equivalently, $M^* = \bigcap_{\mathbf{x} \in M} K(\mathbf{x}, 1)$.

2.3.2 Polytopes

We recall basic definitions and facts about polytopes, which we will use throughout the proofs. The following can be found in textbooks, such as [94] and [193].

A set $P \subseteq \mathbb{R}^p$ is called a polytope if it is the convex hull of a finite set of points $\{\mathbf{v}_1, \dots, \mathbf{v}_k\} \subseteq \mathbb{R}^p$, that is,

$$P = \text{conv}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}.$$

The dimension $\dim(P)$ of a polytope is defined as the dimension of $\text{aff}(P)$, the affine subspace spanned by P . An inequality $\mathbf{a}'\mathbf{x} \leq c$ is called a valid inequality of P if $P \subseteq \{\mathbf{x} \in \mathbb{R}^p : \mathbf{a}'\mathbf{x} \leq c\}$. A face F of P is any subset $F \subseteq P$ that satisfies

$$F = \{\mathbf{x} \in P : \mathbf{a}'\mathbf{x} = c\}, \text{ where } P \subseteq \{\mathbf{x} \in \mathbb{R}^p : \mathbf{a}'\mathbf{x} \leq c\},$$

for some $\mathbf{a} \in \mathbb{R}^p$ and $c \in \mathbb{R}$. Note that $F = \emptyset$ and $F = P$ are faces of P and that any face F is again a polytope. A non-empty face F with $F \neq P$ is called proper. A point $\mathbf{x}_0 \in P$ lies in $\text{ri}(P)$ if \mathbf{x}_0 is not contained in a proper face of P . We state two useful properties about faces in the following lemma.

Lemma 2.3.1. Let $P \subseteq \mathbb{R}^p$ be a polytope given by $P = \text{conv}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$, where $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^p$. The following properties hold.

- (a) If F and \tilde{F} are faces of P , then so is $F \cap \tilde{F}$.
- (b) Let L be an affine line contained in the affine hull of P . If $L \cap \text{ri}(P) \neq \emptyset$, then L intersects a proper face of P .

2.3.3 Pattern equivalence class

Definition 2.3.15 (Pattern equivalence class).

- (a) Let $\mathbf{x} \in \mathbb{R}^p$ and let pen be a polyhedral gauge. The pattern equivalence class $C_{\mathbf{x}}$ is the set of all vectors having the same subdifferential as \mathbf{x} :

$$C_{\mathbf{x}} := \{\mathbf{w} \in \mathbb{R}^p : \partial \text{pen}(\mathbf{w}) = \partial \text{pen}(\mathbf{x})\}.$$

Definition 2.3.16. Let $F \subset \mathbb{R}^p$. As F^\perp we denote the orthogonal complement of $\overrightarrow{\text{aff}(F)}$.

Example 2.3.5. For any $\mathbf{x} \in \mathbb{R}^p$, we have $\{\mathbf{x}\}^\perp = \{\mathbf{0}\}^\perp = \mathbb{R}^p$.

Lemma 2.3.2. Let $\mathbf{x}_0 \in F$. Then

$$\text{lin}(F - \mathbf{x}_0) = \overrightarrow{\text{aff}(F)}.$$

Proof. By [101, Sec. III.5.3.] we have

$$\text{lin}(\mathbb{R}_+(F - \mathbf{x}_0)) = \text{aff}(F) - \mathbf{x}_0.$$

The statement is proved after observing that $\text{lin}(F - \mathbf{x}_0) = \text{lin}(\mathbb{R}_+(F - \mathbf{x}_0))$ and $\text{aff}(F) - \mathbf{x}_0 = \overrightarrow{\text{aff}(F)}$. \square

2.3.4 Normal cones

Definition 2.3.17. [31, p. 15], [101, p.136] The normal cone to a convex set C at a point $\bar{\mathbf{x}} \in C$, written $N_C(\bar{\mathbf{x}})$ is the convex cone of normal vectors, i.e. vectors \mathbf{d} in \mathbb{R}^p such that $\langle \mathbf{d}, \mathbf{x} - \bar{\mathbf{x}} \rangle \leq 0$ for all points \mathbf{x} in C .

Definition 2.3.18. [70, Def. 4.12.] If F is a face of a closed convex set K and $\mathbf{x} \in \text{ri}F$, then, $N_K(\mathbf{x})$ does not depend on $\mathbf{x} \in \text{ri} F$ and is denoted by $N_K(F)$ and is called the cone of normals of K in F .

Lemma 2.3.3. [70, Lemma 3.1.] Let K be a closed convex set in \mathbb{R}^n . To each $\mathbf{x} \in \mathbb{R}^n$ there exists a unique $\mathbf{x}' \in K$ such that

$$\|\mathbf{x} - \mathbf{x}'\| = \inf_{\mathbf{y} \in K} \|\mathbf{x} - \mathbf{y}\|.$$

Definition 2.3.19. [70, Def. 3.2.] The map

$$\begin{aligned} p_K : \mathbb{R}^n &\longrightarrow K \\ \mathbf{x} &\mapsto p_K(\mathbf{x}) = \mathbf{x}' \end{aligned}$$

of Lemma 2.3.3 is called the nearest point map relative to K .

Definition 2.3.20. [70, Def. 4.7.] Let \mathbf{x} be a point of the closed convex set K . We call

$$N_K(\mathbf{x}) := -\mathbf{x} + p_K^{-1}(\mathbf{x})$$

the normal cone of K at \mathbf{x} .

Proposition 2.3.1. The Definitions 2.3.17 and 2.3.20 are equivalent.

Proof. By [101, Proposition III. 5.3.3.], the point $\mathbf{s} \in \mathbb{R}^p$ belongs to $N_K(\mathbf{x})$ from Definition 2.3.17 if and only if $\mathbf{x} = p_K(\mathbf{x} + \mathbf{s})$. It implies that $p_K^{-1}(\mathbf{x}) = \mathbf{x} + N_K(\mathbf{x})$, thus

$$N_K(\mathbf{x}) = -\mathbf{x} + p_K^{-1}(\mathbf{x}).$$

\square

Corollary 2.3.1. The normal cone $N_K(\mathbf{x})$ is invariant under translations of the set K .

Theorem 2.3.6. [70, Theorem 4.13.] Let K be a compact convex set in \mathbb{R}^p and \mathbf{x} a relative interior point of a face $F \neq \emptyset$ of K . Then

$$\{\text{ri } N_K(\mathbf{x}) : F \text{ is a face of } K\} = \{\text{ri } N_K(F) : F \text{ is a face of } K\}$$

is a partition (disjoint covering) of \mathbb{R}^p .

Definition 2.3.21 (Conjugate face). [34, p. 40] Let F be a face of a polytope P . The conjugate face F° of a dual polytope P^* is defined as follows

$$F^\circ := \{\mathbf{y} \in P^* : \langle \mathbf{x}, \mathbf{y} \rangle = 1 \quad \forall \mathbf{x} \in F\}.$$

Equivalently, we get

$$F^\circ = \bigcap_{\mathbf{x} \in P} K(\mathbf{x}, 1) \cap \bigcap_{\mathbf{x}' \in F} H(\mathbf{x}', 1),$$

where the first and second intersection are equal respectively to P^* and $\text{aff}(F^\circ)$.

Consider a convex polytope $P = \{\mathbf{v} \in \mathbb{R}^p : \langle \mathbf{s}_j, \mathbf{v} \rangle \leq r_j, j = 1, \dots, m\}$, cf. [101, p. 138]. We define its active set for $\mathbf{x} \in \mathbb{R}^p$ as $A_P(\mathbf{x}) := \{j = 1, \dots, m : \langle \mathbf{s}_j, \mathbf{x} \rangle = r_j\}$

Proposition 2.3.2. [94, Proposition 14.1] We have

$$N_P(\mathbf{x}) = \text{cone}(\{\mathbf{s}_j : j \in A_P(\mathbf{x})\}).$$

Theorem 2.3.7. Let B^* be the polar set of the polytope B . Let $\mathbf{x} \in \text{ri}(F)$, where F is a face of B^* . Then the normal cone at \mathbf{x} to B^* is given by $N_{B^*}(\mathbf{x}) = \mathbb{R}_+ F^\circ$, where $F^\circ \subset B$ is the conjugate face to F .

Proof. ($\mathbb{R}_+ F^\circ \subset N_{B^*}(\mathbf{x})$):

Let $\mathbf{y} \in \mathbb{R}_+ F^\circ$. Then there exists such $\gamma \geq 0$, that $\mathbf{y} \in \gamma B$ and $\langle \mathbf{x}, \mathbf{y} \rangle = \gamma$ for every $\mathbf{x} \in F$. It implies that

$$\langle \mathbf{v}, \mathbf{y} \rangle \leq \gamma \quad \forall \mathbf{v} \in B^* \text{ and} \quad \langle \mathbf{x}, \mathbf{y} \rangle = \gamma \quad \forall \mathbf{x} \in F.$$

Therefore $\langle \mathbf{v}, \mathbf{y} \rangle \leq \langle \mathbf{x}, \mathbf{y} \rangle$ for all $\mathbf{v} \in B^*$, i.e. $\langle \mathbf{v} - \mathbf{x}, \mathbf{y} \rangle \leq 0$ for all $\mathbf{v} \in B^*$, which means that $\mathbf{y} \in N_{B^*}(\mathbf{x})$.

($N_{B^*}(\mathbf{x}) \subset \mathbb{R}_+ F^\circ$):

Let $\mathbf{y} \in N_{B^*}(\mathbf{x})$, i.e. $\langle \mathbf{y}, \mathbf{v} - \mathbf{x} \rangle \leq 0$ for every $\mathbf{v} \in B^*$. Thus the maximum of the value of $\langle \mathbf{y}, \mathbf{v} - \mathbf{x} \rangle$ over $\mathbf{v} \in B^*$ is attained at $\mathbf{v} = \mathbf{x}$. Moreover, as $\langle \mathbf{y}, \mathbf{v} \rangle \leq \langle \mathbf{y}, \mathbf{x} \rangle$ for every $\mathbf{v} \in B^*$, we get $\mathbf{y} \in \langle \mathbf{y}, \mathbf{x} \rangle B$. Now it suffices to prove that for every $\mathbf{v} \in F$ we have $\langle \mathbf{y}, \mathbf{v} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$.

Take any $\mathbf{v} \in F$. Since $\mathbf{x} \in \text{ri}(F) \subset F$, we have $N_{B^*}(\mathbf{x}) \subset N_{B^*}(\mathbf{v})$. Then $\langle \mathbf{y}, \mathbf{w} - \mathbf{x} \rangle \leq 0$ for every $\mathbf{w} \in B^*$, which implies that $\langle \mathbf{y}, \mathbf{v} - \mathbf{x} \rangle \leq 0$, i.e. $\langle \mathbf{y}, \mathbf{v} \rangle \leq \langle \mathbf{y}, \mathbf{x} \rangle$.

On the other hand, $\mathbf{y} \in N_{B^*}(\mathbf{v})$. Therefore $\langle \mathbf{y}, \mathbf{w} - \mathbf{v} \rangle \leq 0$ for every $\mathbf{w} \in B^*$, which implies that $\langle \mathbf{y}, \mathbf{x} - \mathbf{v} \rangle \leq 0$, i.e. $\langle \mathbf{y}, \mathbf{x} \rangle \leq \langle \mathbf{y}, \mathbf{v} \rangle$.

Thus for every $\mathbf{v} \in F$ we have $\langle \mathbf{y}, \mathbf{v} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$, hence $\mathbf{y} \in \langle \mathbf{y}, \mathbf{x} \rangle F^\circ \subset \mathbb{R}_+ F^\circ$. \square

2.3.5 Polyhedral gauges

Definition 2.3.22 (Gauge). [155, p. 53] Let $K \subset \mathbb{R}^p$ be a closed convex set containing $\mathbf{0}$. The gauge function of K is defined as

$$\inf\{\lambda \geq 0 : \mathbf{x} \in \lambda K\}$$

The gauge function of K is also known as the Minkowski functional of K . When K is a polyhedron, its gauge function is known as a polyhedral gauge pen and can be written as the maximum of linear functions [152, 133]:

$$\forall \mathbf{x} \in \mathbb{R}^p, \text{pen}(\mathbf{x}) = \max\{0, \mathbf{u}'_1 \mathbf{x}, \dots, \mathbf{u}'_l \mathbf{x}\}, \text{ for some } \mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p.$$

Note that a polyhedral gauge, whose unit ball $\{\mathbf{x} \in \mathbb{R}^p : \text{pen}(\mathbf{x}) \leq 1\}$ is a bounded and symmetric with respect to the origin polyhedron, is a polyhedral norm.

2.3.6 Thresholded penalized least squares estimation

Definition 2.3.23 (Thresholded penalized least squares estimator).

Let pen be a polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$ and $\lambda > 0$. Given $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$, we say that $\tilde{\boldsymbol{\beta}}$ is a thresholded estimator of $\hat{\boldsymbol{\beta}}$ if $\partial \text{pen}(\tilde{\boldsymbol{\beta}}) \subseteq \partial \text{pen}(\hat{\boldsymbol{\beta}})$.

Example 2.3.8. In SLOPE optimization problem, a vector $\tilde{\boldsymbol{\beta}} = (5, 5, 5, 0, 0)'$ is a thresholded estimator of $\hat{\boldsymbol{\beta}}^{\text{SLOPE}} = (4, 4, 2, 0, 0)'$.

2.3.7 Permutahedron

Definition 2.3.24. Let $\boldsymbol{\Lambda} \in \mathbb{R}^p$. A permutahedron $P_{\boldsymbol{\Lambda}}$ is a convex hull of all possible permutations of $\boldsymbol{\Lambda}$:

$$P_{\boldsymbol{\Lambda}} := \text{Conv}((\lambda_{\pi(1)}, \dots, \lambda_{\pi(p)})' : \pi \in \mathcal{S}_p). \quad (2.3.1)$$

Definition 2.3.25 (Signed permutahedron). Let $\boldsymbol{\Lambda} \in \mathbb{R}^p$. A signed permutahedron $C_{\boldsymbol{\Lambda}}$ is a convex hull of all possible permutations and sign changes of $\boldsymbol{\Lambda}$:

$$C_{\boldsymbol{\Lambda}} = \{\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_p) \in \mathbb{R}^p : \sum_{j \leq i} |\pi|_{(j)} \leq \sum_{j \leq i} \lambda_j : i = 1, 2, \dots, p\}. \quad (2.3.2)$$

$C_{\boldsymbol{\Lambda}}$ is a unit ball in a norm $J_{\boldsymbol{\Lambda}}^*$ dual to the SLOPE norm $J_{\boldsymbol{\Lambda}}$.

2.3.8 Basics on Moore-Penrose inverse

The notion of the Moore-Penrose inverse is crucial in the SLOPE irrepresentability condition and is outlined below, see [91, 16].

If \mathbf{A} is an $n \times p$ real matrix then a $p \times n$ matrix \mathbf{A}^+ is called a Moore-Penrose inverse of \mathbf{A} if

$$\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}, \quad \mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$$

and if the matrices $\mathbf{A}\mathbf{A}^+$ and $\mathbf{A}^+\mathbf{A}$ are symmetric.

There always exists a unique real Moore-Penrose inverse \mathbf{A}^+ of a real matrix \mathbf{A} . In some cases it may be computed quickly to a more convenient form:

- If $\mathbf{A}'\mathbf{A}$ is an invertible matrix, then $\mathbf{A}^+ = (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'$,
- If $\mathbf{A}\mathbf{A}'$ is an invertible matrix, then $\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^{-1}$,
- $\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^+$,
- $\mathbf{0}^+ = \mathbf{0}$.

In general, the Moore-Penrose inverse \mathbf{A}^+ is computed using the formula $\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^+$ and the Singular Value Decomposition, which allows to do it numerically.

Remark 2.3.9.

- $\mathbf{A}\mathbf{A}^+$ is an orthogonal projector onto $\text{col}(\mathbf{A})$,
- $\text{col}(\mathbf{A}^+) = \text{col}(\mathbf{A}')$.

2.3.9 Functional analysis

Definition 2.3.26 (Dual norm). [33, Appendix A.1.6] Let $\|\cdot\|$ be a norm on \mathbb{R}^p . The dual norm, denoted $\|\cdot\|_*$, is defined as

$$\|\mathbf{x}\|_* = \sup\{\mathbf{x}'\mathbf{b} : \|\mathbf{b}\| \leq 1\}.$$

In particular, the dual sorted ℓ_1 norm J_Λ has an explicit expression given in [136]:

$$J_\Lambda^*(\mathbf{b}) = \max \left\{ \frac{|b|_{(1)}}{\lambda_1}, \frac{\sum_{i=1}^2 |b|_{(i)}}{\sum_{i=1}^2 \lambda_i}, \dots, \frac{\sum_{i=1}^p |b|_{(i)}}{\sum_{i=1}^p \lambda_i} \right\}, \quad \mathbf{b} \in \mathbb{R}^p.$$

2.3.10 Tools from optimization

Definition 2.3.27 (Quasi-convexity). [33, Sec. 3.4.1] A function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is called *quasiconvex* if its domain and all its sublevel sets

$$S_\alpha = \{\mathbf{x} : f(\mathbf{x}) \leq \alpha\},$$

for $\alpha \in \mathbb{R}$, are convex. If a function $-f$ is quasiconvex, then we say that f is *quasiconcave*.

Definition 2.3.28 (Semi-continuity). [101, Definition 3.2.1] We say that the convex function f is *lower semi-continuous*, if

$$\liminf_{\mathbf{x} \rightarrow \mathbf{x}_0} f(\mathbf{x}) \geq f(\mathbf{x}_0)$$

for all $\mathbf{x}_0 \in \mathbb{R}^p$. If a function $-f$ is lower semi-continuous, then we say that f is *upper semi-continuous*.

Theorem 2.3.10 (max-min inequality). [33, Sec. 5.4.1]

$$\sup_{\mathbf{b} \in B} \inf_{\mathbf{a} \in A} f(\mathbf{a}, \mathbf{b}) \leq \inf_{\mathbf{a} \in A} \sup_{\mathbf{b} \in B} f(\mathbf{a}, \mathbf{b})$$

Definition 2.3.29 (Saddle point). [33, Sec. 5.4.2] The pair $(\bar{\mathbf{a}}, \bar{\mathbf{b}}) \in A \times B$ is a *saddle point* for a function f (and A and B) if

$$f(\bar{\mathbf{a}}, \mathbf{b}) \leq f(\bar{\mathbf{a}}, \bar{\mathbf{b}}) \leq f(\mathbf{a}, \bar{\mathbf{b}})$$

for all $\mathbf{a} \in A$ and $\mathbf{b} \in B$. In other words, we have

$$f(\bar{\mathbf{a}}, \bar{\mathbf{b}}) = \inf_{\mathbf{a} \in A} f(\mathbf{a}, \bar{\mathbf{b}}) = \sup_{\mathbf{b} \in B} f(\bar{\mathbf{a}}, \mathbf{b}).$$

The existence of the saddle point implies that it attains the equality in the max-min inequality:

$$\sup_{\mathbf{b} \in B} \inf_{\mathbf{a} \in A} f(\mathbf{a}, \mathbf{b}) = \inf_{\mathbf{a} \in A} \sup_{\mathbf{b} \in B} f(\mathbf{a}, \mathbf{b}) = f(\bar{\mathbf{a}}, \bar{\mathbf{b}}).$$

Theorem 2.3.11 (Sion). [7, Ch. 7, Theorem 7.] Suppose that

- (a) X and Y are convex compact subsets,

(b) for all $\mathbf{y} \in Y$ the function $\mathbf{x} \mapsto f(\mathbf{x}, \mathbf{y})$ is quasiconvex and lower semi-continuous,

(c) for all $\mathbf{x} \in X$ the function $\mathbf{y} \mapsto f(\mathbf{x}, \mathbf{y})$ is quasiconcave and upper semi-continuous,

Then there exists a saddle point $\{\bar{\mathbf{x}}, \bar{\mathbf{y}}\}$.

Theorem 2.3.12 (Hardy-Littlewood-Pólya rearrangement inequality). [98, Theorem 368] For every $x_1 \leq x_2 \leq \dots \leq x_k \in \mathbb{R}$, $y_1 \leq y_2 \leq \dots \leq y_k \in \mathbb{R}$ and for every permutation $\sigma \in \mathcal{S}_k$ we have

$$x_k y_1 + x_{k-1} y_2 + \dots + x_1 y_k \leq x_{\sigma(1)} y_1 + x_{\sigma(2)} y_2 \dots x_{\sigma(k)} y_k \leq x_1 y_1 + x_2 y_2 + \dots + x_k y_k.$$

2.3.11 Tools from probability

Definition 2.3.30 (Convergence in distribution). [19, p. 329] Let X_n and X be random variables with respective distribution functions F_n and F . If $F_n \rightarrow F$, then X_n is said to converge in distribution or in law to X , written $X_n \xrightarrow{d} X$.

Definition 2.3.31 (Convergence in probability). [19, p. 330]

$$\lim_n \mathbb{P}[|X_n - X| > \varepsilon] = 0, \quad \text{for every } \varepsilon > 0.$$

Definition 2.3.32 (Consistency). $\hat{\beta}_n$ is a consistent estimator of β , if β_n converges to β in probability.

Definition 2.3.33 (Almost sure convergence). [19, pp. 59-60]

$$\mathbb{P}\left(\omega : \lim_{n \rightarrow \infty} X_n(\omega) \rightarrow X(\omega)\right) = 1.$$

Definition 2.3.34 (Strong consistency). $\hat{\beta}_n$ is a strong consistent estimator of β , if β_n converges to β almost surely.

Theorem 2.3.13 (Borel-Cantelli Lemma). [19, Theorem 4.3. and 4.4.]

(a) If $\sum_n \mathbb{P}(A_n)$ converges, then $\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 0$.

(b) If $\{A_n\}$ is an independent sequence of events and $\sum_n \mathbb{P}(A_n)$ diverges, then

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} A_n\right) = 1.$$

Theorem 2.3.14 (Strong law of large numbers). [73, VII.8, Theorem 1] Let X_1, X_2, \dots be i.i.d. random variables with $\mathbb{E}(X) = 0$. For each $n \geq 1$ denote $S_n := X_1 + \dots + X_n$. Then $n^{-1}S_n \rightarrow 0$ with probability 1.

Theorem 2.3.15 (Lindeberg-Feller Multivariate Central Limit Theorem). [147, Sec. 4.3.2] Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be independent random variables with multivariate distribution having mean μ and covariance matrix Σ . Then the distribution of $\sqrt{n}\left(\frac{1}{n}\sum_{i=1}^n \mathbf{X}_i - \mu\right)$ tends to $N(\mathbf{0}, \Sigma)$.

Theorem 2.3.16 (CLT for linearly negative quadrant dependent variables). [113, Corollary 1.1] Let $\{\xi_i\}$ be a centered sequence of linearly negative quadrant dependent random variables

such that $\{\xi_i^2\}$ is a uniformly integrable family, and let $\{a_{ni}, 1 \leq i \leq n\}$ be a triangular array of nonnegative numbers such that

$$\sup_n \sum_{i=1}^n \frac{a_{ni}^2}{\sigma_n^2} < \infty$$

and

$$\max_{1 \leq i \leq n} \frac{a_{ni}}{\sigma_n} \rightarrow 0 \quad \text{as} \quad n \rightarrow \infty, \quad (2.3.3)$$

where $\sigma^2 = \text{Var}\left(\sum_{i=1}^n a_{ni}\xi_i\right)$. If

$$\sum_{j:|i-j|\geq u} \text{Cov}(\xi_i, \xi_j) \rightarrow 0 \quad \text{as} \quad u \rightarrow \infty \quad \text{uniformly in} \quad i \geq 1,$$

then, as $n \rightarrow \infty$,

$$\frac{1}{\sigma_n} \sum_{i=1}^n a_{ni}\xi_i \xrightarrow{d} N(0, 1).$$

Theorem 2.3.17 (Tail inequality). *If $X \sim \mathcal{N}(0, 1)$, then*

$$\mathbb{P}(X > t) \leq \frac{1}{t} e^{-t^2/2} \frac{1}{\sqrt{2\pi}}.$$

Proof. [43] Note that for $x \geq t > 0$ we have $\frac{x}{t} \geq 1$. Therefore

$$\mathbb{P}(X \geq t) = \frac{1}{\sqrt{2\pi}} \int_t^\infty 1 \cdot e^{-x^2/2} dx \leq \frac{1}{\sqrt{2\pi}} \int_t^\infty \frac{x}{t} \cdot e^{-x^2/2} dx = \frac{1}{t\sqrt{2\pi}} e^{-t^2/2}.$$

□

Theorem 2.3.18 (Law of Iterated Logarithm for triangular array). [117, Theorem 1 (i)] *Let $\{\varepsilon_i\}_{i \in \mathbb{Z}}$ be independent random variables such that*

$$\mathbb{E}\varepsilon_n = 0 \quad \text{and} \quad \mathbb{E}\varepsilon_n^2 = \sigma^2 \quad \text{for all } n \quad \text{and} \quad \sup_n \mathbb{E}\varepsilon_n^2 < \infty \quad (2.3.4)$$

for some $r > 2$. Let $\{a_{ni}\}$ be a double array of constants satisfying $\sum_{-\infty}^{\infty} a_{ni}^2 < \infty$ for every n .

Define

$$S_n = \sum_{-\infty}^{\infty} a_{ni}\varepsilon_n. \quad (2.3.5)$$

Assume that as $n \rightarrow \infty$,

$$A_n = \sum_{-\infty}^{\infty} a_{ni}^2 \rightarrow \infty,$$

and

$$\sup_i a_{ni}^2 = o(A_n(\log A_n))^{-\rho} \quad \text{for all } \rho > 0. \quad (2.3.6)$$

If there exist constants $c_i \geq 0$ and $d > 2/r$ such that

$$\|a_n - a_m\|^2 \leq \left(\sum_{i=m+1}^n c_i \right)^d \quad \text{for} \quad n > m > m_0 \quad (2.3.7)$$

and

$$\left(\sum_{i=m_0}^n c_i \right)^d = O(A_n) \quad a.s. \quad n \rightarrow \infty, \quad (2.3.8)$$

then

$$\limsup_{n \rightarrow \infty} \frac{|S_n|}{(2A_n \log \log A_n)^{1/2}} \leq \sigma \quad a.s.$$

Chapter 3

Pattern recovery and signal denoising by SLOPE when the design matrix is orthogonal

3.1 Introduction

3.1.1 Introduction and motivations

The content of this chapter may be found in the recently published article of the author of the dissertation, Graczyk, Kołodziejek and Wilczyński [165]. To start our discussion on the pattern recovery by penalized linear regression methods, we start with the SLOPE estimator and the simplest case of the design matrix \mathbf{X} being orthogonal. For that reason let us recall the Linear Multiple Regression. It concerns the model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\mathbf{Y} \in \mathbb{R}^n$ is an output vector, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a fixed design matrix, $\boldsymbol{\beta} \in \mathbb{R}^p$ is an unknown vector of predictors and $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ is a noise vector. The primary goal is to estimate $\boldsymbol{\beta}$. In the low-dimensional setting, that is, when the number of predictors p is not greater than the number of explanatory variables n and \mathbf{X} is of full rank, the ordinary least squares estimator $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ has an exact formula $\hat{\boldsymbol{\beta}}^{\text{OLS}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$. For practical reasons there is an urge to avoid the high-dimensionality curse, therefore we want the estimate to be sparse, i.e., to be describable by a smaller number of parameters. Several solutions were proposed to deal with such problem. One of them, the Least Absolute Shrinkage and Selection Operator (LASSO [47, 176]) involves penalizing the residual sum of squares $\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2$ with an ℓ_1 norm of $\hat{\boldsymbol{\beta}}$ multiplied by a tuning parameter λ :

$$\hat{\boldsymbol{\beta}}^{\text{LASSO}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_1 \right].$$

We turn our focus on one of the extensions of LASSO, which is the Sorted ℓ_1 Penalized Estimator (SLOPE [27, 26, 189]). In addition, SLOPE allows one to clusterize the similar coefficients of $\boldsymbol{\beta}$. In SLOPE, the ℓ_1 norm is replaced by its sorted version $J_{\boldsymbol{\Lambda}}$, which depends on the tuning vector $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_p) \in \mathbb{R}^p$, where $\lambda_1 \geq \dots \geq \lambda_p \geq 0$:

$$J_{\boldsymbol{\Lambda}}(\boldsymbol{\beta}) := \sum_{i=1}^p \lambda_i |\boldsymbol{\beta}|_{(i)},$$

where $\{|\boldsymbol{\beta}|_{(i)}\}_{i=1}^p$ is a decreasing permutation of absolute values of $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_p$:

$$\hat{\boldsymbol{\beta}}^{\text{SLOPE}} := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_{\boldsymbol{\Lambda}}(\mathbf{b}) \right].$$

The case of $\mathbf{\Lambda}$ being an arithmetic sequence was studied by Bondell and Reich [29] and called the Octagonal Shrinkage and Clustering Algorithm for Regression (OSCAR). The special case of SLOPE with $\lambda_1 = \lambda_2 = \dots = \lambda_p > 0$ is LASSO. For $\mathbf{\Lambda} = (0, \dots, 0)$ we obtain the OLS estimator.

Clustering the predictors allows for additional dimension reduction by identifying variables with the same absolute values of the regression coefficients. Recently, interest has increased in methods that cluster highly correlated predictors [30, 85, 128, 137, 139, 168]. SLOPE is ideal for this task, since it is capable to identify the low-dimensional structure, which is called the SLOPE pattern, defined by Schneider and Tardivel with the subdifferential of the SLOPE norm $J_{\mathbf{\Lambda}}$, see Theorem 2.2.4.

The clustering properties of SLOPE have been studied before, cf. [29, 76], but the researchers consider strongly correlated predictors, which are used in financial mathematics to group the assets with respect to their partial correlation with hedge fund return times series [116]. We decided to consider the pattern recovery in more general case. We start with supposing the orthogonal design

$$\mathbf{X}'\mathbf{X} = n\mathbf{I}_p. \quad (3.1.1)$$

This is a classical and natural assumption in the case of experimental data; cf. [176]. Moreover, in the asymptotic case, where $n \rightarrow \infty$ and p is fixed, it is usually supposed that $\mathbf{X}'\mathbf{X}/n \rightarrow \mathbf{C} > 0$, cf. [192, 194]. In (3.1.1) the design matrix \mathbf{X} is orthogonal. Then, the Euclidean norm of each n -dimensional column of \mathbf{X} equals n . If it was 1, the terms of \mathbf{X} would approach zero for large n , which is not natural. This class of matrices is being widely used in signal analysis, [146, 48]. For general \mathbf{X} the problem is considered in Chapter 4.

To study the properties of SLOPE we often use the closed unit ball $C_{\mathbf{\Lambda}}$ in the dual norm of $J_{\mathbf{\Lambda}}$, which was studied, for example, by Zeng and Figueiredo [189]. This dual ball is described explicitly as a signed permutahedron, see, e.g. [136, 156]:

$$C_{\mathbf{\Lambda}} = \left\{ \boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_p) \in \mathbb{R}^p : \sum_{j \leq i} |\pi_{(j)}| \leq \sum_{j \leq i} \lambda_j : i = 1, 2, \dots, p \right\}. \quad (3.1.2)$$

Here we prove novel results on the strong consistency of SLOPE both in estimation and in pattern recovery. We also introduce a new method, based on the minimax approach, to find the relationship between $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$ and $\hat{\boldsymbol{\beta}}^{\text{OLS}}$.

3.1.2 Outline

In Section 3.2 we derive the connections between $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$ and $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ in the orthogonal design. We use the minimax theorem of Sion, cf. [7]. In Section 3.3 we focus on the properties of $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$. We use the geometric interpretation of SLOPE to explain its ability to identify the SLOPE pattern and provide new theoretical results on support recovery and clustering properties using a representation of SLOPE as a function of the ordinary least squares (OLS) estimator. A similar approach for LASSO was used by Ewald and Schneider, cf. [71].

To analyze the asymptotic properties of the SLOPE estimator, e.g., its consistency, we have to assume that the sample size n tends to infinity. Therefore, in Section 3.4 we define a sequence of linear regression models

$$\mathbf{Y}^{(n)} = \mathbf{X}^{(n)}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_n^{(n)}.$$

In this sequence, the response vector $\mathbf{Y}^{(n)} \in \mathbb{R}^n$, the design matrix $\mathbf{X}^{(n)} \in \mathbb{R}^{n \times p}$ and the error term $\boldsymbol{\varepsilon}_n^{(n)} = (\varepsilon_1^{(n)}, \varepsilon_2^{(n)}, \dots, \varepsilon_n^{(n)})' \in \mathbb{R}^n$ varies with n and has the normal distribution $N(0, \sigma^2 \mathbf{I}_n)$. We make no assumptions about the relations between $\boldsymbol{\varepsilon}_n^{(n)}$ and $\boldsymbol{\varepsilon}_m^{(m)}$ for $n \neq m$.

In this chapter we consider the specific, but statistically important, model in which $n \geq p$ and the columns of \mathbf{X} are orthogonal. The orthogonality assumption allows us to derive, by simple techniques, relatively precise results on the SLOPE estimator (e.g., Theorem 3.1), which seem unavailable when columns of \mathbf{X} are not orthogonal. We provide the conditions under which the SLOPE estimator is strongly consistent. Additionally, in case when for each n the design matrix is orthogonal, we provide the conditions on the sequence of tuning parameters such that SLOPE is strongly consistent in the pattern recovery. In Section 3.5 we show the applications of the SLOPE clustering in terms of high frequency signal denoising and illustrate them with simulations. The Appendix covers the proofs of technical results.

Substantially more difficult techniques based on subdifferential calculus are developed in [23]. These techniques are used in [23] to establish the properties of the SLOPE estimator in the general case, where the columns of \mathbf{X} are not orthogonal and p may be much larger than n . However, the asymptotic results of [23] are derived under stronger assumptions than those of this paper. In [23] the sequence of error terms $\varepsilon^{(n)}$ is incremental and the sequence of tuning parameters has the form $\lambda_n = \alpha_n \Lambda$, where α_n is a given sequence of positive numbers and $\Lambda \in \mathbb{R}^p$ is fixed. In this paper we make no assumptions about the relations between $\varepsilon^{(n)}$ and $\varepsilon^{(m)}$ for $n \neq m$, and the sequence of tuning parameters has a general form.

3.2 Approach by minimax theorem

3.2.1 Technical results

Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ be a real-valued matrix. Let r_{SLOPE} denote the minimum value of the SLOPE criterion, attained by $\hat{\boldsymbol{\beta}}^{\text{SLOPE}}$, i.e.

$$r_{\text{SLOPE}} := \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_{\Lambda}(\mathbf{b}) \right] = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_2^2 + J_{\Lambda}(\hat{\boldsymbol{\beta}}^{\text{SLOPE}}).$$

Since

$$\|\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_2 \leq \sqrt{p} \|\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_{\infty} \quad \text{and} \quad \lambda_1 \|\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_{\infty} \leq J_{\Lambda}(\hat{\boldsymbol{\beta}}^{\text{SLOPE}}) \leq r_{\text{SLOPE}},$$

it follows that

$$\lambda_1 \|\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_2 \leq \sqrt{p} r_{\text{SLOPE}} \leq \sqrt{p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{0}\|_2^2 + J_{\Lambda}(\mathbf{0}) \right] = \frac{\sqrt{p}}{2} \|\mathbf{Y}\|_2^2.$$

We immediately get the following result.

Corollary 3.2.1. $\|\hat{\boldsymbol{\beta}}^{\text{SLOPE}}\|_2^2 \leq M_0$, where $M_0 = \left(\frac{p \|\mathbf{Y}\|_2^4}{4\lambda_1^2} \right)$.

From this corollary it is seen that we can clearly limit our search to vectors $\boldsymbol{\beta}$ from the compact set $\mathcal{M} \subset \mathbb{R}^p$ defined by $\mathcal{M} := \{\mathbf{b} \in \mathbb{R}^p : \|\mathbf{b}\|_2^2 \leq M_0\}$. Therefore, we can equivalently define a SLOPE solution by

$$\hat{\boldsymbol{\beta}}^{\text{SLOPE}} = \arg \min_{\mathbf{b} \in \mathcal{M}} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_{\Lambda}(\mathbf{b}) \right]. \quad (3.2.1)$$

Proposition 3.2.1. Let C_{Λ} be the unit ball in the dual SLOPE norm. Then, for each $\mathbf{b} \in \mathbb{R}^p$,

$$J_{\Lambda}(\mathbf{b}) = \max_{\boldsymbol{\pi} \in C_{\Lambda}} \boldsymbol{\pi}'\mathbf{b}. \quad (3.2.2)$$

The proof is a simple application of the definition of the dual norm and the reflexivity of $(\mathbb{R}^p, J_\Lambda) = (\mathbb{R}^p, J_\Lambda^*)^*$. Thus

$$J_\Lambda(\mathbf{b}) = \|\mathbf{b}\|_{(\mathbb{R}^p, J_\Lambda)} = \sup_{\mathbf{x}: J_\Lambda^*(\mathbf{x}) \leq 1} \mathbf{x}'\mathbf{b}.$$

Remark 3.2.1.

- (a) A different, longer proof is given in [27, Proposition 1.1]
- (b) The formula (3.2.2) holds in much greater generality for Lovász extensions in place of the J_Λ norm, see [132].

3.2.2 Saddle point

In this section we continue the assumption that $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a real-valued matrix. Let the function $r : \mathcal{M} \times C_\Lambda \rightarrow \mathbb{R}$ be defined by

$$r(\mathbf{b}, \boldsymbol{\pi}) := \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \boldsymbol{\pi}'\mathbf{b}.$$

As an immediate consequence of (3.2.1) and Proposition 3.2.1 we obtain

$$\begin{aligned} r_{\text{SLOPE}} &= \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_\Lambda(\mathbf{b}) \right] = \min_{\mathbf{b} \in \mathcal{M}} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_\Lambda(\mathbf{b}) \right] \\ &= \min_{\mathbf{b} \in \mathcal{M}} \max_{\boldsymbol{\pi} \in C_\Lambda} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \boldsymbol{\pi}'\mathbf{b} \right] = \min_{\mathbf{b} \in \mathcal{M}} \max_{\boldsymbol{\pi} \in C_\Lambda} r(\mathbf{b}, \boldsymbol{\pi}). \end{aligned}$$

It turns out that the order of maximization over $\boldsymbol{\pi} \in C_\Lambda$ and minimization over $\mathbf{b} \in \mathcal{M}$ can be switched without affecting the result. To see this, note that both C_Λ and \mathcal{M} are convex and compact. Furthermore, for each fixed $\boldsymbol{\pi} \in C_\Lambda$, $r(\mathbf{b}, \boldsymbol{\pi})$ is a convex continuous function with respect to $\mathbf{b} \in \mathcal{M}$ and, for each fixed $\mathbf{b} \in \mathcal{M}$, $r(\mathbf{b}, \boldsymbol{\pi})$ is concave with respect to $\boldsymbol{\pi} \in C_\Lambda$ (in fact, it is linear). Therefore, all assumptions of the Sion's minimax theorem are fulfilled (see [7, p. 218]) and thus there exists a saddle point $(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*) \in \mathcal{M} \times C_\Lambda$ such that

$$\begin{aligned} \max_{\boldsymbol{\pi} \in C_\Lambda} \min_{\mathbf{b} \in \mathcal{M}} r(\mathbf{b}, \boldsymbol{\pi}) &= \min_{\mathbf{b} \in \mathcal{M}} r(\mathbf{b}, \boldsymbol{\pi}^*) = r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*) \\ &= \max_{\boldsymbol{\pi} \in C_\Lambda} r(\boldsymbol{\beta}^*, \boldsymbol{\pi}) = \min_{\mathbf{b} \in \mathcal{M}} \max_{\boldsymbol{\pi} \in C_\Lambda} r(\mathbf{b}, \boldsymbol{\pi}) = r_{\text{SLOPE}}. \end{aligned}$$

In the next section we shall see that the first coordinate of any saddle point $(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*)$ is the SLOPE estimator.

3.2.3 SLOPE solution when \mathbf{X} has full column rank

Since for each fixed $\boldsymbol{\pi} \in C_\Lambda$, the function $r(\mathbf{b}, \boldsymbol{\pi})$ is convex with respect to $\mathbf{b} \in \mathcal{M}$, any point $\mathbf{b}\boldsymbol{\pi} \in \mathcal{M}$, at which the gradient $\frac{\partial r(\mathbf{b}, \boldsymbol{\pi})}{\partial \mathbf{b}}$ is zero, is a global minimum. If we rewrite $r(\mathbf{b}, \boldsymbol{\pi})$ as

$$r(\mathbf{b}, \boldsymbol{\pi}) = \frac{1}{2} \mathbf{Y}'\mathbf{Y} - \mathbf{Y}'\mathbf{X}\mathbf{b} + \frac{1}{2} \mathbf{b}'\mathbf{X}'\mathbf{X}\mathbf{b} + \boldsymbol{\pi}'\mathbf{b}$$

and differentiate with respect to \mathbf{b} , we obtain

$$\frac{\partial r(\mathbf{b}, \boldsymbol{\pi})}{\partial \mathbf{b}} = -\mathbf{X}'(\mathbf{Y} - \mathbf{X}\mathbf{b}) + \boldsymbol{\pi}.$$

Equating this gradient with $\mathbf{0}$ gives the following equation for the optimum point \mathbf{b}_π :

$$\mathbf{X}'\mathbf{X}\mathbf{b}_\pi = \mathbf{X}'\mathbf{Y} - \boldsymbol{\pi}. \quad (3.2.3)$$

Substituting this into the equation for $r(\mathbf{b}_\pi, \boldsymbol{\pi})$ and using the fact that $(\mathbf{X}'\mathbf{X})^{-1}$ exists, we find that

$$\begin{aligned} r(\mathbf{b}_\pi, \boldsymbol{\pi}) &= \frac{1}{2}\mathbf{Y}'\mathbf{Y} - \mathbf{b}'_\pi\mathbf{X}'\mathbf{Y} + \frac{1}{2}\mathbf{b}'_\pi\mathbf{X}'\mathbf{X}\mathbf{b}_\pi + \boldsymbol{\pi}'\mathbf{b}_\pi \\ &= \frac{1}{2}\mathbf{Y}'\mathbf{Y} - \mathbf{b}'_\pi\mathbf{X}'\mathbf{Y} + \mathbf{b}'_\pi\mathbf{X}'\mathbf{X}\mathbf{b}_\pi + \boldsymbol{\pi}'\mathbf{b}_\pi - \frac{1}{2}\mathbf{b}'_\pi\mathbf{X}'\mathbf{X}\mathbf{b}_\pi \\ &= \frac{1}{2}\mathbf{Y}'\mathbf{Y} - \frac{1}{2}\mathbf{b}'_\pi\mathbf{X}'\mathbf{X}\mathbf{b}_\pi = \frac{1}{2}\mathbf{Y}'\mathbf{Y} - \frac{1}{2}\mathbf{b}'_\pi\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{b}_\pi \\ &= \frac{1}{2}\mathbf{Y}'\mathbf{Y} - \frac{1}{2}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi})'(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi}). \end{aligned}$$

Let $p_j = |\{i : |m_i| = k + 1 - j\}|$ be the number of elements of the j^{th} cluster of $\boldsymbol{\beta}$, $P_j = \sum_{i \leq j} p_i$ and $P_{k+1} = p$.

Lemma 3.2.1. *Assume that \mathbf{X} has full column rank. Let $\boldsymbol{\pi}^* = (\pi_1^*, \dots, \pi_p^*)' \in C_\Lambda$ be any solution of*

$$\boldsymbol{\pi}^* = \arg \min_{\boldsymbol{\pi} \in C_\Lambda} \left[(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi})'(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi}) \right]$$

and let $\boldsymbol{\beta}^* = (\beta_1^*, \dots, \beta_p^*)'$ be the corresponding point from \mathcal{M} given by

$$\boldsymbol{\beta}^* = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi}^*).$$

Then, $(\boldsymbol{\pi} - \boldsymbol{\pi}^*)'\boldsymbol{\beta}^* \leq 0$, for all $\boldsymbol{\pi} \in C_\Lambda$ and hence

- (a) $\text{sign}(\beta_i^*) \cdot \text{sign}(\pi_i^*) \geq 0$, $i = 1, 2, \dots, p$,
- (b) $(|\pi_1^*|, \dots, |\pi_p^*|)$ and $(|\beta_1^*|, \dots, |\beta_p^*|)$ are similarly sorted, i.e. if $|(patt(\boldsymbol{\beta}))_i| = k + 1 - j$, then $|\pi_i^*| \in \{|\pi^*|_{(P_{j-1}+1)}, \dots, |\pi^*|_{(P_j)}\}$,
- (c) for any permutation τ satisfying $|\beta_{\tau(1)}^*| \geq \dots \geq |\beta_{\tau(p)}^*|$, if there is a $k \in \{2, \dots, p\}$ such that $\sum_{i=1}^{k-1} |\pi_{\tau(i)}^*| < \sum_{i=1}^{k-1} \lambda_i$ and $|\pi_{\tau(k)}^*| > 0$, then $|\beta_{\tau(k-1)}^*| = |\beta_{\tau(k)}^*|$.

The proof is given in the Appendix. An immediate consequence of the Lemma is the following result.

Lemma 3.2.2. *Assume that \mathbf{X} has full column rank. The point $(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*)$ defined as in Lemma 3.2.1 is the saddle point of the function $r(\mathbf{b}, \boldsymbol{\pi})$.*

The proof is given in the Appendix. We use the last lemma to prove the main result of this section.

Theorem 3.2.2. *Assume that \mathbf{X} has full column rank. Let the point $\boldsymbol{\beta}^*$ be defined as in Lemma 3.2.1. Then $\boldsymbol{\beta}^*$ is the SLOPE estimator of $\boldsymbol{\beta}$.*

Proof. Using the fact that $\max_{\pi \in C_{\Lambda}} r(\beta^*, \pi) = \min_{b \in \mathcal{M}} \max_{\pi \in C_{\Lambda}} r(b, \pi)$ (see previous lemma) we have

$$\begin{aligned} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta^*\|_2^2 + J_{\Lambda}(\beta^*) &= \max_{\pi \in C_{\Lambda}} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta^*\|_2^2 + \pi' \beta^* \right] \\ &= \max_{\pi \in C_{\Lambda}} r(\beta^*, \pi) = \min_{b \in \mathcal{M}} \max_{\pi \in C_{\Lambda}} r(b, \pi) = \min_{b \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y} - \mathbf{X}b\|_2^2 + J_{\Lambda}(b) \right]. \end{aligned}$$

□

Corollary 3.2.2. *In the linear model satisfying $\frac{1}{n} \mathbf{X}'\mathbf{X} = \mathbf{I}_p$ we have*

$$\hat{\beta}^{\text{OLS}} - \hat{\beta}^{\text{SLOPE}} = \frac{1}{n} \pi^* = \frac{1}{n} \arg \min_{\pi \in C_{\Lambda}} \left\| \hat{\beta}^{\text{OLS}} - \pi \right\|_2^2 = \arg \min_{\pi \in C_{\Lambda/n}} \left\| \hat{\beta}^{\text{OLS}} - \pi \right\|_2^2,$$

is the proximal projection of $\hat{\beta}^{\text{OLS}}$ onto $C_{\Lambda/n}$.

Projections onto C_{Λ} are widely used in [132] in the study of the notion of degrees of freedom. However, the Corollary 3.2.2 is not stated there explicitly.

Remark 3.2.3. Assume that \mathbf{X} has full column rank. For each $\pi \in C_{\Lambda}$, the point b_{π} defined in (3.2.3) belongs to

$$\left\{ b \in \mathbb{R}^p : \|b\|_2^2 \leq M \right\},$$

where M is chosen so that $M > \max\{M_0, M_1\}$ with

$$M_1 := \max_{\pi \in C_{\Lambda}} \|(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \pi)\|_2^2 \leq M.$$

3.3 Properties of SLOPE in the orthogonal design

3.3.1 SLOPE vs. OLS

By Theorem 3.2.2 and Corollary 3.2.2, when $\frac{1}{n} \mathbf{X}'\mathbf{X} = \mathbf{I}_p$, the orthogonal projection of the ordinary least squares estimator $\hat{\beta}^{\text{OLS}} = \frac{1}{n} \mathbf{X}'\mathbf{Y}$ onto the unit ball $C_{\Lambda/n}$ is equal to $\hat{\beta}^{\text{OLS}} - \hat{\beta}^{\text{SLOPE}}$. For $\Lambda = (200, 100)'$ and $n = 50$ this property is illustrated in Figure 3.1. The figure represents $\hat{\beta}^{\text{SLOPE}}$ (black arrows) depending on the localization of $\hat{\beta}^{\text{OLS}}$ in the orthogonal design. For $\hat{\beta}^{\text{OLS}}$ being the blue point located in the area labeled by $(1, 0)$ the first component of $\hat{\beta}^{\text{SLOPE}}$ is positive and the second is null. For $\hat{\beta}^{\text{OLS}}$ being the yellow point located on the area labeled by $(-1, 1)$ both components of $\hat{\beta}^{\text{SLOPE}}$ have equal absolute value (clusterization), but their signs are opposite. For $\hat{\beta}^{\text{OLS}}$ being the red point located on the area labeled by $(1, 2)$, both components of $\hat{\beta}^{\text{SLOPE}}$ are positive and the first component is smaller than the second one. The blue polytope is the dual SLOPE unit ball C_{Λ} and labels

$$\mathcal{M}_2 = \{(0, 0), (\pm 1, 0), (0, \pm 1), (\pm 1, \pm 1), (\pm 2, \pm 1), (\pm 1, \pm 2)\}$$

associated to the areas of this figure correspond to all SLOPE patterns for $n = 50$ and $p = 2$. In the orthogonal design, one may also explicitly compute the SLOPE estimator. Indeed, by the Corollary 3.2.2, $\hat{\beta}^{\text{SLOPE}}$ is the image of $\hat{\beta}^{\text{OLS}}$ by the proximal operator of the SLOPE norm. Therefore, this operator has a closed form formula [26, 175, 61]. This explicit expression gives an analytical way to learn that the SLOPE solution is sparse and built of clusters.

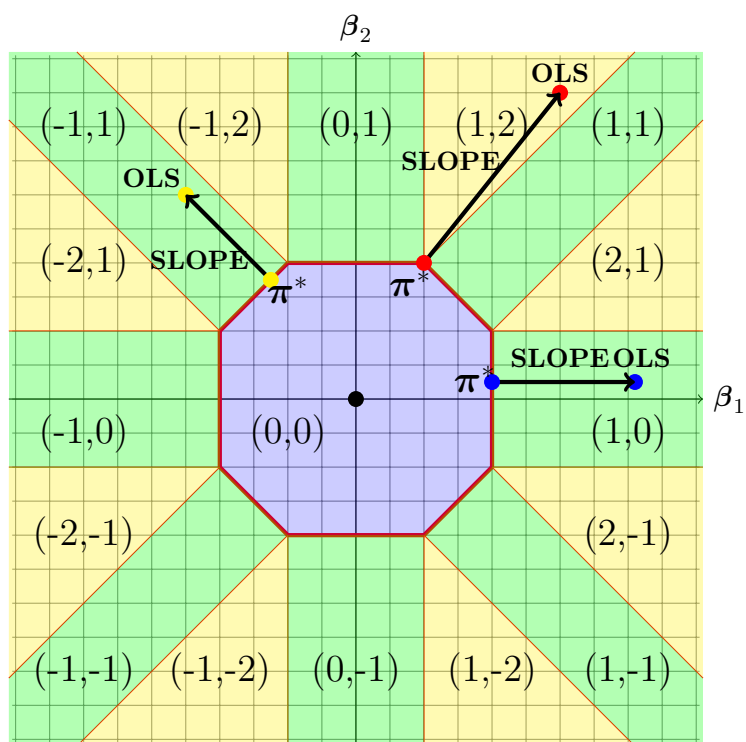


Figure 3.1: The dual unit ball $C_{\Lambda/n}$ for $\Lambda = (200, 100)'$ and examples of $\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$ in the orthogonal design for $n = 50$ and $p = 2$. The labels of each colored set refer to the pattern of $\hat{\beta}^{\text{SLOPE}}$ for $\hat{\beta}^{\text{OLS}}$ lying in this set. The arrows point from $(\hat{\beta}^{\text{OLS}} - \hat{\beta}^{\text{SLOPE}})$ to $\hat{\beta}^{\text{OLS}}$.

Lemma 3.3.1. *In the linear model satisfying $\frac{1}{n}\mathbf{X}'\mathbf{X} = \mathbf{I}_p$ we have*

$$\arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_\Lambda(\mathbf{b}) \right] = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \left\| \widehat{\boldsymbol{\beta}}^{\text{OLS}} - \mathbf{b} \right\|_2^2 + J_\Lambda(\mathbf{b}) \right]. \quad (3.3.1)$$

Our proof of this Lemma is given in the Appendix. The next theorem gives a sufficient condition for the clustering effect of the SLOPE estimator in the orthogonal design.

Theorem 3.3.1. *Consider a linear model with orthogonal design $\frac{1}{n}\mathbf{X}'\mathbf{X} = \mathbf{I}_p$. Let π be a permutation of $(1, 2, \dots, p)$ such that*

$$\left| \widehat{\boldsymbol{\beta}}_{\pi(1)}^{\text{OLS}} \right| \geq \left| \widehat{\boldsymbol{\beta}}_{\pi(2)}^{\text{OLS}} \right| \geq \dots \geq \left| \widehat{\boldsymbol{\beta}}_{\pi(p)}^{\text{OLS}} \right|.$$

For $i \in \{1, 2, \dots, p-1\}$,

$$\text{if } \left| \widehat{\boldsymbol{\beta}}_{\pi(i)}^{\text{OLS}} \right| - \left| \widehat{\boldsymbol{\beta}}_{\pi(i+1)}^{\text{OLS}} \right| \leq \frac{\lambda_i - \lambda_{i+1}}{n}, \text{ then } \left| \widehat{\boldsymbol{\beta}}_{\pi(i)}^{\text{SLOPE}} \right| = \left| \widehat{\boldsymbol{\beta}}_{\pi(i+1)}^{\text{SLOPE}} \right|.$$

Proof. By Lemma 3.3.1, in the orthogonal design, $\widehat{\boldsymbol{\beta}}^{\text{SLOPE}}$ is the proximal map of $J_{\Lambda/n}(\cdot)$ at $\widehat{\boldsymbol{\beta}}^{\text{OLS}}$. The result may be inferred from [26, Lemma 2.3]. \square

In the following theorem we derive the necessary and sufficient conditions under which SLOPE in the orthogonal design recovers the support of the vector

$$\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)', \text{ i.e. } \widehat{\beta}_i^{\text{SLOPE}} = 0 \iff \beta_i = 0.$$

Theorem 3.3.2. *Under orthogonal design $\frac{1}{n}\mathbf{X}'\mathbf{X} = \mathbf{I}_p$, let π be a permutation of $(1, 2, \dots, p)$ that satisfies $\left| \widehat{\boldsymbol{\beta}}_{\pi(1)}^{\text{OLS}} \right| \geq \left| \widehat{\boldsymbol{\beta}}_{\pi(2)}^{\text{OLS}} \right| \geq \dots \geq \left| \widehat{\boldsymbol{\beta}}_{\pi(p)}^{\text{OLS}} \right|$. Without loss of generality suppose that $\text{supp}(\boldsymbol{\beta}) = \{1, 2, \dots, p_0\}$ with $p_0 < p$. The necessary and sufficient condition for SLOPE to identify the set of relevant covariables is:*

- (a) $\min_{1 \leq i \leq p_0} \left| \widehat{\boldsymbol{\beta}}_i^{\text{OLS}} \right| > \max_{p_0+1 \leq i \leq p} \left| \widehat{\boldsymbol{\beta}}_i^{\text{OLS}} \right|,$
- (b) $\sum_{i=k}^{p_0} \left| \widehat{\boldsymbol{\beta}}_{\pi(i)}^{\text{OLS}} \right| > \frac{1}{n} \sum_{i=k}^{p_0} \lambda_i, \quad \text{for } k = 1, 2, \dots, p_0,$
- (c) $\sum_{i=p_0+1}^k \left| \widehat{\boldsymbol{\beta}}_{\pi(i)}^{\text{OLS}} \right| \leq \frac{1}{n} \sum_{i=p_0+1}^k \lambda_i, \quad \text{for } k = p_0 + 1, p_0 + 2, \dots, p.$

Proof. The result may be inferred from the properties of the proximal SLOPE [27, Lemma 2.3 and Lemma 2.4] and from Lemma 3.3.1. \square

3.4 Asymptotic properties of SLOPE

In this section we discuss several asymptotic properties of SLOPE estimators in the low-dimensional regression model in which p is fixed and the sample size n tends to infinity. For each $n \geq 1$ we consider a linear model

$$\mathbf{Y}^{(n)} = \mathbf{X}^{(n)}\boldsymbol{\beta} + \boldsymbol{\varepsilon}^{(n)}, \quad (3.4.1)$$

where $\mathbf{Y}^{(n)} = (y_1^{(n)}, y_2^{(n)}, \dots, y_n^{(n)})' \in \mathbb{R}^n$ is a vector of observations, $\mathbf{X}^{(n)} \in \mathbb{R}^{n \times p}$ is a deterministic design matrix with $\text{rank}(\mathbf{X}^{(n)}) = p$, $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)' \in \mathbb{R}^p$ is a vector of unknown

regression coefficients and $\boldsymbol{\varepsilon}^{(n)} = (\varepsilon_1^{(n)}, \varepsilon_2^{(n)}, \dots, \varepsilon_n^{(n)})' \in \mathbb{R}^n$ is a noise term, which has the normal distribution $N(0, \sigma^2 \mathbf{I}_n)$. We make no assumptions about the dependence between $\boldsymbol{\varepsilon}^{(n)}$ and $\boldsymbol{\varepsilon}^{(m)}$ for $n \neq m$. In particular, $\boldsymbol{\varepsilon}^{(n)}$ does not need to be a subsequence of $\boldsymbol{\varepsilon}^{(m)}$.

When defining the sequence $(\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}})$ of SLOPE estimators, we assume that the tuning vector varies with n . More precisely, for each $n \geq 1$ its coefficients $\lambda_1^{(n)} \geq \lambda_2^{(n)} \geq \dots \geq \lambda_p^{(n)} \geq 0$ are fixed and $\lambda_1^{(n)} > 0$. By $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ we denote the SLOPE estimator corresponding to the tuning vector $\boldsymbol{\Lambda}^{(n)} = (\lambda_1^{(n)}, \dots, \lambda_p^{(n)})'$:

$$\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{Y}^{(n)} - \mathbf{X}^{(n)} \mathbf{b}\|_2^2 + J_{\boldsymbol{\Lambda}^{(n)}}(\mathbf{b}) \right]. \quad (3.4.2)$$

3.4.1 Strong consistency of the SLOPE estimator

Let us recall the definition of a strongly consistent estimator $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ of $\boldsymbol{\beta}$, i.e. $\forall \boldsymbol{\beta} \in \mathbb{R}^p$ we have $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \rightarrow \boldsymbol{\beta}$ almost surely. Below, we discuss strong consistency of the sequence $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ of SLOPE estimators, defined by (3.4.2).

Theorem 3.4.1. *Consider the linear regression model (3.4.1) and assume that*

$$\lim_n n^{-1} (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} = \mathbf{C},$$

where \mathbf{C} is a positive definite matrix. Let $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$, $n \geq 1$, be the SLOPE estimator corresponding to the tuning vector $\boldsymbol{\Lambda}^{(n)} = (\lambda_1^{(n)}, \lambda_2^{(n)}, \dots, \lambda_p^{(n)})'$.

(a) If $\lim_{n \rightarrow \infty} \frac{\lambda_1^{(n)}}{n} = 0$, then $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \boldsymbol{\beta}$.

(b) If $\lim_{n \rightarrow \infty} \frac{\lambda_1^{(n)}}{n} = \lambda_0 > 0$ and if the true parameter $\boldsymbol{\beta}$ satisfies the inequality $\lambda_0 \|\boldsymbol{\beta}\|_\infty > \boldsymbol{\beta}' \mathbf{C} \boldsymbol{\beta} / 2$, then $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ does not converge to $\boldsymbol{\beta}$. Hence, $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ is not strongly consistent for $\boldsymbol{\beta}$.

Before proving the above theorems, we start with stating a simple technical lemma. It follows quickly from the Borel-Cantelli Lemma and the tail inequality:

If $Z \sim N(0, 1)$, then $\mathbb{P}(Z > t) \leq t^{-1} e^{-t^2/2} / \sqrt{2\pi}$, $t > 0$.

Lemma 3.4.1. *Assume that $(Q_n)_{n \in \mathbb{N}}$ is a sequence of Gaussian random variables, defined on the same probability space, which converges in distribution to $N(0, \sigma^2)$ for some $\sigma \in (0, \infty)$. Then, for any $\delta > 0$,*

$$\lim_{n \rightarrow \infty} \frac{Q_n}{(\log(n))^{1/2+\delta}} = 0 \quad \text{a.s.}$$

Our proof of the strong consistency of SLOPE is based on the strong consistency of the OLS estimator. The latter result is a folklore and we prove it in our setting.

Proposition 3.4.1. *Consider the linear regression model (3.4.1).*

If $\lim_n n^{-1} (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} = \mathbf{C}$, where \mathbf{C} is positive definite, then $\hat{\boldsymbol{\beta}}_n^{\text{OLS}} \xrightarrow{\text{a.s.}} \boldsymbol{\beta}$.

Proof. We have

$$\widehat{\boldsymbol{\beta}}_n^{\text{OLS}} - \boldsymbol{\beta} = ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} (\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} - \boldsymbol{\beta} = ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} (\mathbf{X}^{(n)})' \boldsymbol{\varepsilon}^{(n)}.$$

Then $\sqrt{n} \left(\widehat{\boldsymbol{\beta}}_n^{\text{OLS}} - \boldsymbol{\beta} \right)$ has the normal distribution $N(0, \sigma^2 (n^{-1} (\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1})$ and its components satisfy the assumptions of Lemma 3.4.1. Since $\log(n)^{1/2+\delta} = o(\sqrt{n})$, we get the assertion by Lemma 3.4.1. \square

Proof of Theorem 3.4.1. (a) It follows from Theorem 3.2.1 that there exists a vector $\boldsymbol{\pi}_n^* \in C_{\Lambda^{(n)}}$ such that

$$\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} = ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} ((\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} - \boldsymbol{\pi}_n^*).$$

Since $\boldsymbol{\pi}_n^*$ takes values in $C_{\Lambda^{(n)}}$, it follows that $\|\boldsymbol{\pi}_n^*\|_\infty \leq \lambda_1^{(n)}$. Hence,

$$\frac{\boldsymbol{\pi}_n^*}{n} \xrightarrow{a.s.} \mathbf{0}, \quad (3.4.3)$$

because $\left\| \frac{\boldsymbol{\pi}_n^*}{n} \right\|_\infty \leq \frac{\lambda_1^{(n)}}{n} \rightarrow 0$. The assumption that $\text{rank}(\mathbf{X}^{(n)}) = p$ implies that the matrix $(\mathbf{X}^{(n)})' \mathbf{X}^{(n)}$ is invertible and therefore the least squares estimator of $\boldsymbol{\beta}$ is unique and has the form $\widehat{\boldsymbol{\beta}}_n^{\text{OLS}} = ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} (\mathbf{X}^{(n)})' \mathbf{Y}^{(n)}$. Combining with (3.4.3) the fact that $\widehat{\boldsymbol{\beta}}_n^{\text{OLS}} \xrightarrow{a.s.} \boldsymbol{\beta}$, we conclude that

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} &= ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} ((\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} - \boldsymbol{\pi}_n^*) = \widehat{\boldsymbol{\beta}}_n^{\text{OLS}} - ((\mathbf{X}^{(n)})' \mathbf{X}^{(n)})^{-1} \boldsymbol{\pi}_n^* \\ &= \widehat{\boldsymbol{\beta}}_n^{\text{OLS}} - \left(\frac{(\mathbf{X}^{(n)})' \mathbf{X}^{(n)}}{n} \right)^{-1} \frac{\boldsymbol{\pi}_n^*}{n} \xrightarrow{a.s.} \boldsymbol{\beta} - \mathbf{C}^{-1} \mathbf{0} = \boldsymbol{\beta}. \end{aligned}$$

(b) Since $\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ minimizes over $\mathbf{b} \in \mathbb{R}^p$ the function

$$l(\mathbf{b}) := \frac{1}{2} \|\mathbf{Y}^{(n)} - \mathbf{X}^{(n)} \mathbf{b}\|_2^2 + J_{\Lambda^{(n)}}(\mathbf{b})$$

and since $\lambda_1^{(n)} \|\mathbf{b}\|_\infty \leq J_{\Lambda^{(n)}}(\mathbf{b})$, it follows that

$$\begin{aligned} 0 &\leq l(0) - l(\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}}) = (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} \\ &\quad - \frac{1}{2} (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} \widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} - J_{\Lambda^{(n)}}(\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}}) \\ &\leq (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} - \frac{1}{2} (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} \widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \\ &\quad - \lambda_1^{(n)} \|\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}}\|_\infty = (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} \widehat{\boldsymbol{\beta}}_n^{\text{OLS}} \\ &\quad - \frac{1}{2} (\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})' (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} \widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} - \lambda_1^{(n)} \|\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}}\|_\infty. \end{aligned}$$

The last equality follows from the fact that $(\mathbf{X}^{(n)})' \mathbf{Y}^{(n)} = (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} \widehat{\boldsymbol{\beta}}_n^{\text{OLS}}$. Suppose to the contrary that the true parameter $\boldsymbol{\beta}$ satisfies $\lambda_0 \|\boldsymbol{\beta}\|_\infty > \boldsymbol{\beta}' \mathbf{C} \boldsymbol{\beta} / 2$ (which is the case when $\|\boldsymbol{\beta}\|_\infty$ is sufficiently close to 0) and that $\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{a.s.} \boldsymbol{\beta}$. Then, using the facts that $\widehat{\boldsymbol{\beta}}_n^{\text{OLS}} \xrightarrow{a.s.} \boldsymbol{\beta}$ and that $\lim_n n^{-1} (\mathbf{X}^{(n)})' \mathbf{X}^{(n)} = \mathbf{C}$, we have

$$0 \leq \frac{l(0) - l(\widehat{\boldsymbol{\beta}}_n^{\text{SLOPE}})}{n} \xrightarrow{a.s.} \boldsymbol{\beta}' \mathbf{C} \boldsymbol{\beta} - \frac{1}{2} \boldsymbol{\beta}' \mathbf{C} \boldsymbol{\beta} - \lambda_0 \|\boldsymbol{\beta}\|_\infty = \frac{1}{2} \boldsymbol{\beta}' \mathbf{C} \boldsymbol{\beta} - \lambda_0 \|\boldsymbol{\beta}\|_\infty,$$

which provides a contradiction. This proves the first part of Theorem 3.4.1 (b). To prove the second part, note that strong consistency requires convergence for any value of the parameter β . \square

Remark 3.4.2. The proof of Theorem 3.4.1 (b) does not exclude the possibility that $\widehat{\beta}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \beta$ when β satisfies $\lambda_0 \|\beta\|_\infty \leq \beta' C \beta / 2$.

3.4.2 Asymptotic pattern recovery in the orthogonal design

We again consider a sequence of linear models (3.4.1) but this time we assume that for each n the deterministic design matrix $\mathbf{X}^{(n)}$ of size $n \times p$ satisfies

$$(\mathbf{X}^{(n)})' \mathbf{X}^{(n)} = n \mathbf{I}_p. \quad (3.4.4)$$

As usual, we assume Gaussian errors $\varepsilon^{(n)} \sim N(0, \sigma^2 \mathbf{I}_n)$.

Let $\widehat{\beta}_n^{\text{SLOPE}} = \left(\widehat{\beta}_1^{\text{SLOPE}}(n), \dots, \widehat{\beta}_p^{\text{SLOPE}}(n) \right)'$ be the SLOPE estimator defined by (3.4.2). With the above notation we present the main result of this section.

Theorem 3.4.3. *Assume that*

$$\lim_{n \rightarrow \infty} \frac{\lambda_1^{(n)}}{n} = 0$$

and that there exists $\delta > 0$ such that

$$\liminf_{n \rightarrow \infty} \frac{\lambda_i^{(n)} - \lambda_{i+1}^{(n)}}{\sqrt{n} (\log(n))^{1/2 + \delta}} = m > 0 \quad \text{for } i = 1, \dots, p-1. \quad (3.4.5)$$

Then we have

$$\mathbf{patt}(\widehat{\beta}_n^{\text{SLOPE}}) \xrightarrow{\text{a.s.}} \mathbf{patt}(\beta).$$

Note that above conditions are satisfied e.g. by $\lambda_i^{(n)} = c(p+1-i)n^{2/3}$ for any constant $c > 0$.

Proof. Without loss of generality we may assume that $\beta = (\beta_1, \dots, \beta_p)'$ and $\beta_1 \geq \beta_2 \geq \dots \geq \beta_p \geq 0$. Indeed, we can always achieve such condition by permuting the columns of $\mathbf{X}^{(n)}$ and changing their signs. Since the space of patterns is discrete, we have to show that for large n , $\mathbf{patt}(\widehat{\beta}_n^{\text{SLOPE}}) = \mathbf{patt}(\beta)$ a.s. We divide the proof into the following four parts:

- (a) $\beta_i = \beta_j > 0 \implies \widehat{\beta}_i^{\text{SLOPE}}(n) = \widehat{\beta}_j^{\text{SLOPE}}(n)$ a.s. for large n ,
- (b) $\beta_i > \beta_{i+1} \implies \widehat{\beta}_i^{\text{SLOPE}}(n) > \widehat{\beta}_{i+1}^{\text{SLOPE}}(n)$ a.s. for large n ,
- (c) $\beta_i = 0 \implies \widehat{\beta}_i^{\text{SLOPE}}(n) = 0$ a.s. for large n ,
- (d) $\beta_i > 0 \implies \widehat{\beta}_i^{\text{SLOPE}}(n) > 0$ a.s. for large n .

The points (b) and (d) follow quickly by the strong consistency of $\widehat{\beta}^{\text{SLOPE}}(n)$. To prove (a) and (c) we observe that for each n we are in the orthogonal design case.

Let π_n be a permutation of $(1, 2, \dots, p)$ satisfying

$$|\widehat{\beta}_{\pi_n(1)}^{\text{OLS}}(n)| \geq |\widehat{\beta}_{\pi_n(2)}^{\text{OLS}}(n)| \geq \dots \geq |\widehat{\beta}_{\pi_n(p)}^{\text{OLS}}(n)|.$$

By the strong consistency of the OLS estimator, taking n sufficiently large, we may ensure that the clusters of β do not interlace in $\widehat{\beta}_n^{\text{OLS}}$ in the sense that if $\beta_i > \beta_j$, then $\widehat{\beta}_i^{\text{OLS}}(n) > \widehat{\beta}_j^{\text{OLS}}(n)$ a.s. for n sufficiently large.

Let us now consider point (a). Let S_i denote the cluster containing $\beta_i > 0$, that is, the set $S_i = \{j \in \{1, \dots, p\} : \beta_j = \beta_i\}$. In view of the ordering of β , there exists $k_i \in \{1, \dots, p\}$ such that

$$S_i = \{\pi_n(j) : j \in \{k_i, k_i + 1, \dots, k_i + \#S_i - 1\}\}.$$

We will show that if $\pi_n(k), \pi_n(k+1) \in S_i$, then for large n

$$\widehat{\beta}_{\pi_n(k)}^{\text{SLOPE}}(n) = \widehat{\beta}_{\pi_n(k+1)}^{\text{SLOPE}}(n) \quad \text{a.s.}, \quad (3.4.6)$$

thus $\widehat{\beta}_j^{\text{SLOPE}}(n) = \widehat{\beta}_k^{\text{SLOPE}}(n)$ for $j, k \in S_i$, which finishes the proof of (a).

Now assume that $\pi_n(k), \pi_n(k+1) \in S_i$. Then, by Theorem 3.3.1, the condition (3.4.6) is satisfied if

$$\left| \widehat{\beta}_{\pi_n(k)}^{\text{OLS}}(n) \right| - \left| \widehat{\beta}_{\pi_n(k+1)}^{\text{OLS}}(n) \right| \leq \frac{1}{n} \left(\lambda_k^{(n)} - \lambda_{k+1}^{(n)} \right) \quad (3.4.7)$$

holds for large n and both $\widehat{\beta}_{\pi_n(k)}^{\text{OLS}}(n)$ and $\widehat{\beta}_{\pi_n(k+1)}^{\text{OLS}}(n)$ have the same sign. The latter is ensured by the strong consistency of the OLS estimator and the fact that $\beta_i > 0$.

If $\pi_n(k), \pi_n(k+1) \in S_i$, then we have the following bound

$$\left| \widehat{\beta}_{\pi_n(k)}^{\text{OLS}}(n) - \widehat{\beta}_{\pi_n(k+1)}^{\text{OLS}}(n) \right| \leq \sum_{j \in S_i} \left| \widehat{\beta}_j^{\text{OLS}}(n) - \widehat{\beta}_i^{\text{OLS}}(n) \right|. \quad (3.4.8)$$

Take any $j \in S_i$. Since both $\widehat{\beta}_j^{\text{OLS}}(n)$ and $\widehat{\beta}_i^{\text{OLS}}(n)$ have the normal distribution with the same mean, by Lemma 3.4.1, we have

$$\lim_{n \rightarrow \infty} \frac{\sqrt{n} \left(\widehat{\beta}_j^{\text{OLS}}(n) - \widehat{\beta}_i^{\text{OLS}}(n) \right)}{(\log(n))^{1/2+\delta}} = 0 \quad \text{a.s.}$$

In view of (3.4.8) and (3.4.5), this implies that (3.4.7) holds true for large n . Hence, (a) follows. It remains to establish (c). Assume that $\beta_{p_0} > 0 = \beta_{p_0+1} = \dots = \beta_p$. Clearly, condition (a) from Theorem 3.3.2 is satisfied thanks to the strong consistency of the OLS estimator. For (b), we have for $k = 1, 2, \dots, p_0$,

$$\frac{1}{n} \sum_{i=k}^{p_0} \lambda_i^{(n)} \leq p_0 \frac{\lambda_1^{(n)}}{n},$$

which converges to 0. On the other hand, the left-hand side of (b) converges a.s. to $\sum_{i=k}^{p_0} \beta_i$, which is positive. Thus, condition (b) from Theorem 3.3.2 holds for large n . Condition (c) from Theorem 3.3.2 follows from Lemma 3.4.1. Indeed, we have for $\delta > 0$ and $k = p_0 + 1, \dots, p$,

$$\lim_{n \rightarrow \infty} \frac{\sqrt{n}}{(\log(n))^{1/2+\delta}} \sum_{i=p_0+1}^k |\widehat{\beta}_{\pi_n(i)}^{\text{OLS}}(n)| = \sum_{i=p_0+1}^k \lim_{n \rightarrow \infty} \frac{|\sqrt{n} \widehat{\beta}_{\pi_n(i)}^{\text{OLS}}(n)|}{(\log(n))^{1/2+\delta}} = 0 \quad \text{a.s.},$$

while

$$\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}(\log(n))^{1/2+\delta}} \sum_{i=p_0+1}^k \lambda_i^{(n)} \geq \sum_{i=p_0+1}^k \lim_{n \rightarrow \infty} \frac{\lambda_i^{(n)} - \lambda_{i+1}^{(n)}}{\sqrt{n}(\log(n))^{1/2+\delta}} = m > 0$$

Thus, all assumptions of Theorem 3.3.2 are verified and the proof is complete. \square

	OLS	LASSO-CV	LASSO-LS	SLOPE-LS
$MSE(\beta, \cdot)$	613.6797	426.3705	171.7957	20.74967

Table 3.1: Comparison of MSE between different regression methods.

3.5 Numerical experiment

Below we present an application of SLOPE in signal denoising. In our example $\mathbf{X} \in \mathbb{R}^{300 \times 100}$ is an orthogonal system of trigonometric functions, i.e.

$$X_{i,(2*j-1)} = \sin(2\pi ij/n) \text{ and } X_{i,(2*j)} = \cos(2\pi ij/n) \text{ for } i = 1, \dots, 100$$

and $j = 1, \dots, 150$. Here $\beta \in \mathbb{R}^p$ is a vector consisting of two clusters: 20 coordinates with absolute value 100 and 20 coordinates with absolute value 80. The absolute values of coordinates of β are sorted in a decreasing way. The signs of the nonzero coordinates are chosen independently with random uniform distribution. To avoid large bias caused by the shrinkage nature of LASSO and SLOPE, we debias them by combining with the OLS method. For that reason we use the pattern matrix \mathbf{U}_M and the clustered design matrix $\widetilde{\mathbf{X}}_M$, which is based on the SLOPE pattern.

To perform the debiased SLOPE, we begin with recovering the support and clusters of a true vector β with SLOPE. Then, using the obtained SLOPE pattern M , we replace the design matrix with its clustered version $\widetilde{\mathbf{X}}_M = \mathbf{X}\mathbf{U}_M$. Then we perform the Ordinary Least Squares regression for the model $\mathbf{Y} = \widetilde{\mathbf{X}}_M \mathbf{b} + \varepsilon$, where \mathbf{b} consists only of distinct absolute values of $\hat{\beta}^{\text{SLOPE}}$.

Analogously we proceed with the debiased LASSO. However, in this method we use the LASSO pattern matrix defined in a following way:

For LASSO we have the LASSO pattern that is a vector of signs, see Chapter 5. For $\mathbf{S} \in \{-1, 0, 1\}^p$, $\|\mathbf{S}\|_1$ denotes the number of nonzero coordinates. If $\|\mathbf{S}\|_1 = k \geq 1$, then we define the corresponding pattern matrix $\mathbf{U}_S \in \mathbb{R}^{p \times k}$ by

$$\mathbf{U}_S = \text{diag}(\mathbf{S})_{\text{supp}(\mathbf{S})},$$

i.e. the submatrix of $\text{diag}(\mathbf{S})$ obtained by keeping columns corresponding to indices in $\text{supp}(\mathbf{S})$. Then we define the reduced matrix $\widetilde{\mathbf{X}}_S$ by

$$\widetilde{\mathbf{X}}_S = \mathbf{X}\mathbf{U}_S.$$

Equivalently, we have $\widetilde{\mathbf{X}}_S = (S_i X_i)_{i \in \text{supp}(S)}$. For a broader discussion on the pattern matrix, we encourage to see the next Chapter. In our example $\varepsilon \in \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ and $\sigma = 30$.

We compare the Mean Square Error and the signal denoising of the classical OLS estimation, the LASSO with the tuning parameter λ_{cv} minimizing the cross-validated error, the debiased version of LASSO with $\lambda = 5\lambda_{cv}$ and the debiased version of SLOPE with the tuning vector Λ chosen with respect to the sequence proposed below Theorem 3.4.3 ($\lambda_i = 0.1(p+1-i)n^{2/3}$).

We also compare debiased SLOPE with debiased LASSO based on a single trial, as shown in Figure 3.4 and Table 3.1. The horizontal lines correspond to the true values of β . As one may observe, in the presented setting LASSO does not recover the true support, while debiased SLOPE perfectly recovers support, sign and clusters.

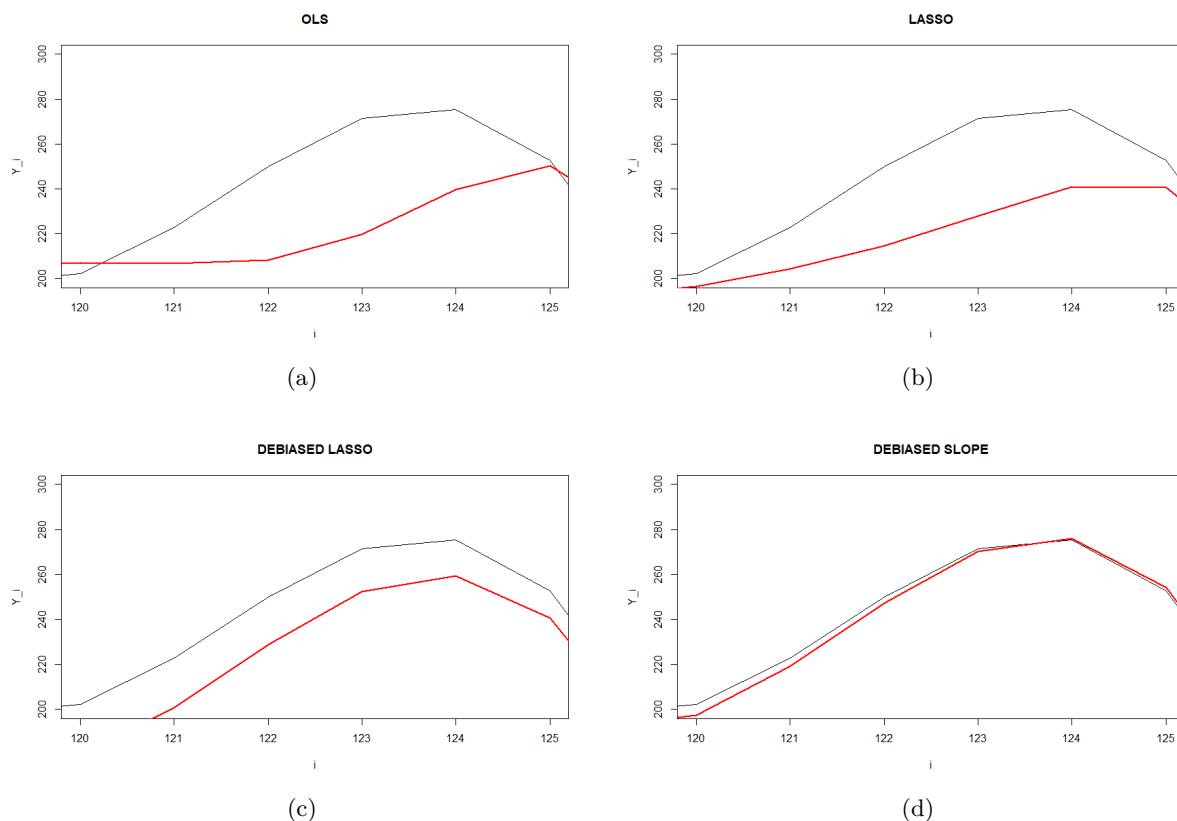


Figure 3.2: Comparison of signal denoising by OLS (a), LASSO (b), debiased LASSO (c) and debiased SLOPE (d) on the coordinates $[120, 125]$ of the regression model $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$. The black lines correspond to the true values of $\mathbf{X}\beta$. The red lines correspond to the estimators $\mathbf{Y} = \mathbf{X}\hat{\beta}$.

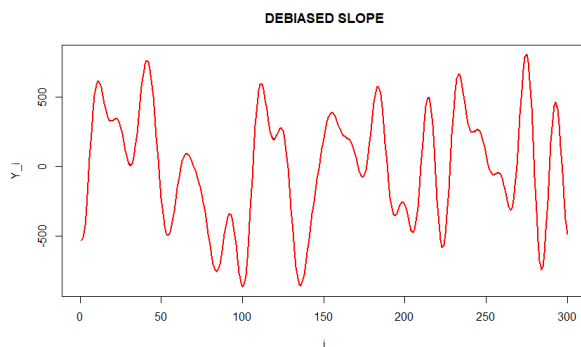


Figure 3.3: Signal denoising by debiased SLOPE on all coordinates of the regression model $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$. The (almost overlapping) black line and the red line correspond respectively to the true values of $\mathbf{X}\beta$ and to $\mathbf{Y} = \mathbf{X}\hat{\beta}^{\text{SLOPE}}$.

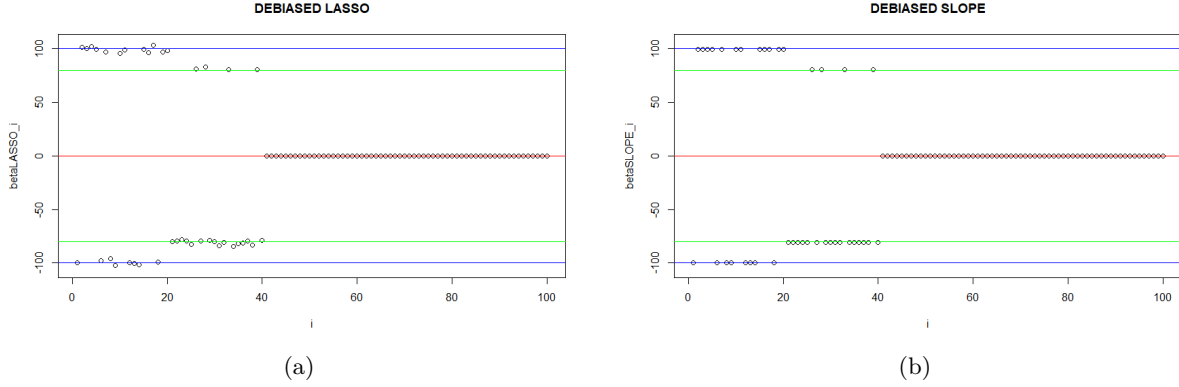


Figure 3.4: Pattern recovery by debiased LASSO (left) and by debiased SLOPE (right) in the same setting as above.

3.6 Appendix

Proof of Lemma 3.2.1. Since the matrix $(\mathbf{X}'\mathbf{X})^{-1}$ is nonnegative definite, it follows that the function $g : C_{\Lambda} \rightarrow [0, \infty)$ defined by

$$g(\boldsymbol{\pi}) := (\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi})'(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi})$$

is convex in $\boldsymbol{\pi}$. Therefore, at the point $\boldsymbol{\pi}^* = (\pi_1^*, \dots, \pi_p^*)'$, where g attains its global minimum over C_{Λ} , the gradient ∇g of g satisfies

$$[\nabla g(\boldsymbol{\pi}^*)]'(\boldsymbol{\pi} - \boldsymbol{\pi}^*) \geq 0, \quad \text{for all } \boldsymbol{\pi} \in C_{\Lambda}.$$

This implies $(\boldsymbol{\pi} - \boldsymbol{\pi}^*)'\boldsymbol{\beta}^* \leq 0$, for all $\boldsymbol{\pi} \in C_{\Lambda}$, because

$$\nabla g(\boldsymbol{\pi}^*) = -2(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y} - \boldsymbol{\pi}^*) = -2\boldsymbol{\beta}^*.$$

In the proof of parts (a), (b) and (c), we use the fact that $\boldsymbol{\pi}^*$ maximizes $\boldsymbol{\pi}'\boldsymbol{\beta}^*$ over $\boldsymbol{\pi} \in C_{\Lambda}$. To prove part (a), suppose that $\text{sign}(\beta_i^*) \cdot \text{sign}(\pi_i^*) < 0$ for some i and define

$$\boldsymbol{\pi} = (\pi_1^*, \dots, \pi_{i-1}^*, -\pi_i^*, \pi_{i+1}^*, \dots, \pi_p^*)'.$$

Then we have $(\boldsymbol{\pi}^*)'\boldsymbol{\beta}^* < \boldsymbol{\pi}'\boldsymbol{\beta}^*$, which is impossible since $\boldsymbol{\pi} \in C_{\Lambda}$.

To prove part (b), consider a permutation τ of $(1, 2, \dots, p)$ such that

$(|\pi_{\tau(1)}^*|, \dots, |\pi_{\tau(p)}^*|)$ and $(|\beta_1^*|, \dots, |\beta_p^*|)$ are similarly sorted. Define the point

$\boldsymbol{\pi} = (s_1 \cdot \pi_{\tau(1)}^*, s_2 \cdot \pi_{\tau(2)}^*, \dots, s_p \cdot \pi_{\tau(p)}^*)$, where $s_i = \text{sign}(\beta_i^*)$, for $i = 1, 2, \dots, p$. If $(|\pi_{\tau(1)}^*|, \dots, |\pi_{\tau(p)}^*|) \neq (|\pi_1^*|, \dots, |\pi_p^*|)$, then, by the Hardy-Littlewood-Pólya rearrangement inequality,

$$\boldsymbol{\pi}'\boldsymbol{\beta}^* = \sum_{i=1}^p |\pi_{\tau(i)}^*| |\beta_i^*| > \sum_{i=1}^p |\pi_i^*| |\beta_i^*| \geq (\boldsymbol{\pi}^*)'\boldsymbol{\beta}^*,$$

which is impossible since $\boldsymbol{\pi} \in C_{\Lambda}$.

Finally, to prove part (c), suppose that $\sum_{i=1}^{k-1} |\pi_{\tau(i)}^*| < \sum_{i=1}^{k-1} \lambda_i$, and that $|\pi_{\tau(k)}^*| > 0$. In this case there is a sufficiently small $\delta > 0$, such that

$$\boldsymbol{\pi} = (\pi_1^*, \dots, \pi_{i-2}^*, \pi_{i-1}^* + \delta s_{i-1}, \pi_i^* - \delta s_i, \pi_{i+1}^*, \dots, \pi_p^*)' \in C_{\Lambda}.$$

If $|\beta_{\tau(k-1)}^*| > |\beta_{\tau(k)}^*|$ then

$$\boldsymbol{\pi}'\boldsymbol{\beta}^* = (\boldsymbol{\pi}^*)'\boldsymbol{\beta}^* + \delta(|\beta_{\tau(k-1)}^*| - |\beta_{\tau(k)}^*|) > (\boldsymbol{\pi}^*)'\boldsymbol{\beta}^*,$$

which is impossible. \square

Proof of Lemma 3.2.2. At first we note that for all $\boldsymbol{\pi} \in C_{\Lambda}$

$$\begin{aligned} r(\boldsymbol{\beta}^*, \boldsymbol{\pi}) &= \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*\|_2^2 + \boldsymbol{\pi}'\boldsymbol{\beta}^* = \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^*\|_2^2 + (\boldsymbol{\pi}^*)'\boldsymbol{\beta}^* \\ &\quad + (\boldsymbol{\pi} - \boldsymbol{\pi}^*)'\boldsymbol{\beta}^* = r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*) + (\boldsymbol{\pi} - \boldsymbol{\pi}^*)'\boldsymbol{\beta}^* \leq r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*), \end{aligned}$$

where the last inequality follows from the fact that $(\boldsymbol{\pi} - \boldsymbol{\pi}^*)'\boldsymbol{\beta}^* \leq 0$, for all $\boldsymbol{\pi} \in C_{\Lambda}$, see the proof of 3.2.1. Therefore, $\max_{\boldsymbol{\pi} \in C_{\Lambda}} r(\boldsymbol{\beta}^*, \boldsymbol{\pi}) = r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*)$. Moreover, from the definition of the point $\boldsymbol{\beta}^*$ it is seen that $r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*) = \min_{\boldsymbol{\beta} \in \mathcal{M}} r(\boldsymbol{\beta}, \boldsymbol{\pi}^*)$. These two facts imply that

$$\begin{aligned} \min_{\boldsymbol{\beta} \in \mathcal{M}} \max_{\boldsymbol{\pi} \in C_{\Lambda}} r(\boldsymbol{\beta}, \boldsymbol{\pi}) &\leq \max_{\boldsymbol{\pi} \in C_{\Lambda}} r(\boldsymbol{\beta}^*, \boldsymbol{\pi}) = r(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*) \\ &= \min_{\boldsymbol{\beta} \in \mathcal{M}} r(\boldsymbol{\beta}, \boldsymbol{\pi}^*) \leq \max_{\boldsymbol{\pi} \in C_{\Lambda}} \min_{\boldsymbol{\beta} \in \mathcal{M}} r(\boldsymbol{\beta}, \boldsymbol{\pi}). \end{aligned}$$

Since $\max_{\boldsymbol{\pi} \in C_{\Lambda}} \min_{\boldsymbol{\beta} \in \mathcal{M}} r(\boldsymbol{\beta}, \boldsymbol{\pi}) \leq \min_{\boldsymbol{\beta} \in \mathcal{M}} \max_{\boldsymbol{\pi} \in C_{\Lambda}} r(\boldsymbol{\beta}, \boldsymbol{\pi})$ (by the max-min inequality), we have the equality throughout. This completes the proof. \square

Proof of Lemma 3.3.1. Observe that

$$\begin{aligned} \frac{1}{n}\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 &= \frac{1}{n}\mathbf{Y}'\mathbf{Y} - \frac{2}{n}\mathbf{Y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{b} \\ \|\widehat{\boldsymbol{\beta}}^{\text{OLS}} - \mathbf{b}\|_2^2 &= \frac{1}{n^2}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y} - \frac{2}{n}\mathbf{Y}'\mathbf{X}\mathbf{b} + \mathbf{b}'\mathbf{b}. \end{aligned}$$

Therefore, both optimization problems differ by $\frac{1}{2n}(\mathbf{Y}'\mathbf{Y} - \frac{1}{n}\mathbf{Y}'\mathbf{X}\mathbf{X}'\mathbf{Y})$, which does not depend on \mathbf{b} , which implies their equivalence. \square

Chapter 4

Pattern recovery by SLOPE

4.1 Introduction

Most of the content of this chapter may be found in the preprint of Bogdan, Dupuis, Graczyk, Kołodziejek, the author of the dissertation, Tardivel and Wilczyński [23]. We decided to enrich the chapter with geometrical explanation of the SLOPE pattern and its recovery to make it more accessible for the general audience and to inform the reader about the connections between the SLOPE estimator and the convex geometry. Such connections exist also for LASSO and may be useful in its study. The geometrical approach was the one, which made us build the intuition, which led us to establish the theoretical results.

In this chapter we give a deeper and more general focus on recovering the pattern of β by SLOPE. Now we abstain from restricting to the case of the orthogonal design matrix and we give novel results on the pattern recovery for a broad generality of situations.

In particular, unlike in the previous chapter, below we concern a model, in which the error is incremental. We also extend the discussion on the properties of the SLOPE estimator to the high-dimensional case, in which p may be much larger than n .

4.1.1 History of SLOPE

SLOPE estimator was introduced by Bogdan et al. [27, 26]. In their seminal paper, SLOPE was motivated by the control of the expected rate of false discoveries (FDR control) in multiple hypotheses testing. Moreover, they proposed a choice of Λ being a Benjamini-Hochberg sequence dependent on a parameter $q \in (0, 1)$, under which, in the orthogonal design, SLOPE controls FDR at a level $q \cdot \frac{p-k}{p}$, where k is the number of non-zero coordinates of β . They also proposed a fast proximal algorithm computing $\hat{\beta}^{SLOPE}$ in the orthogonal design.

However, the first non-trivial example of SLOPE was introduced a few years earlier by Bondell and Reich [29] under the name OSCAR. Their article deals with the tuning vector Λ being an arithmetic sequence. The OSCAR penalty was reformulated in terms of a sorted ℓ_1 norm by Zeng and Figueiredo [190]. Going with the flow, they independently propose the SLOPE estimator [189] under the name OWL. In [76] they propose its representation as a gauge function, compute its dual and propose another formula for the proximal operator algorithm.

Figueiredo and Nowak [76] prove that if the columns of \mathbf{X} are correlated enough, then SLOPE results in their clusterization. Negrinho and Martins [136] connect the SLOPE estimator with the notion of a signed permutahedron C_Λ , noticing that it is a unit ball in the dual of SLOPE norm. Bellec and Tsybakov [13] prove the equality $\mathbf{X}\hat{\beta} = \mathbf{Y} - Proj_{C_{\Lambda/n}}(\mathbf{Y})$ in the linear regression penalized by norm. They also propose novel oracle inequalities on a prediction error,

given $\varepsilon \sim N(0, \sigma^2 \mathbf{I}_n)$ and $\lambda_j = \sigma \sqrt{\log(\frac{2p}{j})/n}$.

Su and Candès [171] prove that with $\mathbf{\Lambda}$ being the Benjamini-Hochberg sequence, SLOPE is an asymptotically minimax estimator and show more detailed results on the asymptotics of its squared error for orthogonal ($\mathbf{X}'\mathbf{X} = \mathbf{I}_p$) and gaussian ($X_{ij} \stackrel{iid}{\sim} N(0, \frac{1}{n})$) designs. Moreover, they point out such sizes of the support of β (ℓ_0 -sparsities), for which the risk for SLOPE is much smaller than for LASSO or SURE. Bellec, Lecué and Tsybakov [14] extend the results of Su and Candès [171] on achieving the minimax optimal prediction rate of SLOPE. In particular, they allow the design matrix \mathbf{X} to be deterministic, leaving only mild assumptions on its eigenvalues. They also derive sharp oracle inequalities and prove that the minimaxity holds for an ℓ_q -estimation error for any $1 \leq q \leq 2$.

Under the gaussian design, Kos and Bogdan [115] use the results of Su and Candès [171] on the asymptotics of SLOPE in order to propose conditions on the cardinality and magnitude of non-zero coordinates of β , under which the SLOPE with $\mathbf{\Lambda}$ being a Benjamini-Hochberg sequence asymptotically achieves $FDR \rightarrow 0$ and the power converging to 1. Hu and Lu [103] put their effort to look for an optimal choice of $\mathbf{\Lambda}$ in terms of type I error and power of variable selection. They consider an asymptotic scenario, in which $\frac{n}{p} \rightarrow \delta > 0$, $\frac{k}{p} \rightarrow \rho > 0$ and the vectors β , ε and $\mathbf{\Lambda}$ converge to limiting measures.

Sankaran, Bach and Bhattacharaya [154] prove that SLOPE belongs to the family of Lovász extensions. They also propose the smoothed version of SLOPE (called SOWL) and derive sufficient condition for its consistency. Minami [132] studies the projections onto the dual norm ball to study the notion of degrees of freedom. SLOPE is used here as one of examples of submodular norm regularizations.

Schneider and Tardivel [156] propose a geometrical approach to the SLOPE estimator. Firstly, they give full characterization of the uniqueness of $\hat{\beta}^{SLOPE}$ that for a given \mathbf{X} holds for any $\mathbf{Y} \in \mathbb{R}^n$. They express this condition both in analytical and geometrical way. This is a generalization of an analogous result for LASSO [71]. They also introduce the notion of SLOPE pattern and highlight its connection with faces of the unit ball in the dual $J_{\mathbf{\Lambda}}^*$ of the SLOPE norm. In their article, the definition of accessibility of the SLOPE pattern is also introduced and enriched with its geometrical full characterization. Tardivel, Servien and Concordet [175] propose an exact formula for SLOPE in the orthogonal case, using the Cesàro summation.

For gaussian designs, random β , random ε and for $\frac{n}{p} \xrightarrow{p \rightarrow \infty} \delta > 0$, Zhang and Bu [191] use the projected gradient descent algorithm to find $\mathbf{\Lambda}$, which reduces the Mean Squared Error (MSE) and compare obtained MSE with ones achieved by SLOPE with Benjamini-Hochberg $\mathbf{\Lambda}$ and by LASSO. To minimize MSE for arbitrary data, they propose finding $\mathbf{\Lambda}$ with help of a coordinate descent algorithm, while imposing the restriction of $\mathbf{\Lambda}$ to contain no more than k distinct values with $k \ll p$. In their recent article, Larsson et al. [120] improved solving numerically the SLOPE optimization problem, using modifications of the proximal gradient descent and the proximal coordinate descent algorithm.

Bu et al. [40] proposed an approximate message passing (AMP) algorithm, to quickly obtain a solution, which they prove to converge in ℓ_2 to $\hat{\beta}^{SLOPE}$. They also refine the asymptotic results of the ℓ_2 -convergence of $\hat{\beta}^{SLOPE}$ to the true β . Larsson, Bogdan and Wallin [119] introduced the strong screening rule for SLOPE, which generalized the analogous rule for LASSO. The safe screening rules were proposed by Bao, Gu and Huang [8] as well as by Elvira and Herzet [65]. Their safety is understood as no detection any non-zero coordinate of β as zero. It is useful in situation, when some prior knowledge of β is available.

Brzyski et al. [38] introduce the group SLOPE method. Lee, Sobczyk and Bogdan [123] apply SLOPE to gaussian graphical models, proposing the Neighborhood Selection SLOPE (nsSLOPE)

algorithm with the FDR control of edges detection. Mazza-Anthony, Mazouze and Coates [129] independently introduce the graphical SLOPE method (under the name gOWL), which applies the SLOPE norm to estimate the precision matrix by clustering its off-diagonal entries. Recently, Riccobello et al. [148] extended the discussion on graphical SLOPE to t-Student (Tslope) data.

Other applications of SLOPE include e.g. the paper of Sepehri [161], who proposes the bayesian approach in which SLOPE is considered as a maximum a posteriori (MAP) distribution in a bayesian regression problem. Its adaptive version was later proposed by Jiang et al. [109]. Stucky and van de Geer [170] introduce the Square Root SLOPE, in which $\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2$ is replaced by $\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2$. It is an analogue of the Square Root LASSO [15].

Recently, Dexheimer and Strauch [56] found a way to apply the results of [14] and used SLOPE to improve the estimation of the drift parameter in Lévy-driven Ornstein-Uhlenbeck processes.

4.1.2 Our contribution

In this chapter, we give the necessary and sufficient conditions for recovering the SLOPE pattern. From a mathematical perspective, our main result is Theorem 4.3.1, in which we propose the positivity and the subdifferential condition for the pattern recovery. To ease the interpretation, we also present the dual norm condition and the SLOPE irrepresentability condition, which is a necessary and sufficient condition for pattern recovery in the noiseless case. The word “irrepresentability” is a tribute to works written a decade ago on sign recovery by LASSO [82, 130, 183, 192, 194]. However, we believe that our mathematical perspective is novel, and paves the path for similar analyses of other penalized estimators. For the noisy case, in Theorem 4.5.1 we give the open SLOPE irrepresentability condition, which implies the consistency of the SLOPE pattern for fixed \mathbf{X} and the difference between non-equal coordinates of $|\beta|$ diverging to ∞ . Later on, in Theorems 4.5.3 and 4.5.5, we give conditions for asymptotic pattern recovery, when p is fixed and n diverges to ∞ . As another main results of this chapter, we provide geometrical interpretations of irrepresentability conditions as well as refined results of the strong consistency of $\hat{\beta}^{SLOPE}$ and of its pattern.

4.1.3 Motivation

While the SLOPE ability to identify the pattern of the vector of regression coefficients β is interesting by itself, the related reduction of model dimension also brings the advantage in terms of the precision of the estimation of β . This phenomenon is illustrated in Figure 4.1, which presents the difference in precision of the LASSO and SLOPE estimators when some of the regression coefficients are equal to each other.

In this example $n = 100$, $p = 200$, and the rows of the design matrix are generated as independent binary Markov chains, with $\mathbb{P}(X_{i1} = 1) = \mathbb{P}(X_{i1} = -1) = 0.5$ and $\mathbb{P}(X_{i(j+1)} \neq X_{ij}) = 1 - \mathbb{P}(X_{i(j+1)} = X_{ij}) \approx 0.0476$. This value corresponds to the probability of the crossover event between genetic markers spaced every 5 centimorgans. To be more specific, it is close to the inverse of the Haldane mapping function [81, pp. 13-14] evaluated at 0.05, which is equal to $e^{-0.05} \sinh(0.05)$. Our design matrix can be viewed as an example of 100 independent haplotypes, each resulting from a single meiosis event. In this example, the correlation between columns of the design matrix decays exponentially, $\rho(\mathbf{X}_i, \mathbf{X}_j) \approx 0.9048^{|i-j|}$. The design matrix is then standardized so that each column has a zero mean and a unit variance, and the response variable is generated according to the linear model $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$, with $\beta_1 = \dots = \beta_{30} = 40$, $\beta_{31} = \dots = \beta_{200} = 0$ and $\sigma = 5$. In this experiment, the data matrix \mathbf{X} and the regression model

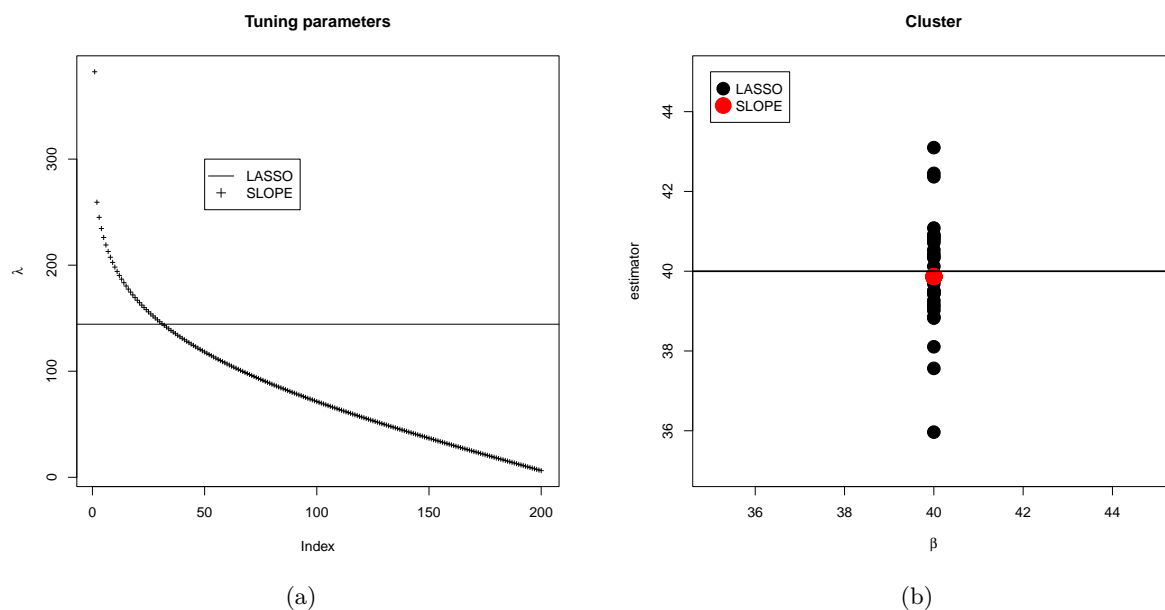


Figure 4.1: Comparison of LASSO and SLOPE when the cluster structure is present in the data. Here $n = 100$, $p = 200$, the rows of the matrix \mathbf{X} are simulated as independent binary Markov chains, with the transition probability 0.0476 (corresponding to the genetic distance of 5 centimorgans). The correlation between i^{th} and the j^{th} column of \mathbf{X} decays exponentially as $0.9048^{|i-j|}$. The first $k = 30$ columns of \mathbf{X} are associated with \mathbf{Y} and their nonzero regression coefficients are all equal to 40 (other details are provided in the text). The left panel represents the value of the tuning parameter for LASSO (solid line) and the sequence of tuning parameters for SLOPE (crosses). The sequences are selected such that both LASSO and SLOPE recover their corresponding patterns with a minimal bias. The right panel represents the LASSO and SLOPE estimates.

are constructed such that the LASSO irrepresentability condition holds. The tuning parameter for LASSO is selected as the smallest value of λ for which LASSO can correctly identify the sign of β . Similarly, the tuning parameter Λ is designed such that the SLOPE irrepresentability condition holds and Λ is multiplied by the smallest constant for which SLOPE properly returns the SLOPE pattern. The selected tuning parameters for LASSO and SLOPE are represented in the left panel of Figure 4.1. Both in the case of LASSO and SLOPE, the proposed tuning parameters are close to the values minimizing the mean squared estimation error. Since in this example both LASSO and SLOPE properly estimate null components of β at $\mathbf{0}$, the right panel in Figure 4.1 illustrates only the accuracy of the estimation of the nonzero coefficients. Here, we can observe that the SLOPE ability to identify the cluster structure leads to superior estimation properties. SLOPE estimates the vector of regression coefficients β virtually without an error, while LASSO estimates are scattered over the interval between 36 and 44. In the result, the squared error of the LASSO estimator is more than 100 times larger than the squared error of SLOPE (63.4 vs 0.53).

4.2 Preliminaries on clustering properties by SLOPE

As the central notion of this chapter is the recovery of the SLOPE pattern, we recall its definition:

Definition 4.2.1. *The SLOPE pattern is a function $\mathbf{patt} : \mathbb{R}^p \rightarrow \mathbb{Z}^p$ such that*

$$\mathbf{patt}(\mathbf{b})_i = \text{sign}(b_i) \text{rank}(|b_i|),$$

where $\text{rank}(|b_i|) \in \{1, 2, \dots, k\}$ is defined to be the number of $|c_j|$'s satisfying $|b_i| \geq |c_j|$, where $|c_1|, |c_2|, \dots, |c_k|$, $k \leq p$, are distinct non-zero values among $|b_1|, \dots, |b_p|$. We adopt the convention that $\text{rank}(0) = 0$.

We also recall that $\mathbb{R}^{k+} = \{\boldsymbol{\kappa} \in \mathbb{R}^k : \kappa_1 > \dots > \kappa_k > 0\}$ and that \mathcal{M}_p is the set of SLOPE patterns of $\mathbf{b} \in \mathbb{R}^p$. Definition 2.2.4 implies that for $\mathbf{0} \neq \mathbf{M} \in \mathcal{M}_p^{\text{SLOPE}}$ and $k = \|\mathbf{M}\|_\infty$, for $\mathbf{b} \in \mathbb{R}^p$ we have

$$\mathbf{patt}(\mathbf{b}) = \mathbf{M} \iff \text{there exists } \boldsymbol{\kappa} \in \mathbb{R}^{k+} \text{ such that } \mathbf{b} = \mathbf{U}_\mathbf{M} \boldsymbol{\kappa}.$$

Example 4.2.1. Let $\mathbf{b} = (2, 5, -2, 0, 8)'$. Then $\mathbf{M} = \mathbf{patt}(\mathbf{b}) = (1, 2, -1, 0, 3)'$ and for $\boldsymbol{\kappa} = (8, 5, 2)'$ we have

$$\mathbf{U}_\mathbf{M} \boldsymbol{\kappa} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 8 \\ 5 \\ 2 \end{bmatrix} = \mathbf{b}.$$

4.2.1 Clustered design matrix and clustered parameter

If $\mathbf{M} = \mathbf{patt}(\beta)$ for $\beta \in \mathbb{R}^p$ satisfies $\|\mathbf{M}\|_\infty < p$, then the pattern $\mathbf{M} = (M_1, \dots, M_p)'$ leads naturally to reduce the dimension of the design matrix \mathbf{X} in the regression problem, by replacing \mathbf{X} by $\widetilde{\mathbf{X}}_\mathbf{M}$. Actually, if $\mathbf{patt}(\beta) = \mathbf{M}$, then $\mathbf{X}\beta = \mathbf{X}\mathbf{U}_\mathbf{M}\boldsymbol{\kappa} = \widetilde{\mathbf{X}}_\mathbf{M}\boldsymbol{\kappa}$ for $\boldsymbol{\kappa} \in \mathbb{R}^{k+}$. In particular,

- (a) null components $m_i = 0$ lead to the elimination of columns \mathbf{X}_i from the design matrix \mathbf{X} ,

- (b) a cluster $K \subset \{1, \dots, p\}$ of \mathbf{M} (a subset of \mathbf{M} having coordinates equal in absolute value) leads to replace the columns $(\mathbf{X}_i)_{i \in K}$ by one column equal to their signed sum: $\sum_{i \in K} \text{sign}(m_i) \mathbf{X}_i$.

Example 4.2.2. Let \mathbf{b} be from Example 4.2.1, $\mathbf{X} = (\mathbf{X}_1 | \mathbf{X}_2 | \mathbf{X}_3 | \mathbf{X}_4 | \mathbf{X}_5)$ and $\mathbf{\Lambda} = (\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5)' \in \mathbb{R}^{5+}$. Then the clustered matrix and the clustered parameter are:

$$\widetilde{\mathbf{X}}_{\mathbf{M}} = (\mathbf{X}_5 | \mathbf{X}_2 | \mathbf{X}_1 - \mathbf{X}_3) \text{ and } \widetilde{\mathbf{\Lambda}}_{\mathbf{M}} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 + \lambda_4 \end{pmatrix}.$$

4.2.2 Sorted ℓ_1 norm, dual sorted ℓ_1 norm and subdifferential

As SLOPE is based on penalizing the residual sum of squared with the sorted ℓ_1 norm, we recall its definition:

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

where $|b|_{(1)} \geq \dots \geq |b|_{(p)}$ are the sorted components of \mathbf{b} with respect to the absolute value. Given a norm $\|\cdot\|$ on \mathbb{R}^p , we recall that the dual norm $\|\cdot\|^*$ is defined by $\|\mathbf{b}\|^* = \max\{\mathbf{v}'\mathbf{b} : \|\mathbf{v}\| \leq 1\}$, for some $\mathbf{b} \in \mathbb{R}^p$. In particular, the dual sorted ℓ_1 norm has an explicit expression given in [136]:

$$J_{\mathbf{\Lambda}}^*(\mathbf{b}) = \max \left\{ \frac{|b|_{(1)}}{\lambda_1}, \frac{\sum_{i=1}^2 |b|_{(i)}}{\sum_{i=1}^2 \lambda_i}, \dots, \frac{\sum_{i=1}^p |b|_{(i)}}{\sum_{i=1}^p \lambda_i} \right\}.$$

We recall the subdifferential of a norm $\|\cdot\|$ at \mathbf{b} (see *e.g.* [101, Def. VI.1.2.1]):

$$\begin{aligned} \partial \|\cdot\|(\mathbf{b}) &= \{\mathbf{v} \in \mathbb{R}^p : \|\mathbf{z}\| \geq \|\mathbf{b}\| + \mathbf{v}'(\mathbf{z} - \mathbf{b}) \quad \forall \mathbf{z} \in \mathbb{R}^p\}, \\ &= \{\mathbf{v} \in \mathbb{R}^p : \|\mathbf{v}\|^* \leq 1 \text{ and } \mathbf{v}'\mathbf{b} = \|\mathbf{b}\|\}. \end{aligned} \quad (4.2.1)$$

For the sorted ℓ_1 norm, geometrical descriptions of the subdifferential at $\mathbf{b} \in \mathbb{R}^p$ have been given in the particular case where $b_1 \geq \dots \geq b_p \geq 0$ [61, 156, 175]. Below, for an arbitrary $\mathbf{b} \in \mathbb{R}^p$, we propose a new and useful formula for the subdifferential of the sorted ℓ_1 norm. This representation is the crux of the mathematical content of the present chapter.

Proposition 4.2.1. *Let $\mathbf{b} \in \mathbb{R}^p$ and $\mathbf{M} = \mathbf{patt}(\mathbf{b})$. Then:*

$$\partial J_{\mathbf{\Lambda}}(\mathbf{b}) = \left\{ \mathbf{v} \in \mathbb{R}^p : J_{\mathbf{\Lambda}}^*(\mathbf{v}) \leq 1 \text{ and } \mathbf{U}'_{\mathbf{M}} \mathbf{v} = \widetilde{\mathbf{\Lambda}}_{\mathbf{M}} \right\}. \quad (4.2.2)$$

In Proposition 4.9.1 we derive a simple characterization of elements of $\partial J_{\mathbf{\Lambda}}(\mathbf{b})$. The notion of SLOPE pattern is related to the subdifferential via the following result.

Proposition 4.2.2. *Let $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ and $\mathbf{\Lambda} = (\lambda_1, \dots, \lambda_p)'$ where $\lambda_1 > \dots > \lambda_p > 0$. We have $\mathbf{patt}(\mathbf{a}) = \mathbf{patt}(\mathbf{b})$ if and only if $\partial J_{\mathbf{\Lambda}}(\mathbf{a}) = \partial J_{\mathbf{\Lambda}}(\mathbf{b})$.*

A proof of Proposition 4.2.2 can be found in [156]. In the Appendix, we provide an independent proof, which is based on Proposition 4.2.1.

From now, to comply with Proposition 4.2.2, we assume that the tuning parameter $\mathbf{\Lambda} = (\lambda_1, \dots, \lambda_p)'$ satisfies

$$\lambda_1 > \dots > \lambda_p > 0.$$

4.2.3 Characterization of SLOPE solutions

SLOPE estimator is a solution of the following optimization problem:

$$S_{\mathbf{X},\Lambda}(\mathbf{Y}) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \left\{ \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + J_\Lambda(\mathbf{b}) \right\}. \quad (4.2.3)$$

We do not assume that $S_{\mathbf{X},\Lambda}(\mathbf{Y})$ is a singleton. However, note that the cases in which the SLOPE estimator is not unique are very rare. Indeed, the family of matrices $\mathbf{X} \in \mathbb{R}^{n \times p}$, for which there exists a $\mathbf{Y} \in \mathbb{R}^n$ such that $S_{\mathbf{X},\Lambda}(\mathbf{Y})$ is not a singleton, has a null Lebesgue measure on $\mathbb{R}^{n \times p}$ [156]. If $\ker(\mathbf{X}) = \{\mathbf{0}\}$, then $S_{\mathbf{X},\Lambda}(\mathbf{Y})$ consists of one element. Recall that a convex function f attains its minimum at a point \mathbf{b} if and only if $\mathbf{0} \in \partial f(\mathbf{b})$. Since $\partial \left(\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 \right) = \{-\mathbf{X}'(\mathbf{Y} - \mathbf{X}\mathbf{b})\}$, the SLOPE estimator admits the following characterization:

$$\hat{\beta} \in S_{\mathbf{X},\Lambda}(\mathbf{Y}) \quad \Leftrightarrow \quad \mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta}) \in \partial J_\Lambda(\hat{\beta}). \quad (4.2.4)$$

4.3 Characterization of pattern recovery by SLOPE

The characterization of pattern recovery by SLOPE given in Theorem 4.3.1 is one of most important results of this chapter. We recall that $\tilde{\mathbf{P}}_{\mathbf{M}} = (\tilde{\mathbf{X}}_{\mathbf{M}}')^+ \tilde{\mathbf{X}}_{\mathbf{M}}' = \tilde{\mathbf{X}}_{\mathbf{M}} \tilde{\mathbf{X}}_{\mathbf{M}}^+$ is the orthogonal projection onto $\text{col}(\tilde{\mathbf{X}}_{\mathbf{M}})$, where \mathbf{A}^+ denotes the Moore-Penrose inverse of the matrix \mathbf{A} .

Theorem 4.3.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{0} \neq \beta \in \mathbb{R}^p$, $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$, $\varepsilon \in \mathbb{R}^n$ and $\Lambda \in \mathbb{R}^{p+}$. Let $\mathbf{M} = \mathbf{patt}(\beta) \in \mathcal{M}_p^{\text{SLOPE}}$ and $k = \|\mathbf{M}\|_\infty$. Denote*

$$\boldsymbol{\pi} = \mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{M}}')^+ \tilde{\Lambda}_{\mathbf{M}} + \mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\mathbf{Y}. \quad (4.3.1)$$

There exists $\hat{\beta} \in S_{\mathbf{X},\Lambda}(\mathbf{Y})$ with $\mathbf{patt}(\hat{\beta}) = \mathbf{patt}(\beta)$ if and only if the two conditions below hold:

$$\begin{cases} \text{there exists } \mathbf{s} \in \mathbb{R}^{k+} \text{ such that } \tilde{\mathbf{X}}_{\mathbf{M}}'\mathbf{Y} - \tilde{\Lambda}_{\mathbf{M}} = \tilde{\mathbf{X}}_{\mathbf{M}}'\tilde{\mathbf{X}}_{\mathbf{M}}\mathbf{s}, & (\text{positivity condition}) \\ \boldsymbol{\pi} \in \partial J_\Lambda(\mathbf{M}). & (\text{subdifferential condition}) \end{cases}$$

If the positivity and subdifferential conditions are satisfied, then $\hat{\beta} = \mathbf{U}_{\mathbf{M}}\mathbf{s} \in S_{\mathbf{X},\Lambda}(\mathbf{Y})$ and $\boldsymbol{\pi} = \mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta})$.

Remark 4.3.2.

- (a) When \mathbf{X} is deterministic and ε has a $N(0, \sigma^2 \mathbf{I}_n)$ distribution, then the pattern recovery by SLOPE is the intersection of statistically independent events:

$$\begin{aligned} A &= \left\{ \omega \in \Omega: \text{there exists } \mathbf{s} \in \mathbb{R}^{k+} \text{ such that } \tilde{\mathbf{X}}_{\mathbf{M}}'\mathbf{Y}(\omega) - \tilde{\Lambda}_{\mathbf{M}} = \tilde{\mathbf{X}}_{\mathbf{M}}'\tilde{\mathbf{X}}_{\mathbf{M}}\mathbf{s} \right\}, \\ B &= \left\{ \omega \in \Omega: \boldsymbol{\pi}(\omega) \in \partial J_\Lambda(\mathbf{M}) \right\}. \end{aligned}$$

Indeed, since $\tilde{\mathbf{X}}_{\mathbf{M}}' = \tilde{\mathbf{X}}_{\mathbf{M}}'\tilde{\mathbf{P}}_{\mathbf{M}}$ then $\tilde{\mathbf{X}}_{\mathbf{M}}'\mathbf{Y}(\omega)$ depends on $\varepsilon_A(\omega) = \tilde{\mathbf{P}}_{\mathbf{M}}\varepsilon(\omega)$. Moreover, $\boldsymbol{\pi}(\omega)$ depends on $\varepsilon_B(\omega) = (\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\varepsilon(\omega)$. Since $\tilde{\mathbf{P}}_{\mathbf{M}}$ is an orthogonal projection, then both ε_A and ε_B have null covariance matrices. But ε is Gaussian and hence ε_A and ε_B are independent. Therefore the events A and B are independent.

- (b) Under the positivity condition, the subdifferential condition is equivalent to

$$J_\Lambda^*(\boldsymbol{\pi}) \leq 1. \quad (\text{dual norm condition})$$

Indeed, observe that $\tilde{\Lambda}_{\mathbf{M}} \in \text{col}(\tilde{\mathbf{X}}_{\mathbf{M}}')$ (or equivalently, $\tilde{\mathbf{X}}_{\mathbf{M}}'(\tilde{\mathbf{X}}_{\mathbf{M}}')^+\tilde{\Lambda}_{\mathbf{M}} = \tilde{\Lambda}_{\mathbf{M}}$) is necessary for the positivity condition. In view of (4.2.2), using the definition of $\boldsymbol{\pi}$, we see that $\mathbf{U}'_{\mathbf{M}}\boldsymbol{\pi} = \tilde{\Lambda}_{\mathbf{M}}$ is equivalent to $\tilde{\mathbf{X}}_{\mathbf{M}}'(\tilde{\mathbf{X}}_{\mathbf{M}}')^+\tilde{\Lambda}_{\mathbf{M}} = \tilde{\Lambda}_{\mathbf{M}}$. This follows from the fact that $\tilde{\mathbf{P}}_{\mathbf{M}}$ is the projection matrix onto the vector subspace $\text{col}(\tilde{\mathbf{X}}_{\mathbf{M}})$, and thus $\mathbf{0}' = [(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\tilde{\mathbf{X}}_{\mathbf{M}}]' = \mathbf{U}'_{\mathbf{M}}\mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})$.

- (c) The assertion of Theorem 4.3.1 cannot be strengthened. Indeed, when the SLOPE solution is not unique, the elements of $S_{\mathbf{X},\Lambda}(\mathbf{Y})$ may have different SLOPE patterns.

Even if many theoretical properties on sign recovery by LASSO are known (see *e.g.* [183]), we believe that it is relevant to give a characterization of sign recovery by LASSO similar as the characterization of pattern recovery by SLOPE given in Theorem 4.3.1.

To show the similarity between proposed irrepresentability condition for SLOPE and the irrepresentability condition for LASSO [82, 192, 194], we propose analogous definitions to the clustered design matrix and parameter.

Remark 4.3.3. Let $\mathbf{0} \neq \mathbf{S} \in \{-1, 0, 1\}^p$ and $k = \|\mathbf{S}\|_1$ (k is the number of nonzero components of \mathbf{S}). The signed matrix $\mathbf{U}_{\mathbf{S}} \in \mathbb{R}^{p \times k}$ is defined by $\mathbf{U}_{\mathbf{S}} = (\mathbf{diag}(\mathbf{S}))_{\text{supp}(\mathbf{S})}$, where $\mathbf{diag}(\mathbf{S}) \in \mathbb{R}^{p \times p}$ is a diagonal matrix and $(\mathbf{diag}(\mathbf{S}))_{\text{supp}(\mathbf{S})}$ denotes the submatrix of $\mathbf{diag}(\mathbf{S})$ obtained by keeping columns corresponding to indices in $\text{supp}(\mathbf{S})$. Observe that for any $\mathbf{0} \neq \boldsymbol{\beta} \in \mathbb{R}^p$ there exists a unique $\tilde{\mathbf{S}} \in \{-1, 0, 1\}^p$ and a unique $\tilde{\boldsymbol{\kappa}}_0 \in (0, \infty)^k$ such that $\boldsymbol{\beta} = \mathbf{U}_{\mathbf{S}}\tilde{\boldsymbol{\kappa}}_0$. Define the reduced matrix $\tilde{\mathbf{X}}_{\mathbf{S}}$ and the reduced parameter $\tilde{\lambda}_{\mathbf{S}}$ by

$$\tilde{\mathbf{X}}_{\mathbf{S}} = \mathbf{X}\mathbf{U}_{\mathbf{S}} \text{ and } \tilde{\lambda}_{\mathbf{S}} = \lambda\mathbf{1}_k, \text{ where } \mathbf{1}_k = (1, \dots, 1)' \in \mathbb{R}^k.$$

Similarly as in the proof of Theorem 4.3.1, one may prove that the necessary and sufficient conditions for the LASSO sign recovery (*i.e.* existence of estimator $\hat{\boldsymbol{\beta}}^{\text{LASSO}}$ such that $\text{sign}(\hat{\boldsymbol{\beta}}^{\text{LASSO}}) = \text{sign}(\boldsymbol{\beta}) = \mathbf{S}$) are the following

$$\begin{cases} \text{there exists } \boldsymbol{\kappa} \in \mathbb{R}_+^k \text{ such that } \tilde{\mathbf{X}}_{\mathbf{S}}'\mathbf{Y} - \tilde{\lambda}_{\mathbf{S}} = \tilde{\mathbf{X}}_{\mathbf{S}}'\tilde{\mathbf{X}}_{\mathbf{S}}\boldsymbol{\kappa}, & (\text{positivity condition}) \\ \mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{S}}')^+\mathbf{1}_k + \frac{1}{\lambda}\mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{X}}_{\mathbf{S}}\tilde{\mathbf{X}}_{\mathbf{S}}^+)\mathbf{Y} \in \partial\|\cdot\|_1(\mathbf{S}). & (\text{subdifferential condition}) \end{cases}$$

In the noiseless case, when $\boldsymbol{\varepsilon} = \mathbf{0}$ and $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$, the subdifferential condition reduces to $\mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{S}}')^+\mathbf{1}_k \in \partial\|\cdot\|_1(\mathbf{S})$ (or equivalently, $\|\mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{S}}')^+\mathbf{1}_k\|_{\infty} \leq 1$ and $\mathbf{1}_k \in \text{col}(\tilde{\mathbf{X}}_{\mathbf{S}}')$). Moreover, if we have $\ker(\mathbf{X}_{\mathbf{S}}) = \{\mathbf{0}\}$, then $\mathbf{1}_k \in \text{col}(\tilde{\mathbf{X}}_{\mathbf{S}}')$ and therefore $\|\mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{S}}')^+\mathbf{1}_k\|_{\infty} \leq 1$ is equivalent to $\|\mathbf{X}'_{\bar{I}}\mathbf{X}_I(\mathbf{X}'_I\mathbf{X}_I)^{-1}\mathbf{S}_I\|_{\infty} \leq 1$, where $I = \text{supp}(\mathbf{S})$, $\bar{I} = \{1, \dots, p\} \setminus I$ and \mathbf{X}_I (resp. $\mathbf{X}_{\bar{I}}$) denotes the submatrix of \mathbf{X} obtained by keeping columns corresponding to indices in I (resp. \bar{I}). This latter expression is known as the irrepresentability condition [82, 192, 194].

From now on, in the definition of SLOPE (4.2.3), we consider the penalty term $J_{\Lambda}(\mathbf{b})$ (with a fixed $\Lambda \in \mathbb{R}^{p+}$) to be multiplied by a scaling parameter $\alpha > 0$ and we denote the set of SLOPE solutions by $S_{\mathbf{X},\alpha\Lambda}(\mathbf{Y})$.

4.3.1 SLOPE irrepresentability condition

As illustrated by Fuchs [82, Theorem 2], Bühlmann and van de Geer [41, Theorem 7.1] and also reminded in Remark 4.3.3, the irrepresentability condition is related to sign recovery by LASSO in the noiseless case. Analogously, studying pattern recovery by SLOPE in the noiseless case allows us to introduce the SLOPE irrepresentability condition. The latter condition will be very useful in the discussed later case when $\boldsymbol{\varepsilon}$ is no longer null. Corollary 4.3.1, which provides a characterization of pattern recovery by SLOPE in the noiseless case, is a consequence of Theorem 4.3.1.

Definition 4.3.1. Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{\Lambda} \in \mathbb{R}^{p+}$ and $\mathbf{0} \neq \boldsymbol{\beta} \in \mathbb{R}^p$. The noiseless pattern recovery by SLOPE is defined as

$$\exists \alpha > 0 \exists \hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \alpha \mathbf{\Lambda}}(\mathbf{X}\boldsymbol{\beta}) \text{ such that } \mathbf{patt}(\hat{\boldsymbol{\beta}}) = \mathbf{patt}(\boldsymbol{\beta}). \quad (4.3.2)$$

Corollary 4.3.1. Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{\Lambda} \in \mathbb{R}^{p+}$ and $\mathbf{0} \neq \boldsymbol{\beta} \in \mathbb{R}^p$. The noiseless pattern recovery by SLOPE is equivalent to

$$J_{\mathbf{\Lambda}}^*(\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ \widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) \leq 1 \text{ and } \widetilde{\mathbf{\Lambda}}_{\mathbf{M}} \in \text{col}(\widetilde{\mathbf{X}}_{\mathbf{M}}')$$

(or equivalently, to $\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ \widetilde{\mathbf{\Lambda}}_{\mathbf{M}} \in \partial J_{\mathbf{\Lambda}}(\mathbf{M})$). Moreover, under this condition, there exists $\alpha_0 > 0$ such that for all $\alpha \in (0, \alpha_0)$ there exists $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \alpha \mathbf{\Lambda}}(\mathbf{X}\boldsymbol{\beta})$ for which $\mathbf{patt}(\hat{\boldsymbol{\beta}}) = \mathbf{patt}(\boldsymbol{\beta})$.

Now we are ready to define the SLOPE irrepresentability condition.

Definition 4.3.2. Let $\mathbf{M} = \mathbf{patt}(\boldsymbol{\beta})$. We define the SLOPE irrepresentability condition as the following inequality and inclusion:

$$J_{\mathbf{\Lambda}}^*(\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ \widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) \leq 1 \text{ and } \widetilde{\mathbf{\Lambda}}_{\mathbf{M}} \in \text{col}(\widetilde{\mathbf{X}}_{\mathbf{M}}'). \quad (4.3.3)$$

Remark 4.3.4. If $\ker(\widetilde{\mathbf{X}}_{\mathbf{M}}) = \{\mathbf{0}\}$, then $\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ = \mathbf{X}'\widetilde{\mathbf{X}}_{\mathbf{M}}(\widetilde{\mathbf{X}}_{\mathbf{M}}'\widetilde{\mathbf{X}}_{\mathbf{M}})^{-1}$ and consequently the SLOPE irrepresentability condition reads $J_{\mathbf{\Lambda}}^*(\mathbf{X}'\widetilde{\mathbf{X}}_{\mathbf{M}}(\widetilde{\mathbf{X}}_{\mathbf{M}}'\widetilde{\mathbf{X}}_{\mathbf{M}})^{-1}\widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) \leq 1$.

Example 4.3.5. Let $p = 2$, $\mathbf{\Lambda} = (4, 2)'$, $\boldsymbol{\beta} = (5, 0)'$ and $\bar{\boldsymbol{\beta}} = (5, 3)'$. Consider a design matrix $\mathbf{X} = (\mathbf{X}_1 | \mathbf{X}_2) \in \mathbb{R}^{n \times 2}$ satisfying

$$\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix}.$$

- The SLOPE irrepresentability condition does not occur when $\boldsymbol{\beta} = (5, 0)'$. Indeed, $\mathbf{M} = \mathbf{patt}(\boldsymbol{\beta}) = (1, 0)'$, $\widetilde{\mathbf{X}}_{\mathbf{M}} = \mathbf{X}_1$ (thus $\widetilde{\mathbf{X}}_{\mathbf{M}}'\widetilde{\mathbf{X}}_{\mathbf{M}} = 1$) and $\widetilde{\mathbf{\Lambda}}_{\mathbf{M}} = \lambda_1 = 4$. Therefore

$$J_{\mathbf{\Lambda}}^*(\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ \widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) = J_{\mathbf{\Lambda}}^*(\mathbf{X}'\widetilde{\mathbf{X}}_{\mathbf{M}}(\widetilde{\mathbf{X}}_{\mathbf{M}}'\widetilde{\mathbf{X}}_{\mathbf{M}})^{-1}\widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) = J_{\mathbf{\Lambda}}^*(4\mathbf{X}'\widetilde{\mathbf{X}}_{\mathbf{M}}) = 6.4/6 > 1.$$

- The SLOPE irrepresentability condition occurs when $\bar{\boldsymbol{\beta}} = (5, 3)'$. Indeed, $\mathbf{M} = \mathbf{patt}(\bar{\boldsymbol{\beta}}) = (2, 1)'$, $\widetilde{\mathbf{X}}_{\mathbf{M}} = \mathbf{X}$ and $\widetilde{\mathbf{\Lambda}}_{\mathbf{M}} = \mathbf{\Lambda}$. Therefore $\ker(\widetilde{\mathbf{X}}_{\mathbf{M}}) = \{\mathbf{0}\}$ and

$$J_{\mathbf{\Lambda}}^*(\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}}')^+ \widetilde{\mathbf{\Lambda}}_{\mathbf{M}}) = J_{\mathbf{\Lambda}}^*(\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{\Lambda}) = J_{\mathbf{\Lambda}}^*(\mathbf{\Lambda}) = 1 \leq 1.$$

Figure 4.2 corroborates graphically that SLOPE irrepresentability condition does not hold for $\boldsymbol{\beta}$, but it is satisfied for $\bar{\boldsymbol{\beta}}$. Note that, in this setup, the SLOPE solution is unique, since $\ker(\mathbf{X}) = \{\mathbf{0}\}$. By $\hat{\boldsymbol{\beta}}(\alpha)$ we denote the unique element of $S_{\mathbf{X}, \alpha \mathbf{\Lambda}}(\mathbf{X}\boldsymbol{\beta})$. Then the SLOPE solution path is the function $(0, \infty) \ni \alpha \mapsto \hat{\boldsymbol{\beta}}(\alpha)$.

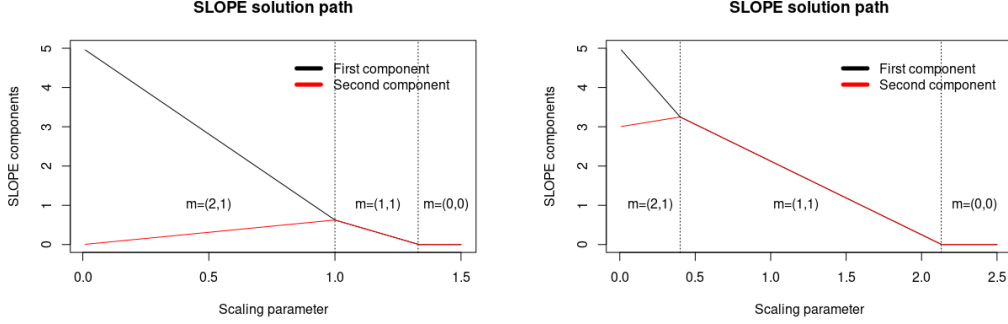


Figure 4.2: SLOPE solution paths for $p = 2$ and $\Lambda = (4, 2)'$ in the noiseless case. β is equal respectively to $(5, 0)'$ (left panel) and $(5, 3)'$ (right panel). As may be seen, the pattern of $\beta = (5, 0)'$ cannot be recovered in the noiseless case for any scaling of Λ . The pattern of $\bar{\beta} = (5, 3)'$ is recovered for $\alpha \in (0, 0.4)$.

4.4 Geometrical interpretation of Irrepresentability Condition

Let $\mathbf{0} \neq \beta \in \mathbb{R}^p$. By (4.2.4), for a SLOPE minimizer $\hat{\beta} \in S_{\mathbf{X}, \alpha \Lambda}(\mathbf{X}\beta)$ the following occurs:

$$\frac{1}{\alpha} \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \in \partial J_{\Lambda}(\hat{\beta}).$$

In addition, for $M = \mathbf{patt}(\hat{\beta})$, the following facts hold:

- $\beta - \hat{\beta} \in \text{col}(U_M)$, thus $\frac{1}{\alpha} \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \in \mathbf{X}' \mathbf{X} \text{col}(U_M)$.
- $\partial J_{\Lambda}(\hat{\beta}) = \partial J_{\Lambda}(M)$.

Therefore, the noiseless pattern recovery by SLOPE implies that the vector space $\mathbf{X}' \mathbf{X} \text{col}(U_M) = \text{col}(\mathbf{X}' \tilde{\mathbf{X}}_M)$ intersects $\partial J_{\Lambda}(M)$. Actually, the vector $\bar{\Pi} = \mathbf{X}'(\tilde{\mathbf{X}}_M')^+ \tilde{\Lambda}_M$ appearing in Corollary 4.3.1 may be interpreted in a geometrical way as we propose below.

Proposition 4.4.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{0} \neq M \in \mathcal{M}_p^{\text{SLOPE}}$ and $\Lambda \in \mathbb{R}^{p+}$. Recall that $\tilde{\mathbf{X}}_M = \mathbf{X} U_M$, $\tilde{\Lambda}_M = U_{|\mathbf{M}|_{\downarrow}} \Lambda$ and $\bar{\Pi} = \mathbf{X}'(\tilde{\mathbf{X}}_M')^+ \tilde{\Lambda}_M$. Then:*

- If $\tilde{\Lambda}_M \notin \text{col}(\tilde{\mathbf{X}}_M')$, then $\text{aff}(\partial J_{\Lambda}(M)) \cap \text{col}(\mathbf{X}' \tilde{\mathbf{X}}_M) = \emptyset$.
- If $\tilde{\Lambda}_M \in \text{col}(\tilde{\mathbf{X}}_M')$, then $\text{aff}(\partial J_{\Lambda}(M)) \cap \text{col}(\mathbf{X}' \tilde{\mathbf{X}}_M) = \{\bar{\Pi}\}$.
- Pattern recovery by SLOPE for $\varepsilon = \mathbf{0}$ is equivalent to

$$\text{col}(\mathbf{X}' \tilde{\mathbf{X}}_M) \cap \partial J_{\Lambda}(M) \neq \emptyset.$$

In other words, the accessibility condition means that the intersection of $\text{col}(\mathbf{X}' \tilde{\mathbf{X}}_M) \cap \text{aff}(\partial J_{\Lambda}(M))$ is not empty. Moreover, then it is equal to the vector $\bar{\Pi} = \mathbf{X}'(\tilde{\mathbf{X}}_M')^+ \tilde{\Lambda}_M$. Moreover, when the accessibility condition holds, by Proposition 4.4.1 iii) the noiseless recovery of the SLOPE pattern is equivalent to the subdifferential condition $\bar{\Pi} \in \partial J_{\Lambda}(M)$.

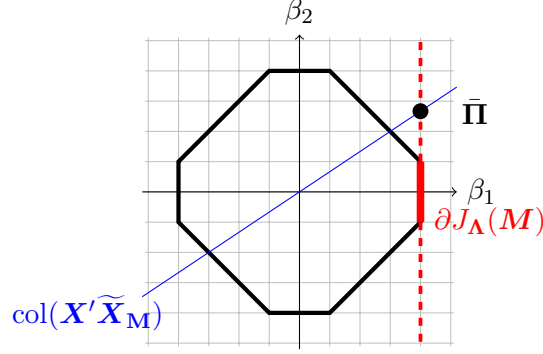


Figure 4.3: Intersection $\bar{\Pi}$ of $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) = \text{col}((1, 2/3)')$ and $\text{aff}(\partial J_\Lambda(\mathbf{M}))$ for $\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 2/3 \\ 2/3 & 1 \end{pmatrix}$, $\Lambda = (4, 1)'$ and $\text{patt}(\beta) = \mathbf{M} = (1, 0)'$. Since $\bar{\Pi} = (4, 8/3)'$ $\notin \partial J_\Lambda(\mathbf{M})$, then in the noiseless case, by Proposition 4.4.1, SLOPE does not recover \mathbf{M} . However, as $\bar{\Pi}$ exists, the pattern \mathbf{M} is accessible, as in the latter examples.

Proof. i) Recall that, according to Lemma 4.9.2, $\text{aff}(\partial J_\Lambda(\mathbf{M})) = \{\mathbf{v} \in \mathbb{R}^p : \mathbf{U}'_M \mathbf{v} = \tilde{\Lambda}_M\}$. If $\text{aff}(\partial J_\Lambda(\mathbf{M})) \cap \text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) \neq \emptyset$, then there exists $\mathbf{z} \in \mathbb{R}^k$, where $k = \|\mathbf{M}\|_\infty$, such that $\mathbf{X}'\tilde{\mathbf{X}}_M \mathbf{z} \in \text{aff}(\partial J_\Lambda(\mathbf{M}))$. Consequently, $\tilde{\Lambda}_M = \mathbf{U}'_M \mathbf{X}'\tilde{\mathbf{X}}_M \mathbf{z} = \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{z}$, thus $\tilde{\Lambda}_M \in \text{col}(\tilde{\mathbf{X}}'_M)$.

ii) If $\tilde{\Lambda}_M \in \text{col}(\tilde{\mathbf{X}}'_M)$, then $\bar{\Pi} \in \text{aff}(\partial J_\Lambda(\mathbf{M}))$. Indeed, since $\tilde{\mathbf{X}}'_M (\tilde{\mathbf{X}}'_M)^+$ is the projection onto $\text{col}(\tilde{\mathbf{X}}'_M)$, we have

$$\mathbf{U}'_M \bar{\Pi} = \tilde{\mathbf{X}}'_M (\tilde{\mathbf{X}}'_M)^+ \tilde{\Lambda}_M = \tilde{\Lambda}_M.$$

Moreover, since $\text{col}((\tilde{\mathbf{X}}'_M)^+) = \text{col}(\tilde{\mathbf{X}}_M)$, we deduce that $\bar{\Pi} \in \text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M)$. To prove that $\bar{\Pi}$ is the unique point in the intersection, we will show that $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) \cap \text{col}(\mathbf{U}_M)^\perp = \{\mathbf{0}\}$. Indeed, if $\mathbf{v} \in \text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) \cap \text{col}(\mathbf{U}_M)^\perp$, then there exists such $\mathbf{z} \in \mathbb{R}^k$ that $\mathbf{v} = \mathbf{X}'\tilde{\mathbf{X}}_M \mathbf{z} = \mathbf{X}'\mathbf{X}\mathbf{U}_M \mathbf{z}$ and $(\mathbf{U}_M \mathbf{z})' \mathbf{v} = \mathbf{0}$. Therefore, $(\mathbf{U}_M \mathbf{z})' \mathbf{X}'\mathbf{X}\mathbf{U}_M \mathbf{z} = \mathbf{0}$, which implies that $\mathbf{X}\mathbf{U}_M \mathbf{z} = \tilde{\mathbf{X}}_M \mathbf{z} = \mathbf{0}$ and thus $\mathbf{v} = \mathbf{0}$. Finally, if $\bar{\Pi} \in \text{aff}(\partial J_\Lambda(\mathbf{M})) \cap \text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M)$, then $\bar{\Pi} - \bar{\Pi} \in \text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M)$ and $\mathbf{U}'_M (\bar{\Pi} - \bar{\Pi}) = \mathbf{0}$, which implies that $\bar{\Pi} = \bar{\Pi}$.

iii) According to Corollary 4.3.1, pattern recovery by SLOPE in the noiseless case is equivalent to $\bar{\Pi} \in \partial J_\Lambda(\mathbf{M})$ which is equivalent, by i) and ii), to $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) \cap \partial J_\Lambda(\mathbf{M}) \neq \emptyset$. \square

Proposition 4.4.1 may be interpreted with a notion of the accessibility condition for SLOPE, which was introduced in [156].

Definition 4.4.1 (Accessible pattern). [156] Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\lambda > 0$ and pen be a polyhedral gauge. We say that the pattern of $\beta \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} and λpen , if there exist $\mathbf{y} \in \mathbb{R}^n$ and $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ such that $\hat{\beta} \stackrel{\text{pen}}{\approx} \beta$.

The accessibility of a pattern can be characterized in a geometric and an analytic way.

Proposition 4.4.2 (Characterization of accessible patterns). Let $\mathbf{X} \in \mathbb{R}^{n \times p}$.

(a) *Geometric characterization:* The SLOPE pattern of $\beta \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} if and only if

$$\text{row}(\mathbf{X}) \cap \partial \text{pen}(\beta) \neq \emptyset.$$

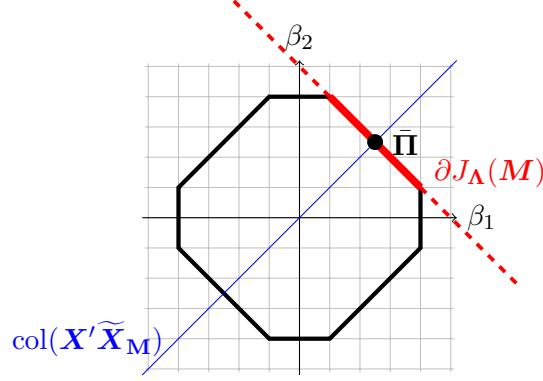


Figure 4.4: Intersection $\bar{\Pi}$ of $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) = \text{col}((1,1)')$ and $\text{aff}(\partial J_\Lambda(M))$ for $\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 2/3 \\ 2/3 & 1 \end{pmatrix}$, $\Lambda = (4,1)'$ and $\mathbf{patt}(\beta) = M = (1,1)'$. Since $\bar{\Pi} = (2.5, 2.5)' \in \partial J_\Lambda(M)$, then in the noiseless case, by Proposition 4.4.1, SLOPE recovers the pattern M of β .

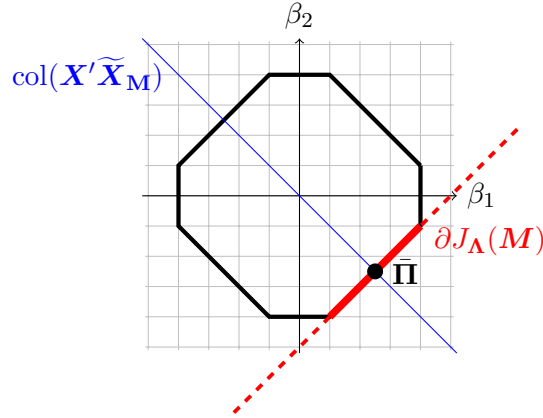


Figure 4.5: Intersection $\bar{\Pi}$ of $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_M) = \text{col}((1,-1)')$ and $\text{aff}(\partial J_\Lambda(M))$ for $\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 2/3 \\ 2/3 & 1 \end{pmatrix}$, $\Lambda = (4,1)'$ and $\mathbf{patt}(\beta) = M = (1,-1)'$. Since $\bar{\Pi} \notin \partial J_\Lambda(M)$, then in the noiseless case, by Proposition 4.4.1, SLOPE recovers the pattern M of β .

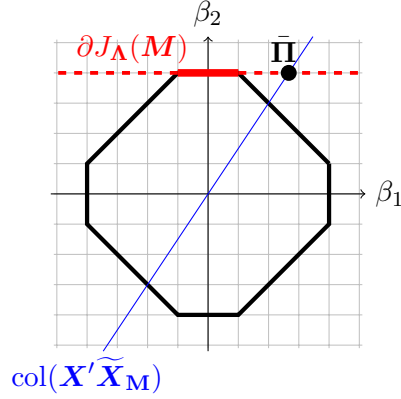


Figure 4.6: Intersection $\bar{\Pi}$ of $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_{\mathbf{M}}) = \text{col}((2/3, 1)')$ and $\text{aff}(\partial J_{\Lambda}(\mathbf{M}))$ for $\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 2/3 \\ 2/3 & 1 \end{pmatrix}$, $\Lambda = (4, 1)'$ and $\text{patt}(\beta) = \mathbf{M} = (0, 1)'$. Since $\bar{\Pi} = (8/3, 4)' \notin \partial J_{\Lambda}(\mathbf{M})$, then in the noiseless case, by Proposition 4.4.1, SLOPE does not recover \mathbf{M} .

(b) *Analytic characterization: The SLOPE pattern of $\beta \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} if and only if for any $\mathbf{b} \in \mathbb{R}^p$ we have*

$$\mathbf{X}\beta = \mathbf{X}\mathbf{b} \implies \text{pen}(\beta) \leq \text{pen}(\mathbf{b}).$$

For $p = 2, n \geq 2$, a full rank matrix \mathbf{X} and the models with zeros or with 2-element clusters, the Proposition 4.4.1 and Corollary 4.3.1 are illustrated on Figures 4.3 – 4.6. They present the SLOPE dual unit ball C_{Λ} , linear spaces $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_{\mathbf{M}})$ (blue lines) and, if existing, their intersection points $\bar{\Pi} = \mathbf{X}'(\tilde{\mathbf{X}}_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}}$ (black dots) with the affine spaces $\text{aff}(F_{\mathbf{M}})$ (red dashed lines) of the corresponding pattern faces $F_{\mathbf{M}} = \partial J_{\Lambda}(\mathbf{M})$ (red segments). The existence of $\bar{\Pi}$ is equivalent to the accessibility condition. The SLOPE irrepresentability condition holds if and only if $\bar{\Pi} \in F_{\mathbf{M}}$. It holds true in Figures 4.4 and 4.5 and fails in Figures 4.3 and 4.6.

In the low-dimensional setting, when the design matrix \mathbf{X} is of a full column rank, another geometrical interpretation of $\hat{\beta}^{\text{SLOPE}}$ might be used. As a corollary of Theorem 3.2.2, one may deduce that the difference between $\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$ is equal to

$$(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\pi}^* \in (\mathbf{X}'\mathbf{X})^{-1}C_{\Lambda},$$

where C_{Λ} is the unit ball in J_{Λ}^* and $\boldsymbol{\pi}^*$ is defined as in Lemma 3.2.1. On Figure 4.7 we present graphically this relation between $\hat{\beta}^{\text{SLOPE}}$ and $\hat{\beta}^{\text{OLS}}$.

The illustration of the SLOPE dual norm unit ball C_{Λ} in \mathbb{R}^p for $p = 3$ may be found in the recently published article of Schneider and Tardivel [156]. We recall it to present that with growing number of explanatory variables p the complexity of C_{Λ} rises drastically.

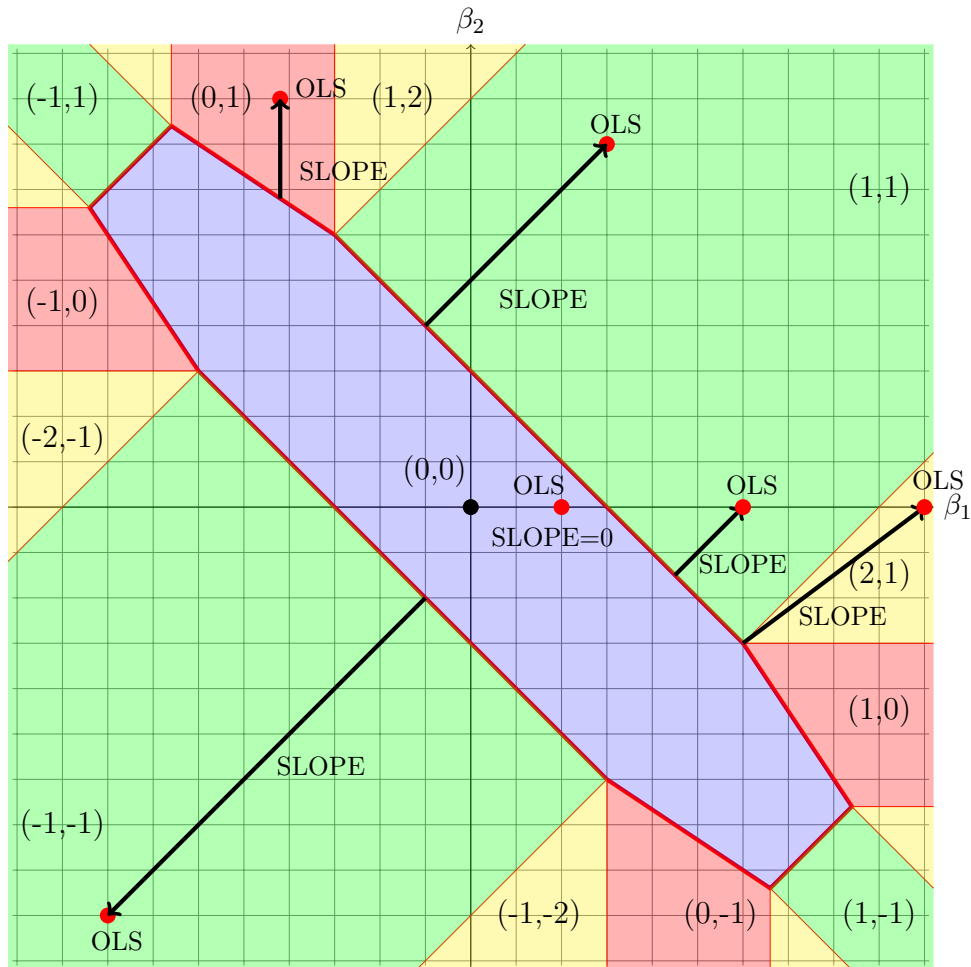


Figure 4.7: Comparison of $\hat{\beta}^{SLOPE}$, $\hat{\beta}^{OLS}$ and their difference $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}\varepsilon \in (\mathbf{X}'\mathbf{X})^{-1}C_{\Lambda}$ for $p = 2$, $\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 2/3 \\ 2/3 & 1 \end{bmatrix}$ and $\Lambda = (4, 1)'$.

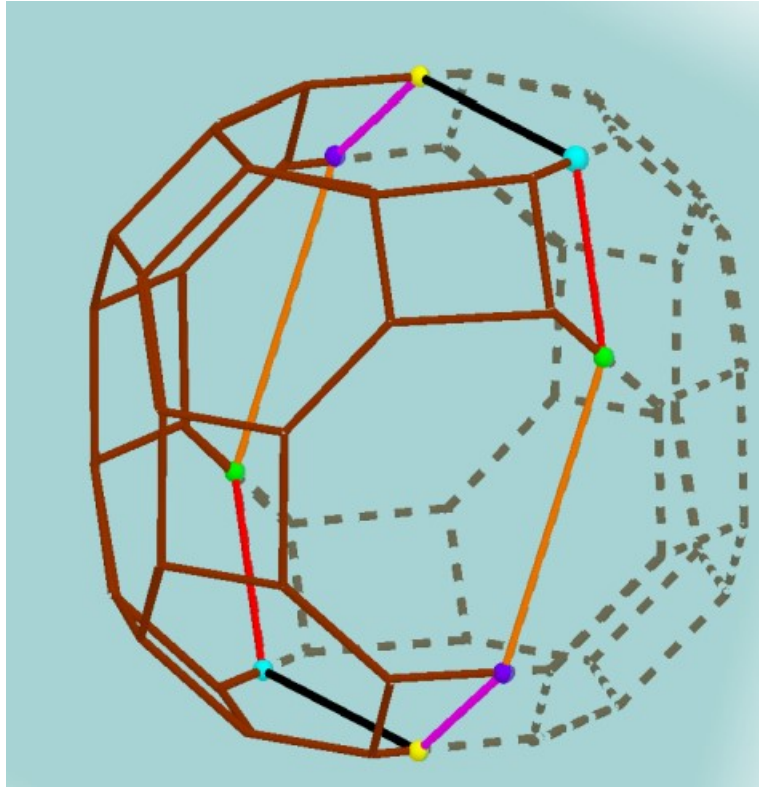


Figure 4.8: C_{Λ} in \mathbb{R}^p for $p = 3$ with its two-dimensional subset corresponding to the intersection of C_{Λ} with $\text{row}(\mathbf{X})$ from [156, Example 4]. One may observe that in this example there are only 17 accessible patterns out of all 147 patterns for $p = 3$.

While C_{Λ} consists of 17 faces for $p = 2$, for $p = 3$ the number of faces increases to 147. For larger values of p , the exact number of faces of C_{Λ} may be found at the On-line Encyclopedia of Integer Sequences [138, A080253]. The growth of the number of faces is super-exponential with respect to p .

4.5 Asymptotics of pattern recovery and pattern consistency

In this section we consider two asymptotic scenarios and establish conditions on tuning parameters for which the pattern of β is recovered. In Section 4.5.1 we consider the case where gaps between distinct absolute values of β diverge and in Section 4.5.2 the case where the sample size n diverges. The proofs rely on Theorem 4.3.1. We show that the positivity and subdifferential conditions are satisfied under our settings. It turns out that for the positivity condition the tuning parameter cannot be too large, while for the subdifferential condition it cannot be too small. In this way we consider a tuning parameter of the form $\alpha\Lambda$, where $\Lambda \in \mathbb{R}^{p+}$ is fixed and α varies. We determine the assumptions for the sequence (α) for which both positivity and subdifferential conditions hold true, *i.e.* for which the pattern is recovered.

4.5.1 \mathbf{X} is a fixed matrix

The subdifferential condition, given in Theorem 4.3.1, says that a vector π defined in (4.3.1) belongs to $\partial J_{\alpha\Lambda}(\mathbf{M})$, where α is a scaling parameter. This condition is equivalent to requiring

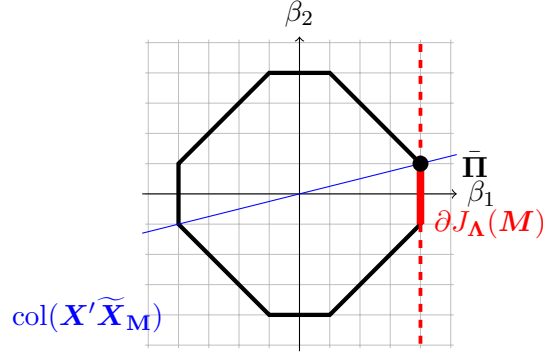


Figure 4.9: Intersection $\bar{\Pi}$ of $\text{col}(\mathbf{X}'\tilde{\mathbf{X}}_{\mathbf{M}}) = \text{col}((1, 1/4)')$ and $\text{aff}(\partial J_{\Lambda}(\mathbf{M}))$ for $\mathbf{X}'\mathbf{X} = \begin{pmatrix} 1 & 1/4 \\ 1/4 & 1 \end{pmatrix}$, $\Lambda = (4, 1)'$ and $\mathbf{patt}(\beta) = \mathbf{M} = (1, 0)'$. Since $\bar{\Pi} = (4, 1)' \in F_{\mathbf{M}} = \partial J_{\Lambda}(\mathbf{M})$, but it lies on a boundary of $F_{\mathbf{M}}$. Therefore the open SLOPE irrepresentability condition does not hold, even though the irrepresentability condition is satisfied.

that a vector π/α belongs to $\partial J_{\Lambda}(\mathbf{M})$. We denote the vector π/α by

$$\pi_{\alpha} = \mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^{\dagger} \tilde{\Lambda}_{\mathbf{M}} + \frac{1}{\alpha} \mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\mathbf{Y} = \mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^{\dagger} \tilde{\Lambda}_{\mathbf{M}} + \frac{1}{\alpha} \mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\boldsymbol{\varepsilon}, \quad (4.5.1)$$

where in the latter equality we use the fact that $(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})$ is an orthogonal projection onto $\text{col}(\tilde{\mathbf{X}}_{\mathbf{M}})^{\perp}$ and therefore $(\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\mathbf{X}\beta = (\mathbf{I}_n - \tilde{\mathbf{P}}_{\mathbf{M}})\tilde{\mathbf{X}}_{\mathbf{M}}\mathbf{s} = \mathbf{0}$, where $\beta = \mathbf{U}_{\mathbf{M}}\mathbf{s}$ and $\mathbf{s} \in \mathbb{R}^{\|\mathbf{M}\|_{\infty}+}$.

By Theorem 4.3.1, the probability of the pattern recovery by SLOPE is upper bounded by

$$\mathbb{P}\left(\exists \hat{\beta} \in S_{\mathbf{X}, \alpha \Lambda}(\mathbf{Y}) \text{ such that } \mathbf{patt}(\hat{\beta}) = \mathbf{patt}(\beta)\right) \leq \begin{cases} \mathbb{P}(J_{\Lambda}^*(\pi_{\alpha}) \leq 1), \\ 0 \text{ if } \tilde{\Lambda}_{\mathbf{M}} \notin \text{col}(\tilde{\mathbf{X}}'_{\mathbf{M}}). \end{cases} \quad (4.5.2)$$

The first point in Theorem 4.5.1 shows that the probability of pattern recovery matches with the upper bound (4.5.2) when gaps between different absolute values of terms of β are large enough. The last point provides the pattern consistency by SLOPE.

Before stating this Theorem, we introduce the open SLOPE irrepresentability condition.

Definition 4.5.1 (Open SLOPE irrepresentability condition).

$$\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^{\dagger} \tilde{\Lambda}_{\mathbf{M}} \in \text{ri}(\partial J_{\Lambda}(\mathbf{M})).$$

We illustrate the difference between the SLOPE irrepresentability condition and its open version on Figure 4.9.

Theorem 4.5.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{0} \neq \mathbf{M} \in \mathcal{M}_p^{\text{SLOPE}}$, and $\Lambda = (\lambda_1, \dots, \lambda_p)' \in \mathbb{R}^{p+}$. Consider a sequence of signals $(\beta^{(r)})_{r \geq 1}$ with pattern \mathbf{M} :*

$$\beta^{(r)} = \mathbf{U}_{\mathbf{M}}\mathbf{s}^{(r)} \quad \text{with} \quad \mathbf{s}^{(r)} \in \mathbb{R}^{k+} \text{ and } k = \|\mathbf{M}\|_{\infty},$$

whose strength is increasing in the following sense:

$$\Delta_r = \min_{1 \leq i < k} (s_i^{(r)} - s_{i+1}^{(r)}) \xrightarrow{r \rightarrow \infty} \infty, \text{ with the convention } s_{k+1}^{(r)} = 0$$

and let $\mathbf{Y}^{(r)} = \mathbf{X}\beta^{(r)} + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon}$ is a vector in \mathbb{R}^n .

- (a) *Sharpness of the upper bound:* Let $\alpha > 0$. If ε is random, then the upper bound (4.5.2) is asymptotically reached:

$$\lim_{r \rightarrow \infty} \mathbb{P} \left(\exists \hat{\beta} \in S_{\mathbf{X}, \alpha \Lambda}(\mathbf{Y}^{(r)}) \text{ such that } \mathbf{patt}(\hat{\beta}) = \mathbf{M} \right) = \begin{cases} \mathbb{P}(J_{\Lambda}^*(\boldsymbol{\pi}_{\alpha}) \leq 1), \\ 0 \text{ if } \tilde{\Lambda}_{\mathbf{M}} \notin \text{col}(\tilde{\mathbf{X}}'_{\mathbf{M}}). \end{cases}$$

- (b) *Pattern consistency:* If $\alpha_r \rightarrow \infty$, $\alpha_r/\Delta_r \rightarrow 0$ as $r \rightarrow \infty$ and

$$\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}} \in \text{ri}(\partial J_{\Lambda}(\mathbf{M})),$$

then for any $\varepsilon \in \mathbb{R}^n$ we have

$$\exists r_0 > 0 \quad \forall r \geq r_0 \quad \exists \hat{\beta} \in S_{\mathbf{X}, \alpha_r \Lambda}(\mathbf{Y}^{(r)}) \text{ such that } \mathbf{patt}(\hat{\beta}) = \mathbf{M}.$$

Remark 4.5.2.

- (a) The open SLOPE irrepresentability condition $\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}} \in \text{ri}(\partial J_{\Lambda}(\mathbf{M}))$ is equivalent to the following computationally verifiable conditions:

$$\left\{ \begin{array}{l} J_{\Lambda}^*(\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}}) \leq 1 \text{ and } \tilde{\Lambda}_{\mathbf{M}} \in \text{col}(\tilde{\mathbf{X}}'_{\mathbf{M}}), \\ \left| \left\{ i \in \{1, \dots, p\} : \sum_{j=1}^i |\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}}|_{(j)} = \sum_{j=1}^i \lambda_j \right\} \right| = \|\mathbf{M}\|_{\infty}. \end{array} \right.$$

This equivalence follows from Proposition 4.9.1.

- (b) Let us assume that the distribution of ε and $-\varepsilon$ is equal. Because the unit ball of the dual sorted ℓ_1 norm is convex, when $J_{\Lambda}^*(\mathbf{X}'(\tilde{\mathbf{X}}'_{\mathbf{M}})^+ \tilde{\Lambda}_{\mathbf{M}}) > 1$ then, independently on $\alpha > 0$, the probability of pattern recovery is smaller than 1/2, namely

$$\mathbb{P} \left(\exists \hat{\beta} \in S_{\mathbf{X}, \alpha \Lambda}(\mathbf{Y}) \text{ such that } \mathbf{patt}(\hat{\beta}) = \mathbf{M} \right) \leq 1/2.$$

For LASSO, a similar inequality on the probability of sign recovery is given in [183].

- (c) In Section 4.7, we illustrate that, under the open irrepresentability condition, one may select $\alpha > 0$ to fix the asymptotic probability of pattern recovery at a level arbitrarily close to 1 (a similar result for LASSO is given in [174]).

4.5.2 X is random, p is fixed, n tends to infinity

In this section we discuss asymptotic properties of the SLOPE estimator in the low-dimensional regression model in which p is fixed and the sample size n tends to infinity.

For each $n \geq p$ we consider a linear regression problem

$$\mathbf{Y}_n = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}_n, \tag{4.5.3}$$

where $\mathbf{X}_n \in \mathbb{R}^{n \times p}$ is a random design matrix. We will use the following assumptions:

- A. $\boldsymbol{\varepsilon}_n = (\varepsilon_1, \dots, \varepsilon_n)'$, where $(\varepsilon_i)_i$ are i.i.d. centered with finite variance.
 B1. A sequence of design matrices $\mathbf{X}_1, \mathbf{X}_2, \dots$ satisfies the condition

$$\frac{1}{n} \mathbf{X}'_n \mathbf{X}_n \xrightarrow{\mathbb{P}} \mathbf{C}, \tag{4.5.4}$$

where \mathbf{C} is a deterministic positive definite symmetric $p \times p$ matrix.

B2. For each $j = 1, \dots, p$,

$$\frac{\max_{i=1, \dots, n} |X_{ij}^{(n)}|}{\sqrt{\sum_{i=1}^n (X_{ij}^{(n)})^2}} \xrightarrow{\mathbb{P}} 0.$$

C. $(\mathbf{X}_n)_n$ and $(\varepsilon_n)_n$ are independent.

We will consider a sequence of tuning parameters $(\Lambda_n)_n$ defined by

$$\Lambda_n = \alpha_n \Lambda,$$

where $\Lambda \in \mathbb{R}^{p+}$ is fixed and $(\alpha_n)_n$ is a sequence of positive numbers.

Let $\hat{\beta}_n^{\text{SLOPE}}$ be an element of the set $S_{\mathbf{X}_n, \Lambda_n}(\mathbf{Y}_n)$ of SLOPE minimizers. Under Assumption B1, for large n with high probability, the set $S_{\mathbf{X}_n, \Lambda_n}(\mathbf{Y}_n)$ consists of one element. Indeed, we have

$$\mathbb{P}(\ker(\mathbf{X}_n) = \{\mathbf{0}\}) = \mathbb{P}(\mathbf{X}'_n \mathbf{X}_n \text{ is positive definite}) \xrightarrow{n \rightarrow \infty} 1$$

and $\ker(\mathbf{X}_n) = \{\mathbf{0}\}$ ensures existence of the unique SLOPE minimizer. In a natural setting, the strong consistency of $\hat{\beta}_n^{\text{SLOPE}}$ can be characterized in terms of behaviour of the tuning parameter. We use such approach in the Theorem 3.4.1 and Theorem 4.6.2. At this point we note that if (4.6.8) holds almost surely, then the condition $\alpha_n/n \rightarrow 0$ ensures that $\hat{\beta}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \beta$. Thus, if β does not have any non-trivial clusters nor zeros, *i.e.* $\|\mathbf{patt}(\beta)\|_\infty = p$, then $\alpha_n/n \rightarrow 0$ suffices for $\mathbf{patt}(\hat{\beta}_n^{\text{SLOPE}}) \xrightarrow{\text{a.s.}} \mathbf{patt}(\beta)$. However, if $\|\mathbf{patt}(\beta)\| < p$, then the situation is more complex as we shall show below.

The first of our asymptotic results concerns the consistency of the pattern recovery by the SLOPE estimator. We note that Assumption B2 is not necessary for the SLOPE pattern recovery. This assumption was introduced to ensure the existence of a Gaussian vector in the Theorem 4.5.3 (a).

Theorem 4.5.3. *Under the assumptions A, B1, C, the following statements hold true.*

(a) *If B2 is additionally satisfied and moreover $\alpha_n = \sqrt{n}$, then*

$$\lim_{n \rightarrow \infty} \mathbb{P}(\mathbf{patt}(\hat{\beta}_n^{\text{SLOPE}}) = \mathbf{patt}(\beta)) = \mathbb{P}(J_\Lambda^*(\mathbf{V}) \leq 1),$$

where $\mathbf{V} \sim \mathcal{N}(\mathbf{C} \mathbf{U}_M (\mathbf{U}'_M \mathbf{C} \mathbf{U}_M)^{-1} \tilde{\Lambda}_M, \sigma^2 [\mathbf{C} - \mathbf{C} \mathbf{U}_M (\mathbf{U}'_M \mathbf{C} \mathbf{U}_M)^{-1} \mathbf{U}'_M \mathbf{C}])$.

(b) *Let*

$$\mathbf{C} \mathbf{U}_M (\mathbf{U}'_M \mathbf{C} \mathbf{U}_M)^{-1} \tilde{\Lambda}_M \in \text{ri}(\partial J_\Lambda(\mathbf{M})). \quad (4.5.5)$$

The pattern of SLOPE estimator is consistent, i.e.

$$\mathbf{patt}(\hat{\beta}_n^{\text{SLOPE}}) \xrightarrow{\mathbb{P}} \mathbf{patt}(\beta),$$

if and only if

$$\lim_{n \rightarrow \infty} \frac{\alpha_n}{n} = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\alpha_n}{\sqrt{n}} = \infty.$$

(c) *The condition*

$$J_\Lambda^*(\mathbf{C} \mathbf{U}_M (\mathbf{U}'_M \mathbf{C} \mathbf{U}_M)^{-1} \tilde{\Lambda}_M) \leq 1 \quad (4.5.6)$$

is necessary for pattern consistency of SLOPE estimator.

The random vector \mathbf{V} belongs to the smallest affine space containing $\partial J_{\Lambda}(\mathbf{b})$, *i.e.* $\text{aff}(\partial J_{\Lambda}(\mathbf{b})) = \{\mathbf{v} \in \mathbb{R}^p: \mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}_{\mathbf{M}}\}$, see Lemma 4.9.2.

Condition (4.5.5) is the open SLOPE irrepresentability condition in the $n \rightarrow \infty$ regime. The above result should be compared with [192, Theorem 1], where the same conditions on the LASSO tuning parameter ensure consistency of sign recovery by LASSO estimator. Below we make one step further and consider the strong consistency of SLOPE pattern recovery by $\hat{\beta}_n^{\text{SLOPE}}$. Although this was not the main focus of Zhao and Yu, it can be deduced from [192, Theorem 1] that if for $c \in (0, 1)$ the LASSO tuning parameter λ_n satisfies $\lambda_n/n \rightarrow 0$ and $\lambda_n/n^{\frac{1+c}{2}} \rightarrow \infty$, then under the strong LASSO irrepresentability condition, one has $\text{sign}(\hat{\beta}_n^{\text{LASSO}}) \xrightarrow{\text{a.s.}} \text{sign}(\beta)$. Even though the patterns are discrete objects, as the underlying probability space is uncountable, the convergence in probability does not imply the almost sure convergence. Below we show that if $\alpha_n/n \rightarrow 0$ and $\alpha_n/\sqrt{n} \rightarrow \infty$, then $\text{patt}(\hat{\beta}_n^{\text{SLOPE}})$ is not strongly consistent and one actually has to impose a slightly stronger condition (4.5.7).

For the purpose of the a.s. convergence, we strengthen the assumption on design matrices:

- B'. Assume that the rows of \mathbf{X}_n are independent and that each row of \mathbf{X}_n has the same law as Ξ , where Ξ is a random vector whose components are linearly independent almost surely and that $\mathbb{E}[\Xi_i^2] < \infty$ for $i = 1, \dots, p$.

Remark 4.5.4. Under B', by the strong law of large numbers, we have $n^{-1}\mathbf{X}'_n\mathbf{X}_n \xrightarrow{\text{a.s.}} \mathbf{C}$, where $\mathbf{C} = (C_{ij})_{ij}$ with $C_{ij} = \mathbb{E}[\Xi_i\Xi_j]$. Moreover, \mathbf{C} is positive definite if and only if the random variables (Ξ_1, \dots, Ξ_p) are linearly independent almost surely. Indeed, for $\mathbf{t} \in \mathbb{R}^p$ we have $\mathbf{t}'\mathbf{C}\mathbf{t} = \mathbb{E}[(\sum_{i=1}^p t_i\Xi_i)^2] > 0$ if and only if $\sum_{i=1}^p t_i\Xi_i \neq 0$ a.s. for all $\mathbf{t} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$.

Since B' ensures that (4.6.8) holds a.s., it also implies that for large n , there almost surely exists a unique SLOPE minimizer. We denote this element by $\hat{\beta}_n^{\text{SLOPE}}$.

Theorem 4.5.5. Under A, B' and C, assume that a sequence $(\alpha_n)_n$ satisfies

$$\lim_{n \rightarrow \infty} \frac{\alpha_n}{n} = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\alpha_n}{\sqrt{n \log \log n}} = \infty. \quad (4.5.7)$$

If (4.5.5) holds, then the sequence $(\hat{\beta}_n^{\text{SLOPE}})_n$ recovers almost surely the pattern of β asymptotically, *i.e.*

$$\text{patt}(\hat{\beta}_n^{\text{SLOPE}}) \xrightarrow{\text{a.s.}} \text{patt}(\beta). \quad (4.5.8)$$

Remark 4.5.6. Assume that (4.5.5) is satisfied and set $\alpha_n = c\sqrt{n \log \log n}$ for $c > 0$. Then (4.5.7) is not satisfied and the probability that the correct SLOPE pattern is not recovered is greater than zero. In Section 4.6.1 we present more refined results on the strong consistency of the SLOPE pattern. The $\log \log n$ correction in (4.5.7) comes from the law of iterated logarithm.

4.6 Strong consistency of SLOPE and its pattern

4.6.1 Refined results on strong consistency of the SLOPE pattern

In this section we aim to give assumptions on the design matrix that are weaker than condition B', but they ensure the almost sure convergence of the pattern of $\hat{\beta}_n^{\text{SLOPE}}$.

- A'. $\varepsilon_n = (\varepsilon_1, \dots, \varepsilon_n)'$, where $(\varepsilon_i)_i$ are independent random variables such that

$$\mathbb{E}[\varepsilon_n] = 0 \quad \text{and} \quad \text{Var}(\varepsilon_n) = \sigma^2 \quad \text{for all } n, \quad \text{and} \quad \sup_n \mathbb{E}[|\varepsilon_n|^r] < \infty \quad (4.6.1)$$

for some $r > 2$.

B". A sequence of design matrices $\mathbf{X}_1, \mathbf{X}_2, \dots$ satisfies the condition

$$\frac{1}{n} \mathbf{X}'_n \mathbf{X}_n \xrightarrow{a.s.} \mathbf{C}, \quad (4.6.2)$$

where \mathbf{C} is a deterministic positive definite symmetric $p \times p$ matrix.

With $\mathbf{X}_n = (X_{ij}^{(n)})_{ij}$,

$$\lim_{n \rightarrow \infty} \frac{(\log n)^\rho}{\sqrt{n}} \sup_{i,j} |X_{ij}^{(n)}| = 0 \quad \text{a.s. for all } \rho > 0 \quad (4.6.3)$$

and there exist nonnegative random variables $(c_i)_i$, constants $d > 2/r$ and $m_0 \in \mathbb{N}$ such that for $n > m \geq m_0$,

$$\sup_j \left[\sum_{i=1}^m (X_{ij}^{(n)} - X_{ij}^{(m)})^2 + \sum_{i=m+1}^n (X_{ij}^{(n)})^2 \right] \leq \left(\sum_{i=m+1}^n c_i \right)^d \quad \text{a.s.}, \quad (4.6.4)$$

$$\left(\sum_{i=m_0}^n c_i \right)^d = O(n) \quad \text{a.s.} \quad (4.6.5)$$

C. $(\mathbf{X}_n)_n$ and $(\varepsilon_n)_n$ are independent.

We note that conditions (4.6.3) and (4.6.4) are trivially satisfied in the i.i.d. rows setting of Remark 4.5.4 or Assumption B'. The main ingredient of the proof of the strong pattern consistency is the law of iterated logarithm (4.9.9) which holds trivially under B'. Below we establish the same result under more general B". The technical assumption (4.6.4) is a kind of weak continuity assumption on the rows of \mathbf{X}_n as it says that the ℓ_2 distance between j^{th} rows of \mathbf{X}_n and \mathbf{X}_m should not be too large.

Lemma 4.6.1. *Assume A', B" and C. Then*

$$\limsup_{n \rightarrow \infty} \frac{\|\mathbf{X}'_n \varepsilon_n\|_\infty}{\sqrt{n \log \log n}} < \infty \quad \text{a.s.} \quad (4.6.6)$$

Proof. In view of (4.6.8) we have for $j = 1, \dots, p$,

$$n^{-1} A_n^{(j)} := n^{-1} \sum_{i=1}^n (X_{ij}^{(n)})^2 = (n^{-1} \mathbf{X}'_n \mathbf{X}_n)_{jj} \xrightarrow{a.s.} C_{jj} > 0. \quad (4.6.7)$$

We apply the general law of iterated logarithm for weights forming a triangular array from [117]. The result follows directly from [117, Theorem 1], which we recall in Theorem 2.3.18. Defining $a_{ni}^{(j)} := X_{ij}^{(n)}$ for $i = 1, \dots, n$, $j = 1, \dots, p$, $n \geq 1$ and 0 otherwise, we have

$$(\mathbf{X}'_n \varepsilon_n)_j = \sum_{i=-\infty}^{\infty} a_{ni}^{(j)} \varepsilon_i$$

and therefore we fall within the framework of (2.3.5). Then, (4.6.1), (4.6.3), (4.6.4) and (4.6.5) coincide with (2.3.4), (2.3.6), (2.3.7) and (2.3.8), respectively. Let $\mathbb{P}(\cdot | (\mathbf{X}_n)_n)$ be a regular

conditional probability. Then, applying Theorem 2.3.18 on the probability space $(\Omega, \mathcal{F}, \mathbb{P}_{\mathbf{X}})$ to our sequence, we obtain that for $j = 1, \dots, n$,

$$\mathbb{P} \left(\limsup_{n \rightarrow \infty} \frac{|(\mathbf{X}'_n \boldsymbol{\varepsilon}_n)_j|}{\sqrt{2A_n^{(j)} \log \log A_n^{(j)}}} \leq \sigma \mid (\mathbf{X}_n)_n \right) = 1 \quad \text{a.s.}$$

Averaging over $(\mathbf{X}_n)_n$ and using (4.6.7) again, we obtain the assertion. \square

Theorem 4.6.1. *Assume A', B'' and C. Suppose that $(\alpha_n)_n$ satisfies*

$$\lim_{n \rightarrow \infty} \frac{\alpha_n}{n} = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\alpha_n}{\sqrt{n \log \log n}} = \infty.$$

If (4.5.5) is satisfied, then $\mathbf{patt}(\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}) \xrightarrow{\text{a.s.}} \mathbf{patt}(\boldsymbol{\beta})$.

Comments:

- (a) Under reasonable assumptions (see *e.g.* [117, Theorem 1 (iii)]), one can show that

$$\limsup_{n \rightarrow \infty} \frac{\|\mathbf{X}'_n \boldsymbol{\varepsilon}_n\|_{\infty}}{\sqrt{n \log \log n}} > 0 \quad \text{a.s.}$$

Since $\alpha_n^{-1} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \xrightarrow{\text{a.s.}} \mathbf{0}$ is necessary for the a.s. pattern recovery, we can show that the condition $\alpha_n / \sqrt{n \log \log n} \rightarrow \infty$ cannot be weakened. Thus, the gap between the convergence in probability and the a.s. convergence is integral to the problem and in general cannot be reduced.

- (b) One can relax assumption B'' by imposing stronger conditions on the error $\boldsymbol{\varepsilon}_n$. For example, if $\boldsymbol{\varepsilon}_n$ is Gaussian, then one can use results from [167]. We note that [167] offers a very similar result as [117], but their assumptions are not quite comparable, see [167, Section 3 i)] for detailed discussion.
- (c) For Gaussian errors, one can consider a more general setting where one does not assume any relation between $\boldsymbol{\varepsilon}_n$ and $\boldsymbol{\varepsilon}_{n+1}$, *i.e.* the error need not be incremental. For orthogonal design such approach was taken in [165]. It is proved there that one obtains the a.s. SLOPE pattern consistency with the second limit condition of Theorem 4.6.1 replaced by $\lim_{n \rightarrow \infty} \alpha_n / \sqrt{n \log n} = \infty$. This result can be generalized to non-orthogonal designs.

4.6.2 Strong consistency of the SLOPE estimator

Lemma 4.6.2. *Let $\boldsymbol{\varepsilon}_n = (\varepsilon_1, \dots, \varepsilon_n)'$ with ε_i i.i.d., centered and having finite variance. Suppose that*

$$\frac{1}{n} \mathbf{X}'_n \mathbf{X}_n \xrightarrow{\text{a.s.}} \mathbf{C} > 0. \quad (4.6.8)$$

and that $(\boldsymbol{\varepsilon}_n)_n$ and $(\mathbf{X}_n)_n$ are independent. Then $n^{-1} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \xrightarrow{\text{a.s.}} \mathbf{0}$.

Proof. Denote by $\mathbb{P}(\cdot \mid (\mathbf{X}_n)_n)$ the regular conditional probability. By [52, Th. 1.1] applied to the sequence $(n^{-1} \mathbf{X}'_n \boldsymbol{\varepsilon}_n)_j$ on the probability space $(\Omega, \mathcal{F}, \mathbb{P}(\cdot \mid (\mathbf{X}_n)_n))$, we obtain

$$\mathbb{P} \left(\lim_{n \rightarrow \infty} n^{-1} (\mathbf{X}'_n \boldsymbol{\varepsilon}_n)_j = 0 \mid (\mathbf{X}_n)_n \right) = 1, \quad j = 1, \dots, p, \quad \text{a.s.}$$

Thus, applying the expectation to both sides above we obtain the assertion. \square

Theorem 4.6.2. Assume that $\mathbf{Y}_n = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}_n$, where $\boldsymbol{\beta} \in \mathbb{R}^p$, $\boldsymbol{\varepsilon}_n = (\varepsilon_1, \dots, \varepsilon_n)'$ with ε_i i.i.d., centered and having finite variance. Suppose (4.6.8) and that $(\boldsymbol{\varepsilon}_n)_n$ and $(\mathbf{X}_n)_n$ are independent. Let $\boldsymbol{\Lambda}_n = (\lambda_1^{(n)}, \dots, \lambda_p^{(n)})'$. Then, for large n , $S_{\mathbf{X}_n, \boldsymbol{\Lambda}_n}(\mathbf{Y}_n) = \{\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}\}$ almost surely. If $\boldsymbol{\beta} \neq \mathbf{0}$, then

$$\left(\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \boldsymbol{\beta} \right) \iff \left(\lim_{n \rightarrow \infty} \frac{\lambda_1^{(n)}}{n} = 0 \right). \quad (4.6.9)$$

If $\boldsymbol{\beta} = \mathbf{0}$ and $\lim_{n \rightarrow \infty} \frac{\lambda_1^{(n)}}{n} = 0$, then $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \mathbf{0}$.

Proof of Theorem 4.6.2. The assumption (4.6.8) implies that the matrix $\mathbf{X}'_n \mathbf{X}_n$ is positive definite for large n almost surely and hence ensures that $\ker(\mathbf{X}_n) = \{\mathbf{0}\}$. It implies the uniqueness of the SLOPE solution.

By Proposition 4.2.1, $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ is the SLOPE estimator of $\boldsymbol{\beta}$ in a linear regression model $\mathbf{Y}_n = \mathbf{X}_n \boldsymbol{\beta} + \boldsymbol{\varepsilon}_n$ if and only if for $\boldsymbol{\pi}_n = \mathbf{X}'_n (\mathbf{Y}_n - \mathbf{X}_n \hat{\boldsymbol{\beta}}_n^{\text{SLOPE}})$ we have

$$J_{\boldsymbol{\Lambda}}^*(\boldsymbol{\pi}_n) \leq 1 \quad (4.6.10)$$

and

$$\mathbf{U}'_{\mathbf{M}_n} \boldsymbol{\pi}_n = \tilde{\boldsymbol{\Lambda}}_n, \quad (4.6.11)$$

where $\mathbf{M}_n = \text{patt}(\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}})$ and $\tilde{\boldsymbol{\Lambda}}_n = \mathbf{U}'_{|\mathbf{M}_n|} \boldsymbol{\Lambda}_n$. By the definition of $\boldsymbol{\pi}_n$ we have

$$\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} = (\mathbf{X}'_n \mathbf{X}_n)^{-1} \mathbf{X}'_n \mathbf{Y}_n - (\mathbf{X}'_n \mathbf{X}_n)^{-1} \boldsymbol{\pi}_n = \hat{\boldsymbol{\beta}}_n^{\text{OLS}} - \left(\frac{1}{n} \mathbf{X}'_n \mathbf{X}_n \right)^{-1} \left(\frac{1}{n} \boldsymbol{\pi}_n \right).$$

Since in our setting $\hat{\boldsymbol{\beta}}_n^{\text{OLS}}$ is strongly consistent, $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \boldsymbol{\beta}$ if and only if $(n^{-1} \mathbf{X}'_n \mathbf{X}_n)^{-1} (n^{-1} \boldsymbol{\pi}_n) \xrightarrow{\text{a.s.}} \mathbf{0}$. In view of (4.6.8), we have $(n^{-1} \mathbf{X}'_n \mathbf{X}_n)^{-1} (n^{-1} \boldsymbol{\pi}_n) \xrightarrow{\text{a.s.}} \mathbf{0}$ if and only if $n^{-1} \boldsymbol{\pi}_n \xrightarrow{\text{a.s.}} \mathbf{0}$.

Assume $n^{-1} \lambda_1^{(n)} \rightarrow 0$. By (4.6.10) we have $\|\boldsymbol{\pi}_n\|_\infty \leq \lambda_1^{(n)}$, which gives

$$\left\| \frac{\boldsymbol{\pi}_n}{n} \right\|_\infty \leq \frac{\lambda_1^{(n)}}{n} \rightarrow \mathbf{0}.$$

Therefore, (4.6.9) implies that $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} \xrightarrow{\text{a.s.}} \boldsymbol{\beta}$.

Now assume that $\boldsymbol{\beta} \neq \mathbf{0}$ and $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}$ is strongly consistent, i.e. $n^{-1} \boldsymbol{\pi}_n \xrightarrow{\text{a.s.}} \mathbf{0}$. Then, (4.6.11) gives

$$p \|\boldsymbol{\pi}_n\|_\infty \geq \|\mathbf{U}'_{\mathbf{M}_n} \boldsymbol{\pi}_n\|_\infty = \|\tilde{\boldsymbol{\Lambda}}_n\|_\infty \geq \lambda_1^{(n)} \quad (4.6.12)$$

provided $\mathbf{M}_n \neq \mathbf{0}$. Applying (4.6.10) for $\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}} = \mathbf{0}$, we note that $\mathbf{M}_n(\omega) = \mathbf{0}$ if and only if

$$J_{n^{-1} \boldsymbol{\Lambda}_n}^* \left(n^{-1} \mathbf{X}_n(\omega)' \mathbf{Y}_n(\omega) \right) \leq 1.$$

In view of Lemma 4.6.2, it can be easily verified that $n^{-1} \mathbf{X}'_n \mathbf{Y}_n \xrightarrow{\text{a.s.}} \mathbf{C} \boldsymbol{\beta}$. Since

$$\left\| \frac{1}{n} \boldsymbol{\pi}_n \right\|_\infty \geq \left\| \frac{1}{n} \boldsymbol{\pi}_n \right\|_\infty \mathbb{1}_{\{\mathbf{M}_n = \mathbf{0}\}} = \left\| \frac{1}{n} \mathbf{X}'_n \mathbf{Y}_n \right\|_\infty \mathbb{1}_{\{\mathbf{M}_n = \mathbf{0}\}},$$

we see that for $\boldsymbol{\beta} \neq \mathbf{0}$, we have $\mathbf{M}_n \neq \mathbf{0}$ for large n almost surely. Thus, for $\boldsymbol{\beta} \neq \mathbf{0}$ we eventually obtain for large n

$$\frac{\lambda_1^{(n)}}{n} \leq p \left\| \frac{\boldsymbol{\pi}_n}{n} \right\|_\infty \quad \text{a.s.}$$

□

4.7 Simulation study

Our simulations aim at illustrating Theorems 4.5.1 and 4.5.3 and at showing that the results provided in these theorems are somehow unified. We consider the linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ has i.i.d. $N(0, 1)$ entries. Up to a constant, we choose components of $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_p)'$ as expected values of ordered standard Gaussian statistics. Let $V_{(1)} \geq \dots \geq V_{(p)}$ be ordered statistics of i.i.d. $N(0, 1)$ random variables. An approximation of $\mathbb{E}[V_{(i)}]$ for some $i \in \{1, \dots, p\}$, denoted $E(i, p)$, is given below (see [99] and references therein)

$$E(i, p) = -\Phi^{-1} \left(\frac{i - 0.375}{p + 1 - 0.750} \right),$$

where Φ is the cumulative distribution function of a $N(0, 1)$ random variable. We set $\lambda_i = E(i, p) + E(p - 1, p) - 2E(p, p)$ for $i = 1, \dots, p$. Note that since $E(1, p) > \dots > E(p, p)$, then $\boldsymbol{\Lambda} = (\lambda_1, \dots, \lambda_p)' \in \mathbb{R}^{p+}$.

For the design matrix \mathbf{X} and the vector of regression coefficients $\boldsymbol{\beta}$ we consider two cases:

- \mathbf{X} is orthogonal, and the components of $\boldsymbol{\beta}$ are all equal with a magnitude that tends to infinity.
- \mathbf{X} is asymptotically orthogonal, the sample size diverges to infinity and the components of $\boldsymbol{\beta}$ are equal to 1.

4.7.1 Sharp upper bound when \mathbf{X} is orthogonal

In Figure 4.10 we have $p = 100$, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is orthogonal (*i.e.* $\mathbf{X}'\mathbf{X} = \mathbf{I}_{100}$) and $\beta_1 = \dots = \beta_p = c > 0$. To compute the value $\alpha_{0.95}$ of the scaling parameter for which the upper bound equals 0.95, we note that $\boldsymbol{\pi}_\alpha$ is a Gaussian vector of $N(\mathbf{X}'(\widetilde{\mathbf{X}}_M')^+ \widetilde{\boldsymbol{\Lambda}}_M, \alpha^{-2} \mathbf{X}'(\mathbf{I} - \widetilde{\mathbf{X}}_M \widetilde{\mathbf{X}}_M^+) \mathbf{X})$ distribution. Moreover, since $\mathbf{M} = \mathbf{patt}(\boldsymbol{\beta}) = (1, \dots, 1)'$, we have

$$\begin{aligned} \mathbf{X}'(\widetilde{\mathbf{X}}_M')^+ \widetilde{\boldsymbol{\Lambda}}_M &= \left(\frac{1}{p} \sum_{i=1}^p \lambda_i, \dots, \frac{1}{p} \sum_{i=1}^p \lambda_i \right) \text{ and} \\ \mathbf{X}'(\mathbf{I}_n - \widetilde{\mathbf{X}}_M \widetilde{\mathbf{X}}_M^+) \mathbf{X} &= \begin{pmatrix} 1 - 1/p & -1/p & \dots & -1/p \\ -1/p & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & -1/p \\ -1/p & \dots & -1/p & 1 - 1/p \end{pmatrix}. \end{aligned} \quad (4.7.1)$$

Since the distribution of $\boldsymbol{\pi}_\alpha$ is given and the open SLOPE irrepresentability condition occurs, for arbitrary $q \in (0, 1)$ one can choose α_q for which $\mathbb{P}(J_{\boldsymbol{\Lambda}}^*(\boldsymbol{\pi}_{\alpha_q}) \leq 1) = q$.

In the following graph, $q = 0.95$ and $\alpha_{0.95} \approx 1.391$.

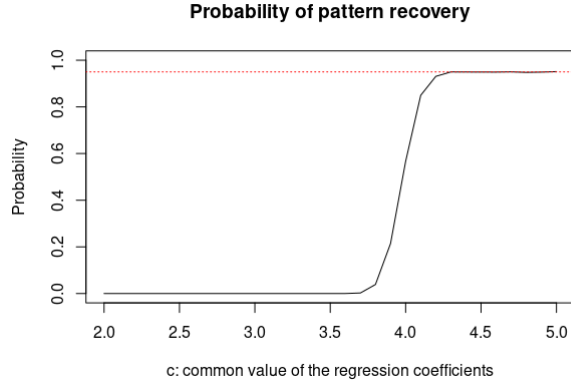


Figure 4.10: The estimates of probability of pattern recovery by SLOPE as a function of c where $c = \beta_1 = \dots = \beta_{100} > 0$. For each point the probability is computed via 10^5 Monte-Carlo experiments. The scaling parameter $\alpha_{0.95} \approx 1.391$ is chosen to fix the upper bound at 0.95. Note that when c is large, the probability of pattern recovery is approximately equal to 0.95.

4.7.2 Limiting probability when \mathbf{X} is asymptotically orthogonal

In Figure 4.11, $\mathbf{X} \in \mathbb{R}^{n \times 100}$ has i.i.d. $N(0, 1)$ entries, $\beta_1 = \dots = \beta_{100} = 1$ and $\alpha_{0.95} \approx 1.391$. Actually, since $n^{-1} \mathbf{X}' \mathbf{X}$ converges to \mathbf{I}_{100} , when $\mathbf{patt}(\boldsymbol{\beta}) = (1, \dots, 1)'$ the Gaussian vector involved in the limiting probability has the same mean and covariance as (4.7.1).

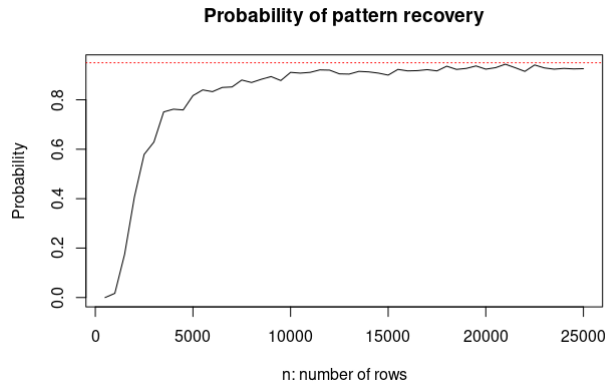


Figure 4.11: The estimates of probability of pattern recovery by SLOPE as a function of n . For each point the probability is computed via 10^3 Monte-Carlo experiments.

4.8 Discussion

This chapter makes an important step in understanding the clustering properties of SLOPE and we have shown that the irrepresentability condition provides theoretical guarantees for SLOPE pattern recovery. However, this by no means closes the topic of the SLOPE pattern recovery. Similarly to the irrepresentability condition for LASSO, SLOPE irrepresentability condition is rather stringent and imposes a strict restriction on the number of nonzero clusters in $\boldsymbol{\beta}$. On the other hand, in [174] it is shown that a much weaker condition for LASSO is required to separate

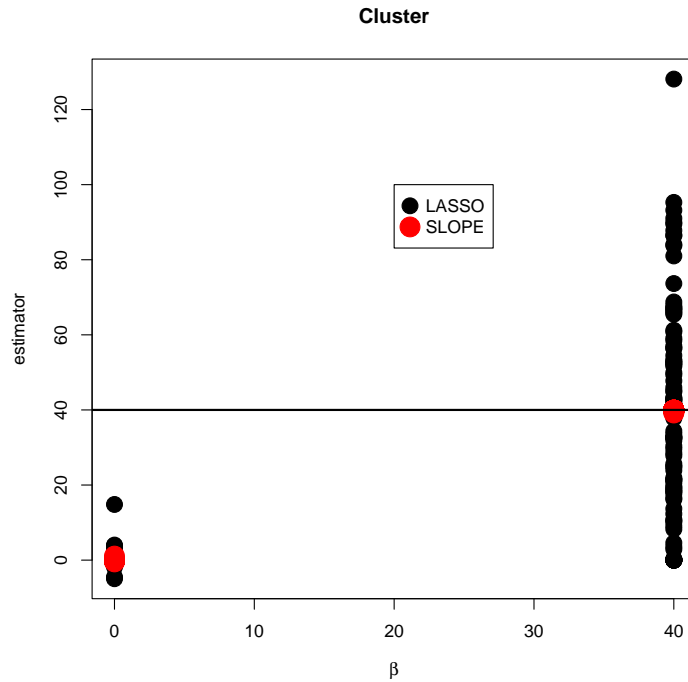


Figure 4.12: Comparison of LASSO and SLOPE when the cluster structure is present in the data. Here $n = 100$, $p = 200$, and the correlation between i^{th} and j^{th} column of \mathbf{X} is equal to $0.9048^{|i-j|}$. First $k = 100$ columns of \mathbf{X} are associated with \mathbf{Y} and their nonzero regression coefficient are all equal to 40. The SLOPE and LASSO irrepresentability conditions are not satisfied, but SLOPE, contrary to LASSO, satisfies the accessibility condition.

the estimators of the null components of β from the estimators of nonzero regression coefficients. This condition, called accessibility (also called identifiability), requires that the vector β has a minimal ℓ_1 norm among all vectors γ such that $\mathbf{X}\beta = \mathbf{X}\gamma$. Thus, when the accessibility condition is satisfied, one can recover the sign of β by thresholding LASSO estimates. Empirical results from [174] suggest that this weaker condition is also sufficient for the sign recovery by the adaptive LASSO [194]. In Chapter 5 it is shown that a similar result holds for SLOPE. In this case, the accessibility condition is satisfied if the vector β has the smallest sorted ℓ_1 norm among all vectors γ such that $\mathbf{X}\beta = \mathbf{X}\gamma$. In Chapter ?? it is shown that when the accessibility condition is satisfied then SLOPE properly ranks the estimators of regression coefficients and the SLOPE pattern can be recovered by shrinking similar estimates towards the cluster centers. Figure 4.12 illustrates this phenomenon and shows that the accessibility condition for SLOPE can be much less restrictive than the accessibility condition for LASSO. In this example, the matrix \mathbf{X} and the vector \mathbf{Y} are generated as in the example illustrated in Figure 4.1 and the only difference is that now first $k = 100 = n$ regression coefficients are all equal to 40. In this situation, the accessibility condition for LASSO is not satisfied and LASSO can not properly separate the null and nonzero regression coefficients. Also, despite the selection of the tuning parameter so as to minimize the squared estimation error, the precision of LASSO estimates is very poor. As far as SLOPE is concerned, the irrepresentability condition is not satisfied, but the accessibility condition holds. Thus, while SLOPE can not properly identify the pattern, it estimates β with such a good precision that the difference between the estimated and the true

pattern is hardly visible on the graph. These nice ranking and estimation properties of SLOPE bring a promise for efficient pattern recovery by appropriate thresholded and adaptive SLOPE versions.

4.9 Appendix — Proofs

4.9.1 Proof of Proposition 4.2.1

Note that if $\mathbf{M} = \mathbf{0}$, then the statement is valid by (4.2.1). Thus, we may later assume that $\mathbf{M} \neq \mathbf{0}$. To simplify the notation, we write $\tilde{\Lambda}$ instead of $\tilde{\Lambda}_{\mathbf{M}}$. The elements of $\tilde{\Lambda}$ are denoted by $\tilde{\Lambda}_l$, $l = 1, \dots, k$. Let $k = \|\mathbf{M}\|_{\infty}$. Before proving Proposition 4.2.1, note that, by assumption, there exists $\mathbf{s} \in \mathbb{R}^{k+}$ such that $\mathbf{b} = \mathbf{U}_{\mathbf{M}}\mathbf{s}$. Consequently, $|\mathbf{b}|_{\downarrow} = \mathbf{U}_{|\mathbf{M}|_{\downarrow}}\mathbf{s}$ and thus

$$J_{\Lambda}(\mathbf{b}) = \lambda_1|b|_{(1)} + \dots + \lambda_p|b|_{(p)} = \Lambda'\mathbf{U}_{|\mathbf{M}|_{\downarrow}}\mathbf{s} = \tilde{\Lambda}'\mathbf{s} = s_1\tilde{\Lambda}_1 + \dots + s_k\tilde{\Lambda}_k.$$

Moreover, with $p_l = |\{i: |m_i| \geq k+1-l\}|$, we have $\tilde{\Lambda}_l = \lambda_{p_{l-1}+1} + \dots + \lambda_{p_l}$, $l = 1, \dots, k$.

Proof of Proposition 4.2.1.

First, we prove the inclusion $\partial J_{\Lambda}(\mathbf{b}) \subset \{\mathbf{v} \in \mathbb{R}^p: J_{\Lambda}^*(\mathbf{v}) \leq 1 \text{ and } \mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}\}$. Let $\mathbf{v} \in \partial J_{\Lambda}(\mathbf{b})$. Since $J_{\Lambda}^*(\mathbf{v}) \leq 1$ (see (4.2.1)) then, by definition of the dual sorted ℓ_1 norm, for all $j \in \{1, 2, \dots, p\}$ we have $\sum_{i=1}^j |v|_{(i)} \leq \sum_{i=1}^j \lambda_i$. It remains to prove that $\mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}$. For all $l \in \{1, \dots, k\}$ we have the following inequality

$$\sum_{i=1}^l [\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_i = \sum_{i: |m_i| \geq k+1-l} \text{sign}(m_i)v_i \leq \sum_{i: |m_i| \geq k+1-l} |v_i| \leq \sum_{i=1}^{p_l} |v|_{(i)} \leq \sum_{i=1}^{p_l} \lambda_i = \sum_{i=1}^l \tilde{\Lambda}_i. \quad (4.9.1)$$

Note that

$$\begin{aligned} \mathbf{b}'\mathbf{v} &= (\mathbf{U}_{\mathbf{M}}\mathbf{s})'\mathbf{v} = \sum_{i=1}^k s_i [\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_i = \sum_{l=1}^{k-1} (s_l - s_{l+1}) \sum_{i=1}^l [\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_i + s_k \sum_{i=1}^k [\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_i \\ &\leq \sum_{l=1}^{k-1} (s_l - s_{l+1}) \sum_{i=1}^l \tilde{\Lambda}_i + s_k \sum_{i=1}^k \tilde{\Lambda}_i = \sum_{l=1}^k s_l \tilde{\Lambda}_l = J_{\Lambda}(\mathbf{b}). \end{aligned}$$

Moreover, since $\mathbf{v} \in \partial J_{\Lambda}(\mathbf{b})$, we have $\mathbf{b}'\mathbf{v} = J_{\Lambda}(\mathbf{b})$ (see (4.2.1)). Therefore,

$$\sum_{i=1}^l [\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_i = \sum_{i=1}^l \tilde{\Lambda}_i \quad \text{for } l = 1, \dots, k$$

and thus the inequalities given in (4.9.1) are the equalities. Thus

$$[\mathbf{U}'_{\mathbf{M}}\mathbf{v}]_l = \tilde{\Lambda}_l \quad \text{for } l = 1, \dots, k$$

and hence that $\mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}$.

Now we prove the other inclusion, $\partial J_{\Lambda}(\mathbf{b}) \supset \{\mathbf{v} \in \mathbb{R}^p: J_{\Lambda}^*(\mathbf{v}) \leq 1 \text{ and } \mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}\}$. Assume that $\mathbf{v} \in \mathbb{R}^p$ satisfies $J_{\Lambda}^*(\mathbf{v}) \leq 1$ and $\mathbf{U}'_{\mathbf{M}}\mathbf{v} = \tilde{\Lambda}$. To prove that $\mathbf{v} \in \partial J_{\Lambda}(\mathbf{b})$, it remains to establish that $\mathbf{b}'\mathbf{v} = J_{\Lambda}(\mathbf{b})$ (see (4.2.1)). Since $\mathbf{b} = \mathbf{U}_{\mathbf{M}}\mathbf{s}$, we have

$$\mathbf{b}'\mathbf{v} = (\mathbf{U}_{\mathbf{M}}\mathbf{s})'\mathbf{v} = \mathbf{s}'\mathbf{U}'_{\mathbf{M}}\mathbf{v} = \mathbf{s}'\tilde{\Lambda} = J_{\Lambda}(\mathbf{b}).$$

□

4.9.2 Proof of Proposition 4.2.2

Lemma 4.9.1. *Let $\mathbf{\Lambda} \in \mathbb{R}^{p+}$ and $\mathbf{b} \in \mathbb{R}^p$. If $\mathbf{\Lambda} \in \partial J_{\mathbf{\Lambda}}(\mathbf{b})$, then $b_1 \geq \dots \geq b_p \geq 0$.*

Proof. Let us assume that $b_i < 0$ for some $i \in \{1, \dots, p\}$.

For $\check{\boldsymbol{\pi}} = (\lambda_1, \dots, \lambda_{i-1}, -\lambda_i, \lambda_{i+1}, \dots, \lambda_p)$ we have $J_{\mathbf{\Lambda}}^*(\check{\boldsymbol{\pi}}) \leq 1$ and one may deduce that

$$\mathbf{\Lambda}'\mathbf{b} < \check{\boldsymbol{\pi}}'\mathbf{b} \leq \max\{\boldsymbol{\pi}'\mathbf{b} : J_{\mathbf{\Lambda}}^*(\boldsymbol{\pi}) \leq 1\} = J_{\mathbf{\Lambda}}(\mathbf{b}).$$

Consequently $\mathbf{\Lambda} \notin \partial J_{\mathbf{\Lambda}}(\mathbf{b})$, which leads to a contradiction.

Now, let us assume that $b_i < b_j$ for some $1 \leq i < j \leq p$. We define $\check{\boldsymbol{\pi}}$ as the following:

$$\check{\boldsymbol{\pi}}_k = \begin{cases} \lambda_k & \text{if } k \neq i, k \neq j, \\ \lambda_j & \text{if } k = i, \\ \lambda_i & \text{if } k = j, \end{cases} \quad k = 1, \dots, p.$$

Note that $J_{\mathbf{\Lambda}}^*(\check{\boldsymbol{\pi}}) \leq 1$. Since $\lambda_i > \lambda_j$, by the rearrangement inequality we have $\lambda_i b_i + \lambda_j b_j < \lambda_j b_i + \lambda_i b_j$. Thus, one may deduce the following inequality

$$\mathbf{\Lambda}'\mathbf{b} < \check{\boldsymbol{\pi}}'\mathbf{b} \leq \max\{\boldsymbol{\pi}'\mathbf{b} : \boldsymbol{\pi} \in \mathbb{R}^p, J_{\mathbf{\Lambda}}^*(\boldsymbol{\pi}) \leq 1\} = J_{\mathbf{\Lambda}}(\mathbf{b}).$$

Consequently, $\mathbf{\Lambda} \notin \partial J_{\mathbf{\Lambda}}(\mathbf{b})$, which again leads to a contradiction. \square

Let ψ be an orthogonal transformation defined by $\psi: \mathbf{x} \in \mathbb{R}^p \mapsto (v_1 b_{r(1)}, \dots, v_p b_{r(p)})$ where $v_1, \dots, v_p \in \{-1, 1\}$ and let r be a permutation on $\{1, \dots, p\}$. Before proving Proposition 4.2.2, let us recall that for any $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ we have $J_{\mathbf{\Lambda}}(\mathbf{b}) = J_{\mathbf{\Lambda}}(\psi(\mathbf{b}))$, $J_{\mathbf{\Lambda}}^*(\mathbf{b}) = J_{\mathbf{\Lambda}}^*(\psi(\mathbf{b}))$ and $\mathbf{b}'\mathbf{a} = \psi(\mathbf{b})'\psi(\mathbf{a})$, therefore $\partial J_{\mathbf{\Lambda}}(\psi(\mathbf{b})) = \psi(\partial J_{\mathbf{\Lambda}}(\mathbf{b}))$.

Proof of Proposition 4.2.2. If $\mathbf{patt}(\mathbf{a}) = \mathbf{patt}(\mathbf{b})$, then, according to Proposition 4.2.1, $\partial J_{\mathbf{\Lambda}}(\mathbf{a}) = \partial J_{\mathbf{\Lambda}}(\mathbf{b})$. Set $\mathbf{M} = \mathbf{patt}(\mathbf{a})$ and $\widetilde{\mathbf{M}} = \mathbf{patt}(\mathbf{b})$. It remains to prove that if $\partial J_{\mathbf{\Lambda}}(\mathbf{a}) = \partial J_{\mathbf{\Lambda}}(\mathbf{b})$, then $\mathbf{M} = \widetilde{\mathbf{M}}$. Since the subdifferential $\partial J_{\mathbf{\Lambda}}(\mathbf{a})$ depends on \mathbf{a} only through its pattern, then by Proposition 4.2.1 we have $\partial J_{\mathbf{\Lambda}}(\mathbf{a}) = \partial J_{\mathbf{\Lambda}}(\mathbf{M})$ and similarly $\partial J_{\mathbf{\Lambda}}(\mathbf{b}) = \partial J_{\mathbf{\Lambda}}(\widetilde{\mathbf{M}})$.

First let us assume that $\mathbf{M} = |\mathbf{M}|_{\downarrow}$ namely $M_1 \geq M_2 \geq \dots \geq M_p \geq 0$. In this case, $\mathbf{M}'\mathbf{\Lambda} = J_{\mathbf{\Lambda}}(\mathbf{M})$ and hence $\mathbf{\Lambda} = (\lambda_1, \dots, \lambda_p)' \in \partial J_{\mathbf{\Lambda}}(\mathbf{M})$. Since $\partial J_{\mathbf{\Lambda}}(\mathbf{M}) = \partial J_{\mathbf{\Lambda}}(\widetilde{\mathbf{M}})$, it follows from Lemma 4.9.1 that $\widetilde{M}_1 \geq \dots \geq \widetilde{M}_p \geq 0$, because $\mathbf{\Lambda} \in \partial J_{\mathbf{\Lambda}}(\widetilde{\mathbf{M}})$. To prove that $\mathbf{M} = \widetilde{\mathbf{M}}$, first let us establish that $M_p = \widetilde{M}_p = 0$ or $M_p = \widetilde{M}_p = 1$. If $M_p = 0$ and $\widetilde{M}_p = 1$, then we set $\check{\boldsymbol{\pi}} = (\lambda_1, \dots, \lambda_{p-1}, 0)'$, where $J_{\mathbf{\Lambda}}^*(\check{\boldsymbol{\pi}}) \leq 1$. Because

$$J_{\mathbf{\Lambda}}(\mathbf{M}) = \mathbf{\Lambda}'\mathbf{M} = \check{\boldsymbol{\pi}}'\mathbf{M} \text{ and } J_{\mathbf{\Lambda}}(\widetilde{\mathbf{M}}) = \mathbf{\Lambda}'\widetilde{\mathbf{M}} > \check{\boldsymbol{\pi}}'\widetilde{\mathbf{M}},$$

we have $\check{\boldsymbol{\pi}} \in \partial J_{\mathbf{\Lambda}}(\mathbf{M})$ and $\check{\boldsymbol{\pi}} \notin \partial J_{\mathbf{\Lambda}}(\widetilde{\mathbf{M}})$ which provides a contradiction. We proceed analogously for $M_p = 1$ and $\widetilde{M}_p = 0$. To achieve proving that $\mathbf{M} = \widetilde{\mathbf{M}}$, let us establish that $m_i = m_{i+1}$ and $\widetilde{m}_i = \widetilde{m}_{i+1}$ or $m_i > m_{i+1}$ and $\widetilde{m}_i > \widetilde{m}_{i+1}$. If $m_i = m_{i+1}$ and $\widetilde{m}_i > \widetilde{m}_{i+1}$ then we define $\check{\boldsymbol{\pi}}$ satisfying $J_{\mathbf{\Lambda}}^*(\check{\boldsymbol{\pi}}) \leq 1$, as follows

$$\check{\boldsymbol{\pi}}_k = \begin{cases} \lambda_k & \text{if } k \neq i, k \neq i+1, \\ \lambda_{i+1} & \text{if } k = i, \\ \lambda_i & \text{if } k = i+1, \end{cases} \quad k = 1, \dots, p.$$

Since $\lambda_i m_i + \lambda_{i+1} m_{i+1} = \lambda_{i+1} m_i + \lambda_i m_{i+1}$ and $\lambda_i \tilde{m}_i + \lambda_{i+1} \tilde{m}_{i+1} > \lambda_{i+1} \tilde{m}_i + \lambda_i \tilde{m}_{i+1}$ then

$$J_\Lambda(\mathbf{M}) = \Lambda' \mathbf{M} = \check{\pi}' \mathbf{M} \text{ and } J_\Lambda(\tilde{\mathbf{M}}) = \Lambda' \tilde{\mathbf{M}} > \check{\pi}' \tilde{\mathbf{M}}.$$

Consequently $\check{\pi} \in \partial J_\Lambda(\mathbf{M})$ and $\check{\pi} \notin \partial J_\Lambda(\tilde{\mathbf{M}})$, which provides a contradiction. We proceed analogously for $m_i > m_{i+1}$ and $\tilde{m}_i = \tilde{m}_{i+1}$. Finally, if $\mathbf{M} \neq |\mathbf{M}|_\downarrow$, then we may pick an orthogonal transformation ψ as defined above for which $\psi(\mathbf{M}) = |\mathbf{M}|_\downarrow$. Since $\partial J_\Lambda(\mathbf{M}) = \partial J_\Lambda(\tilde{\mathbf{M}})$ implies that $\partial J_\Lambda(\psi(\mathbf{M})) = \partial J_\Lambda(\psi(\tilde{\mathbf{M}}))$, then the first part of the proof establishes that $\psi(\tilde{\mathbf{M}}) = \psi(\mathbf{M})$ and thus $\mathbf{M} = \tilde{\mathbf{M}}$. \square

Recall that $J_\Lambda^*(\mathbf{x}) \leq 1$ if and only if

$$|x|_{(1)} + \dots + |x|_{(j)} \leq \lambda_1 + \dots + \lambda_j, \quad j = 1, \dots, p. \quad (4.9.2)$$

The following result follows from the proof of Proposition 4.2.1.

Proposition 4.9.1. *Assume $\mathbf{x} \in \mathbb{R}^p$ satisfies $J_\Lambda^*(\mathbf{x}) \leq 1$ and let $\mathbf{b} \in \mathbb{R}^p$. Then \mathbf{x} belongs to $\partial J_\Lambda(\mathbf{b})$ if and only if the following three conditions hold true:*

- (a) *If $b_i \neq 0$, then $\text{sign}(x_i) = \text{sign}(b_i)$,*
- (b) *If $|b_i| > |b_j|$ then $|x_i| \geq |x_j|$,*
- (c) *The equalities hold in (4.9.2) for $j \in \{n_1, n_2, \dots, n_k\}$, where $n_j = |\{i: |m_i| \geq k + 1 - j\}|$ with $(M_1, \dots, M_p)' = \mathbf{patt}(\mathbf{b})$.*

4.9.3 Proof of Theorem 4.3.1

Proof of Theorem 4.3.1. Necessity. Let us assume that there exists $\hat{\beta} \in S_{\mathbf{X}, \Lambda}(\mathbf{Y})$ with $\mathbf{patt}(\hat{\beta}) = \mathbf{M}$. Consequently, $\hat{\beta} = \mathbf{U}_\mathbf{M} \mathbf{s}$ for some $\mathbf{s} \in \mathbb{R}^{k+}$. By Proposition 4.2.2, $\mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta}) \in \partial J_\Lambda(\hat{\beta}) = \partial J_\Lambda(\mathbf{M})$. Multiplying this inclusion by $\mathbf{U}'_\mathbf{M}$, due to (4.2.2), we get $\tilde{\mathbf{X}}'_\mathbf{M}(\mathbf{Y} - \mathbf{X}\hat{\beta}) = \tilde{\Lambda}_\mathbf{M}$ and so

$$\tilde{\mathbf{X}}'_\mathbf{M} \mathbf{Y} - \tilde{\Lambda}_\mathbf{M} = \tilde{\mathbf{X}}'_\mathbf{M} \mathbf{X} \hat{\beta} = \tilde{\mathbf{X}}'_\mathbf{M} \tilde{\mathbf{X}}_\mathbf{M} \mathbf{s}. \quad (4.9.3)$$

The positivity condition is proven.

We apply $(\tilde{\mathbf{X}}'_\mathbf{M})^+$ from the left to (4.9.3) and use the fact that $\tilde{\mathbf{P}}_\mathbf{M} = (\tilde{\mathbf{X}}'_\mathbf{M})^+ \tilde{\mathbf{X}}'_\mathbf{M}$ is the projection onto $\text{col}(\tilde{\mathbf{X}}_\mathbf{M})$. Since $\mathbf{X}\hat{\beta} \in \text{col}(\tilde{\mathbf{X}}_\mathbf{M})$, we have $\tilde{\mathbf{P}}_\mathbf{M} \mathbf{X}\hat{\beta} = \mathbf{X}\hat{\beta}$. Thus,

$$\tilde{\mathbf{P}}_\mathbf{M} \mathbf{Y} - (\tilde{\mathbf{X}}'_\mathbf{M})^+ \tilde{\Lambda}_\mathbf{M} = \mathbf{X}\hat{\beta}.$$

The above equality gives the subdifferential condition:

$$\begin{aligned} \partial J_\Lambda(\mathbf{M}) \ni \mathbf{X}'(\mathbf{Y} - \mathbf{X}\hat{\beta}) &= \mathbf{X}'(\mathbf{Y} - (\tilde{\mathbf{P}}_\mathbf{M} \mathbf{Y} - (\tilde{\mathbf{X}}'_\mathbf{M})^+ \tilde{\Lambda}_\mathbf{M})) \\ &= \mathbf{X}'(\tilde{\mathbf{X}}'_\mathbf{M})^+ \tilde{\Lambda}_\mathbf{M} + \mathbf{X}'(\mathbf{I}_n - \tilde{\mathbf{P}}_\mathbf{M}) \mathbf{Y} = \boldsymbol{\pi}. \end{aligned} \quad (4.9.4)$$

Sufficiency. Assume that the positivity condition and the subdifferential conditions hold true. Then, by the positivity condition, one may pick $\mathbf{s} \in \mathbb{R}^{k+}$ for which

$$\tilde{\Lambda}_\mathbf{M} = \tilde{\mathbf{X}}'_\mathbf{M} \mathbf{Y} - \tilde{\mathbf{X}}'_\mathbf{M} \tilde{\mathbf{X}}_\mathbf{M} \mathbf{s}. \quad (4.9.5)$$

Let us show that $\mathbf{U}_M \mathbf{s} \in S_{\mathbf{X}, \Lambda}(\mathbf{Y})$. By definition of \mathbf{U}_M , we have $\mathbf{patt}(\mathbf{U}_M \mathbf{s}) = \mathbf{M}$, thus $\partial J_\Lambda(\mathbf{U}_M \mathbf{s}) = \partial J_\Lambda(\mathbf{M})$. Moreover, using (4.9.4) and (4.9.5) one may deduce that

$$\begin{aligned} \partial J_\Lambda(\mathbf{U}_M \mathbf{s}) \ni \boldsymbol{\pi} &= \mathbf{X}'(\mathbf{Y} - \tilde{\mathbf{P}}_M \mathbf{Y} + (\tilde{\mathbf{X}}'_M)^+ \tilde{\boldsymbol{\Lambda}}_M) \\ &= \mathbf{X}'(\mathbf{Y} - \tilde{\mathbf{P}}_M \mathbf{Y} + (\tilde{\mathbf{X}}'_M)^+ (\tilde{\mathbf{X}}_M \mathbf{Y} - \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{s})) \\ &= \mathbf{X}'(\mathbf{Y} - \mathbf{X} \mathbf{U}_M \mathbf{s}). \end{aligned}$$

Consequently, $\mathbf{U}_M \mathbf{s} \in S_{\mathbf{X}, \Lambda}(\mathbf{Y})$. \square

4.9.4 Proof of Corollary 4.3.1

Proof of Corollary 4.3.1. If SLOPE recovers the pattern of $\boldsymbol{\beta}$ in the noiseless case, then, by Theorem 4.3.1, the subdifferential condition reads as: $\mathbf{X}'(\tilde{\mathbf{X}}'_M)^+ \tilde{\boldsymbol{\Lambda}}_M \in \partial J_\Lambda(\mathbf{M})$.

Conversely, if $\mathbf{X}'(\tilde{\mathbf{X}}'_M)^+ \tilde{\boldsymbol{\Lambda}}_M \in \partial J_\Lambda(\mathbf{M})$ then, by Theorem 4.3.1, it remains to show that the positivity condition occurs for $\alpha > 0$ small enough. Since $\boldsymbol{\beta} = \mathbf{U}_M \mathbf{s}$ for some $\mathbf{s} \in \mathbb{R}^{k+}$, where $k = \|\mathbf{M}\|_\infty$, we have

$$\tilde{\mathbf{X}}'_M \mathbf{Y} - \alpha \tilde{\boldsymbol{\Lambda}}_M = \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{s} - \alpha \tilde{\boldsymbol{\Lambda}}_M.$$

Therefore for $\alpha > 0$ small enough, $\tilde{\mathbf{X}}'_M \mathbf{Y} - \alpha \tilde{\boldsymbol{\Lambda}}_M \in \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbb{R}^{k+}$ and thus, the positivity condition is proven. \square

4.9.5 Proof of Theorem 4.5.1

Lemma 4.9.2. *Let $\mathbf{0} \neq \mathbf{b} \in \mathbb{R}^p$ and $\mathbf{M} = \mathbf{patt}(\mathbf{b})$. Then the smallest affine space containing $\partial J_\Lambda(\mathbf{b})$ is $\text{aff}(\partial J_\Lambda(\mathbf{b})) = \{\mathbf{v} \in \mathbb{R}^p : \mathbf{U}'_M \mathbf{v} = \tilde{\boldsymbol{\Lambda}}_M\}$.*

Proof. According to Proposition 4.2.1 we have

$$\text{aff}(\partial J_\Lambda(\mathbf{b})) \subset \{\mathbf{v} \in \mathbb{R}^p : \mathbf{U}'_M \mathbf{v} = \tilde{\boldsymbol{\Lambda}}_M\}.$$

Moreover, according to Theorem 4 in [156] we have

$$\dim(\text{aff}(\partial J_\Lambda(\mathbf{b}))) = \|\mathbf{M}\|_\infty = \dim(\{\mathbf{v} \in \mathbb{R}^p : \mathbf{U}'_M \mathbf{v} = \tilde{\boldsymbol{\Lambda}}_M\}),$$

which achieves the proof. \square

Proof of Theorem 4.5.1. (i) Sharpness of the upper bound. According to Theorem 4.3.1, the pattern recovery by SLOPE is equivalent to have simultaneously the positivity condition and the subdifferential condition satisfied. The upper bound (4.5.2) coincides with the probability of the subdifferential condition. Thus to prove that this upper bound is sharp, it remains to show that the probability of the positivity condition tends to 1 when r tends to ∞ . Clearly the upper bound is reached when $\tilde{\boldsymbol{\Lambda}}_M \notin \text{col}(\tilde{\mathbf{X}}'_M)$, thus we may assume that $\tilde{\boldsymbol{\Lambda}}_M \in \text{col}(\tilde{\mathbf{X}}'_M)$. Recall that $\boldsymbol{\beta}^{(r)} = \mathbf{U}_M \mathbf{s}^{(r)}$ for $\mathbf{s}^{(r)} \in \mathbb{R}^{k+}$ and thus $\tilde{\mathbf{X}}'_M \mathbf{Y}^{(r)} = \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{s}^{(r)} + \tilde{\mathbf{X}}'_M \boldsymbol{\varepsilon}$. As $\tilde{\mathbf{X}}'_M (\tilde{\mathbf{X}}'_M)^+ = \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M (\tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M)^+$ is the projection on $\text{col}(\tilde{\mathbf{X}}'_M)$, we obtain

$$\begin{aligned} \tilde{\mathbf{X}}'_M \mathbf{Y}^{(r)} - \alpha_r \tilde{\boldsymbol{\Lambda}}_M &= \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{s}^{(r)} - \alpha_r \tilde{\boldsymbol{\Lambda}}_M + \tilde{\mathbf{X}}'_M \boldsymbol{\varepsilon} \\ &= \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \mathbf{s}^{(r)} - \alpha_r \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M (\tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M)^+ \tilde{\boldsymbol{\Lambda}}_M + \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M (\tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M)^+ \tilde{\boldsymbol{\Lambda}}_M + \tilde{\mathbf{X}}'_M \boldsymbol{\varepsilon} \\ &= \tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M \Delta_r \left(\frac{1}{\Delta_r} \mathbf{s}^{(r)} - \frac{\alpha_r}{\Delta_r} (\tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M)^+ \tilde{\boldsymbol{\Lambda}}_M + \frac{1}{\Delta_r} (\tilde{\mathbf{X}}'_M \tilde{\mathbf{X}}_M)^+ \tilde{\boldsymbol{\Lambda}}_M + \tilde{\mathbf{X}}'_M \boldsymbol{\varepsilon} \right). \end{aligned}$$

Note that by the assumption on Δ_r :

- the vector $\mathbf{s}^{(r)}/\Delta_r \in \mathbb{R}^{k+}$ is component-wise larger than or equal to $(k, \dots, 1)$,
- $\lim_{r \rightarrow \infty} \alpha_r/\Delta_r = 0$ and $\lim_{r \rightarrow \infty} 1/\Delta_r = 0$.

Consequently, for r large enough we have

$$\widetilde{\mathbf{X}}'_M \mathbf{Y}^{(r)} - \alpha_r \widetilde{\Lambda}_M \in \widetilde{\mathbf{X}}'_M \widetilde{\mathbf{X}}_M \mathbb{R}^{k+}.$$

Since this fact is true for any realization of $\boldsymbol{\varepsilon}$, we get

$$\lim_{r \rightarrow \infty} \mathbb{P} \left(\widetilde{\mathbf{X}}'_M \mathbf{Y}^{(r)} - \alpha_r \widetilde{\Lambda}_M \in \widetilde{\mathbf{X}}'_M \widetilde{\mathbf{X}}_M \mathbb{R}^{k+} \right) = 1.$$

(ii) *Pattern consistency.* In the proof of the previous part, we see that positivity condition occurs when r is sufficiently large. Thus it remains to prove that subdifferential condition occurs as $r \rightarrow \infty$ when $\mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M \in \text{ri}(\partial J_\Lambda(\mathbf{M}))$. First we observe that

$$\mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M + \frac{1}{\alpha_r} \mathbf{X}'(\mathbf{I}_n - \widetilde{P}_M) \boldsymbol{\varepsilon} \xrightarrow{r \rightarrow \infty} \mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M. \quad (4.9.6)$$

Note by Lemma 4.9.2 that $\mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M + \alpha_r^{-1} \mathbf{X}'(\mathbf{I}_n - \widetilde{P}_M) \boldsymbol{\varepsilon} \in \text{aff}(\partial J_\Lambda(\mathbf{M}))$. Indeed, since $\widetilde{\Lambda}_M \in \text{col}(\widetilde{\mathbf{X}}'_M)$, we have

$$\underbrace{U'_M \mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M}_{=\widetilde{\Lambda}_M} + \frac{1}{\alpha_r} \underbrace{U'_M \mathbf{X}'(\mathbf{I}_n - \widetilde{P}_M) \boldsymbol{\varepsilon}(\omega)}_{=0} = \widetilde{\Lambda}_M.$$

The second term above is zero due to the fact that $(\mathbf{I}_n - \widetilde{P}_M)$ is an orthogonal projection onto $\text{col}(\widetilde{\mathbf{X}}'_M)^\perp$. When $\mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M \in \text{ri}(\partial J_\Lambda(\mathbf{M}))$, due to (4.9.6), one may deduce that for sufficiently large r we have

$$\mathbf{X}'(\widetilde{\mathbf{X}}'_M)^+ \widetilde{\Lambda}_M + \frac{1}{\alpha_r} \mathbf{X}'(\mathbf{I}_n - \widetilde{P}_M) \boldsymbol{\varepsilon} \in \partial J_\Lambda(\mathbf{M}).$$

Consequently, when r is sufficiently large, both the positivity and the subdifferential conditions occur which, by Theorem 4.3.1, concludes the proof. \square

4.9.6 Proofs from Section 4.5.2

In this part we give proofs of Theorem 4.5.3 and Theorem 4.5.5. They are preceded by a series of simple lemmas. For reader's convenience we recall the setting of Section 4.5.2.

- A. $\boldsymbol{\varepsilon}_n = (\varepsilon_1, \dots, \varepsilon_n)'$, where $(\varepsilon_i)_i$ are i.i.d. centered with finite variance σ^2 .
- B1. $n^{-1} \mathbf{X}'_n \mathbf{X}_n \xrightarrow{\mathbb{P}} \mathbf{C} > 0$.
- B2. $\frac{\max_{i=1, \dots, n} |X_{ij}^{(n)}|}{\sqrt{\sum_{i=1}^n (X_{ij}^{(n)})^2}} \xrightarrow{\mathbb{P}} 0$, where $\mathbf{X}_n = (X_{ij}^{(n)})_{ij}$, for each $j = 1, \dots, p$.
- B'. Rows of \mathbf{X}_n are i.i.d. distributed as $\boldsymbol{\Xi}$, where $\boldsymbol{\Xi}$ is a random vector whose components are linearly independent a.s. and such that $\mathbb{E}[\boldsymbol{\Xi}_i^2] < \infty$ for $i = 1, \dots, p$.
- C. $(\mathbf{X}_n)_n$ and $(\varepsilon_n)_n$ are independent.

We consider a sequence of tuning parameters $(\mathbf{\Lambda}_n)_n$ defined by $\mathbf{\Lambda}_n = \alpha_n \mathbf{\Lambda}$, where $\mathbf{\Lambda} \in \mathbb{R}^{p+}$ is fixed and $(\alpha_n)_n$ is a sequence of positive numbers.

To ease the notation, we write the clustered matrices and clustered parameters without the subscript indicating the pattern \mathbf{M} , *i.e.* $\tilde{\mathbf{\Lambda}} := \mathbf{U}'_{|\mathbf{M}|_{\downarrow}} \mathbf{\Lambda}$, $\tilde{\mathbf{\Lambda}}_n := \alpha_n \tilde{\mathbf{\Lambda}}$ and $\tilde{\mathbf{X}}_n := \mathbf{X}_n \mathbf{U}_{\mathbf{M}}$.

Lemma 4.9.3.

(a) Under A, B1, B2 and C,

$$\frac{1}{\sqrt{n}} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \xrightarrow{d} \mathbf{V} \sim \mathbf{N}(0, \sigma^2 \mathbf{C}). \quad (4.9.7)$$

(b) Under A, B1 and C,

$$\frac{1}{n} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \xrightarrow{\mathbb{P}} 0. \quad (4.9.8)$$

(c) Under A, B' and C,

$$0 < \limsup_{n \rightarrow \infty} \frac{\|\mathbf{X}'_n \boldsymbol{\varepsilon}_n\|_{\infty}}{\sqrt{n \log \log n}} < \infty \quad a.s. \quad (4.9.9)$$

Proof of (4.9.7). It is enough to show that for any Borel subset $A \subset \mathbb{R}^p$ one has

$$\mathbb{P} \left(\frac{1}{\sqrt{n}} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \in A \mid (\mathbf{X}_n)_n \right) \xrightarrow{\mathbb{P}} \mathbb{P}(\mathbf{V} \in A). \quad (4.9.10)$$

Since both sides above are bounded, the convergence in probability implies convergence in L^1 and therefore establishes (4.9.7). To show (4.9.10), we will prove that for any subsequence $(n_k)_k$, there exists a sub-subsequence $(n_{k_l})_l$ for which, as $l \rightarrow \infty$,

$$\mathbb{P} \left(\frac{1}{\sqrt{n_{k_l}}} \mathbf{X}'_{n_{k_l}} \boldsymbol{\varepsilon}_{n_{k_l}} \in A \mid (X_n)_n \right) \xrightarrow{a.s.} \mathbb{P}(\mathbf{V} \in A). \quad (4.9.11)$$

Let $\mathbb{P}_{\mathbf{X}}$ denote the regular conditional probability $\mathbb{P}(\cdot \mid (\mathbf{X}_n)_n)$ on (Ω, \mathcal{F}) . By assumptions B1 and B2, from sequences $(n_k)_k$ one can choose a subsequence $(n_{k_l})_l$ for which

$$\frac{1}{n_{k_l}} \mathbf{X}'_{n_{k_l}} \mathbf{X}_{n_{k_l}} \xrightarrow{a.s.} \mathbf{C} > 0 \quad \text{and} \quad \frac{\max_{i=1, \dots, n_{k_l}} |X_{ij}^{(n_{k_l})}|}{\sqrt{\sum_{i=1}^{n_{k_l}} (X_{ij}^{(n_{k_l})})^2}} \xrightarrow{a.s.} 0, \quad j = 1, \dots, p.$$

We have

$$\begin{aligned} \text{Var}_{\mathbf{X}} \left(\frac{1}{\sqrt{n_{k_l}}} \mathbf{X}'_{n_{k_l}} \boldsymbol{\varepsilon}_{n_{k_l}} \right) &= \frac{1}{n_{k_l}} \mathbb{E} \left[\mathbf{X}'_{n_{k_l}} \boldsymbol{\varepsilon}_{n_{k_l}} \boldsymbol{\varepsilon}'_{n_{k_l}} \mathbf{X}_{n_{k_l}} \mid (\mathbf{X}_n)_n \right] \\ &= \frac{1}{n_{k_l}} \mathbf{X}'_{n_{k_l}} \mathbb{E} \left[\boldsymbol{\varepsilon}_{n_{k_l}} \boldsymbol{\varepsilon}'_{n_{k_l}} \right] \mathbf{X}_{n_{k_l}} = \frac{\sigma^2}{n_{k_l}} \mathbf{X}'_{n_{k_l}} \mathbf{X}_{n_{k_l}} \xrightarrow{a.s.} \sigma^2 \mathbf{C} > 0, \end{aligned}$$

and one can apply multivariate Lindeberg-Feller CLT on the space $(\Omega, \mathcal{F}, \mathbb{P}_{\mathbf{X}})$ to prove (4.9.11). Alternatively, the same result follows from Theorem 2.3.16, which concerns more general Central Limit Theorem for linearly negative quadrant dependent variables with weights forming a triangular array (in particular assumption B2 coincides with (2.3.3)). For our application, the assumption of nonnegative weights in Theorem 2.3.16 is not essential.

For (ii) we observe that previous derivations imply that $\text{Var}_{\mathbf{X}}(n^{-1} \mathbf{X}'_n \boldsymbol{\varepsilon}_n) \xrightarrow{\mathbb{P}} 0$. We deduce that $\mathbb{P}_{\mathbf{X}}(n^{-1} \mathbf{X}'_n \boldsymbol{\varepsilon}_n) \xrightarrow{\mathbb{P}} 0$ and hence (ii) follows upon averaging over $(\mathbf{X}_n)_n$.

Eq. (4.9.9) is the law of iterated logarithm for an i.i.d. sequence $(\Xi_i \varepsilon_i)_i$. \square

Lemma 4.9.4. *Let $\mathbf{M} = \mathbf{patt}(\boldsymbol{\beta})$. Assume $\alpha_n/n \rightarrow 0$.*

- (a) *Under A, B1 and C, the positivity condition is satisfied for large n with high probability.*
 (b) *Under A, B' and C, the positivity condition is almost surely satisfied for large n .*

Proof. If $\mathbf{M} = \mathbf{0}$, then the positivity condition is trivially satisfied. Thus, we consider $\mathbf{M} \neq \mathbf{0}$.
 (i) Since $\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n$ is invertible for large n with high probability, the positivity condition is equivalent to

$$s_n := (\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} [\widetilde{\mathbf{X}}_n' \mathbf{Y}_n - \widetilde{\boldsymbol{\Lambda}}_n] \in \mathbb{R}^{k+}.$$

Let $\mathbf{s}_0 \in \mathbb{R}^{k+}$ be defined through $\boldsymbol{\beta} = \mathbf{U}_M \mathbf{s}_0$, where $k = \|\mathbf{M}\|_\infty$. We will show that if $\alpha_n/n \rightarrow 0$, then $s_n \xrightarrow{\mathbb{P}} \mathbf{s}_0$. Since \mathbb{R}^{k+} is an open set, this will imply that for large n with high probability, the positivity condition is satisfied.

First we rewrite s_n as

$$s_n = (\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\mathbf{X}}_n' \mathbf{Y}_n - \alpha_n (\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\boldsymbol{\Lambda}}_n.$$

Since $\boldsymbol{\beta} = \mathbf{U}_M \mathbf{s}_0$, we conclude $\mathbf{X}_n \boldsymbol{\beta} = \mathbf{X}_n \mathbf{U}_M \mathbf{s}_0 = \widetilde{\mathbf{X}}_n \mathbf{s}_0$, so the linear regression model takes the form $\mathbf{Y}_n = \widetilde{\mathbf{X}}_n \mathbf{s}_0 + \boldsymbol{\varepsilon}_n$. Thus, $(\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\mathbf{X}}_n' \mathbf{Y}_n$ is the OLS estimator of \mathbf{s}_0 .
 By assumption B and Lemma 4.9.3, we deduce that

$$(\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\mathbf{X}}_n' \mathbf{Y}_n = \mathbf{s}_0 + (n^{-1} \widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \mathbf{U}_M \frac{1}{n} \mathbf{X}_n' \boldsymbol{\varepsilon}_n \xrightarrow{\mathbb{P}} \mathbf{s}_0 + [(\mathbf{U}_M' \mathbf{C} \mathbf{U}_M)^{-1} \mathbf{U}_M] \mathbf{0} = \mathbf{s}_0.$$

To complete the proof, we note that

$$\alpha_n (\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\boldsymbol{\Lambda}}_n = \frac{\alpha_n}{n} \left[n (\widetilde{\mathbf{X}}_n' \widetilde{\mathbf{X}}_n)^{-1} \widetilde{\boldsymbol{\Lambda}}_n \right] \xrightarrow{\mathbb{P}} \mathbf{0} \left[(\mathbf{U}_M' \mathbf{C} \mathbf{U}_M)^{-1} \widetilde{\boldsymbol{\Lambda}}_n \right] = \mathbf{0}.$$

(ii) If one assumes B' instead of B1, then $n^{-1} \mathbf{X}_n' \mathbf{X}_n \xrightarrow{a.s.} \mathbf{C}$ and by (4.9.9), $n^{-1} \mathbf{X}_n' \boldsymbol{\varepsilon}_n \xrightarrow{a.s.} \mathbf{0}$. The result follows along the same lines as (i). \square

For $\mathbf{M} \neq \mathbf{0}$ we denote

$$\begin{aligned} \boldsymbol{\pi}_n^{(1)} &= \mathbf{X}_n' (\widetilde{\mathbf{X}}_n')^+ \widetilde{\boldsymbol{\Lambda}}_n, & \boldsymbol{\pi}_n^{(2)} &= \mathbf{X}_n' (\mathbf{I}_n - \widetilde{\mathbf{P}}_n) \mathbf{Y}_n, \\ \boldsymbol{\pi}_n &= \boldsymbol{\pi}_n^{(1)} + \boldsymbol{\pi}_n^{(2)}, \end{aligned}$$

which simplifies in the $\mathbf{M} = \mathbf{0}$ case to $\boldsymbol{\pi}_n = \boldsymbol{\pi}_n^{(2)} = \mathbf{X}_n' \mathbf{Y}_n$.

Recall that the subdifferential condition is equivalent to $J_{\boldsymbol{\Lambda}_n}^*(\boldsymbol{\pi}_n) \leq 1$ and $\widetilde{\boldsymbol{\Lambda}}_n \in \text{col}(\widetilde{\mathbf{X}}_M')$ and the latter is satisfied in our setting. Since $\alpha J_{\boldsymbol{\Lambda}} = J_{\alpha \boldsymbol{\Lambda}}$, the subdifferential condition is satisfied if and only if

$$1 \geq J_{\boldsymbol{\Lambda}}^* (\alpha_n^{-1} \boldsymbol{\pi}_n) = J_{\boldsymbol{\Lambda}}^* \left(\alpha_n^{-1} \boldsymbol{\pi}_n^{(1)} + \frac{\sqrt{n}}{\alpha_n} n^{-1/2} \boldsymbol{\pi}_n^{(2)} \right).$$

In view of results shown below, $\alpha_n^{-1} \boldsymbol{\pi}_n^{(1)}$ converges almost surely, while $n^{-1/2} \boldsymbol{\pi}_n^{(2)}$ converges in distribution to a Gaussian vector. Thus, the pattern recovery properties of SLOPE estimator strongly depend on the behavior of the sequence $(\alpha_n/\sqrt{n})_n$.

Lemma 4.9.5. (a)

(i) *Assume A, B1 and C. If $\mathbf{M} \neq \mathbf{0}$, then*

$$\frac{1}{\alpha_n} \boldsymbol{\pi}_n^{(1)} \xrightarrow{\mathbb{P}} \mathbf{C} \mathbf{U}_M (\mathbf{U}_M' \mathbf{C} \mathbf{U}_M)^{-1} \widetilde{\boldsymbol{\Lambda}}_n.$$

(ii) Assume A, B1, B2 and C. The sequence $(n^{-1/2}\boldsymbol{\pi}_n^{(2)})_n$ converges in distribution to a Gaussian vector \mathbf{V} with

$$\mathbf{V} \sim \mathcal{N}\left(0, \sigma^2 \left[\mathbf{C} - \mathbf{C}\mathbf{U}_M(\mathbf{U}'_M\mathbf{C}\mathbf{U}_M)^{-1}\mathbf{U}'_M\mathbf{C} \right]\right).$$

(iii) Assume A, B1 and C. If $\lim_{n \rightarrow \infty} \alpha_n/\sqrt{n} = \infty$, then $\alpha_n^{-1}\boldsymbol{\pi}_n^{(2)} \xrightarrow{\mathbb{P}} \mathbf{0}$.

(b) Assume A, B' and C.

(i') If $\mathbf{M} \neq \mathbf{0}$, then

$$\frac{1}{\alpha_n}\boldsymbol{\pi}_n^{(1)} \xrightarrow{a.s.} \mathbf{C}\mathbf{U}_M(\mathbf{U}'_M\mathbf{C}\mathbf{U}_M)^{-1}\tilde{\boldsymbol{\Lambda}}.$$

(ii') If $\lim_{n \rightarrow \infty} \alpha_n/\sqrt{n \log \log n} = \infty$, then $\alpha_n^{-1}\boldsymbol{\pi}_n^{(2)} \xrightarrow{a.s.} 0$.

Proof. (i) Assumption B implies that

$$\mathbf{X}'_n \tilde{\mathbf{X}}_n (\tilde{\mathbf{X}}'_n \tilde{\mathbf{X}}_n)^{-1} = \frac{1}{n} \mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M (\mathbf{U}'_M n^{-1} \mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M)^{-1} \xrightarrow{\mathbb{P}} \mathbf{C}\mathbf{U}_M (\mathbf{U}'_M \mathbf{C}\mathbf{U}_M)^{-1}.$$

(ii) When $\boldsymbol{\beta} = \mathbf{U}_M \mathbf{s}_0$, then the linear regression model takes the form $\mathbf{Y}_n = \tilde{\mathbf{X}}_n \mathbf{s}_0 + \boldsymbol{\varepsilon}_n$. Since $\tilde{\mathbf{P}}_n$ is the projection matrix onto $\text{col}(\tilde{\mathbf{X}}_n)$, we have $(\mathbf{I}_n - \tilde{\mathbf{P}}_n)\tilde{\mathbf{X}}_n = \mathbf{0}$. Thus,

$$\begin{aligned} n^{-1/2}\boldsymbol{\pi}_n^{(2)} &= n^{-1/2}\mathbf{X}'_n(\mathbf{I}_n - \tilde{\mathbf{P}}_n)\mathbf{Y}_n = n^{-1/2}\mathbf{X}'_n(\mathbf{I}_n - \tilde{\mathbf{P}}_n)\boldsymbol{\varepsilon}_n \\ &= \left[\mathbf{I}_p - \mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M (\mathbf{U}'_M \mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M)^{-1} \mathbf{U}'_M \right] \left[n^{-1/2} \mathbf{X}'_n \boldsymbol{\varepsilon}_n \right]. \end{aligned}$$

By assumption B we have,

$$n^{-1}\mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M (\mathbf{U}'_M n^{-1} \mathbf{X}'_n \mathbf{X}_n \mathbf{U}_M)^{-1} \mathbf{U}'_M \xrightarrow{\mathbb{P}} \mathbf{C}\mathbf{U}_M (\mathbf{U}'_M \mathbf{C}\mathbf{U}_M)^{-1} \mathbf{U}'_M.$$

Thus, by Lemma 4.9.3 and Slutsky's theorem, we obtain (ii).

(iii) follows directly from (ii). Assumption B' implies that $n^{-1}\mathbf{X}'_n \mathbf{X}_n \xrightarrow{a.s.} \mathbf{C}$ and thus (i') is proven in the same way as (i). (ii') follows from (4.9.9). \square

Proof of Theorem 4.5.3. (a) is a direct consequence of Lemmas 4.9.4 and 4.9.5. Since positivity condition is satisfied for large n with high probability, for (b) we have with $\mathbf{M} = \mathbf{patt}(\boldsymbol{\beta})$,

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}\left(\mathbf{patt}(\hat{\boldsymbol{\beta}}_n^{\text{SLOPE}}) = \mathbf{M}\right) &= \lim_{n \rightarrow \infty} \mathbb{P}\left(\boldsymbol{\pi}_n \in \partial J_{\alpha_n \boldsymbol{\Lambda}}(\mathbf{M})\right) = \lim_{n \rightarrow \infty} \mathbb{P}\left(\alpha_n^{-1}\boldsymbol{\pi}_n \in \partial J_{\boldsymbol{\Lambda}}(\mathbf{M})\right) \\ &\geq \lim_{n \rightarrow \infty} \mathbb{P}\left(\alpha_n^{-1}\boldsymbol{\pi}_n \in \text{ri}(\partial J_{\boldsymbol{\Lambda}}(\mathbf{M}))\right) = 1, \end{aligned} \tag{4.9.12}$$

where in the last inequality we use Portmanteau Theorem, assumption (4.5.5) and the fact that the sequence $(\alpha_n^{-1}\boldsymbol{\pi}_n)_n$ converges in distribution to $\mathbf{C}\mathbf{U}_M(\mathbf{U}'_M\mathbf{C}\mathbf{U}_M)^{-1}$ if and only if $\alpha_n/\sqrt{n} \rightarrow \infty$.

Condition (4.5.6) implies that $\mathbf{C}\mathbf{U}_M(\mathbf{U}'_M\mathbf{C}\mathbf{U}_M)^{-1} \in \partial J_{\boldsymbol{\Lambda}}(\mathbf{M})$. Since $(\alpha_n^{-1}\boldsymbol{\pi}_n)_n$ converges in probability to $\mathbf{C}\mathbf{U}_M(\mathbf{U}'_M\mathbf{C}\mathbf{U}_M)^{-1}$, the necessity of this condition is explained by (4.9.12). \square

Proof of Theorem 4.5.5. By Lemma 4.9.4, the positivity condition is satisfied for large n almost surely. By Lemma 4.9.5 (i) and (iii), we have

$$\mathbf{a}_n := \frac{1}{\alpha_n} \boldsymbol{\pi}_n \xrightarrow{a.s.} \mathbf{C} \mathbf{U}_M (\mathbf{U}'_M \mathbf{C} \mathbf{U}_M)^{-1} \tilde{\boldsymbol{\Lambda}} =: \mathbf{a}_0.$$

It is easy to see that $\mathbf{U}'_M \mathbf{a}_n = \tilde{\boldsymbol{\Lambda}}$. By the condition $\mathbf{a}_0 \in \text{ri}(J_{\boldsymbol{\Lambda}}(\mathbf{M}))$ it follows that $\mathbf{a}_n \in J_{\boldsymbol{\Lambda}}(\mathbf{M})$ almost surely for sufficiently large n . Therefore $\boldsymbol{\pi}_n \in J_{\boldsymbol{\Lambda}_n}(\mathbf{M})$ for large n almost surely and thus the subdifferential condition is also satisfied. \square

Chapter 5

Geometry of Pattern Recovery by Penalized and Thresholded Estimators

5.1 Introduction

The content of this chapter is based on the preprint of Graczyk, Schneider, the author of the dissertation and Tardivel [93].

As in previous chapters, we consider the linear regression model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix, $\boldsymbol{\beta} \in \mathbb{R}^p$ is an unknown vector of regression coefficients and $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ is a random noise term. We assume that the distribution of $\boldsymbol{\varepsilon}$ is symmetric, continuous and its density is strictly positive.

Many penalized estimators of $\boldsymbol{\beta}$ have been broadly studied in literature, e.g., LASSO [47, 176], SLOPE [189, 26, 136], OSCAR [29], fused LASSO [177], fused adaptive LASSO [149], clustered LASSO [163], PACS [162] and generalized LASSO [179]. When the loss function is the residual sum of squares, these estimators minimize, with respect to $\mathbf{b} \in \mathbb{R}^p$, a function:

$$f(\mathbf{b}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \text{pen}(\mathbf{b}),$$

where $\lambda > 0$ is the tuning parameter and the penalty term pen is a real-valued polyhedral gauge, i.e. a nonnegative convex function, homogeneous, vanishing at $\mathbf{0}$ and whose unit ball is a polyhedron. Note that in this definition we do not require boundedness of a unit ball.

The literature related to penalized least squares estimators is vast and many of these estimators have interesting and relevant structures as illustrated, e.g., in [182]. For instance, LASSO is sparse, i.e. some components of this estimator are null. The fused LASSO is a sparse and piecewise constant estimator [177], the supremum norm promotes a cluster of components maximal in absolute value [107], and SLOPE and OSCAR estimators promote clusters made of those components of these estimators, which are equal in absolute value [29, 76, 156, 24, 165].

5.1.1 Pattern recovery by penalized least squares estimators

We define a polyhedral gauge pen as a nonnegative convex function, homogeneous, vanishing at $\mathbf{0}$ and whose unit ball is a polyhedron. Every polyhedral gauge can be written as the maximum

of a finite number of linear functions [152, 133]

$$\forall \mathbf{x} \in \mathbb{R}^p, \text{pen}(\mathbf{x}) = \max\{0, \mathbf{u}'_1 \mathbf{x}, \dots, \mathbf{u}'_l \mathbf{x}\}, \text{ for some } \mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p.$$

Note that a polyhedral gauge with a bounded and symmetric unit ball $\{\mathbf{x} \in \mathbb{R}^p : \text{pen}(\mathbf{x}) \leq 1\}$ is a polyhedral norm.

Given $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$ and $\lambda > 0$, the set $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$ of solutions of a penalized least squares optimization problem is defined as follows:

$$S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y}) := \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \text{pen}(\mathbf{b}). \quad (5.1.1)$$

It is important to know that the set $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$ is not empty:

Theorem 5.1.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$, $\lambda > 0$ and $\text{pen}(\mathbf{x}) = \max\{0, \mathbf{u}'_1 \mathbf{x}, \dots, \mathbf{u}'_l \mathbf{x}\}$, where $\mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p$. Then the function*

$$f(\mathbf{b}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}(\mathbf{b})$$

has at least one minimizer.

The proof is given in the Appendix.

Note that, potentially, the set $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$ might be not a singleton, i.e. the penalized least squares estimator might be not unique. Below we recall the definition of the pattern equivalence class, which is one of the most important notions in this chapter.

Definition 2.2.2 (Equality of patterns). Let $\text{pen} : \mathbb{R}^p \mapsto \mathbb{R}$ be a polyhedral gauge. We say that $\mathbf{x} \in \mathbb{R}^p$ and $\mathbf{z} \in \mathbb{R}^p$ have the same pattern with respect to pen , denoted $\mathbf{x} \stackrel{\text{pen}}{\sim} \mathbf{z}$, when

$$\partial \text{pen}(\mathbf{x}) = \partial \text{pen}(\mathbf{z}),$$

where ∂pen is the subdifferential of pen . The set of all vectors having the same subdifferential as \mathbf{x} , denoted $C_{\mathbf{x}}$, is called the pattern equivalence class.

In Theorem 5.2.2 we prove that pattern equivalence classes, illustrated in Section 2.2, are given by normal cones of B^* , where B^* is the polar set of the unit ball of pen .

For the ℓ_1 norm penalization, two vectors $\mathbf{x}, \mathbf{z} \in \mathbb{R}^p$ have the same pattern if and only if $\text{sign}(\mathbf{x}) = \text{sign}(\mathbf{z})$. More generally, two vectors having the same pattern with respect to a polyhedral gauge penalty have a specific structure as illustrated on many examples in Section 5.3. Given \mathbf{X} and \mathbf{Y} , we aim at recovering the pattern of β . For LASSO this means the recovery of $\text{sign}(\beta)$.

In Theorem 5.4.4 we give a necessary condition for pattern recovery by penalized least squares estimators, called noiseless recovery condition. Later, in Section 5.5, we introduce penalized estimators relaxing this condition. Beforehand, we recall some of the already known conditions for pattern recovery.

Conditions for pattern recovery — examples

LASSO. Below we assume the uniqueness of the LASSO estimator. Then we note $\hat{\beta}^{\text{LASSO}}$ as a unique element of $S_{\mathbf{X}, \lambda \|\cdot\|_1}(\mathbf{Y})$. As mentioned above, LASSO estimation is a sparse method that nullifies some of the components with positive probability, entailing that the estimator also performs variable selection. Instigated by this sparsity property, an abundant literature

has arisen to deal with the recovery of the location of the non-null components of β , or, more specifically, the recovery of the sign vector of β [83, 130, 183, 192, 194]. An evident necessary condition for sign recovery by LASSO is for $\mathbf{sign}(\beta)$ to be accessible by LASSO, i.e. for a fixed $\lambda > 0$, there has to exist such $\mathbf{Y} \in \mathbb{R}^n$ that $\mathbf{sign}(\hat{\beta}^{\text{LASSO}}(\mathbf{Y})) = \mathbf{sign}(\beta)$. Otherwise, the sign recovery is impossible. A geometrical characterization of accessible sign vectors is given in [160, 156]. However, the accessibility of $\mathbf{sign}(\beta)$ by LASSO does not mean that the probability of sign recovery by LASSO is close to 1 even if the non-null components of β are extremely large. Actually, for sign recovery with a probability larger than 1/2 a stronger condition is needed, called the irrepresentability condition [183], which is satisfied when

$$\|\mathbf{X}'_I \mathbf{X}_I (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{sign}(\beta_I)\|_\infty \leq 1,$$

where $I := \{i \in [p] : \beta_i \neq 0\}$ and $\bar{I} := \{i \in [p] : \beta_i = 0\}$.

SLOPE. The notions of accessibility condition and irrepresentability condition for SLOPE have been recently introduced respectively in [156] and in Chapter 4. In particular, in Chapter 4, similarly as for LASSO, when the SLOPE irrepresentability condition does not occur, the probability of pattern recovery is smaller than 1/2.

Generalized LASSO. By substituting the ℓ_1 norm with a polyhedral gauge $\text{pen} = \|\mathbf{D}\cdot\|_1$ for a fixed matrix \mathbf{D} of a linear map on \mathbb{R}^p , one constructs an estimator $\hat{\beta} \in S_{\mathbf{X}, \lambda \|\mathbf{D}\cdot\|_1}(\mathbf{Y})$, where $\mathbf{D}\hat{\beta}$ has some null components. It is a reason why the generalized LASSO is frequently used for structure recovery. One should be aware that the structure induced by generalized LASSO depends on the matrix \mathbf{D} .

For instance, when \mathbf{D} is a matrix such that $\mathbf{D}\mathbf{b} = (b_2 - b_1, \dots, b_p - b_{p-1})'$ (denoted \mathbf{D}^{tv} below), then the penalty term $\|\mathbf{D}\cdot\|_1$ promotes equality between neighbouring components of $\hat{\beta}$. Moreover, this estimator can recover the jump set: $\{i \in [p-1] : \beta_i \neq \beta_{i+1}\}$ [104]. Actually, articles [145, 140] provide theoretical properties for the jump set recovery under the irrepresentability condition.

The noiseless recovery condition and the irrepresentability condition can be relaxed using thresholded estimators as explained below.

5.1.2 Pattern recovery by a thresholded estimator

Theorem 5.5.1 generalizes results known for LASSO to a wide class of penalized estimators. Specifically, we prove that thresholding penalized least squares estimators allows the recovery of the pattern of β with large probability under a weaker condition than before. We recall the definition of thresholded estimator below.

Definition 2.3.23.

Let pen be a polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$ and $\lambda > 0$. Given $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$, we say that $\tilde{\beta}$ is a thresholded estimator of $\hat{\beta}$ if $\partial \text{pen}(\tilde{\beta}) \subseteq \partial \text{pen}(\hat{\beta})$.

One of the examples of a thresholded penalized estimator is the thresholded LASSO 5.5.1.

Given a threshold $\tau \geq 0$, the thresholded LASSO $\hat{\beta}^{\text{LASSO}, \tau}$ is defined as follows

$$\hat{\beta}_i^{\text{LASSO}, \tau} = \begin{cases} \hat{\beta}_i^{\text{LASSO}} & \text{if } |\hat{\beta}_i^{\text{LASSO}}| > \tau, \\ 0 & \text{otherwise.} \end{cases}$$

Note that for every $\tau \geq 0$ we have $\partial\|\cdot\|_1(\hat{\boldsymbol{\beta}}^{\text{LASSO}}) \subseteq \partial\|\cdot\|_1(\hat{\boldsymbol{\beta}}^{\text{LASSO},\tau})$ and thus $\hat{\boldsymbol{\beta}}^{\text{LASSO},\tau}$ is a thresholded estimator of $\hat{\boldsymbol{\beta}}^{\text{LASSO}}$ in the sense of Definition 2.3.23.

The thresholded LASSO does not have the same statistical properties as LASSO, cf. [131, 186]. Concerning sign recovery by thresholded LASSO, the accessibility condition is necessary and sufficient, which was proven by Tardivel and Bogdan [174]. Moreover, they also proved that, contrarily to LASSO, thresholded LASSO can recover the sign of $\boldsymbol{\beta}$ with a large probability under the accessibility condition, even if the irrepresentability condition is not satisfied, as soon as non-null components of $\boldsymbol{\beta}$ are sufficiently large. This nice property for sign recovery by thresholded LASSO remains true for thresholded basis pursuit [153, 55, 54].

5.2 Geometry of pattern equivalence

Let F be a face of a polytope P . We propose the following relation between normal cones $N_P(\boldsymbol{x})$ and $N_F(\boldsymbol{x})$ for $\boldsymbol{x} \in \text{ri}(F)$, that seems unknown.

Proposition 5.2.1. *Let $\boldsymbol{x} \in \text{ri}(F) \subset P$. Then*

$$\text{lin}(N_P(\boldsymbol{x})) = N_F(\boldsymbol{x}).$$

Proof. A polytope P can be represented as $\{\boldsymbol{v} \in \mathbb{R}^p : \langle \boldsymbol{s}_j, \boldsymbol{v} \rangle \leq r_j, j = 1, \dots, m\}$, cf. [101, p. 138]. Then one may describe a face $F \subset P$ as

$$F = P \cap \{\boldsymbol{v} \in \mathbb{R}^p : \langle -\boldsymbol{s}_j, \boldsymbol{v} \rangle \leq -r_j, j \in A_P(\boldsymbol{x})\},$$

where $A_P(\boldsymbol{x}) := \{j = 1, \dots, m : \langle \boldsymbol{s}_j, \boldsymbol{v} \rangle = r_j\}$ is the set of active constraints for \boldsymbol{x} . Without loss of generality one may rearrange the sequence $\{\boldsymbol{s}_j\}_{j=1}^m$ such that $A_P(\boldsymbol{x}) = \{1, \dots, a\}$, where $a \leq m$. Then F is equal to

$$F = \{\boldsymbol{v} \in \mathbb{R}^p : \langle \widetilde{\boldsymbol{s}}_j, \boldsymbol{v} \rangle \leq \widetilde{r}_j, j = 1, \dots, m+a\},$$

where

$$\widetilde{\boldsymbol{s}}_j = \begin{cases} \boldsymbol{s}_j, & j \leq m, \\ -\boldsymbol{s}_{j-m}, & j > m \end{cases} \quad \text{and} \quad \widetilde{r}_j = \begin{cases} r_j, & 1 \leq m, \\ -r_j, & j > m. \end{cases}$$

The corresponding set of active constraints $A_F(\boldsymbol{x})$ is equal to

$$A_F(\boldsymbol{x}) = A_P(\boldsymbol{x}) \cup \{m+1, \dots, m+a\}.$$

By Proposition 2.3.2 we obtain

$$N_F(\boldsymbol{x}) = \text{cone}\{\widetilde{\boldsymbol{s}}_j : j \in A_F(\boldsymbol{x})\} = \text{cone}\{\boldsymbol{s}_j, -\boldsymbol{s}_j : j \in A_P(\boldsymbol{x})\} = \text{lin}\{\boldsymbol{s}_j : j \in A_P(\boldsymbol{x})\} = \text{lin}(N_P(\boldsymbol{x})).$$

□

We will also use the following property of normal cones.

Proposition 5.2.2. *Let $\boldsymbol{x} \in \text{ri}(F) \subset P$. Then the normal cone $N_P(\boldsymbol{x})$ has the property*

$$\text{lin}(N_P(\boldsymbol{x})) = \left(\overrightarrow{\text{aff}(F)}\right)^\perp.$$

Proof. (C) We start by noticing that for any $\mathbf{d} \in \text{lin}(N_P(\mathbf{x})) = N_F(\mathbf{x})$ and for any $\mathbf{s} \in F$ we have $\langle \mathbf{d}, \mathbf{s} - \mathbf{x} \rangle = 0$. Indeed, since $\mathbf{s} = \mathbf{x} - (\mathbf{x} - \mathbf{s})$ and $\mathbf{x} \in \text{ri}(F)$, then there exists such $\delta > 0$ that $\mathbf{x} + \delta(\mathbf{x} - \mathbf{s}) \in F$. Therefore,

$$0 = \langle \mathbf{d}, \mathbf{x} - \mathbf{x} \rangle = \frac{\delta}{\delta + 1} \langle \mathbf{d}, \mathbf{x} - (\mathbf{x} - (\mathbf{x} - \mathbf{s})) \rangle + \frac{1}{\delta + 1} \langle \mathbf{d}, \mathbf{x} - (\mathbf{x} + \delta(\mathbf{x} - \mathbf{s})) \rangle,$$

while both summands are not larger than zero. Thus, $\langle \mathbf{d}, \mathbf{x} - \mathbf{s} \rangle = \langle \mathbf{d}, \mathbf{x} - (\mathbf{x} - (\mathbf{x} - \mathbf{s})) \rangle = 0$. It means that for every $\mathbf{s} \in \text{aff}(F)$ we have $\langle \mathbf{d}, \mathbf{s} \rangle = \langle \mathbf{d}, \mathbf{x} \rangle$. It implies that for every $\mathbf{s} \in \overrightarrow{\text{aff}(F)}$ we have $\langle \mathbf{d}, \mathbf{s} \rangle = 0$, thus $\mathbf{d} \in \left(\overrightarrow{\text{aff}(F)}\right)^\perp$. By the arbitrariness of \mathbf{d} , we obtain $N_F(\mathbf{x}) \subset \left(\overrightarrow{\text{aff}(F)}\right)^\perp$ and then $\text{lin}(N_P(\mathbf{x})) \subset \left(\overrightarrow{\text{aff}(F)}\right)^\perp$.

(D) Let $\mathbf{v} \in \left(\overrightarrow{\text{aff}(F)}\right)^\perp \stackrel{\text{Lemma 2.3.2}}{=} \text{lin}\{\mathbf{w} - \mathbf{x} : \mathbf{w} \in F\}$. Then $\langle \mathbf{v}, \mathbf{x} - \mathbf{w} \rangle = 0$ for every $\mathbf{w} \in F$. It implies that $\mathbf{v} \in N_F(\mathbf{x})$, which by Proposition 5.2.1 equals $\text{lin}(N_P(\mathbf{x}))$. \square

Remark 5.2.1. In addition to the property of normal cones from Proposition 2.3.2, we deduce from Propositions 5.2.2 and 2.3.2 that

$$\text{lin}(\{\mathbf{s}_j : j \in A_P(\mathbf{x})\}) = \left(\overrightarrow{\text{aff}(F)}\right)^\perp.$$

5.2.1 Pattern equivalence classes and normal cones

Definition 2.3.15 (Pattern equivalence class). [Pattern equivalence class] Let $\mathbf{x} \in \mathbb{R}^p$ and let pen be a polyhedral gauge. The pattern equivalence class $C_{\mathbf{x}}$ is the set of all vectors having the same subdifferential as \mathbf{x} :

$$C_{\mathbf{x}} := \{\mathbf{w} \in \mathbb{R}^p : \partial \text{pen}(\mathbf{w}) = \partial \text{pen}(\mathbf{x})\}.$$

Lemma 5.2.1. *Let pen be a polyhedral gauge and $\mathbf{x} \neq \mathbf{0}$. Then $C_{\mathbf{x}} \subset N(F_{\mathbf{x}})$.*

Proof. Let $\mathbf{w} \in C_{\mathbf{x}}$ and $\mathbf{s} \in \text{ri}(F_{\mathbf{x}})$. Since $\mathbf{s} \in F_{\mathbf{x}} = F_{\mathbf{w}}$ then, for all $\mathbf{z} \in B^*$, we have

$$\mathbf{w}'(\mathbf{z} - \mathbf{s}) = \underbrace{\mathbf{w}'\mathbf{z}}_{\leq \text{pen}(\mathbf{w})} - \underbrace{\mathbf{w}'\mathbf{s}}_{=\text{pen}(\mathbf{w})} \leq \text{pen}(\mathbf{w}) - \text{pen}(\mathbf{w}) = 0$$

Consequently, $\mathbf{w} \in N_{B^*}(\mathbf{s}) = N_{B^*}(F_{\mathbf{x}})$. \square

Theorem 5.2.2. *The pattern cone $C_{\mathbf{x}}$ is the relative interior of the normal cone of the face $F_{\mathbf{x}}$ of the dual unit ball B^**

$$C_{\mathbf{x}} = \text{ri}(N_{B^*}(F_{\mathbf{x}})). \quad (5.2.1)$$

Proof. If $\partial \text{pen}(\mathbf{x}) = F$, we denote $C_{\mathbf{x}} = C_F$. \mathbb{R}^p can be partitioned both into pairwise disjoint pattern sets and into pairwise disjoint relative interiors of corresponding normal cones:

$$\bigsqcup_F C_F = \mathbb{R}^p = \bigsqcup_F \text{ri}(N(F)).$$

Thus it is sufficient to prove the inclusion $C_F \subset \text{ri}(N(F))$, or equivalently

$$C_{\mathbf{x}} \subset \text{ri}(N_{B^*}(F_{\mathbf{x}})) \text{ for } \mathbf{x} \in \mathbb{R}^p. \quad (5.2.2)$$

Let $\mathbf{w} \in C_x$. We want to show that $\mathbf{w} \in \text{ri}(N_{B^*}(F_x))$, i.e. there exists $\delta > 0$ such that any point $\mathbf{z} \in B(\mathbf{w}, \delta) \cap \text{aff}(N_{B^*}(F_x))$ belongs to C_x .

According to Lemma 5.8.9, for $\delta > 0$ small enough we obtain $F_z \subseteq F_w$, which equals F_x , since $\mathbf{w} \in C_x$. Moreover, if $F_z \neq F_x$, then one may pick $\mathbf{u} \in F_z \subset F_x$ and $\mathbf{v} \in F_x \setminus F_z$. Observe that $\mathbf{u} - \mathbf{v} \in \overrightarrow{\text{aff}(F_x)}$. Then, by Lemma 2.3.1, we have $\mathbf{x} \in C_x \subset N_{B^*}(F_x)$, which by Lemma 5.8.1 belongs to $(\overrightarrow{\text{aff}(F_x)})^\perp$. Therefore $\langle \mathbf{x}, \mathbf{u} - \mathbf{v} \rangle = 0$.

We also have $\mathbf{z} \in \text{aff}(N_{B^*}(F_x)) = \text{lin}(N_{B^*}(F_x))$, which by Lemma 5.8.1 belongs to $(\overrightarrow{\text{aff}(F_x)})^\perp$. Therefore $\langle \mathbf{z}, \mathbf{u} - \mathbf{v} \rangle = 0$, too.

Consequently, $\mathbf{v}'\mathbf{z} = \mathbf{u}'\mathbf{z} = \text{pen}(\mathbf{z})$ and thus $\mathbf{v} \in F_z$, which leads to a contradiction. Therefore $F_z = F_x$, which means that $\mathbf{z} \in C_x$. \square

5.2.2 Model subspace recovery

More generally, for a wide class of penalty terms including polyhedral gauges, Vaiter et al. [182] showed that the irrepresentability condition is sufficient for the model subspace recovery by penalized least squares estimators. The notion of model subspace is related to the notion of pattern. Specifically, the model subspace of $\mathbf{x} \in \mathbb{R}^p$ is a vector subspace of \mathbb{R}^p perpendicular to $\partial \text{pen}(\mathbf{x})$. For the ℓ_1 norm two vectors $\mathbf{x}, \mathbf{z} \in \mathbb{R}^p$ have the same model subspace when $\text{supp}(\mathbf{x}) = \text{supp}(\mathbf{z})$. In the particular case of LASSO, Theorem 6 in [182] shows that $\|\mathbf{X}'_I \mathbf{X}_I (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{sign}(\boldsymbol{\beta}_I)\|_\infty < 1$ is a sufficient condition for model subspace recovery, i.e. the recovery of $\text{supp}(\boldsymbol{\beta})$. Whereas correct, this statement is not optimal. Indeed, when $\|\mathbf{X}'_I \mathbf{X}_I (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{sign}(\boldsymbol{\beta}_I)\|_\infty < 1$, it is well known that LASSO actually can recover $\mathbf{sign}(\boldsymbol{\beta})$ and a fortiori $\text{supp}(\boldsymbol{\beta})$ [183]. Whereas we do not retain the notion of model subspace, in Theorem 5.2.3 we prove that the model subspace coincides with the linear span of a pattern equivalence class.

Definition 5.2.1 (Model subspace). ([182]) *The model subspace S_x of \mathbf{x} is the orthogonal complement of $\overrightarrow{\text{aff}(F_x)}$:*

$$S_x := (\overrightarrow{\text{aff}(F_x)})^\perp,$$

Our objective in this section is to prove that $S_x = \text{lin}(C_x)$, i.e. the linear space generated by a pattern equivalence class coincides with the model subspace from [182].

Theorem 5.2.3. $S_x = \text{lin}(C_x)$.

Proof.

$$S_x = (\overrightarrow{\text{aff}(F_x)})^\perp \stackrel{\text{Prop. 5.2.2}}{=} \text{lin}(N(F_x)) = \text{lin}(\text{ri}(N(F_x))) \stackrel{\text{Thm. 5.2.2}}{=} \text{lin}(C_x).$$

\square

Remark 5.2.4. The equality $S_x = \text{lin}(C_x)$ does not hold for general penalizers. Indeed, when pen is a strictly convex function, then $S_x = \mathbb{R}^p$, but $\text{lin}(C_x)$ does not need to be equal to \mathbb{R}^p . For example, $\text{pen}(\mathbf{x}) = \|\mathbf{x}\|_2^2$ gives $\text{lin}(C_0) = \{\mathbf{0}\} \neq \mathbb{R}^p$.

5.3 Examples of polyhedral gauges and their patterns

Pattern for the ℓ_1 norm (LASSO):

Subdifferentials $\partial\|\cdot\|_1(\mathbf{x}) = \partial\|\cdot\|_1(\mathbf{z})$ are equal if and only if $\mathbf{sign}(\mathbf{x}) = \mathbf{sign}(\mathbf{z})$. For instance, if $\mathbf{x} = (1.45, -0.38, 1.56, 0, -2.76)'$, then $\mathbf{sign}(\mathbf{x}) = (1, -1, 1, 0, -1)'$. Figure 5.1 illustrates the comparison between projections onto B^* and the signs of $\hat{\beta}^{\text{LASSO}}$ for $p = 2$ and orthogonal \mathbf{X} .

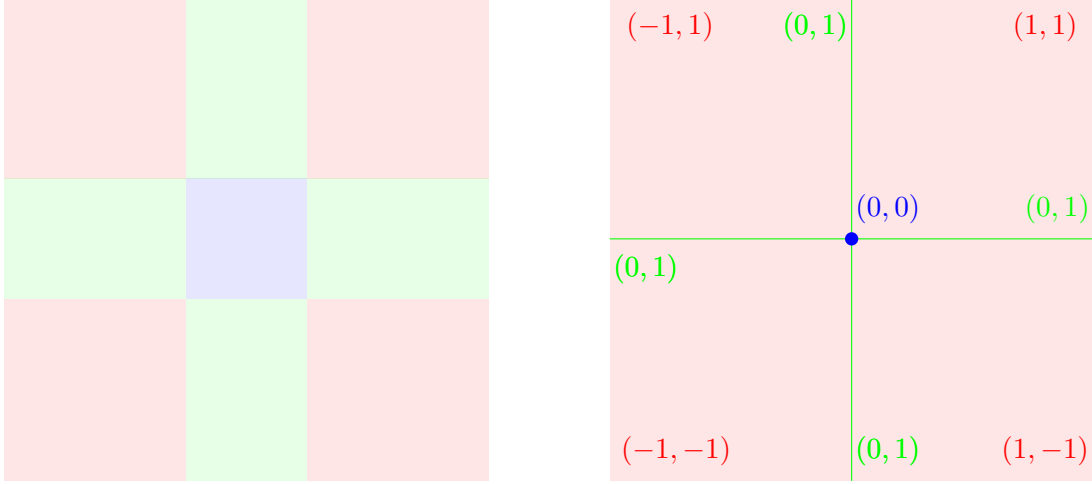


Figure 5.1: Pattern (sign) equivalence classes for LASSO in orthogonal design with $p = 2$: $\text{pen}(\mathbf{b}) = \|\mathbf{b}\|_1 = |b_1| + |b_2|$. On the left the blue polytope is B^* . Red and green (unbounded) sets are the preimages, with respect to the projection onto B^* , of its vertices and edges, respectively. The picture on the right presents sign equivalence classes of vectors in \mathbb{R}^2 .

Pattern for the ℓ_∞ norm:

The vector $\mathbf{sign}^\infty(\mathbf{x}) \in \{-1, *, 1\}^p$ is defined as follows

$$\forall i \in [p] \text{ sign}^\infty(\mathbf{x})_i := \begin{cases} 1 & \text{if } x_i > 0 \text{ and if } x_i = \|\mathbf{x}\|_\infty, \\ * & \text{if } x_i = 0 \text{ or if } |x_i| < \|\mathbf{x}\|_\infty, \\ -1 & \text{if } x_i < 0 \text{ and if } x_i = -\|\mathbf{x}\|_\infty. \end{cases}$$

The notation $*$ represents components that are not maximal in absolute value. Subdifferentials $\partial\|\cdot\|_\infty(\mathbf{x}) = \partial\|\cdot\|_\infty(\mathbf{z})$ are equal if and only if $\mathbf{sign}^\infty(\mathbf{x}) = \mathbf{sign}^\infty(\mathbf{z})$. For instance, if $\mathbf{x} = (1.45, 1.45, 0.56, 0, -1.45)'$ then $\mathbf{sign}^\infty(\mathbf{x}) = (1, 1, *, *, -1)'$. Figure 5.2 illustrates the comparison between projections onto B^* and the pattern equivalence classes of $\hat{\beta}$ for $\text{pen}(\mathbf{x}) = \|\mathbf{x}\|_\infty$, $p = 2$ and orthogonal \mathbf{X} .

Pattern for the SLOPE norm:

Let $\mathbf{x} \in \mathbb{R}^p$. The SLOPE pattern of \mathbf{x} , $\text{patt}(\mathbf{x})$, is defined by

$$\text{patt}(\mathbf{x})_i = \text{sign}(x_i) \text{rank}(|x|)_i, \quad \forall i \in [p]$$

where $\text{rank}(|x|)_i \in \{0, 1, \dots, k\}$, k is the number of non-zero distinct values in $\{|x_1|, \dots, |x_p|\}$, see 1.2.1. Let $\mathbf{\Lambda} \in \mathbb{R}^p$ where $\lambda_1 > \dots > \lambda_p > 0$. Then the subdifferentials $\partial J_{\mathbf{\Lambda}}(\mathbf{x}) = \partial J_{\mathbf{\Lambda}}(\mathbf{z})$

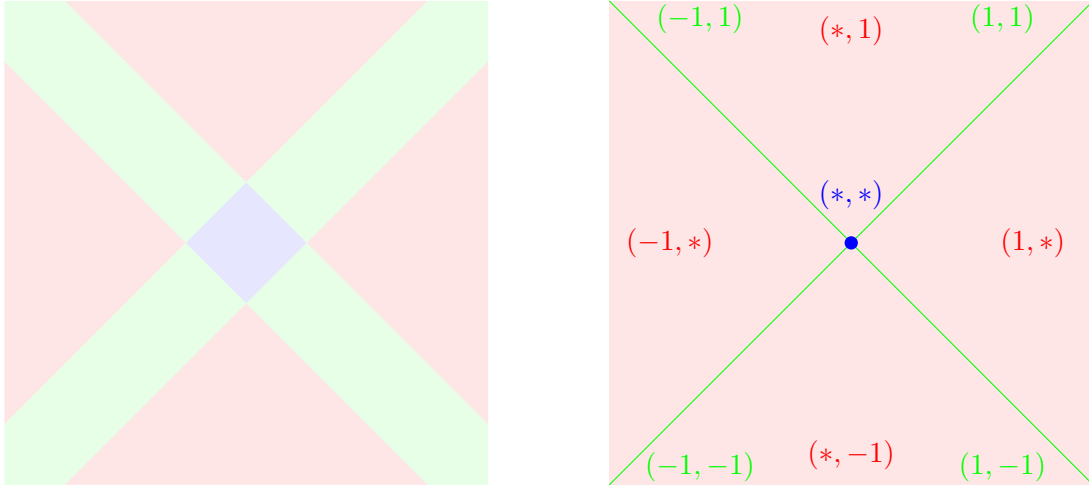


Figure 5.2: Pattern equivalence classes for ℓ_∞ norm, orthogonal \mathbf{X} and $p = 2$: $\text{pen}(\mathbf{b}) = \|\mathbf{b}\|_\infty = \max\{|b_1|, |b_2|\}$. On the left the blue polytope is B^* . Red and green (unbounded) sets are the preimages, with respect to the projection onto B^* , of its vertices and edges. The picture on the right presents pattern equivalence classes of vectors in \mathbb{R}^2 .

are equal if and only if $\text{patt}(\mathbf{x}) = \text{patt}(\mathbf{z})$. For instance, if $\mathbf{x} = (3.1, -1.2, 0.5, 0, 1.2, -3.1)'$, then $\text{patt}(\mathbf{x}) = (3, -2, 1, 0, 2, -3)'$. Figure 5.3 illustrates the comparison between projections onto B^* and the patterns of $\widehat{\beta}^{\text{SLOPE}}$ for $p = 2$ and orthogonal \mathbf{X} .

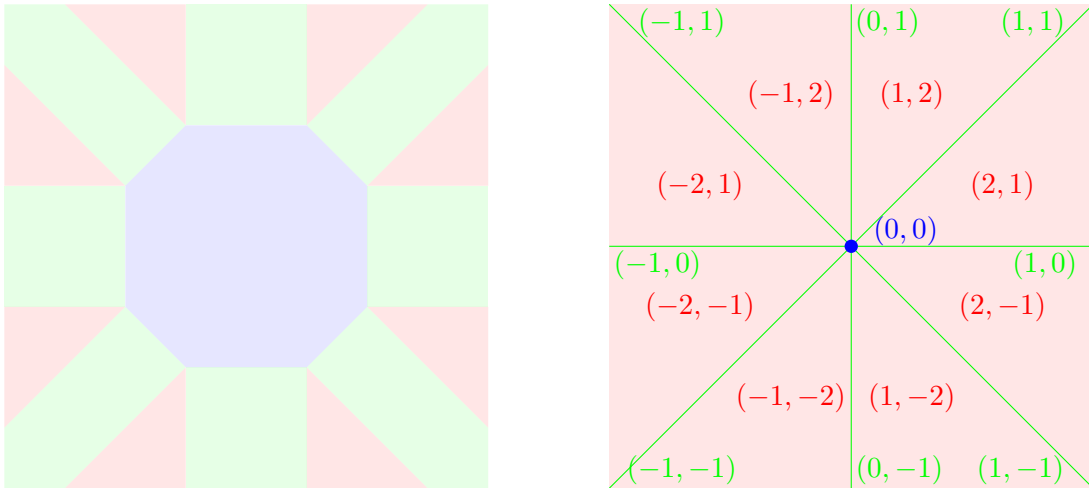


Figure 5.3: Pattern equivalence classes for SLOPE in orthogonal design with $p = 2$ and $\mathbf{\Lambda} = (2, 1)'$: $\text{pen}(\mathbf{b}) = J_{\mathbf{\Lambda}}(\mathbf{b}) = \lambda_1|b_{(1)}| + \lambda_2|b_{(2)}|$. On the left the blue polytope is B^* . Red and green (unbounded) sets are the preimages, with respect to the projection onto B^* , of its vertices and edges. The picture on the right presents pattern equivalence classes of vectors in \mathbb{R}^2 .

Furthermore, the composition of a polyhedral gauge with a linear map is still a polyhedral gauge. For example, for generalized LASSO, the penalty term is the polyhedral gauge $\mathbf{x} \in \mathbb{R}^p \mapsto \|\mathbf{D}\mathbf{x}\|_1$ where $\mathbf{D} \in \mathbb{R}^{m \times p}$. Note that, when $\ker(\mathbf{D}) \neq \mathbf{0}$, the function $\mathbf{x} \in \mathbb{R}^p \mapsto \|\mathbf{D}\mathbf{x}\|_1$ is not a norm but only a semi-norm. Below we present two examples of generalized LASSO.

- Total variation:

Let $p \geq 2$ and let $\mathbf{D}^{\text{tv}} \in \mathbb{R}^{(p-1) \times p}$ be the first order difference matrix defined as follows

$$\mathbf{D}^{\text{tv}} = \begin{pmatrix} -1 & 1 & 0 & \dots & 0 \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -1 & 1 \end{pmatrix}.$$

- ℓ_1 trend filtering:

Let $p \geq 3$ and let $\mathbf{D}^{\text{tf}} \in \mathbb{R}^{(p-2) \times p}$ be the second order difference matrix defined as follows

$$\mathbf{D}^{\text{tf}} = \begin{pmatrix} -1 & 2 & -1 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -1 & 2 & -1 \end{pmatrix}.$$

The total variation and the ℓ_1 trend filtering [111] are examples of the generalized LASSO with the penalty term $\|\mathbf{D}^{\text{tv}}\|_1$ and $\|\mathbf{D}^{\text{tf}}\|_1$, respectively.

Pattern for the total variation $\|\mathbf{D}^{\text{tv}}\|_1$:

Let $p \geq 2$. The vector $\mathbf{jump}(\mathbf{x}) \in \{\nearrow, \rightarrow, \searrow\}^{p-1}$ is defined as follows

$$\forall i \in [p-1], \mathbf{jump}(\mathbf{x})_i := \begin{cases} \nearrow & \text{if } x_{i+1} > x_i \\ \rightarrow & \text{if } x_{i+1} = x_i \\ \searrow & \text{if } x_{i+1} < x_i \end{cases}$$

Subdifferentials $\partial\|\mathbf{D}^{\text{tv}}\|_1(\mathbf{x}) = \partial\|\mathbf{D}^{\text{tv}}\|_1(\mathbf{z})$ are equal if and only if $\mathbf{jump}(\mathbf{x}) = \mathbf{jump}(\mathbf{z})$. For instance, if $\mathbf{x} = (1.45, 1.45, 0.56, 0.56, -0.45, 0.35)'$ then $\mathbf{jump}(\mathbf{x}) = (\rightarrow, \searrow, \rightarrow, \searrow, \nearrow)'$. Figure 5.4 compares the projections onto B^* with the jump sets of $\hat{\beta}$ for $p = 2$ and orthogonal \mathbf{X} .

Pattern for the trend filtering $\|\mathbf{D}^{\text{tf}}\|_1$:

Let $p \geq 3$. The vector $\mathbf{knot}(\mathbf{x}) \in \{l, cx, cv\}^{p-2}$ is defined as follows

$$\forall i \in [2 : p-1], \mathbf{knot}(\mathbf{x})_i := \begin{cases} cx & \text{if } x_i < (x_{i+1} - x_{i-1})/2, \\ l & \text{if } x_i = (x_{i+1} - x_{i-1})/2, \\ cv & \text{if } x_i > (x_{i+1} - x_{i-1})/2. \end{cases}$$

Consider the piecewise linear curve $L_{\mathbf{x}} := \bigcup_{i=1}^{p-1} [(i, x_i), (i+1, x_{i+1})]$. Note that $\mathbf{knot}(\mathbf{x})_i$ is equal to l (resp. cx or cv) when, in the neighborhood of i , the curve $L_{\mathbf{x}}$ is linear (resp. convex or concave). Subdifferentials $\partial\|\mathbf{D}^{\text{tf}}\|_1(\mathbf{x}) = \partial\|\mathbf{D}^{\text{tf}}\|_1(\mathbf{z})$ are equal if and only if $\mathbf{knot}(\mathbf{x}) = \mathbf{knot}(\mathbf{z})$. Figure 5.5 provides an illustration of $\mathbf{knot}(\mathbf{x})$ for $\mathbf{x} = (1, 3, 5, 7, 6, 5, 4, 6, 5)' \in \mathbb{R}^9$.

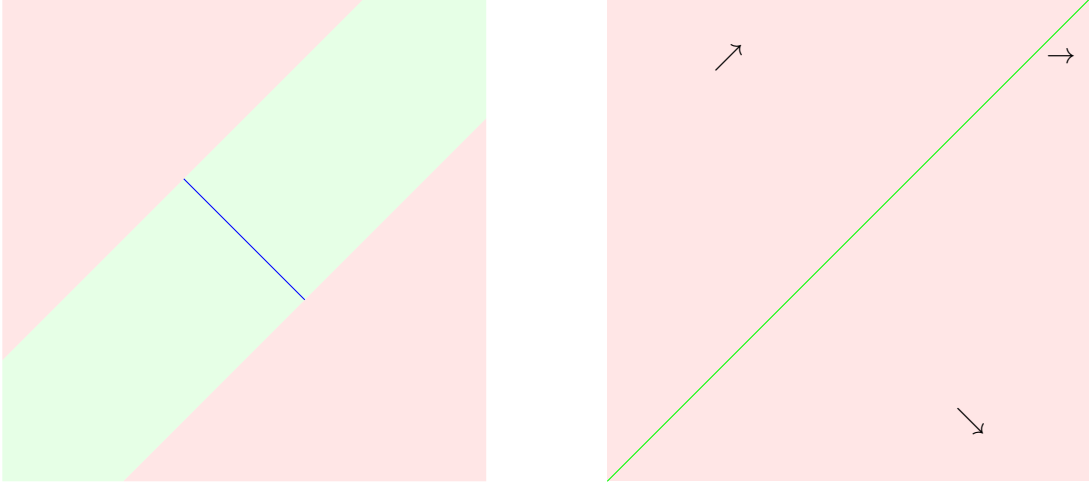


Figure 5.4: Pattern equivalence classes for the total variation in orthogonal design with $p = 2$: $\text{pen}(\mathbf{b}) = |b_1 - b_2|$. On the left the blue polytope is B^* . Red and green (unbounded) sets are the preimages, with respect to the projection onto B^* , of its vertices and edges. The picture on the right presents pattern equivalence classes of vectors in \mathbb{R}^2 .

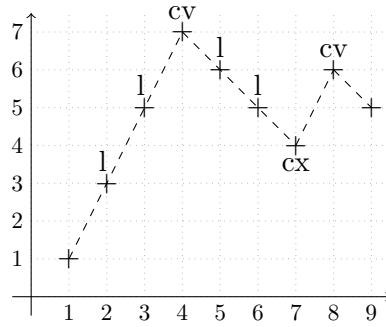


Figure 5.5: In this figure the dotted curve represents the piecewise linear curve L_x for $\mathbf{x} = (1, 3, 5, 7, 6, 5, 4, 6, 5)'$. Here $\mathbf{knot}(\mathbf{x}) = (l, l, cv, l, l, cx, cv)'$.

5.4 Pattern recovery in penalized estimation

5.4.1 Accessibility condition

Below we recall the notion of accessible patterns in the following definition, which generalizes the notion of accessible sign vectors [160, 156] and accessible SLOPE patterns [156] to a broad class of estimators penalized with polyhedral gauges.

Definition 5.4.1 (Accessible pattern). *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\lambda > 0$ and pen be a polyhedral gauge. We say that the pattern of $\beta \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} and λpen , if there exist $\mathbf{y} \in \mathbb{R}^n$ and $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ such that $\hat{\beta} \stackrel{\text{pen}}{\sim} \beta$.*

The accessibility of a pattern can be characterized in a geometric and an analytic way.

Proposition 5.4.1 (Characterization of accessible patterns). *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\text{pen} : \mathbb{R}^p \rightarrow \mathbb{R}$ be a real-valued polyhedral gauge.*

- (a) *Geometric characterization: The pattern of $\beta \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} and*

λ_{pen} if and only if

$$\text{row}(\mathbf{X}) \cap \partial \text{pen}(\boldsymbol{\beta}) \neq \emptyset.$$

(b) *Analytic characterization: The pattern of $\boldsymbol{\beta} \in \mathbb{R}^p$ is accessible with respect to \mathbf{X} and λ_{pen} if and only if for any $\mathbf{b} \in \mathbb{R}^p$ we have*

$$\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\mathbf{b} \implies \text{pen}(\boldsymbol{\beta}) \leq \text{pen}(\mathbf{b}).$$

Based on Proposition 4.4.2, one may see that the notion of accessibility does not depend on the tuning parameter λ .

5.4.2 Noiseless recovery condition

The solution path for a penalized estimator is defined as the curve $0 < \lambda \mapsto \widehat{\boldsymbol{\beta}}(\lambda)$ where $\widehat{\boldsymbol{\beta}}(\lambda)$ is the unique element of $S_{\mathbf{X}, \lambda_{\text{pen}}}(\mathbf{Y})$ for fixed $\mathbf{Y} \in \mathbb{R}^n$ and $\mathbf{X} \in \mathbb{R}^{n \times p}$. The solution paths for the generalized LASSO or OSCAR and the Clustered LASSO are studied in [179] or [173], respectively. Based on this notion, below we define the noiseless recovery condition. Note that the following definition does not require the uniqueness of an estimator.

Definition 5.4.2 (Noiseless recovery condition). *Let pen be a polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$. We say that the pattern of $\boldsymbol{\beta}$ satisfies the noiseless recovery condition with respect to \mathbf{X} and pen if*

$$\exists \lambda > 0, \exists \widehat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda_{\text{pen}}}(\mathbf{X}\boldsymbol{\beta}) \text{ such that } \widehat{\boldsymbol{\beta}} \stackrel{\text{pen}}{\sim} \boldsymbol{\beta}.$$

This condition generalizes the noiseless pattern recovery condition for SLOPE (4.3.2).

Remark 5.4.1. $\boldsymbol{\beta} = \mathbf{0}$ satisfies the noiseless recovery condition with respect to \mathbf{X} and pen , because $\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ and $\mathbf{0} \in S_{\mathbf{X}, \lambda_{\text{pen}}}(\mathbf{0})$. In other words, the noiseless recovery condition means that in the noiseless case when $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}$, in the solution path, one may pick such tuning parameter λ , that the minimizer has the same pattern as $\boldsymbol{\beta}$.

Example 5.4.2. (Solution path for LASSO):
Consider the LASSO optimization problem for

$$\mathbf{X} = \begin{pmatrix} 5/6 & 1 & 0 \\ 1/3 & 0 & 1 \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \mathbf{0} \text{ and } \boldsymbol{\beta} = (10, 0, 0)'$$

Then the LASSO solution path $0 < \lambda \mapsto \widehat{\boldsymbol{\beta}}^{\text{LASSO}}(\lambda)$ is the following curve:

5.4.3 Irrepresentability Condition for polyhedral gauges

The following theorem is one of the main results of this chapter. Its proof is based on Theorem 5.2.2.

Theorem 5.4.3. *Let pen be a polyhedral norm. Let $\boldsymbol{\varepsilon} = \mathbf{0}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$, $\boldsymbol{\beta} \neq \mathbf{0}$. Then the existence of $\widehat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda_{\text{pen}}}(\mathbf{X}\boldsymbol{\beta})$ recovering the pattern of $\boldsymbol{\beta}$ is equivalent to*

$$\mathbf{X}'\mathbf{X}\text{lin}(C_{\boldsymbol{\beta}}) \cap F_{\boldsymbol{\beta}} \neq \emptyset.$$

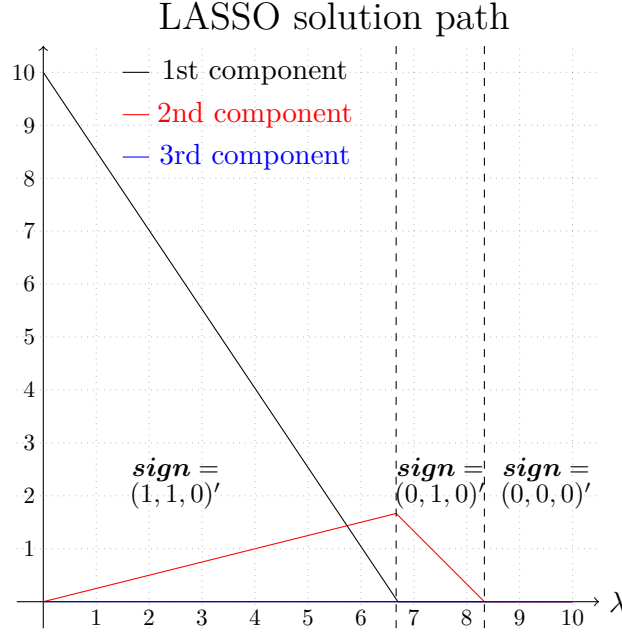


Figure 5.6: Coordinates and the sign of $\hat{\beta}^{\text{LASSO}}$ as the functions of $\lambda > 0$: $\hat{\beta}_1^{\text{LASSO}}(\lambda)$ (black curve), $\hat{\beta}_2^{\text{LASSO}}(\lambda)$ (red) and $\hat{\beta}_3^{\text{LASSO}}(\lambda)$ (blue). Note that $\text{sign}(\beta) = (1, 0, 0)'$ is not recovered for any $\lambda > 0$.

Proof. (\implies):

Let $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{X}\beta)$ and $F_{\hat{\beta}} = F_{\beta}$. Then, by the subdifferential properties, $\frac{1}{\lambda}\mathbf{X}'\mathbf{X}(\beta - \hat{\beta}) \in F_{\beta}$.

$\hat{\beta}$ recovers the pattern of β , thus we have $\beta, \hat{\beta} \in \text{lin}(C_{\beta})$. Therefore $\frac{\beta - \hat{\beta}}{\lambda} \in \text{lin}(C_{\beta})$ and finally $\frac{1}{\lambda}\mathbf{X}'\mathbf{X}(\beta - \hat{\beta}) \in \mathbf{X}'\mathbf{X}\text{lin}(C_{\beta})$.

(\impliedby):

By Theorem 5.2.2, the pattern set C_{β} is relatively open.

Let $\mathbf{X}'\mathbf{X}\text{lin}(C_{\beta}) \cap F_{\beta} \neq \emptyset$, i.e. there exists such $\mathbf{z} \in V_{\beta}$ that $\mathbf{X}'\mathbf{X}\mathbf{z} \in F_{\beta}$. Consider $\hat{\beta} = \beta - \lambda\mathbf{z}$. Since $\beta \in C_{\beta} = \text{ri}(C_{\beta})$ and $\mathbf{z} \in \text{lin}(C_{\beta})$, then for λ small enough we have $F_{\hat{\beta}} = F_{\beta}$.

Moreover, we obtain $\frac{1}{\lambda}\mathbf{X}'(\mathbf{X}\beta - \mathbf{X}\hat{\beta}) = \mathbf{X}'\mathbf{X}\mathbf{z} \in F_{\beta}$. Therefore $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{X}\beta)$, which ends the proof. \square

Below we provide a geometrical characterization of the noiseless recovery condition. Neither the above definition nor the geometrical characterization provide an analytic expression for checking the noiseless recovery condition, but for some of the penalized estimators analogous formulas have been already given. For example, when $\text{pen} = \|\cdot\|_1$, the noiseless recovery condition is equivalent to

$$\|\mathbf{X}'(\mathbf{X}'_I)^+ \text{sign}(\beta_I)\|_{\infty} \leq 1 \text{ and } \text{sign}(\beta_I) \in \text{row}(\mathbf{X}_I), \quad (5.4.1)$$

where $I = \{i \in [p] : \beta_i \neq 0\}$ and \mathbf{X}_I is the matrix whose columns are $(\mathbf{X}_j)_{j \in I}$. Note that under the assumption that $\ker(\mathbf{X}_I) = \{\mathbf{0}\}$ we obtain $\text{sign}(\beta_I) \in \text{row}(\mathbf{X}_I)$ and the expression (5.4.1) coincides with the irrepresentability condition: $\|\mathbf{X}'_{\bar{I}}\mathbf{X}_I(\mathbf{X}'_I\mathbf{X}_I)^{-1}\text{sign}(\beta_I)\|_{\infty} \leq 1$ where $\mathbf{X}'_{\bar{I}}$ is a matrix whose columns are $(\mathbf{X}_j)_{j \notin I}$ [41, 183, 194, 192]. Thus, the well known irrepresentability condition for LASSO can be thought of as an analytical shortcut for checking the noiseless recovery condition, see Figure 5.6. Indeed, in the above example, we have $\|\mathbf{X}'_{\bar{I}}\mathbf{X}_I(\mathbf{X}'_I\mathbf{X}_I)^{-1}\text{sign}(\beta_I)\|_{\infty} = 30/29 > 1$ and based on Figure 5.6, one may observe that

the noiseless recovery condition does not hold for β . For SLOPE with $\mathbf{M} = \text{patt}(\beta)$, the noiseless recovery condition is equivalent to

$$J_{\Lambda}^*(\mathbf{X}'(\widetilde{\mathbf{X}}_{\mathbf{M}})^+\widetilde{\Lambda}_{\mathbf{M}}) \leq 1 \text{ and } \widetilde{\Lambda}_{\mathbf{M}} \in \text{row}(\widetilde{\mathbf{X}}_{\mathbf{M}}),$$

where J_{Λ} is the sorted ℓ_1 norm, $\widetilde{\mathbf{X}}_{\mathbf{M}}$ is the clustered matrix and $\widetilde{\Lambda}_{\mathbf{M}}$ is the clustered parameter, cf. Definition 2.2.5. In Appendix we also provide an analytic characterization of noiseless pattern recovery when the penalty term is the supremum norm. Below we show that

- (a) The noiseless recovery condition is necessary to recover the pattern with a probability larger than 1/2, see Theorem 5.4.4.
- (b) Thresholded penalized estimators can recover the pattern of β under much weaker condition than the noiseless recovery condition, see Section 5.5.

Theorem 5.4.4. *Let $\mathbf{Y} = \mathbf{X}\beta + \varepsilon$ where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a fixed matrix, $\beta \in \mathbb{R}^p$ and ε follows a symmetric distribution. Let pen be a polyhedral gauge. If β does not satisfy the noiseless recovery condition with respect to \mathbf{X} and pen , then*

$$\mathbb{P}\left(\exists \lambda > 0 \exists \widehat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y}) \text{ such that } \widehat{\beta} \stackrel{\text{pen}}{\approx} \beta\right) \leq 1/2.$$

Corollary 5.4.1. *If the noiseless recovery condition does not hold for the LASSO (for example, when $\|\mathbf{X}'_I \mathbf{X}_I (\mathbf{X}'_I \mathbf{X}_I)^{-1} \mathbf{sign}(\beta_I)\|_{\infty} > 1$), the following holds*

$$\mathbb{P}\left(\exists \lambda > 0 \exists \widehat{\beta} \in S_{\mathbf{X}, \lambda \|\cdot\|_1}(\mathbf{Y}) \text{ such that } \mathbf{sign}(\widehat{\beta}) = \mathbf{sign}(\beta)\right) \leq 1/2.$$

The above result extends the Theorem 2 from [183], which shows that $\mathbb{P}(\mathbf{sign}(\widehat{\beta}^{\text{LASSO}}(\lambda)) = \mathbf{sign}(\beta)) \leq 1/2$ for fixed $\lambda > 0$.

If β satisfies the noiseless recovery condition with respect to \mathbf{X} and pen , then β is accessible with respect to \mathbf{X} and pen , by taking $\mathbf{y} = \mathbf{X}\beta$ in the definition of accessibility. In the following section, we show that thresholded penalized least-squares estimators recover the pattern of β under the accessibility condition.

5.5 Pattern recovery by thresholded estimators

In practice, some additional information about β may be priorly known, e.g. its sparsity. Therefore it is quite natural to threshold small components of $\widehat{\beta}^{\text{LASSO}}$ and to consider the thresholded LASSO estimator $\widehat{\beta}^{\text{LASSO}, \tau}$ for some threshold $\tau \geq 0$:

Definition 5.5.1 (Thresholded LASSO). *The thresholded LASSO [174] is defined in a following way [41, Sec. 2.9]:*

$$\widehat{\beta}_i^{\text{LASSO}, \tau} = \begin{cases} \widehat{\beta}_i^{\text{LASSO}} & \text{if } |\widehat{\beta}_i^{\text{LASSO}}| > \tau, \\ 0 & \text{otherwise.} \end{cases} \quad (5.5.1)$$

Moreover, if the threshold is appropriately selected, the estimator allows to recover $\mathbf{sign}(\beta)$ under weaker conditions than LASSO itself [174]. We aim at generalizing this property to the class of penalized estimators with polyhedral penalty. Before introducing the notion of a thresholded estimator, recall that for any threshold $\tau \geq 0$, the inclusion $\partial \|\cdot\|_1(\widehat{\beta}^{\text{LASSO}}) \subseteq \partial \|\cdot\|_1(\widehat{\beta}^{\text{LASSO}, \tau})$ occurs. This last inclusion is the keystone concept to introduce the notion of a thresholded estimator as defined in Definition 2.3.23.

The notion of accessibility introduced for penalized estimators in Section 5.4 also covers the thresholded estimators as can be seen below.

Proposition 5.5.1. *Let pen be a real-valued polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$. Then*

$$\begin{aligned} \exists \mathbf{y} \in \mathbb{R}^n, \exists \hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y}) \text{ such that } \hat{\boldsymbol{\beta}} \stackrel{\text{pen}}{\approx} \boldsymbol{\beta} \\ \iff \exists \mathbf{y} \in \mathbb{R}^n, \exists \hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y}) \text{ such that } \partial \text{pen}(\hat{\boldsymbol{\beta}}) \subseteq \partial \text{pen}(\boldsymbol{\beta}). \end{aligned}$$

According to Propositions 4.4.2 and 5.5.1, if the pattern of $\boldsymbol{\beta}$ is not accessible with respect to \mathbf{X} and λpen , i.e. there exists $\mathbf{b} \in \mathbb{R}^p$ such that $\mathbf{X}\mathbf{b} = \mathbf{X}\boldsymbol{\beta}$ and $\text{pen}(\mathbf{b}) < \text{pen}(\boldsymbol{\beta})$, then for any $\mathbf{y} \in \mathbb{R}^n$, $\lambda > 0$, and $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ we have $\partial \text{pen}(\hat{\boldsymbol{\beta}}) \not\subseteq \partial \text{pen}(\boldsymbol{\beta})$. Consequently, no penalized nor thresholded penalized estimator can recover the pattern of $\boldsymbol{\beta}$.

On the other hand, if the accessibility condition is satisfied, then both penalized and thresholded penalized estimator can recover the pattern of $\boldsymbol{\beta}$ with different choices of \mathbf{y} . However, in practice, instead of choosing an appropriate \mathbf{y} to recover the pattern of $\boldsymbol{\beta}$, the response of a linear regression model is being used to infer this pattern.

In this direction, by Theorem 5.4.4, if $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, then the recovery of the pattern of $\boldsymbol{\beta}$ with probability larger than 1/2 requires the noiseless recovery condition, which is stronger than the accessibility condition. This result remains true for any symmetric and continuous noise. In Theorem 5.5.1, we relax the stringent noiseless recovery condition by considering a thresholded estimator. Before stating this theorem, we introduce the following class of thresholded estimators.

Definition 5.5.2 (τ -thresholded penalized estimator). *Let pen be a real-valued polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$ and $\lambda \geq 0$. Given $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$, we say that $\tilde{\boldsymbol{\beta}}^\tau$ is a τ -thresholded estimator of $\hat{\boldsymbol{\beta}}$ if*

- (a) $\partial \text{pen}(\hat{\boldsymbol{\beta}}) \subseteq \partial \text{pen}(\tilde{\boldsymbol{\beta}}^\tau)$,
- (a) $\|\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}^\tau\|_\infty \leq \tau$,
- (a) $\dim(\partial \text{pen}(\mathbf{b})) \leq \dim(\partial \text{pen}(\tilde{\boldsymbol{\beta}}^\tau))$ for all \mathbf{b} with $\|\hat{\boldsymbol{\beta}} - \mathbf{b}\|_\infty \leq \tau$.

The thresholded LASSO is, in fact, an example of a τ -thresholded estimator with the threshold τ . Another example of a τ -thresholded estimator when the penalty term is the supremum norm, is given in Algorithm 1. Theorem 5.5.1 shows that a thresholded estimator recovers the pattern of $\boldsymbol{\beta}$ under the accessibility condition and the assumption that the signal is large enough, as formalized in the following theorem.

Theorem 5.5.1. *Let pen be a real-valued polyhedral gauge, $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\boldsymbol{\beta} \in \mathbb{R}^p$, and $\lambda > 0$. Assume that the uniform uniqueness holds, i.e. for any $\mathbf{y} \in \mathbb{R}^n$, the set $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ consists of one element $\hat{\boldsymbol{\beta}}(\mathbf{y})$. For $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ and for $r \in \mathbb{N}_+$ set $\mathbf{y}^{(r)} = \mathbf{X}(r\boldsymbol{\beta}) + \boldsymbol{\varepsilon}$. If $\text{pen}(\mathbf{b}) \geq \text{pen}(\boldsymbol{\beta})$ for any $\mathbf{b} \in \mathbb{R}^p$ with $\mathbf{X}\mathbf{b} = \mathbf{X}\boldsymbol{\beta}$, then there exists $r_0 \in \mathbb{N}_+$ and $\tau \geq 0$ such that for all $r \geq r_0$*

$$\begin{cases} \partial \text{pen}(\mathbf{b}) \subseteq \partial \text{pen}(\boldsymbol{\beta}) \text{ for any } \mathbf{b} \in \overline{B}_\infty(\hat{\boldsymbol{\beta}}^{(r)}, \tau) \\ \exists \mathbf{b}_0 \in \overline{B}_\infty(\hat{\boldsymbol{\beta}}^{(r)}, \tau) \text{ such that } \mathbf{b}_0 \stackrel{\text{pen}}{\approx} \boldsymbol{\beta} \end{cases}$$

Consequently, a τ -thresholded penalized estimator $\tilde{\boldsymbol{\beta}}^\tau(\mathbf{y}^{(r)})$ recovers the pattern of $\boldsymbol{\beta}$.

Similar results in which non-null components are large enough (i.e., $r \geq r_0$ in Theorem 5.5.1) are given in [174] and [54]. In particular, Theorem 5.5.1 corroborates Theorem 1 in [174], which proves that the thresholded LASSO estimator recovers the sign of $\boldsymbol{\beta}$ once the accessibility condition holds and non-null components of $\boldsymbol{\beta}$ are large enough. Similarly as thresholded LASSO, while $\text{pen} = \|\cdot\|_\infty$, a τ -estimator can be explicitly computed by Algorithm 1.

Algorithm 1 Thresholded penalized estimator when the penalty term is the ℓ_∞ norm:

Require: estimate: $\widehat{\boldsymbol{\beta}}$, threshold $\tau \geq 0$.

if $\|\widehat{\boldsymbol{\beta}}\|_\infty \leq \tau$ **then**
 $\widetilde{\boldsymbol{\beta}}^\tau \leftarrow \mathbf{0}$.

else

$$\forall j \in [p] \quad \widetilde{\beta}_j^\tau \leftarrow \begin{cases} \|\widehat{\boldsymbol{\beta}}\|_\infty - \tau & \text{if } \|\widehat{\boldsymbol{\beta}}\|_\infty - 2\tau \leq \widehat{\beta}_j \text{ and } \widehat{\beta}_j \geq 0, \\ -\|\widehat{\boldsymbol{\beta}}\|_\infty + \tau & \text{if } \widehat{\beta}_j \leq -\|\widehat{\boldsymbol{\beta}}\|_\infty + 2\tau \text{ and } \widehat{\beta}_j < 0, \\ \widehat{\beta}_j & \text{otherwise.} \end{cases}$$

end if

return $\widetilde{\boldsymbol{\beta}}^\tau$

5.6 Full characterization of the uniform uniqueness

For the pattern recovery by the τ -thresholded penalized estimator, in Theorem 5.5.1 we assume the uniform uniqueness. For that reason, below we provide a necessary and sufficient condition for uniform uniqueness of the penalized optimization problem (5.1.1) in Theorem 5.6.1. This theorem extends the Theorem 1 from [156] to all polyhedral gauges.

Theorem 5.6.1 (Necessary and sufficient condition for uniform uniqueness). *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\lambda > 0$. Let pen be a polyhedral gauge, i.e., $\text{pen}(\mathbf{x}) = \max\{0, \mathbf{u}'_1 \mathbf{x}, \dots, \mathbf{u}'_l \mathbf{x}\}$ for some $\mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p$. For $\mathbf{y} \in \mathbb{R}^n$ consider the following optimization problem*

$$S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y}) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \text{pen}(\mathbf{b}). \quad (5.1.1)$$

Then the solution to the above problem is unique, i.e., $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ is a singleton for all $\mathbf{y} \in \mathbb{R}^n$ if and only if $\text{row}(\mathbf{X})$ does not intersect a face of the polytope $B^ = \text{conv}\{\mathbf{0}, \mathbf{u}_1, \dots, \mathbf{u}_l\}$ whose dimension is smaller than $\dim(\ker(\mathbf{X}))$.*

Note that every face F of B^* satisfies

$$\dim(F) < \dim(\ker(\mathbf{X})) \iff \text{codim}(F) > \text{rk}(\mathbf{X}),$$

where $\text{codim}(F) = p - \dim(F)$. For a better explanation of the non-uniqueness of the estimator, below we give an example for the generalized LASSO with $\text{pen}(\mathbf{b}) = \|\mathbf{D}\mathbf{b}\|_1$. Note that if $\ker(\mathbf{X}) \cap \ker(\mathbf{D}) = \{\mathbf{0}\}$, then for every $\widehat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \|\mathbf{D}\cdot\|_1}(\mathbf{y})$ and $\mathbf{g} \in \ker(\mathbf{X}) \cap \ker(\mathbf{D})$ we have $\widehat{\boldsymbol{\beta}} + \mathbf{g} \in S_{\mathbf{X}, \lambda \|\mathbf{D}\cdot\|_1}(\mathbf{y})$, too. Thus for every $\mathbf{y} \in \mathbb{R}^n$ the minimizer is not unique. Consequently, $\ker(\mathbf{X}) \cap \ker(\mathbf{D}) = \{\mathbf{0}\}$ is a necessary condition for uniform uniqueness, yet, it is not sufficient, as illustrated in the example given in [9], which we revisit below.

Example 5.6.2. [9, p. 19] Consider the following optimization problem

$$\arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \frac{1}{2} \|\mathbf{D}\mathbf{b}\|_1,$$

where

$$\mathbf{X} = \begin{pmatrix} 1 & 1 & 1 \\ 3 & 1 & 1 \\ \sqrt{2} & 0 & 0 \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 2 & 1 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}.$$

We have $S_{\mathbf{X}, \frac{1}{2}\|\mathbf{D}\cdot\|_1}(\mathbf{Y}) = \text{conv}\{(0, 1/2, 0)', (0, 0, 1/2)'\}$ [9, p. 19]. Since

$$\|\mathbf{D}\mathbf{b}\|_1 = \max\{\pm(4b_1 + 2b_2 + 2b_3), \pm(2b_1 + 2b_2), \pm(2b_1 + 2b_3)\}$$

then $B^* = \text{conv}\{\pm(4, 2, 2)', \pm(2, 2, 0)', \pm(2, 0, 2)'\}$. Because the vertex $F = (4, 2, 2)'$ is a face of B^* which lies in $\text{row}(\mathbf{X})$ and satisfies $\dim(F) = 0 < 1 = \dim(\ker(\mathbf{X}))$ then, according to Theorem 5.6.1, the uniform uniqueness cannot hold. This complies with the fact that $S_{\mathbf{X}, \frac{1}{2}\|\mathbf{D}\cdot\|_1}(\mathbf{Y})$ is not a singleton.

When $\ker(\mathbf{X}) \cap \ker(\mathbf{D}) = \{\mathbf{0}\}$, in broad generality, the set of generalized LASSO solutions is a polytope (a bounded polyhedron) [9, Proposition 4.3.] and its extremal points can be explicitly computed [62]. This description is relevant when this set is not a singleton.

Lemma 5.6.1. *Consider the function*

$$f(\mathbf{b}) = \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 + \text{pen}(\mathbf{b}),$$

where $\text{pen}(\mathbf{b}) = \max\{0, \mathbf{u}'_1\mathbf{b}, \dots, \mathbf{u}'_l\mathbf{b}\}$ is a polyhedral gauge. Then the condition $\ker(\mathbf{X}) \cap \ker(\mathbf{D}) = \{\mathbf{0}\}$ is necessary for the uniform uniqueness of the minimizer of f .

Proof. Let $\mathbf{0} \neq \mathbf{h} \in \ker(\mathbf{X}) \cap \ker(\text{pen})$ and let $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda\text{pen}}(\mathbf{Y})$. Then

$$f(\hat{\boldsymbol{\beta}}) \leq f(\hat{\boldsymbol{\beta}} + \mathbf{h}) = \frac{1}{2}\|\mathbf{Y} - \mathbf{X}(\hat{\boldsymbol{\beta}} + \mathbf{h})\|_2^2 + \text{pen}(\hat{\boldsymbol{\beta}} + \mathbf{h}) = \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \max\{0, \mathbf{u}'_1(\hat{\boldsymbol{\beta}} + \mathbf{h}), \dots, \mathbf{u}'_l(\hat{\boldsymbol{\beta}} + \mathbf{h})\}.$$

Since $\text{pen}(\mathbf{h}) = 0$, we have $\mathbf{u}'_i\mathbf{h} \leq 0$ for every $i \in [l]$, hence

$$f(\hat{\boldsymbol{\beta}} + \mathbf{h}) \leq \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \max\{0, \mathbf{u}'_1\hat{\boldsymbol{\beta}}, \dots, \mathbf{u}'_l\hat{\boldsymbol{\beta}}\} = f(\hat{\boldsymbol{\beta}}).$$

Therefore $\hat{\boldsymbol{\beta}}$ is not a unique minimizer of f . □

5.7 Numerical experiments

Below, in our simulations, we consider the linear regression model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where:

- The matrix $\mathbf{X} = (\mathbf{X}_1 | \dots | \mathbf{X}_{150}) \in \mathbb{R}^{100 \times 150}$ has i.i.d. $\mathcal{N}(0, 1/100)$ entries.
- The random noise $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ has i.i.d. $\mathcal{N}(0, 1)$ entries.

5.7.1 Numerical experiments for LASSO

For LASSO, the noiseless recovery condition and the accessibility condition depend on $\boldsymbol{\beta}$ through $\text{sign}(\boldsymbol{\beta}) \in \{-1, 0, 1\}^p$. Moreover, since the distribution of \mathbf{X} is invariant by permutations and sign changes of the columns of \mathbf{X} , then the probability that a k -sparse vector satisfies the noiseless recovery condition is equivalent to the irrepresentability condition:

$$\mathbb{P}_{\mathbf{X}}\left(\|\mathbf{X}'(\mathbf{X}'_I)^+\mathbf{1}_k\|_\infty \leq 1 \text{ and } \mathbf{1}_k \in \text{row}(\mathbf{X}_I)\right),$$

where $I = [k]$ and $\mathbf{1}_k = (1, \dots, 1)' \in \mathbb{R}^k$. Moreover, the accessibility condition is satisfied with probability

$$\mathbb{P}_{\mathbf{X}}(\min\{\|\boldsymbol{\gamma}\|_1 : \mathbf{X}\boldsymbol{\gamma} = \mathbf{X}_I\mathbf{1}_k\} = k).$$

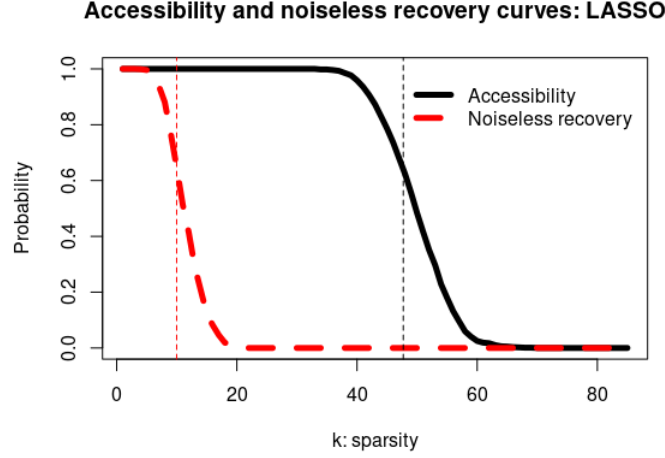


Figure 5.7: Probability of the noiseless recovery condition and the accessibility condition as functions of the support size k . The highlighted values of k are $k_1 = 50/\log(150) \approx 9.9$ and $k_2 = 100\rho_{DT}(2/3) \approx 47.8$, which are the preimages of 0.5 for the noiseless recovery curve [183] and the accessibility curve [59, 60], respectively.

Figure 5.7 provides these probabilities as functions of k .

Figure 5.8 illustrates sign recovery properties by LASSO and thresholded LASSO for a particular observation of $\mathbf{X} \in \mathbb{R}^{100 \times 150}$, a particular observation of $\mathbf{Y} \in \mathbb{R}^{100}$ and for k -sparse $\boldsymbol{\beta} \in \mathbb{R}^{150}$ with $\beta_1 = \dots = \beta_{k/2} = 20$ and $\beta_{k/2+1} = \dots = \beta_k = -20$. Our examples are given for $k = 4$ and $k = 30$. For the LASSO estimator, we consider the following setting:

- LASSO with a large tuning parameter $\lambda = 2\sqrt{2\log(150)}$ (as suggested by Candès and Plan [42]).
- LASSO with a small tuning parameter; the one provided by SURE formula, which for a given \mathbf{X} and \mathbf{Y} minimizes the function $0 < \lambda \mapsto \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}(\lambda)\|_2^2 + |\{i \in [p] : \hat{\beta}_i(\lambda) \neq 0\}|$ where $\hat{\boldsymbol{\beta}}(\lambda)$ is the LASSO estimator (see *e.g.* [179] or [181]).

5.7.2 Numerical experiments when the penalty term is the supremum norm

For the ℓ_∞ regularization, the noiseless recovery condition and the accessibility condition depend on $\boldsymbol{\beta}$ through $\text{sign}^\infty(\boldsymbol{\beta}) \in \{-1, *, 1\}^p$. Same as for LASSO, since the distribution of \mathbf{X} is invariant by permutations and sign changes of the columns of \mathbf{X} , then the probability that a non-zero vector having k non-maximal components in absolute value satisfies the noiseless recovery condition is given by

$$\mathbb{P}_{\mathbf{X}}(\widetilde{\mathbf{X}}'(\widetilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1) \text{ where } \widetilde{\mathbf{X}} = (\widetilde{\mathbf{X}}_1 | \mathbf{X}_I) \text{ with } \widetilde{\mathbf{X}}_1 = \sum_{i=1}^{p-k} \mathbf{X}_i \text{ and } I = \{p-k+1, \dots, p\}.$$

An explicit formula for checking the noiseless recovery condition is given in the Appendix. Moreover, the accessibility condition is satisfied with probability

$$\mathbb{P}_{\mathbf{X}}(\min\{\|\boldsymbol{\gamma}\|_\infty : \mathbf{X}\boldsymbol{\gamma} = \widetilde{\mathbf{X}}_1\} = 1).$$

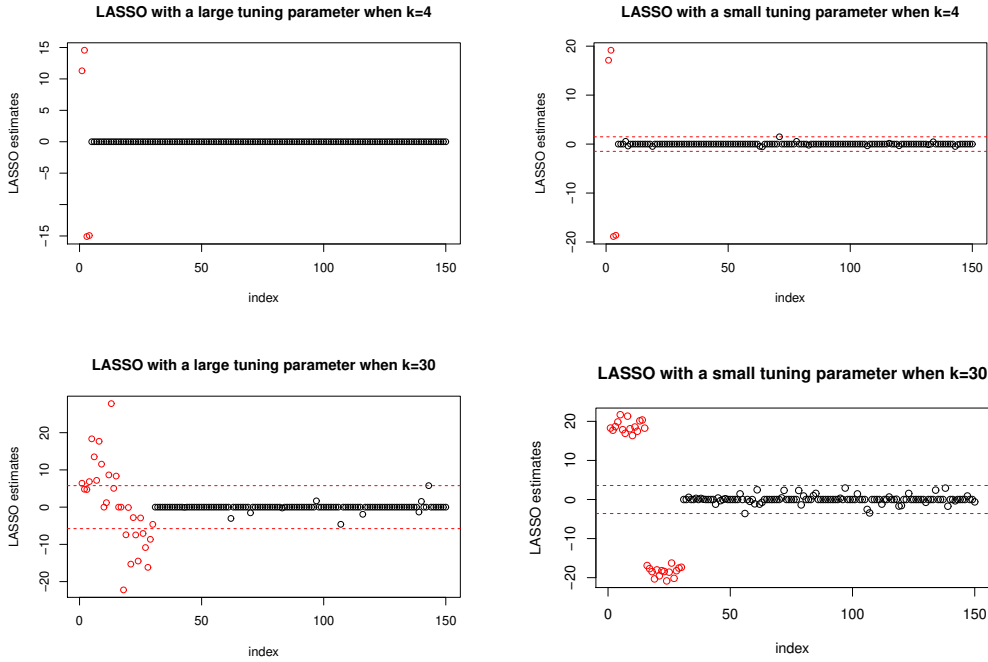


Figure 5.8: Illustrations of sign recovery by LASSO and thresholded LASSO. On the top, when $k = 4$, both the noiseless recovery condition and the accessibility condition hold. Thus, both LASSO and thresholded LASSO can recover the sign of β . With the large tuning parameter $\lambda = 2\sqrt{2\log(150)}$ the sign of β is recovered both by LASSO and thresholded LASSO (top left). When the tuning parameter is small (computed by SURE), some null components of β are not correctly estimated at 0 (black points outside the x-axis), but there exists a threshold, for which the thresholded LASSO recovers the sign of β (top right). On the bottom, when $k = 30$, the accessibility condition holds but the noiseless recovery condition does not hold, thus thresholded LASSO can recover the sign of β but LASSO cannot. When the tuning parameter is large: $\lambda = 2\sqrt{2\log(150)}$, both LASSO and thresholded LASSO fail to recover the sign of β (bottom left). When the tuning parameter is small, some null components of β are not correctly estimated at 0, but there exists a threshold, for which the thresholded LASSO recovers the sign of β (bottom right).

Figure 5.9 provides both the probability of the accessibility condition and the probability of the noiseless recovery condition as functions of k .

In Figure 5.10 we illustrate the pattern recovery properties by a penalized estimator and a thresholded penalized estimator for the supremum norm. Specifically, $\beta \in \mathbb{R}^{150}$ satisfies $\beta_1 = \dots = \beta_{60} = 20$, $\beta_{61} = \dots = \beta_{120} = -20$ and $\beta_{121} = \dots = \beta_{150} = 0$. The tuning parameter is given by the SURE formula, which for a given \mathbf{X} and \mathbf{Y} minimizes the function $0 < \lambda \mapsto \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\hat{\beta}(\lambda)\|_2^2 + |\{i \in [p] : |\hat{\beta}_i| < \|\hat{\beta}(\lambda)\|_\infty\}|$, where $\hat{\beta}(\lambda)$ is the unique element of $S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$ (see *e.g.* [132] or [181]).

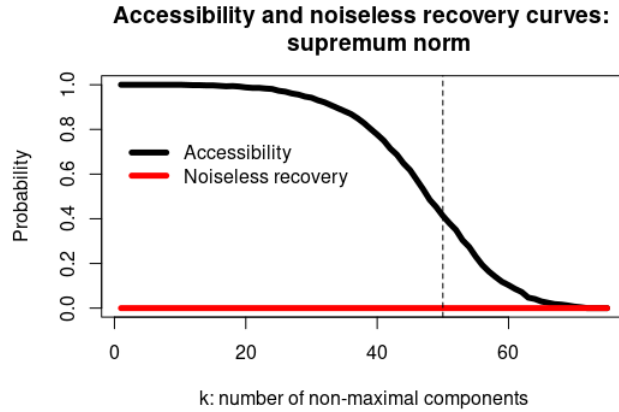


Figure 5.9: Probability of the noiseless recovery and the probability of the accessibility condition as functions of the number k of non-maximal components in absolute value. The value $k = 50$ [4] provides, approximately, the preimage of 0.5 for the accessibility curve.

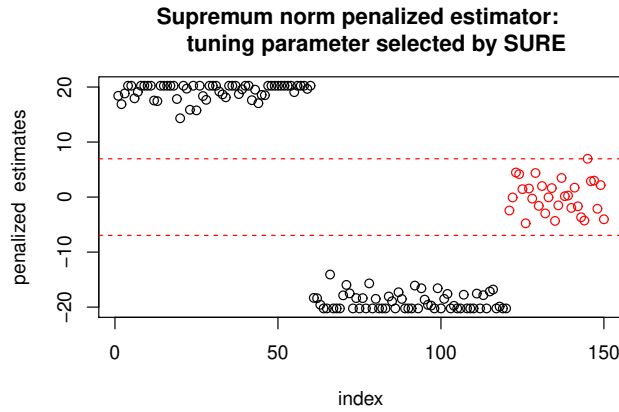


Figure 5.10: Illustrations of the pattern recovery by a penalized estimator and a thresholded penalized for the supremum norm. When $k = 30$, the accessibility condition holds, but the noiseless recovery condition does not hold. Thus, as illustrated on this picture, the recovery of the pattern of β requires thresholding the estimator.

5.8 Appendix

5.8.1 Facts about real-valued polyhedral gauges

Lemma 5.8.1 characterizes the connection between elements of a certain class of convex functions (which encompasses polyhedral gauges) and the faces of a related polytope. It is needed to prove Theorem 5.6.1.

Lemma 5.8.1. *Let $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathbb{R}^p$, $P = \text{conv}\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ and*

$$\phi(\mathbf{x}) = \max\{\mathbf{v}'_1 \mathbf{x}, \dots, \mathbf{v}'_k \mathbf{x}\} \text{ for } \mathbf{x} \in \mathbb{R}^p.$$

Then the subdifferential of ϕ at \mathbf{x} is a face of P and is given by

$$\partial\phi(\mathbf{x}) = \text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\} = \{\mathbf{s} \in P : \mathbf{s}'\mathbf{x} = \phi(\mathbf{x})\}, \text{ where } I_\phi(\mathbf{x}) = \{l \in [k] : \mathbf{v}'_l \mathbf{x} = \phi(\mathbf{x})\}.$$

Conversely, let F be a non-empty face of P . Then $F = \partial\phi(\mathbf{x})$ for some $\mathbf{x} \in \mathbb{R}^p$.

Proof. The fact that $\partial\phi(\mathbf{x}) = \text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\}$ can be found in [101, p. 259]. To prove the second equality, we consider the following. If $l \in I_\phi(\mathbf{x})$, then by definition of $I_\phi(\mathbf{x})$, $\mathbf{v}'_l \mathbf{x} = \phi(\mathbf{x})$ and thus $\mathbf{v}_l \in \{\mathbf{s} \in P : \mathbf{s}'\mathbf{x} = \phi(\mathbf{x})\}$. Since the latter set is convex, one may deduce that

$$\text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\} \subseteq \{\mathbf{s} \in P : \mathbf{s}'\mathbf{x} = \phi(\mathbf{x})\}.$$

For the other inclusion, assume that $\mathbf{s} \in P$ is such that $\mathbf{s} \notin \text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\}$. Then we have $\mathbf{s} = \sum_{l=1}^k \alpha_l \mathbf{v}_l$ where $\alpha_1, \dots, \alpha_k \geq 0$, $\sum_{l=1}^k \alpha_l = 1$ and $\alpha_{l_0} > 0$ for some $l_0 \notin I_\phi(\mathbf{x})$. Since $\mathbf{v}'_l \mathbf{x} \leq \phi(\mathbf{x})$ for all $l \in [k]$ and $\mathbf{v}'_{l_0} \mathbf{x} < \phi(\mathbf{x})$, we also get

$$\mathbf{s}'\mathbf{x} = \sum_{l=1}^k \alpha_l \mathbf{v}'_l \mathbf{x} < \phi(\mathbf{x}).$$

Consequently, $\mathbf{s}'\mathbf{x} \neq \phi(\mathbf{x})$ and thus

$$\{\mathbf{s} \in P : \mathbf{s}'\mathbf{x} = \phi(\mathbf{x})\} \subseteq \text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\}.$$

Therefore, $\partial\phi(\mathbf{x}) = \text{conv}\{\mathbf{v}_l : l \in I_\phi(\mathbf{x})\} = \{\mathbf{s} \in P : \mathbf{s}'\mathbf{x} = \phi(\mathbf{x})\}$. Now we show that the subdifferentials of ϕ are the non-empty faces of P . Let $\mathbf{x} \in \mathbb{R}^p$. By definition of ϕ , $\mathbf{v}'_l \mathbf{x} \leq \phi(\mathbf{x})$ for every $l \in [k]$, so the inequality $\mathbf{x}'\mathbf{s} \leq \phi(\mathbf{x})$ holds for all $\mathbf{s} \in P$. This implies that $\partial\phi(\mathbf{x})$ is a non-empty face of P . Now we show that any non-empty face F of P is equal to $\partial\phi(\mathbf{a})$ for some $\mathbf{a} \in \mathbb{R}^p$. Let $F = \{\mathbf{s} \in P : \mathbf{a}'\mathbf{s} = c\}$ be a non-empty face of P where $\mathbf{a} \in \mathbb{R}^p$, $c \in \mathbb{R}$ and $\mathbf{a}'\mathbf{s} \leq c$ is a valid inequality for all $\mathbf{s} \in P$. We prove that $F = \partial\phi(\mathbf{a})$. Indeed, take any $\mathbf{s} \in F$. We get $\mathbf{a}'\mathbf{s} = c$ as well as $\mathbf{a}'\mathbf{s} \leq \phi(\mathbf{a})$ as shown above, implying that $c \leq \phi(\mathbf{a})$. Analogously, for any $\mathbf{s} \in \partial\phi(\mathbf{a})$, $\mathbf{a}'\mathbf{s} = \phi(\mathbf{a})$ as well as $\mathbf{a}'\mathbf{s} \leq c$ since $\partial\phi(\mathbf{a}) \subseteq P$, yielding $\phi(\mathbf{a}) \leq c$. Therefore one may deduce that $\phi(\mathbf{a}) = c$ and thus $F = \partial\phi(\mathbf{a})$. \square

Corollary 5.8.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\lambda > 0$. Let pen be a polyhedral gauge, i.e., $\text{pen}(\mathbf{x}) = \max\{0, \mathbf{u}'_1 \mathbf{x}, \dots, \mathbf{u}'_l \mathbf{x}\}$ for some $\mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p$. Then the polar set of the unit ball of pen equals*

$$B^* = \text{conv}\{0, \mathbf{u}_1, \dots, \mathbf{u}_l\}.$$

The following lemma, also needed to show Theorem 5.6.1, states that the fitted values are unique over all non-unique solutions of the penalized problem for a given \mathbf{y} . It is a generalization of Lemma 1 in [178], which shows this fact for the special case of the LASSO.

Lemma 5.8.2. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{y} \in \mathbb{R}^n$, $\lambda > 0$ and pen be a polyhedral gauge. Then $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\tilde{\boldsymbol{\beta}}$ and $\text{pen}(\hat{\boldsymbol{\beta}}) = \text{pen}(\tilde{\boldsymbol{\beta}})$ for all $\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{\beta}} \in S_{\mathbf{X}, \text{pen}}(\mathbf{Y})$.*

Proof. Assume that $\mathbf{X}\hat{\boldsymbol{\beta}} \neq \mathbf{X}\tilde{\boldsymbol{\beta}}$ for some $\hat{\boldsymbol{\beta}}, \tilde{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y})$ and let $\check{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}} + \tilde{\boldsymbol{\beta}})/2$. Because the function $\boldsymbol{\mu} \in \mathbb{R}^n \mapsto \|\mathbf{Y} - \boldsymbol{\mu}\|_2^2$ is strictly convex, one may deduce that

$$\|\mathbf{Y} - \mathbf{X}\check{\boldsymbol{\beta}}\|_2^2 < \frac{1}{2}\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}\|_2^2.$$

Moreover, by convexity of pen , we have $\text{pen}(\check{\boldsymbol{\beta}}) \leq \frac{1}{2}(\text{pen}(\hat{\boldsymbol{\beta}}) + \text{pen}(\tilde{\boldsymbol{\beta}}))$. Consequently,

$$\frac{1}{2}\|\mathbf{Y} - \mathbf{X}\check{\boldsymbol{\beta}}\|_2^2 + \lambda \text{pen}(\check{\boldsymbol{\beta}}) < \frac{1}{2}\left(\frac{1}{2}\|\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 + \lambda \text{pen}(\hat{\boldsymbol{\beta}}) + \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}\|_2^2 + \lambda \text{pen}(\tilde{\boldsymbol{\beta}})\right),$$

which contradicts both $\hat{\boldsymbol{\beta}}$ and $\tilde{\boldsymbol{\beta}}$ being minimizers. Finally, $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}\tilde{\boldsymbol{\beta}}$ implies that $\text{pen}(\hat{\boldsymbol{\beta}}) = \text{pen}(\tilde{\boldsymbol{\beta}})$. \square

5.8.2 Proofs

Proof of Theorem 5.1.1

For $\gamma \geq 0$, as P_γ we denote the polyhedron:

$$P_\gamma = \{\mathbf{b} \in \mathbb{R}^p : \text{pen}(\mathbf{b}) \leq \gamma\} = \{\mathbf{b} \in \mathbb{R}^p : \mathbf{u}'_1 \mathbf{b} \leq \gamma, \dots, \mathbf{u}'_l \mathbf{b} \leq \gamma\}.$$

Lemma 5.8.3. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\mathbf{Y} \in \mathbb{R}^n$, $\mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p$ and $R \geq 0$. Then the optimization problem*

$$\min \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|_2^2 \text{ subject to the constraint } \underbrace{\mathbf{u}'_1 \mathbf{b} \leq R, \dots, \mathbf{u}'_l \mathbf{b} \leq R}_{\mathbf{b} \in P_R}. \quad (5.8.1)$$

has at least one solution.

Proof. Let us set $\mathbf{z} = \mathbf{X}\mathbf{b}$ in the optimization problem (5.8.1). Since the set $\mathbf{X}P_R$ is a closed convex set, the problem of minimizing $\|\mathbf{Y} - \mathbf{z}\|_2^2$ subject to the constraint $\mathbf{z} \in \mathbf{X}P_R$ has a unique solution $\bar{\mathbf{z}} \in \mathbf{X}P_R$. Consequently, $\bar{\mathbf{z}} = \mathbf{X}\bar{\mathbf{b}}$ for some $\bar{\mathbf{b}} \in P_R$. Finally, $\bar{\mathbf{b}}$ is a solution of the optimization problem (5.8.1). \square

Let us revisit the polyhedron P_γ . As every polyhedron, P_γ can be decomposed as the sum of its recession cone and a bounded polyhedron [193, Theorem 1.2. and Proposition 1.12.]. Therefore, for $\gamma = 1$ we have:

$$P_1 = \{\mathbf{b} \in \mathbb{R}^p : \mathbf{u}'_1 \mathbf{b} \leq 0, \dots, \mathbf{u}'_l \mathbf{b} \leq 0\} + E,$$

where E is a polytope. Note that, for an arbitrary $\gamma \geq 0$, we have $P_\gamma = P_0 + \gamma E$.

Proof of 5.1.1. Let $m := \inf_{\mathbf{b} \in \mathbb{R}^p} f(\mathbf{b})$. Note that $0 \leq m \leq f(\mathbf{0}) = \frac{1}{2}\|\mathbf{Y}\|_2^2$. We want to show the existence of such $\beta^* \in \mathbb{R}^p$ that $f(\beta^*) = m$. If $f(\mathbf{0}) = m$, then we are done. Assume then that $f(\mathbf{0}) = m + \delta$ for some $\delta > 0$. By the definition of infimum, one may consider such sequence $(\beta_k)_{k \geq 1}$ that $f(\beta_k) \leq m + \frac{\delta}{k} \leq f(\mathbf{0})$ for every $k \geq 1$.

Step 1: Convergence of $(\mathbf{X}\beta_k)_{k \geq 1}$ and $(\text{pen}(\beta_k))_{k \geq 1}$:

At first we show that $(\mathbf{X}\beta_k)_{k \geq 1}$ is bounded. Indeed, we have

$$0 \leq \frac{1}{2}\|\mathbf{Y} - \mathbf{X}\beta_k\|_2^2 \leq f(\beta_k) \leq \frac{1}{2}\|\mathbf{Y}\|_2^2.$$

Thus $\|\mathbf{X}\beta_k\|_2 \leq \|\mathbf{Y}\|_2 + \|\mathbf{X}\beta_k - \mathbf{Y}\|_2 \leq 2\|\mathbf{Y}\|_2$ and the sequence $(\mathbf{X}\beta_k)_{k \geq 1}$ is bounded. Thus, if it is not convergent, then there exist two subsequences $(\mathbf{X}\beta_{r_k})_{k \geq 1} \xrightarrow{k \rightarrow \infty} \mathbf{r}$ and $(\mathbf{X}\beta_{s_k})_{k \geq 1} \xrightarrow{k \rightarrow \infty} \mathbf{s}$ converging to two different limits. By the strict convexity of $\varphi(\mathbf{t}) := \|\mathbf{Y} - \mathbf{t}\|_2^2$, one may deduce the following inequality:

$$\lim_{k \rightarrow \infty} \left\| \mathbf{Y} - \mathbf{X} \frac{\beta_{r_k} + \beta_{s_k}}{2} \right\|_2^2 = \left\| \mathbf{Y} - \frac{\mathbf{r} + \mathbf{s}}{2} \right\|_2^2 < \frac{1}{2} \|\mathbf{Y} - \mathbf{r}\|_2^2 + \frac{1}{2} \|\mathbf{Y} - \mathbf{s}\|_2^2.$$

Moreover, the convexity of pen yields

$$\text{pen} \left(\frac{\beta_{r_k} + \beta_{s_k}}{2} \right) \leq \frac{1}{2} \left(\text{pen}(\beta_{r_k}) + \text{pen}(\beta_{s_k}) \right).$$

Consequently,

$$\limsup_{k \rightarrow \infty} f \left(\frac{\beta_{r_k} + \beta_{s_k}}{2} \right) < \frac{1}{2} \left(f(\beta_{r_k}) + f(\beta_{s_k}) \right) = m,$$

which contradicts $m = \inf_{\mathbf{b} \in \mathbb{R}^p} f(\mathbf{b})$. Therefore the sequence $(\mathbf{X}\boldsymbol{\beta}_k)_{k \geq 1}$ is convergent and thus $(\text{pen}(\boldsymbol{\beta}_k))_{k \geq 1}$ converges, too.

Step 2: Existence of the minimizer:

Denote $\mathbf{g} := \lim_{k \rightarrow \infty} \mathbf{X}\boldsymbol{\beta}_k$ and $\gamma := \lim_{k \rightarrow \infty} \text{pen}(\boldsymbol{\beta}_k)$. Then we have

$$m = \inf_{\mathbf{b} \in \mathbb{R}^p} f(\mathbf{b}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{g}\|_2^2 + \lambda\gamma.$$

Let $\widehat{\boldsymbol{\beta}}$ be an arbitrary solution of (5.8.1) with $R = \gamma$. We are going to prove that $f(\widehat{\boldsymbol{\beta}}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{g}\|_2^2 + \lambda\gamma = m$.

Case 1:

If $\gamma > 0$, then for k large enough such that $\text{pen}(\boldsymbol{\beta}_k) > 0$ one may set $\mathbf{v}_k := \frac{\gamma}{\text{pen}(\boldsymbol{\beta}_k)} \boldsymbol{\beta}_k$. We have $\text{pen}(\mathbf{v}_k) = \gamma$ and thus $\mathbf{v}_k \in P_\gamma$. Consequently, by definition of $\widehat{\boldsymbol{\beta}}$, we have $\|\mathbf{Y} - \mathbf{X}\mathbf{v}_k\|_2^2 \geq \|\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|_2^2$. Therefore

$$f(\widehat{\boldsymbol{\beta}}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|_2^2 + \lambda \text{pen}(\widehat{\boldsymbol{\beta}}) \leq \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{v}_k\|_2^2 + \lambda\gamma \xrightarrow{k \rightarrow \infty} \frac{1}{2} \|\mathbf{Y} - \mathbf{g}\|_2^2 + \lambda\gamma = m.$$

Case 2:

Let $\gamma = 0$. Because $P_1 = P_0 + E$, where E is a bounded polyhedron, one may write $\boldsymbol{\beta}_k = \mathbf{v}_k + \text{pen}(\boldsymbol{\beta}_k)\mathbf{w}_k$ where $\mathbf{v}_k \in P_0$ and $\mathbf{w}_k \in E$. Because $\mathbf{X}\boldsymbol{\beta}_k \rightarrow \mathbf{g}$ and $\text{pen}(\boldsymbol{\beta}_k)\mathbf{w}_k \rightarrow 0$, one may deduce that $\mathbf{X}\mathbf{v}_k \rightarrow \mathbf{g}$. Moreover, since $\text{pen}(\boldsymbol{\beta}_k)\mathbf{v}_k \in P_0$, we have the following inequality:

$$f(\widehat{\boldsymbol{\beta}}) = \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|_2^2 \leq \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\mathbf{v}_k\|_2^2 \xrightarrow{k \rightarrow +\infty} \frac{1}{2} \|\mathbf{Y} - \mathbf{g}\|_2^2 = m,$$

which achieves the proof. \square

Proof of Theorem 5.6.1

Proof of Theorem 5.6.1. (\implies) Assume that there exists a face F of $B^* = \text{conv}\{\mathbf{0}, \mathbf{u}_1, \dots, \mathbf{u}_l\}$ that intersects $\text{row}(\mathbf{X})$ and satisfies $\dim(F) < \dim(\ker(\mathbf{X}))$. By Lemma 5.8.1, $F = \partial \text{pen}(\widehat{\boldsymbol{\beta}})$ for some $\widehat{\boldsymbol{\beta}} \in \mathbb{R}^p$. Let $\mathbf{z} \in \mathbb{R}^n$ with $\mathbf{X}'\mathbf{z} \in F$, which exists by assumption. Now let $\mathbf{y} = \mathbf{X}\widehat{\boldsymbol{\beta}} + \lambda\mathbf{z}$. Note that $\widehat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ since

$$\mathbf{0} \in \partial \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}\|_2^2 + \lambda \text{pen}(\widehat{\boldsymbol{\beta}}) \right) = \mathbf{X}'\mathbf{X}\widehat{\boldsymbol{\beta}} - \mathbf{X}'\mathbf{y} + \lambda \partial \text{pen}(\widehat{\boldsymbol{\beta}}) \quad (5.8.2)$$

$$\iff \frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) = \mathbf{X}'\mathbf{z} \in \partial \text{pen}(\widehat{\boldsymbol{\beta}}). \quad (5.8.3)$$

Now we construct $\widetilde{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ with $\widetilde{\boldsymbol{\beta}} \neq \widehat{\boldsymbol{\beta}}$. According to Lemma 5.8.1, $\partial \text{pen}(\widehat{\boldsymbol{\beta}}) = \text{conv}\{\mathbf{u}_l : l \in I\}$ where $I = I_{\text{pen}}(\widehat{\boldsymbol{\beta}}) = \{l \in [k] : \mathbf{u}'_l \widehat{\boldsymbol{\beta}} = \text{pen}(\widehat{\boldsymbol{\beta}})\}$ and thus $\mathbf{u}'_l \widehat{\boldsymbol{\beta}} < \text{pen}(\widehat{\boldsymbol{\beta}})$ whenever $l \notin I$. Now we show that it is possible to pick $\mathbf{h} \in \ker(\mathbf{X})$ with $\mathbf{h} \neq \mathbf{0}$, but $\mathbf{u}'_l \mathbf{h} = 0$ for all $l \in I$. Then we can make \mathbf{h} small enough such that $\mathbf{u}'_l(\widehat{\boldsymbol{\beta}} + \mathbf{h}) \leq \text{pen}(\widehat{\boldsymbol{\beta}})$ still holds for all $l \notin I$, which in turn implies that $\text{pen}(\widehat{\boldsymbol{\beta}} + \mathbf{h}) = \max\{\mathbf{u}'_l \widehat{\boldsymbol{\beta}} : l \in I\} = \text{pen}(\widehat{\boldsymbol{\beta}})$. This, together with $\mathbf{X}\widehat{\boldsymbol{\beta}} = \mathbf{X}(\widehat{\boldsymbol{\beta}} + \mathbf{h})$, yields $\widehat{\boldsymbol{\beta}} \neq \widetilde{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}} + \mathbf{h} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$. The existence of $\mathbf{0} \neq \mathbf{h} \in \ker(\mathbf{X})$ such that $\mathbf{u}'_l \mathbf{h} = 0$ for every $l \in I$ is equivalent to $\ker(\mathbf{X}) \cap (\text{col}(\mathbf{U}))^\perp \neq \{\mathbf{0}\}$, where $\mathbf{U} = (\mathbf{u}_l)_{l \in I} \in \mathbb{R}^{p \times |I|}$. To prove it, we distinguish two cases:

1) Assume that $\mathbf{0} \in \text{aff}\{\mathbf{u}_l : l \in I\}$. Then $\text{aff}\{\mathbf{u}_l : l \in I\} = \text{col}(\mathbf{U})$ and $\text{rk}(\mathbf{U}) = \dim(F) < \dim(\ker(\mathbf{X}))$. This implies that

$$\dim(\ker(\mathbf{X})) + \dim((\text{col}(\mathbf{U}))^\perp) > p,$$

which proves what was claimed.

2) Assume that $\mathbf{0} \notin \text{aff}\{\mathbf{u}_l : l \in I\}$. It implies that $\mathbf{v} = \mathbf{X}'\mathbf{z} \in \text{row}(\mathbf{X}) \cap \text{conv}\{\mathbf{u}_l : l \in I\}$ satisfies $\mathbf{X}'\mathbf{z} \neq \mathbf{0}$. We also have $\text{rk}(\mathbf{U}) = \dim(\text{aff}\{\mathbf{u}_l : l \in I\}) + 1 = \dim(F) + 1 \leq \dim(\ker(\mathbf{X}))$, which implies that

$$\dim(\ker(\mathbf{X})) + \dim((\text{col}(\mathbf{U}))^\perp) \geq p.$$

If $\ker(\mathbf{X}) \cap (\text{col}(\mathbf{U}))^\perp = \{\mathbf{0}\}$, then $\mathbb{R}^p = \ker(\mathbf{X}) \oplus (\text{col}(\mathbf{U}))^\perp$. But, since $\mathbf{v} \in \text{row}(\mathbf{X}) \cap \text{conv}\{\mathbf{u}_l : l \in I\}$, we also have $\ker(\mathbf{X}) \subseteq \mathbf{v}^\perp$ and $(\text{col}(\mathbf{U}))^\perp \subseteq \mathbf{v}^\perp$, yielding a contradiction and proving the claim.

(\Leftarrow) Assume that there exists $\mathbf{y} \in \mathbb{R}^n$ such that $\hat{\beta}, \tilde{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ with $\hat{\beta} \neq \tilde{\beta}$. Then

$$\frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) \in \partial \text{pen}(\hat{\beta}) \quad \text{and} \quad \frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\tilde{\beta}) \in \partial \text{pen}(\tilde{\beta}).$$

According to Lemma 5.8.2, $\mathbf{X}\hat{\beta} = \mathbf{X}\tilde{\beta}$, thus $\frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) = \frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\tilde{\beta}) \in \text{row}(\mathbf{X})$. Consequently, one may deduce that $\text{row}(\mathbf{X})$ intersects the face $\partial \text{pen}(\hat{\beta}) \cap \partial \text{pen}(\tilde{\beta})$. Let $F^* = \text{conv}\{\mathbf{u}_l : l \in I^*\}$ be a face of $\partial \text{pen}(\hat{\beta}) \cap \partial \text{pen}(\tilde{\beta})$ of the smallest dimension among all faces of $\partial \text{pen}(\hat{\beta}) \cap \partial \text{pen}(\tilde{\beta})$ intersecting $\text{row}(\mathbf{X})$. By minimality of $\dim(F^*)$, $\text{row}(\mathbf{X})$ intersects the relative interior of F^* , namely, there exists $\mathbf{z} \in \mathbb{R}^n$ such that $\mathbf{v} = \mathbf{X}'\mathbf{z}$ lies in F^* , but not on a proper face of F^* . Now we will show that if $\dim(F^*) \geq \dim(\ker(\mathbf{X}))$, then $\text{row}(\mathbf{X})$ intersects a proper face of F^* , yielding a contradiction.

For this, we start with observing that $\dim(F^*) = \dim(\text{aff}\{\mathbf{u}_l : l \in I^*\})$ and that we can write the affine space $\text{aff}\{\mathbf{u}_l : l \in I^*\} = \mathbf{u}_{l_0} + \text{col}(\tilde{\mathbf{U}}^*)$ where $l_0 \in I^*$ and $\tilde{\mathbf{U}}^* = (\mathbf{u}_l - \mathbf{u}_{l_0})_{l \in I^* \setminus \{l_0\}} \in \mathbb{R}^{p \times |I^*| - 1}$. It implies that $\dim(F^*) = \text{rk}(\tilde{\mathbf{U}}^*)$.

Now, let $\mathbf{h} = \hat{\beta} - \tilde{\beta} \neq \mathbf{0}$. By Lemma 5.8.2, $\mathbf{h} \in \ker(\mathbf{X})$ and $\text{pen}(\hat{\beta}) = \text{pen}(\tilde{\beta})$ by Lemma 5.8.2, and. Moreover, since $\mathbf{u}_l \in \partial \text{pen}(\hat{\beta}) \cap \partial \text{pen}(\tilde{\beta})$ for all $l \in I^*$, then by Lemma 5.8.1 we get

$$\mathbf{u}'_l \mathbf{h} = \mathbf{u}'_l \hat{\beta} - \mathbf{u}'_l \tilde{\beta} = \text{pen}(\hat{\beta}) - \text{pen}(\tilde{\beta}) = 0 \quad \forall l \in I^*.$$

Therefore, $\mathbf{h} \in \ker(\mathbf{X}) \cap \text{col}(\mathbf{U}^*)^\perp$, where $\mathbf{U}^* = (\mathbf{u}_l)_{l \in I^*} \in \mathbb{R}^{p \times |I^*|}$.

Assume that $\dim(F^*) \geq \dim(\ker(\mathbf{X}))$. Then

$$\dim(\text{row}(\mathbf{X})) + \dim(\text{col}(\tilde{\mathbf{U}}^*)) \geq \text{rk}(\mathbf{X}) + \dim(\ker(\mathbf{X})) = p.$$

If $\text{row}(\mathbf{X}) \cap \text{col}(\tilde{\mathbf{U}}^*) = \{\mathbf{0}\}$, then $\mathbb{R}^p = \text{row}(\mathbf{X}) \oplus \text{col}(\tilde{\mathbf{U}}^*)$. However, the last equation cannot hold since $\text{row}(\mathbf{X}) = \ker(\mathbf{X})^\perp \subseteq \mathbf{h}^\perp$ and $\text{col}(\tilde{\mathbf{U}}^*) \subseteq \text{col}(\mathbf{U}^*) \subseteq \mathbf{h}^\perp$, where $\mathbf{h} \neq \mathbf{0}$. Consequently, there exists $\mathbf{0} \neq \tilde{\mathbf{v}} \in \text{row}(\mathbf{X}) \cap \text{col}(\tilde{\mathbf{U}}^*)$. The affine line $L = \{\mathbf{X}'\mathbf{z} + t\tilde{\mathbf{v}} : t \in \mathbb{R}\} \subseteq \text{row}(\mathbf{X})$ intersects the relative interior of F^* at $t = 0$. It is contained in $\text{aff}(F^*) = \mathbf{u}_{l_0} + \text{col}(\tilde{\mathbf{U}}^*)$, since $\mathbf{X}'\mathbf{z} \in F^*$ and $\tilde{\mathbf{v}} \in \text{col}(\tilde{\mathbf{U}}^*)$. Therefore, L must intersect a proper face of F^* by Lemma 2.3.1 (ii). But then also $\text{row}(\mathbf{X})$ intersects a proper face of F^* , which yields the required contradiction. \square

Proof of Proposition 4.4.2

The following lemma can be seen as a generalization of Proposition 4.1 from [88] from the ℓ_1 norm to all convex functions.

Lemma 5.8.4. *Let $\beta \in \mathbb{R}^p$ and ϕ be a convex function on \mathbb{R}^p . Then $\text{row}(\mathbf{X})$ intersects $\partial\phi(\beta)$ if and only if, for any $\mathbf{b} \in \mathbb{R}^p$, the following implication holds*

$$\mathbf{X}\beta = \mathbf{X}\mathbf{b} \implies \phi(\beta) \leq \phi(\mathbf{b}). \quad (5.8.4)$$

Proof. Consider the function $\iota_\beta : \mathbb{R}^p \rightarrow \{0, \infty\}$ given by

$$\iota_\beta(\mathbf{b}) = \begin{cases} 0, & \text{when } \mathbf{X}\mathbf{b} = \mathbf{X}\beta, \\ \infty, & \text{else.} \end{cases}$$

Then (5.8.4) holds for any $\mathbf{b} \in \mathbb{R}^p$ if and only if β is a minimizer of the function $\phi(\mathbf{b}) + \iota_\beta(\mathbf{b})$. By definition of the subdifferential it may be shown that $\partial\iota_\beta(\beta) = \text{row}(\mathbf{X})$. Therefore we can deduce that the implication (5.8.4) holds for any $\mathbf{b} \in \mathbb{R}^p$ if and only if

$$\mathbf{0} \in \text{row}(\mathbf{X}) + \partial\phi(\beta) \iff \text{row}(\mathbf{X}) \cap \partial\phi(\beta) \neq \emptyset.$$

□

Proof of Proposition 4.4.2. By Lemma 5.8.4, the geometric characterization of accessible patterns is equivalent to the analytic one. We show the geometric characterization.

(\implies) If the pattern of β is accessible with respect to \mathbf{X} and λpen , then there exists $\mathbf{y} \in \mathbb{R}^n$ and $\hat{\beta} \in S_{\mathbf{X}, \lambda\text{pen}}(\mathbf{y})$ such that $\hat{\beta} \stackrel{\text{pen}}{\approx} \beta$. Because $\hat{\beta} \in S_{\mathbf{X}, \lambda\text{pen}}(\mathbf{y})$, we get $\frac{1}{\lambda}\mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) \in \partial\text{pen}(\hat{\beta}) = \partial\text{pen}(\beta)$, thus $\text{row}(\mathbf{X})$ intersects $\partial\text{pen}(\beta)$.

(\impliedby) If $\text{row}(\mathbf{X})$ intersects the face $\partial\text{pen}(\beta)$, then there exists $\mathbf{z} \in \mathbb{R}^n$ such that $\mathbf{X}'\mathbf{z} \in \partial\text{pen}(\beta)$. For $\mathbf{y} = \mathbf{X}\beta + \lambda\mathbf{z}$, we have $\frac{1}{\lambda}\mathbf{X}'(\mathbf{y} - \mathbf{X}\beta) = \mathbf{X}'\mathbf{z} \in \partial\text{pen}(\beta)$, so that $\beta \in S_{\mathbf{X}, \lambda\text{pen}}(\mathbf{y})$, and the pattern of β is accessible with respect to \mathbf{X} and λpen . □

Proof of Theorem 5.4.4

Lemma 5.8.5. *Let $\phi : \mathbb{R}^p \rightarrow \mathbb{R}$ be the polyhedral gauge defined as*

$$\phi(\mathbf{x}) = \max\{0, \mathbf{u}'_1\mathbf{x}, \dots, \mathbf{u}'_l\mathbf{x}\} \text{ for some } \mathbf{u}_1, \dots, \mathbf{u}_l \in \mathbb{R}^p.$$

If $\partial\phi(\mathbf{x}) = \partial\phi(\mathbf{v})$, then $\partial\phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) = \partial\phi(\mathbf{v})$ for all $\alpha \in [0, 1]$.

Proof. Let $\mathbf{s} \in \partial\phi(\mathbf{x}) = \partial\phi(\mathbf{v})$. Since \mathbf{s} is both a subgradient at \mathbf{x} and at \mathbf{v} , the following inequalities hold

$$\begin{aligned} \phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) &\geq \phi(\mathbf{x}) - (1 - \alpha)\mathbf{s}'(\mathbf{x} - \mathbf{v}) \\ \phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) &\geq \phi(\mathbf{v}) + \alpha\mathbf{s}'(\mathbf{x} - \mathbf{v}). \end{aligned}$$

Multiplying the first inequality by α , the second by $(1 - \alpha)$ and adding them, we get

$$\phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) \geq \alpha\phi(\mathbf{x}) + (1 - \alpha)\phi(\mathbf{v}).$$

Using the convexity of ϕ , we arrive at

$$\phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) = \alpha\phi(\mathbf{x}) + (1 - \alpha)\phi(\mathbf{v}).$$

By Lemma 5.8.1 we have $\partial\phi(\mathbf{x}) = \text{conv}\{\mathbf{u}_l : l \in I\}$, where $I_\phi(\mathbf{x}) = \{l \in [k] : \mathbf{u}'_l\mathbf{x} = \phi(\mathbf{x})\}$. Therefore, if $\mathbf{u}_l \in \partial\phi(\mathbf{x}) = \partial\phi(\mathbf{v})$, then $\mathbf{u}'_l\mathbf{x} = \phi(\mathbf{x})$ and $\mathbf{u}'_l\mathbf{v} = \phi(\mathbf{v})$, thus

$$\mathbf{u}'_l(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) = \alpha\phi(\mathbf{x}) + (1 - \alpha)\phi(\mathbf{v}) = \phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}).$$

Consequently, $\mathbf{u}_l \in \partial\phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v})$. On the other hand, if $\mathbf{u}_l \notin \partial\phi(\mathbf{x})$, then $\mathbf{u}'_l\mathbf{x} < \phi(\mathbf{x})$ and $\mathbf{u}'_l\mathbf{v} < \phi(\mathbf{v})$, thus

$$\mathbf{u}'_l(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}) < \alpha\phi(\mathbf{x}) + (1 - \alpha)\phi(\mathbf{v}) = \phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v}).$$

Consequently, $\mathbf{u}_l \notin \partial\phi(\alpha\mathbf{x} + (1 - \alpha)\mathbf{v})$ and the claim follows. □

Lemma 5.8.6. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\beta \in \mathbb{R}^p$. Then the following set is convex*

$$V_\beta = \{\mathbf{y} \in \mathbb{R}^n : \exists \lambda > 0 \exists \hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y}) \text{ such that } \hat{\beta} \stackrel{\text{pen}}{\sim} \beta\}.$$

Note that V_β may be an empty set.

Proof. Assume that $V_\beta \neq \emptyset$. Let $\mathbf{y}, \tilde{\mathbf{y}} \in V_\beta$. Then there exist $\lambda > 0$ and $\tilde{\lambda} > 0$ such that $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ and $\tilde{\beta} \in S_{\mathbf{X}, \tilde{\lambda} \text{pen}}(\tilde{\mathbf{y}})$ with $\partial \text{pen}(\hat{\beta}) = \partial \text{pen}(\tilde{\beta}) = \partial \text{pen}(\beta)$. Consequently,

$$\mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) \in \lambda \partial \text{pen}(\beta) \text{ and } \mathbf{X}'(\tilde{\mathbf{y}} - \mathbf{X}\tilde{\beta}) \in \tilde{\lambda} \partial \text{pen}(\beta).$$

Let $\alpha \in (0, 1)$ and $\check{\mathbf{y}} = \alpha \mathbf{y} + (1 - \alpha)\tilde{\mathbf{y}}$. Define $\check{\lambda} = \alpha \lambda + (1 - \alpha)\tilde{\lambda}$ and $\check{\beta} = \alpha \hat{\beta} + (1 - \alpha)\tilde{\beta}$. Now we show that $\check{\mathbf{y}} \in V_\beta$. Indeed, observe that

$$\mathbf{X}'(\check{\mathbf{y}} - \mathbf{X}\check{\beta}) = \alpha \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) + (1 - \alpha) \mathbf{X}'(\tilde{\mathbf{y}} - \mathbf{X}\tilde{\beta}) \in \alpha \lambda \partial \text{pen}(\beta) + (1 - \alpha) \tilde{\lambda} \partial \text{pen}(\beta) = \check{\lambda} \partial \text{pen}(\beta).$$

By Lemma 5.8.5, $\partial \text{pen}(\check{\beta}) = \partial \text{pen}(\alpha \hat{\beta} + (1 - \alpha)\tilde{\beta}) = \partial \text{pen}(\beta)$, then also $\check{\beta} \in S_{\mathbf{X}, \check{\lambda} \text{pen}}(\check{\mathbf{y}})$, which proves the claim. \square

Proof of Theorem 5.4.4. Assume that the noiseless recovery condition does not hold for β . Then $\mathbf{X}\beta \notin V_\beta$, where V_β is defined as in Lemma 5.8.6. Consequently, by convexity of V_β , for any realization of $\varepsilon \in \mathbb{R}^n$ we have $\mathbf{X}\beta + \varepsilon \notin V_\beta$ or $\mathbf{X}\beta - \varepsilon \notin V_\beta$. Therefore

$$\begin{aligned} 1 &= \mathbb{P}_\varepsilon(\{\mathbf{X}\beta + \varepsilon \notin V_\beta\} \cup \{\mathbf{X}\beta - \varepsilon \notin V_\beta\}) \\ &\leq \mathbb{P}_\varepsilon(\{\mathbf{X}\beta + \varepsilon \notin V_\beta\}) + \mathbb{P}_\varepsilon(\{\mathbf{X}\beta - \varepsilon \notin V_\beta\}) = 2\mathbb{P}_\varepsilon(\{\mathbf{X}\beta + \varepsilon \notin V_\beta\}). \end{aligned}$$

Consequently,

$$\frac{1}{2} \geq \mathbb{P}_\varepsilon(\{\mathbf{X}\beta + \varepsilon \in V_\beta\}) = \mathbb{P}_\varepsilon(\exists \lambda > 0 \exists \hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{Y}) \text{ such that } \hat{\beta} \stackrel{\text{pen}}{\sim} \beta).$$

\square

Proof of Proposition 5.5.1

Proof of Proposition 5.5.1. We only need to prove the implication (\Leftarrow), as the other implication is obvious. Assume that $\partial \text{pen}(\hat{\beta}) \subseteq \partial \text{pen}(\beta)$. Since $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$, we have $\frac{1}{\lambda} \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\beta}) \in \partial \text{pen}(\hat{\beta}) \subseteq \partial \text{pen}(\beta)$. Consequently, $\text{row}(\mathbf{X})$ intersects $\partial \text{pen}(\beta)$ which implies that the pattern of β is accessible with respect to \mathbf{X} and pen by Proposition 4.4.2. Consequently, there exists $\mathbf{y} \in \mathbb{R}^n$ and there exists $\hat{\beta} \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y})$ for which $\hat{\beta} \stackrel{\text{pen}}{\sim} \beta$. \square

Proof of Theorem 5.5.1

Lemma 5.8.7. *Let pen be a polyhedral gauge on \mathbb{R}^p , $\mathbf{X} \in \mathbb{R}^{n \times p}$, and $\mathbf{v} \in \text{col}(\mathbf{X})$. Let $K_1, K_2 \geq 0$ be large enough such that the set $C = \{\mathbf{b} \in \mathbb{R}^p : \text{pen}(\mathbf{b}) \leq K_1, \|\mathbf{X}\mathbf{b} - \mathbf{v}\|_2 \leq K_2\}$ is non-empty. If $\ker(\mathbf{X}) \cap \ker(\text{pen}) = \{\mathbf{0}\}$, then C is compact.*

Proof. Clearly, C is closed and convex. If $\text{pen}(\mathbf{d}) > 0$ or $\mathbf{X}\mathbf{d} \neq \mathbf{0}$, then $\mathbf{d} \notin \text{rec}(C)$. Consequently, $\text{rec}(C) \subset \ker(\mathbf{X}) \cap \ker(\text{pen}) = \{\mathbf{0}\}$ and thus C is compact. \square

Lemma 5.8.8. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, $\lambda > 0$ and pen be a polyhedral gauge on \mathbb{R}^p . Assume that the uniform uniqueness holds for (5.1.1). Let $\boldsymbol{\beta} \in \mathbb{R}^p$, $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ and set $\mathbf{y}^{(r)} = \mathbf{X}(r\boldsymbol{\beta}) + \boldsymbol{\varepsilon}$. If $\boldsymbol{\beta}$ is accessible with respect to \mathbf{X} and pen , then*

$$\lim_{r \rightarrow \infty} \widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r = \boldsymbol{\beta}.$$

Proof. Since $\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)}) \in S_{\mathbf{X}, \lambda \text{pen}}(\mathbf{y}^{(r)})$, the following inequality holds

$$\frac{1}{2} \|\mathbf{y}^{(r)} - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})\|_2^2 + \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})) \leq \frac{1}{2} \|\mathbf{y}^{(r)} - \mathbf{X}(r\boldsymbol{\beta})\|_2^2 + \lambda \text{pen}(r\boldsymbol{\beta}).$$

$\mathbf{y}^{(r)} - \mathbf{X}(r\boldsymbol{\beta}) = \boldsymbol{\varepsilon}$, therefore one may deduce that

$$\frac{1}{2} \|\boldsymbol{\varepsilon} + \mathbf{X}(r\boldsymbol{\beta}) - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})\|_2^2 + \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})) \leq \frac{1}{2} \|\boldsymbol{\varepsilon}\|_2^2 + \lambda \text{pen}(r\boldsymbol{\beta}) \quad (5.8.5)$$

and in particular

$$\begin{aligned} \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})) &\leq \frac{1}{2} \|\boldsymbol{\varepsilon}\|_2^2 + \lambda \text{pen}(r\boldsymbol{\beta}) \\ \implies \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r) &\leq \frac{\|\boldsymbol{\varepsilon}\|_2^2}{2\lambda r} + \text{pen}(\boldsymbol{\beta}) \\ \implies \limsup_{r \rightarrow \infty} \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r) &\leq \text{pen}(\boldsymbol{\beta}). \end{aligned} \quad (5.8.6)$$

Consequently, the sequence $(\text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r))_{r \in \mathbb{N}_+}$ is bounded. By the Cauchy-Schwarz inequality, the inequality (5.8.5) implies that

$$-\|\boldsymbol{\varepsilon}\|_2 \|\mathbf{X}(r\boldsymbol{\beta}) - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})\|_2 + \frac{1}{2} \|\mathbf{X}(r\boldsymbol{\beta}) - \mathbf{X}\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})\|_2^2 \leq \lambda \text{pen}(r\boldsymbol{\beta}) - \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})) \quad (5.8.7)$$

$$\implies -\frac{\|\boldsymbol{\varepsilon}\|_2}{r} \left\| \mathbf{X} \left(\frac{\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})}{r} - \boldsymbol{\beta} \right) \right\|_2 + \frac{1}{2} \left\| \mathbf{X} \left(\frac{\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})}{r} - \boldsymbol{\beta} \right) \right\|_2^2 \leq \frac{\lambda \text{pen}(\boldsymbol{\beta})}{r} - \frac{\lambda}{r} \text{pen} \left(\frac{\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})}{r} \right). \quad (5.8.8)$$

Let $\alpha \in [0, \infty]$ be the limes superior of the sequence

$$\left(\left\| \mathbf{X} \left(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r - \boldsymbol{\beta} \right) \right\|_2 \right)_{r \in \mathbb{N}}. \quad (5.8.9)$$

By (5.8.8) we get

$$\limsup_{r \rightarrow \infty} \frac{\lambda \text{pen}(\boldsymbol{\beta}) - \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r)}{r} \geq \begin{cases} \alpha^2/2 & \text{if } \alpha < \infty \\ \infty & \text{if } \alpha = \infty. \end{cases}$$

Moreover, since the sequence $(\text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r))_{r \in \mathbb{N}_+}$ is bounded, we obtain

$$\limsup_{r \rightarrow \infty} \frac{\lambda \text{pen}(\boldsymbol{\beta}) - \lambda \text{pen}(\widehat{\boldsymbol{\beta}}(\mathbf{y}^{(r)})/r)}{r} = 0.$$

We conclude that $\alpha = 0$ and the sequence (5.8.9) converges to 0.

Due to uniform uniqueness, by Lemma 5.6.1 we have $\ker(\text{pen}) \cap \ker(\mathbf{X}) = \{\mathbf{0}\}$ and thus, by Lemma 5.8.7, the sequence $(\widehat{\beta}(\mathbf{y}^{(r)})/r)_{r \in \mathbb{N}}$ is bounded. Therefore, to prove that $\lim_{r \rightarrow \infty} \widehat{\beta}(\mathbf{y}^{(r)})/r = \beta$, it suffices to show that β is the unique accumulation point of this sequence. We extract a subsequence $(\widehat{\beta}(\mathbf{y}^{q(r)})/q(r))_{r \in \mathbb{N}}$ converging to $\gamma \in \mathbb{R}^p$. By (5.8.6), one may deduce that $\text{pen}(\gamma) \leq \text{pen}(\beta)$. Moreover, we have

$$0 = \lim_{r \rightarrow \infty} \left\| \mathbf{X} \left(\widehat{\beta}(\mathbf{y}^{q(r)})/q(r) - \beta \right) \right\|_2^2 = \|\mathbf{X}(\gamma - \beta)\|_2^2.$$

Finally, γ satisfies

$$\mathbf{X}\gamma = \mathbf{X}\beta \text{ and } \text{pen}(\gamma) \leq \text{pen}(\beta).$$

Because the pattern of β is accessible, by Lemma 5.8.4 we also have $\text{pen}(\beta) \leq \text{pen}(\gamma)$. Then, Lemma 5.6.1 implies that $\gamma = \beta$. Therefore,

$$\lim_{r \rightarrow \infty} \frac{\widehat{\beta}(\mathbf{y}^{(r)})}{r} = \beta.$$

□

Lemma 5.8.9. *Let pen be a polyhedral gauge on \mathbb{R}^p . Then there exists $\tau_0 > 0$ depending on β such that*

$$\partial \text{pen}(\mathbf{b}) \subseteq \partial \text{pen}(\beta) \text{ for all } \mathbf{b} \in \overline{B}_\infty(\beta, \tau_0).$$

Proof. Let $I = \{l \in [k] : \mathbf{u}'_l \beta = \text{pen}(\beta)\}$. By Lemma 5.8.1, $\partial \text{pen}(\beta) = \text{conv}\{\mathbf{u}_l\}_{l \in I}$. Since

$$\mathbf{u}'_l \beta < \text{pen}(\beta) \quad \forall l \notin I,$$

and by the continuity of pen, one may pick $\tau_0 > 0$ small enough such that

$$\mathbf{u}'_l \mathbf{b} < \text{pen}(\mathbf{b}) \quad \forall l \notin I, \quad \forall \mathbf{b} \in \overline{B}_\infty(\beta, \tau_0).$$

Consequently, for any $\mathbf{b} \in \overline{B}_\infty(\beta, \tau)$, we have $J = \{l \in [k] : \mathbf{u}'_l \mathbf{b} = \text{pen}(\mathbf{b})\} \subseteq I$ and thus

$$\partial \text{pen}(\mathbf{b}) = \text{conv}\{\mathbf{u}_l\}_{l \in J} \subseteq \text{conv}\{\mathbf{u}_l\}_{l \in I} = \partial \text{pen}(\beta).$$

□

Proof of Theorem 5.5.1. By Lemma 5.8.9, there exists $\tau_0 > 0$ such that for any $\mathbf{b} \in \overline{B}_\infty(\beta, \tau_0)$ we have $\partial \text{pen}(\mathbf{b}) \subseteq \partial \text{pen}(\beta)$. By Lemma 5.8.8, $\widehat{\beta}(\mathbf{y}^{(r)})/r$ converges to β when r tends to ∞ . Consequently, we have

$$\exists r_0 \in \mathbb{N} \text{ such that } \forall r \geq r_0, \|\widehat{\beta}(\mathbf{y}^{(r)})/r - \beta\|_\infty \leq \tau_0/2.$$

Consequently, for $r \geq r_0$ we have

$$\begin{cases} \forall \mathbf{b} \in \overline{B}_\infty(\widehat{\beta}(\mathbf{y}^{(r)})/r, \tau_0/2), \partial \text{pen}(\mathbf{b}) \subseteq \partial \text{pen}(\beta) \\ \exists \tilde{\mathbf{b}} \in \overline{B}_\infty(\widehat{\beta}(\mathbf{y}^{(r)})/r, \tau_0/2), \partial \text{pen}(\tilde{\mathbf{b}}) = \partial \text{pen}(\beta). \end{cases}$$

Since for any $t > 0$ and any $\mathbf{x} \in \mathbb{R}^p$, we have $\partial \text{pen}(\mathbf{x}) = \partial \text{pen}(t\mathbf{x})$, one may deduce that

$$\begin{cases} \forall \mathbf{b} \in \overline{B}_\infty(\widehat{\beta}(\mathbf{y}^{(r)}), r\tau_0/2), \partial \text{pen}(\mathbf{b}) \subseteq \partial \text{pen}(\beta) \\ \exists \tilde{\mathbf{b}} \in \overline{B}_\infty(\widehat{\beta}(\mathbf{y}^{(r)}), r\tau_0/2), \partial \text{pen}(\tilde{\mathbf{b}}) = \partial \text{pen}(\beta). \end{cases}$$

Consequently, the claim follows by taking $\tau = r\tau_0/2$. □

Noiseless recovery condition for the supremum norm

Recall the noiseless pattern recovery for the supremum norm:

$$\exists \lambda > 0 \exists \hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \|\cdot\|_\infty}(\mathbf{X}\boldsymbol{\beta}) \text{ such that } \hat{\boldsymbol{\beta}} \stackrel{\|\cdot\|_\infty}{\sim} \boldsymbol{\beta}.$$

Remark 5.8.1. Let $\text{pen} \geq 0$ be minimized at $\mathbf{0}$. Then $\boldsymbol{\beta} = \mathbf{0}$ minimizes the function

$$f(\mathbf{b}) = \frac{1}{2} \|\mathbf{X}\mathbf{b}\|_2^2 + \lambda \text{pen}(\mathbf{b}).$$

Proof. The proof is a direct consequence of $\mathbf{0}$ minimizing both $\|\mathbf{X}\mathbf{b}\|_2$ and $\text{pen}(\mathbf{b})$. \square

Corollary 5.8.2. *The noiseless pattern recovery is satisfied by $\boldsymbol{\beta} = \mathbf{0}$ for every polyhedral gauge pen .*

Proposition 5.8.1. *Let $\mathbf{0} \neq \boldsymbol{\beta} \in \mathbb{R}^p$ and $I := \{i \in [p] : |\beta_i| < \|\boldsymbol{\beta}\|_\infty\}$. Let $\tilde{\mathbf{X}} = (\tilde{\mathbf{X}}_1 | \mathbf{X}_I)$, where*

$$\tilde{\mathbf{X}}_1 = \sum_{i \notin I} \text{sign}(\beta_i) \mathbf{X}_i.$$

Then the noiseless pattern recovery occurs if and only if

$$\tilde{\mathbf{X}}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1,$$

where $\mathbf{e}_1 = (1, 0, \dots, 0)'$.

Before giving the proof, we recall that the subdifferential of the ℓ_∞ norm at $\mathbf{0}$ is the unit ball of the ℓ_1 norm and for $\mathbf{x} \neq \mathbf{0}$ this subdifferential is equal to

$$\begin{aligned} \partial \|\cdot\|_\infty(\mathbf{x}) &= \{\mathbf{s} \in \mathbb{R}^p : \|\mathbf{s}\|_1 \leq 1 \text{ and } \mathbf{s}'\mathbf{x} = \|\mathbf{x}\|_\infty\} \\ &= \left\{ \mathbf{s} \in \mathbb{R}^p : \|\mathbf{s}\|_1 = 1 \text{ and } \begin{cases} s_i x_i \geq 0 & \text{if } |x_i| = \|\mathbf{x}\|_\infty \\ s_i = 0 & \text{otherwise} \end{cases} \right\}. \end{aligned} \quad (5.8.10)$$

Proof. (\implies) Let us assume that there exists $\lambda > 0$ and $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \|\cdot\|_\infty}(\mathbf{X}\boldsymbol{\beta})$ such that $\hat{\boldsymbol{\beta}} \stackrel{\|\cdot\|_\infty}{\sim} \boldsymbol{\beta}$. Then the following property holds

$$\frac{1}{\lambda} \mathbf{X}'\mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \in \partial \|\cdot\|_\infty(\boldsymbol{\beta}) = \partial \|\cdot\|_\infty(\hat{\boldsymbol{\beta}}). \quad (5.8.11)$$

Let us set $\mathbf{c} = (\|\boldsymbol{\beta}\|_\infty, \boldsymbol{\beta}'_I)'$ and $\hat{\mathbf{c}} = (\|\hat{\boldsymbol{\beta}}\|_\infty, \hat{\boldsymbol{\beta}}'_I)'$. By construction, $\mathbf{X}\boldsymbol{\beta} = \tilde{\mathbf{X}}\mathbf{c}$ and $\mathbf{X}\hat{\boldsymbol{\beta}} = \tilde{\mathbf{X}}\hat{\mathbf{c}}$. Consequently, according to (5.8.11), we have

$$\frac{1}{\lambda} \mathbf{X}'\tilde{\mathbf{X}}(\mathbf{c} - \hat{\mathbf{c}}) \in \partial \|\cdot\|_\infty(\boldsymbol{\beta}). \quad (5.8.12)$$

By the representation (5.8.10) of $\partial \|\cdot\|_\infty(\boldsymbol{\beta})$, the relation above implies that

$$\begin{cases} \forall i \in I, \mathbf{X}'_i \tilde{\mathbf{X}}(\mathbf{c} - \hat{\mathbf{c}}) = 0, \\ \frac{1}{\lambda} \sum_{i \notin I} \text{sign}(\beta_i) \mathbf{X}'_i \tilde{\mathbf{X}}(\mathbf{c} - \hat{\mathbf{c}}) = 1, \end{cases}$$

therefore

$$\frac{1}{\lambda} \tilde{\mathbf{X}}' \tilde{\mathbf{X}}(\mathbf{c} - \hat{\mathbf{c}}) = \mathbf{e}_1 \Rightarrow \tilde{\mathbf{X}}(\mathbf{c} - \hat{\mathbf{c}}) = \lambda (\tilde{\mathbf{X}}')^+ \mathbf{e}_1.$$

Using the last implication and (5.8.12), one may deduce that

$$\mathbf{X}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \frac{1}{\lambda} \mathbf{X}' \tilde{\mathbf{X}} (\mathbf{c} - \hat{\mathbf{c}}) \in \partial \|\cdot\|_\infty(\boldsymbol{\beta}).$$

Finally, we prove that $\tilde{\mathbf{X}}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1$:

$$\mathbf{X}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 \in \partial \|\cdot\|_\infty(\boldsymbol{\beta}) \Rightarrow \begin{cases} \forall i \in I \mathbf{X}'_i(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 0 \\ \sum_{i \notin I} \text{sign}(\beta_i) \mathbf{X}'_i(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{X}'_1(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 1 \end{cases} \Rightarrow \tilde{\mathbf{X}}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1.$$

(\Leftarrow) Again, let $\mathbf{c} = (\|\boldsymbol{\beta}\|_\infty, \boldsymbol{\beta}'_I)'$. Denote $\hat{\mathbf{c}} = \mathbf{c} - \lambda \tilde{\mathbf{X}}^+(\tilde{\mathbf{X}}')^+ \mathbf{e}_1$ and set $\hat{\boldsymbol{\beta}}$ as follows:

$$\forall i \in [p], \hat{\beta}_i := \begin{cases} \hat{c}_i & \text{when } i \in I, \\ \text{sign}(\beta_i) \hat{c}_1 & \text{when } i \notin I. \end{cases}$$

Note that $\|\hat{\boldsymbol{\beta}}\|_\infty = c_1 = \|\boldsymbol{\beta}\|_\infty$. Moreover, by definition of I , c_1 is the unique coordinate of \mathbf{c} having the maximal absolute value. Therefore, for λ small enough, we have $\hat{c}_1 = \|\hat{\mathbf{c}}\|_\infty$. It implies that $\|\hat{\boldsymbol{\beta}}\|_\infty = \hat{c}_1$ and this value is attained exactly at coordinates $\hat{\beta}_i : i \notin I$. Therefore $\{i \in [p] : |\hat{\beta}_i| < \|\hat{\boldsymbol{\beta}}\|_\infty\} = \{i \in [p] : |\beta_i| < \|\boldsymbol{\beta}\|_\infty\} = I$. Moreover, for $i \notin I$ we have $\beta_i \hat{\beta}_i = \beta_i \text{sign}(\beta_i) \hat{c}_1 = \|\boldsymbol{\beta}\|_\infty \hat{c}_1 > 0$. Consequently, for λ small enough we have $\text{sign}^\infty(\hat{\boldsymbol{\beta}}) = \text{sign}^\infty(\boldsymbol{\beta})$. To conclude the proof it is enough to show that $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \|\cdot\|_\infty}(\mathbf{X}\boldsymbol{\beta})$, i.e. $\frac{1}{\lambda} \mathbf{X}'(\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}}) \in \partial \|\cdot\|_\infty(\hat{\boldsymbol{\beta}})$. By Remark 2.3.9 $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^+$ is the projection onto $\text{col}(\tilde{\mathbf{X}})$ and $\text{col}((\tilde{\mathbf{X}}')^+) = \text{col}(\tilde{\mathbf{X}})$. Therefore

$$\frac{1}{\lambda} \mathbf{X}' \mathbf{X}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) = \frac{1}{\lambda} \mathbf{X}'(\tilde{\mathbf{X}}\mathbf{c} - \tilde{\mathbf{X}}\hat{\mathbf{c}}) = \mathbf{X}' \tilde{\mathbf{X}} \tilde{\mathbf{X}}^+(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{X}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1.$$

Then it suffices to prove that the latter term belongs to $\partial \|\cdot\|_\infty(\boldsymbol{\beta})$. As we assumed that $\mathbf{X}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1$, we have

$$\begin{aligned} \left(\tilde{\mathbf{X}}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = \mathbf{e}_1 \right) &\Rightarrow \begin{cases} \mathbf{X}'_1(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 1 \\ \forall i \in I \mathbf{X}'_i(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 0 \end{cases} \\ &\Rightarrow \begin{cases} \sum_{i \notin I} \text{sign}(\beta_i) \mathbf{X}'_i(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 1 \\ \forall i \in I \mathbf{X}'_i(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 = 0 \end{cases} \Rightarrow \left(\mathbf{X}'(\tilde{\mathbf{X}}')^+ \mathbf{e}_1 \in \partial \|\cdot\|_\infty(\boldsymbol{\beta}) \right). \end{aligned}$$

Consequently, for $\lambda > 0$ small enough, $\hat{\boldsymbol{\beta}} \in S_{\mathbf{X}, \lambda \|\cdot\|_\infty}(\mathbf{X}\boldsymbol{\beta})$ and $\text{sign}^\infty(\hat{\boldsymbol{\beta}}) = \text{sign}^\infty(\boldsymbol{\beta})$. □

Chapter 6

Maximum likelihood estimation for discrete exponential families and random graphs

6.1 Introduction and preliminaries

The content of this chapter comes from the published paper of Bogdan, Bosy and the author of the dissertation [22].

Exponential families are of paramount importance in probability and statistics. They were introduced by Fisher, Pitman, Darmois and Koopman in 1934-36 and have many properties that make them indispensable in theory and applications, see [125] (Section 2.7), [11] (Chapter 9), [5], [57] (Chapter 9.E), [58], and [121]. In this paper we study *discrete* exponential families, more specifically, exponential families on *finite* sets, and give a new characterization of the existence of the maximum likelihood estimator (MLE) for exponential family and the data at hand. We also present applications, in particular for specific exponential families we give a threshold of the sample size sufficient for the existence of MLE with high probability for *i.i.d.* samples.

The computation of MLE is in general difficult with the number of variables increasing. On the other hand, for given data and an exponential family, MLE may fail to exist. In particular, [50, 51] pointed out to problems with the maximum likelihood estimation when the number of parameters is too large for the sample size. He also gave a sufficient condition for MLE to exist almost surely – the Haar condition.

A complete characterization of the existence of MLE for rather general exponential families was given by Barndorff-Nielsen. Namely, by [11] (Theorem 9.13), MLE for a sample and an exponential family exists if and only if the vector of the sample means calculated for a basis of the linear space of exponents belongs to the interior of the convex hull of the pointwise range of the basis.

This beautiful criterion is alas cumbersome to apply. Therefore, [105] gives an alternative condition for discrete exponential families, together with applications to Cox regression, logistic regression and multiplicative Poisson models. Similar condition is presented by [3] for log-linear model. [96] gives a characterization of the existence of MLE for hierarchical log-linear models. His conditions can be interpreted in terms of polytope geometry, see also [69], and [75]. [36] characterizes the existence of MLE when the log-partition function is steep and regularly convex, and interprets the problem of finding MLE as the optimization of the Kullback-Leibler

divergence. [53] connect the properties of MLE in decomposable models with graph-theoretical notions, thus starting the theory of graphical models in statistics. Sufficient conditions for the existence of MLE in specific exponential families are also given by [169] and [25]. [86] looks for MLE in the closure of convex exponential families, relates the existence of MLE with the linear programming feasibility problem, and in the case of nonexistent MLE, reduces the considered exponential family until MLE exists for the family. He also applies MCMC algorithms to calculate MLE. A comparison between the conditions of Barndorff-Nielsen and Jacobsen is discussed in [114]. In addition, Konis presents an implementation of Jacobsen's test using linear programming. A broad survey of the history of log-linear models and further motivation for the study of the existence of MLE can be found in [74, 75].

The main inspiration for our work is [21] (Theorem 2.3) on the existence of MLE for exponential families of continuous functions on finite interval. In Theorem 6.2.2 below we propose a similar characterization, which is new in the setting of discrete exponential families. We obtain the result by a straightforward, self-contained approach, which does not depend on the delicate convex analysis of [11].

The paper is composed as follows. In Section 6.2 we state and prove our criterion, using the notion of *set of uniqueness*. The criterion is restated in Section 6.2.2 as a linear programming problem. In Section 6.3 we give applications to exponential families spanned by Rademacher and Walsh functions, and to exponential families of random graphs. In particular we give sharp or plain thresholds for the sample size to secure the existence of MLE with high probability. In Appendix 6.6 we give auxiliary results and reformulations of our criterion and pin down its connections with the criterion of Barndorff-Nielsen.

Acknowledgments: We are grateful to Małgorzata Bogdan, Piotr Ciołek, Persi Diaconis, Hélène Massam, Sumit Mukherjee, Krzysztof Oleszkiewicz, Krzysztof Samotij, Maciej Wilczyński and anonymous referees for comments, corrections, references and discussion.

6.1.1 Discrete exponential family

Consider a finite set $\mathcal{X} \neq \emptyset$ and weight function $\mu : \mathcal{X} \rightarrow (0, \infty)$. As usual, $\mathbb{R}^{\mathcal{X}}$ is the family of all the real-valued functions on \mathcal{X} . For $\phi \in \mathbb{R}^{\mathcal{X}}$ we define the *partition* and the *log-partition* functions,

$$Z(\phi) = \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x), \quad \psi(\phi) = \log Z(\phi), \quad (6.1.1)$$

respectively, and the *exponential density*

$$p = e(\phi) = e^{\phi - \psi(\phi)} = e^{\phi} / Z(\phi). \quad (6.1.2)$$

Clearly, $p > 0$ and $\sum_{x \in \mathcal{X}} p(x) \mu(x) = 1$. For arbitrary real number c we have $\psi(\phi + c) = \psi(\phi) + c$, hence

$$e(\phi + c) = e(\phi). \quad (6.1.3)$$

Moreover, for $\phi_1, \phi_2 \in \mathbb{R}^{\mathcal{X}}$ we have $e(\phi_1) = e(\phi_2)$ if and only if $\phi_1 - \phi_2$ is constant. Consider $x_1, \dots, x_n \in \mathcal{X}$, a *sample*. For $\phi \in \mathbb{R}^{\mathcal{X}}$ we denote, as usual,

$$\bar{\phi} = \frac{1}{n} \sum_{i=1}^n \phi(x_i).$$

The *likelihood function* of $p = e(\phi)$ is defined as

$$L_{e(\phi)}(x_1, \dots, x_n) = L_p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i),$$

and the *log-likelihood function* is

$$l_{e(\phi)}(x_1, \dots, x_n) := \log L_{e(\phi)}(x_1, \dots, x_n) = n(\bar{\phi} - \psi(\phi)). \quad (6.1.4)$$

Of course, for every $c \in \mathbb{R}$ we have

$$l_{e(\phi+c)}(x_1, \dots, x_n) = l_{e(\phi)}(x_1, \dots, x_n). \quad (6.1.5)$$

We note that the likelihood functions are uniformly bounded. Indeed, for every $\phi \in \mathbb{R}^{\mathcal{X}}$,

$$\psi(\phi) = \log \sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x) \geq \max_{\mathcal{X}} \phi + \min_{\mathcal{X}} \log \mu, \quad (6.1.6)$$

and so by (6.1.4) and (6.1.6),

$$l_{e(\phi)}(x_1, \dots, x_n) \leq -n \min_{\mathcal{X}} \log \mu \quad \text{and} \quad L_{e(\phi)}(x_1, \dots, x_n) \leq (\min_{\mathcal{X}} \mu)^{-n}. \quad (6.1.7)$$

We fix a linear subspace $\mathcal{B} \subset \mathbb{R}^{\mathcal{X}}$. The *exponential family* spanned by \mathcal{B} is

$$e(\mathcal{B}) := \{p = e(\phi) : \phi \in \mathcal{B}\}. \quad (6.1.8)$$

Since \mathcal{X} is a finite set, $e(\mathcal{B})$ will be called *discrete exponential family* (we do not consider infinite countable sets, for which see [105]).

We call $\hat{p} \in e(\mathcal{B})$ an MLE for x_1, \dots, x_n and $e(\mathcal{B})$ if

$$L_{\hat{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} L_p(x_1, \dots, x_n),$$

or, equivalently,

$$l_{\hat{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n).$$

The following result is well known (see, e.g., [110] (Theorem 2.1) or [57] (p. 177)), but for the reader's convenience we give a proof in Appendix 6.6.1.

Lemma 6.1.1. *If MLE exists, then it is unique.*

Despite the boundedness (6.1.7), MLE may fail to exist, as shown by the following example.

Example 6.1.1. Let $\mathcal{X} = \{0, 1\}$, $\mu \equiv 1$, $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$, $n = 1$ and $x_1 = 1$. Let $a, b \in \mathbb{R}$ and $\phi = a + b\mathbb{1}_{\{1\}}$. Then $Z(\phi) = e^a(1 + e^b)$, $e(\phi) = e^{b\mathbb{1}_{\{1\}}}/(1 + e^b)$, and $L_{e(\phi)}(x_1) = e(\phi)(1) = e^b/(1 + e^b)$. Thus, $\sup L_{e(\phi)}(x_1) = 1$, but the supremum is not attained for any $a, b \in \mathbb{R}$, so MLE does not exist in this case. On the other hand, if $n = 3$, $x_1 = x_2 = 0$, and $x_3 = 1$, then $L_{e(\phi)}(x_1, x_2, x_3) = e^b/(1 + e^b)^3$. By calculus, the maximum is attained when $e^b = 1/2$, therefore $\hat{p} = (2 - \mathbb{1}_{\{1\}})/3$ is the MLE in this case.

We note that the first supremum in Example 6.1.1 is approached when $b \rightarrow \infty$, or for the density $p = \mathbb{1}_{\{1\}}$, which, however, is not in $e(\mathbb{R}^{\mathcal{X}})$ but rather in $e(\mathbb{R}^{\{1\}})$. Below in Theorem 6.2.2 we characterize the situation when the genuine MLE exists, and in Theorem 6.2.5 we treat, by a suitable reduction of \mathcal{X} , the case when the supremum of the likelihood function is ‘‘attained at infinity’’. Before we proceed, we owe the reader some comments on the notation used in this paper and in the literature.

6.1.2 Alternative setting

Let d be a natural number. Consider a nonempty finite set $S \subset \mathbb{R}^d$, weight m on S and the linear space spanned by the coordinate functions on \mathbb{R}^d . The corresponding exponential densities have the form

$$\pi_\theta(y) = e^{\theta \cdot y} / \zeta(\theta), \quad y \in S, \quad (6.1.9)$$

where $\theta \in \mathbb{R}^d$, \cdot is the scalar product in \mathbb{R}^d and $\zeta(\theta) = \sum_{y \in S} e^{\theta \cdot y} m(y)$. Thus, (6.1.9) is a *natural*, or *standard*, exponential family, see [126] or [36]. Since the range of the vector of parameters θ is the whole of \mathbb{R}^d , which is open, the exponential family (6.1.9) is *regular*, see [122] (Appendix D.1). The setting is actually generic, as we explain momentarily. If functions ϕ_1, \dots, ϕ_d span the linear space \mathcal{B} in the general discussion above and we let $T(x) = (\phi_1(x), \dots, \phi_d(x))$ for $x \in \mathcal{X}$, then for every $\phi \in \mathcal{B}$ there is $\theta \in \mathbb{R}^d$ such that $\phi(x) = \theta \cdot T(x)$ for $x \in \mathcal{X}$, and

$$e(\phi) = e^{\theta \cdot T} / Z(\theta \cdot T). \quad (6.1.10)$$

This is the form used by most authors, see [122] or [110], and T is called the *canonical* statistics. Furthermore, we let $S = T(\mathcal{X}) \subset \mathbb{R}^d$ and $m(y) = \sum_{x: T(x)=y} \mu(x)$ for $y \in S$. With the notation of (6.1.9) and (6.1.10) we have

$$\pi_\theta(y) = e(\phi)(x) \quad \text{if} \quad T(x) = y. \quad (6.1.11)$$

If $x_1, \dots, x_n \in \mathcal{X}$ is the sample and we denote $y_1 = T(x_1), \dots, y_n = T(x_n)$, then the corresponding likelihoods are equal, too. Therefore $\pi_{\hat{\theta}}$ is the maximum likelihood estimator for y_1, \dots, y_n and $\{\pi_\theta : \theta \in \mathbb{R}^d\}$ if and only if $e(\hat{\theta} \cdot T)$ is the maximum likelihood estimator for x_1, \dots, x_n and $\{e(\phi) : \phi \in \mathcal{B}\}$. This makes a complete connection between our setting and the setting of natural exponential families with finite support S . The same setting of discrete exponential families on finite set is described, using slightly different language, in [172] (§6.2). We also recall that if ϕ_1, \dots, ϕ_d are affinely independent, then the representation (6.1.10) is *minimal*, see [110] (Chapter 1) or [122], where the affine independence means that $\theta \cdot T = \text{const}$ implies $\theta = 0$. In general, one allows the representation to be nonminimal because over-parametrization is often natural in applications. We shall return to this discussion again in Section 6.6.6, but for now we get back to the setting of \mathcal{B} and (6.1.8). The latter allows to work without coordinates and benefit from properties of specific linear spaces \mathcal{B} , which could otherwise be obscured by an arbitrary choice of T and S .

6.2 Main results

Let $\mathbb{1}$ denote the function on \mathcal{X} identically equal to 1. Assume that $\mathbb{1} \in \mathcal{B}$. This entails no restriction on the considered exponential families $e(\mathcal{B})$, but allows an elegant formulation of the criterion of existence of MLE in terms of \mathcal{B} , in fact in terms of the cone of nonnegative functions in \mathcal{B} :

$$\mathcal{B}_+ := \{\phi \in \mathcal{B} : \phi \geq 0\}.$$

We note in passing that Appendix 6.6.6 gives a reformulation of our criterion for the existence of MLE without requiring that $\mathbb{1} \in \mathcal{B}$.

Let $U \subset \mathcal{X}$. We say that U is a *set of uniqueness* for \mathcal{B} if $\phi = 0$ is the only function in \mathcal{B} such that $\phi = 0$ on U . Similarly, we say that U is a *set of uniqueness* for \mathcal{B}_+ if $\phi = 0$ is the only function in \mathcal{B}_+ such that $\phi = 0$ on U . Put differently, U is of uniqueness for \mathcal{B}_+ if the conditions $\phi \in \mathcal{B}_+$ and $\phi = 0$ on U imply that $\phi = 0$ on \mathcal{X} . Of course, if U is a set of uniqueness for \mathcal{B} , then U is a set of uniqueness for \mathcal{B}_+ .

Example 6.2.1. Let $\mathcal{X} = \{-2, -1, 0, 1, 2\} \subset \mathbb{R}$. Let \mathcal{B} denote the class of all real functions on \mathcal{X} that are of the form $a + bx$ on $\{-2, -1, 0\}$ and $a + cx$ on $\{0, 1, 2\}$ with some $a, b, c \in \mathbb{R}$. Then $\{-1, 2\}$ is a set of uniqueness for \mathcal{B}_+ but $\{-2, 2\}$ is not. We also observe that $\{-1, 2\}$ is not a set of uniqueness for \mathcal{B} , so the nonnegativity of functions in \mathcal{B}_+ plays a role here.

Being a set of uniqueness is a monotone property in the sense that every set larger than a set of uniqueness is also of uniqueness. Furthermore, if U is a set of uniqueness for \mathcal{B}_+ and \mathcal{A} is a linear subspace of \mathcal{B} , then U is of uniqueness for \mathcal{A}_+ .

The following is a crucial definition: For $U \subset \mathcal{X}$ and $\phi \in \mathcal{B}$ we let

$$\lambda_U(\phi) = \max_{\mathcal{X}} \phi - \min_U \phi.$$

Here is our characterization of the existence of MLE for discrete exponential families.

Theorem 6.2.2. *MLE for $e(\mathcal{B})$ and $x_1, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ is of uniqueness for \mathcal{B}_+ .*

Proof. Let us start with the “only if” part. If $U = \{x_1, \dots, x_n\}$ is not a set of uniqueness for \mathcal{B}_+ , then there is a nonzero function $f \in \mathcal{B}_+$ such that $f(x_1) = \dots = f(x_n) = 0$. Let $\phi \in \mathcal{B}$ be arbitrary. Let $\varphi = \phi - f$. We have $\bar{\varphi} = \bar{\phi}$, but $\psi(\varphi) < \psi(\phi)$, where ψ is defined in (6.1.1). So, by (6.1.4), $l_{e(\phi)}(x_1, \dots, x_n) < l_{e(\varphi)}(x_1, \dots, x_n)$. Therefore no $\phi \in \mathcal{B}$ is MLE for x_1, \dots, x_n . To prove the other implication, we let U be a set of uniqueness for \mathcal{B}_+ . By (6.1.4) for $\varphi \in \mathcal{B}$,

$$l_{e(\varphi)}(x_1, \dots, x_n) = n(\bar{\varphi} - \psi(\varphi)) \leq n \left(\frac{1}{n} \left(\min_U \varphi + (n-1) \max_{\mathcal{X}} \varphi \right) - \psi(\varphi) \right).$$

Let $C = \min_{x \in \mathcal{X}} \log \mu(x)$. By (6.1.6), (6.1.5),

$$\begin{aligned} l_{e(\varphi)}(x_1, \dots, x_n) &\leq \min_U \varphi + (n-1) \max_{\mathcal{X}} \varphi - n \max_{\mathcal{X}} \varphi - nC \\ &= -\lambda_U(\varphi) - nC \rightarrow -\infty, \end{aligned}$$

as $\lambda_U(\varphi) \rightarrow \infty$. By Lemma 6.6.1, $\lambda_U(\varphi) \rightarrow \infty$ if $\lambda_{\mathcal{X}}(\varphi) \rightarrow \infty$. In particular, there exists $M > 0$ such that if $\lambda_{\mathcal{X}}(\varphi) > M$, then

$$l_{e(\varphi)}(x_1, \dots, x_n) < l_{e(0)}(x_1, \dots, x_n) = -n \log \mu(\mathcal{X}).$$

By (6.1.5) and continuity, the maximum of $l_{e(\varphi)}(x_1, \dots, x_n)$ is attained on the compact set $\{\varphi \in \mathcal{B} : 0 \leq \varphi \leq M\}$. \square

The above proof is different from that of [21] (Theorem 2.3), [11] (Theorem 9.13) and [172] (Theorem 8.2.1); the use of λ_U makes our arguments more direct.

Remark 6.2.3. By Theorem 6.2.2 we see that the existence of MLE depends on the sequence (x_1, \dots, x_n) only through the set $\{x_1, \dots, x_n\}$. Furthermore, the existence of MLE does not depend on μ , i.e., we may take constant μ without losing generality. Summarizing, the existence of MLE depends only on \mathcal{B} and the set $\{x_1, \dots, x_n\}$. Of course, the actual MLE, say \hat{p} , does depend on the sequence (x_1, \dots, x_n) , the weight μ and \mathcal{B} .

6.2.1 Nonexistence of MLE

In this section we elaborate on the case of nonexistence of MLE in the spirit of [86]. To this end we fix $x_1, \dots, x_n \in \mathcal{X}$ and assume that there is a nontrivial $\delta \in \mathcal{B}_+$ such that $\delta(x_1) = \dots = \delta(x_n) = 0$. By Theorem 6.2.2, $\sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n)$ is not attained at any $p \in e(\mathcal{B})$. However, the supremum is “attained at infinity”, in fact for an exponential density on a proper subset of the state space \mathcal{X} . Indeed, fix δ as above. If $\varphi \in \mathcal{B}$ and $k \in (0, \infty)$, then

$$l_{e(\varphi)}(x_1, \dots, x_n) \leq l_{e(\varphi - k\delta)}(x_1, \dots, x_n),$$

see the first part of the proof of Theorem 6.2.2. Furthermore,

$$\psi(\varphi - k\delta) \rightarrow \log \sum_{x \in \mathcal{X}: \delta(x)=0} e^{\varphi(x)} \mu(x), \quad \text{as } k \rightarrow \infty. \quad (6.2.1)$$

We let $\tilde{\mathcal{X}} = \{x \in \mathcal{X} : \delta(x) = 0\}$ and carrying on with the notation for $\tilde{\mathcal{X}}$ we obtain measure $\tilde{\mu}$, linear space $\tilde{\mathcal{B}}$ with cone $\tilde{\mathcal{B}}_+$, log-partition function $\tilde{\psi}$, likelihood function \tilde{L} , log-likelihood function \tilde{l} and exponential family $e(\tilde{\mathcal{B}})$. Put simpler, we discard $\{x \in \mathcal{X} : \delta(x) > 0\}$ and achieve the following reduction.

Lemma 6.2.1. $\sup_{\tilde{p} \in e(\tilde{\mathcal{B}})} \tilde{l}_{\tilde{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n)$.

Proof. For $\phi \in \mathcal{B}$ we let $\tilde{\phi} = \phi|_{\tilde{\mathcal{X}}}$. Since $\{x_1, \dots, x_n\} \subset \tilde{\mathcal{X}}$,

$$\bar{\phi} = \frac{1}{n} \sum_{i=1}^n \tilde{\phi}(x_i) = \frac{1}{n} \sum_{i=1}^n \phi(x_i) = \bar{\phi}. \quad (6.2.2)$$

Furthermore,

$$\psi(\phi) = \log \left(\sum_{x \in \mathcal{X}} e^{\phi(x)} \mu(x) \right) \geq \log \left(\sum_{x \in \tilde{\mathcal{X}}} e^{\phi(x)} \mu(x) \right) = \tilde{\psi}(\tilde{\phi}).$$

Thus $\bar{\phi} - \psi(\phi) \leq \bar{\phi} - \tilde{\psi}(\tilde{\phi})$, and so

$$\sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n) \leq \sup_{\tilde{p} \in e(\tilde{\mathcal{B}})} \tilde{l}_{\tilde{p}}(x_1, \dots, x_n).$$

Let $\delta \in \mathcal{B}_+$ and k be as in (6.2.1). Using (6.2.1) and (6.2.2),

$$l_{e(\phi - k\delta)}(x_1, \dots, x_n) \rightarrow \tilde{l}_{e(\tilde{\phi})}(x_1, \dots, x_n), \quad \text{as } k \rightarrow \infty.$$

Therefore,

$$\sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n) \geq \sup_{\tilde{p} \in e(\tilde{\mathcal{B}})} \tilde{l}_{\tilde{p}}(x_1, \dots, x_n).$$

□

Motivated by Lemma 6.2.1, we define

$$\{x_1, \dots, x_n\}_{\mathcal{B}_+} = \bigcap \phi^{-1}(\{0\}), \quad (6.2.3)$$

where the intersection is taken over all $\phi \in \mathcal{B}_+$ such that $\phi(x_1) = \dots = \phi(x_n) = 0$. Thus for all $\phi \in \mathcal{B}_+$, if ϕ vanishes on $\{x_1, \dots, x_n\}$, then it vanishes on $\{x_1, \dots, x_n\}_{\mathcal{B}_+}$, and the latter is the largest such set. Put differently, if there is $\delta \in \mathcal{B}_+$ such that $\delta(x_1) = \dots = \delta(x_n) = 0$ but $\delta(x) > 0$, then $x \notin \{x_1, \dots, x_n\}_{\mathcal{B}_+}$, and conversely. In particular, $U \subset \mathcal{X}$ is set of uniqueness for \mathcal{B}_+ if and only if $U_{\mathcal{B}_+} = \mathcal{X}$.

Example 6.2.4. In the setting of Example 6.2.1 we have $\{-2\}_{\mathcal{B}_+} = \{-2\}$ and $\{-1\}_{\mathcal{B}_+} = \{-2, -1, 0\}$.

We note that if $x \notin \{x_1, \dots, x_n\}_{\mathcal{B}_+}$, then there is $\phi \in \mathcal{B}_+$ such that $\phi = 0$ on $\{x_1, \dots, x_n\}$ but $\phi(x) > 0$. Since \mathcal{X} is finite, by adding such functions we can construct $\delta \in \mathcal{B}_+$ that vanishes precisely on $\{x_1, \dots, x_n\}_{\mathcal{B}_+}$, i.e., $\delta^{-1}(\{0\}) = \{x_1, \dots, x_n\}_{\mathcal{B}_+}$. We adopt the setting of Lemma 6.2.1 with this δ , in particular with $\tilde{\mathcal{X}} = \{x_1, \dots, x_n\}_{\mathcal{B}_+}$, and we get the following result.

Theorem 6.2.5. *There is a unique $\tilde{p} \in e(\tilde{\mathcal{B}})$ such that $\tilde{l}_{\tilde{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n)$.*

Proof. By the definition of $\{x_1, \dots, x_n\}_{\mathcal{B}_+}$ and by Theorem 6.2.2, Lemma 6.1.1 and 6.2.1, there is a unique $\tilde{p} \in e(\tilde{\mathcal{B}})$ such that

$$\tilde{l}_{\tilde{p}}(x_1, \dots, x_n) = \sup_{\hat{p} \in e(\tilde{\mathcal{B}})} \tilde{l}_{\hat{p}}(x_1, \dots, x_n) = \sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n).$$

□

Example 6.2.6. For the first sample in Example 6.1.1 we get $\tilde{\mathcal{X}} = \{x_1\}_{\mathcal{B}_+} = \{1\}$, and $\tilde{p} = 1$ on $\tilde{\mathcal{X}}$.

For more substantial applications of Theorem 6.2.5 we refer to Example 6.3.1 and Example 6.3.5.

6.2.2 Linear programming

Before we address special spaces \mathcal{B} , we offer the reader a down-to-earth perspective. To start with, by a comment at the beginning of Section 6.2, we get the following simple result.

Corollary 6.2.1. *If $\{x_1, \dots, x_n\}$ is of uniqueness for \mathcal{B} then MLE exists for $e(\mathcal{B})$ and x_1, \dots, x_n .*

Notably, the condition in Corollary 6.2.1 may be verified by solving the following linear problem:

$$\begin{aligned} \phi &\in \mathcal{B}, \\ \phi(x_1) &= \dots = \phi(x_n) = 0. \end{aligned}$$

Indeed, $\{x_1, \dots, x_n\}$ is of uniqueness for \mathcal{B} if and only if the homogeneous linear system has only the trivial solution. In contrast, Theorem 6.2.2 is a linear programming problem. Indeed, $\{x_1, \dots, x_n\}$ is of uniqueness for \mathcal{B}_+ if and only if the supremum of the (objective) function $\sum_{x \in \mathcal{X}} \phi(x)$ is zero for the class of functions satisfying

$$\begin{aligned} \phi &\in \mathcal{B}, \\ \phi(x_1) &= \dots = \phi(x_n) = 0, \\ \phi &\geq 0. \end{aligned}$$

In this vein [150] (Appendix C) observe that the condition of Barndorff-Nielsen is actually a linear programming problem and make connections to the geometry (of the convex hull of the set S in Section 6.1.2). The linear programming also occurs in the study of the closures of convex exponential families [86] or binary logistic regression models [114]. Furthermore, [185] consider the linear programming in the case when MLE fails to exist. See also [172] for further information on linear programming and cases of nonexistence of MLE for discrete exponential families. Since the linear programming in general runs in polynomial time, see [157], it should be the method of choice when verifying the existence of MLE for discrete exponential families

and data at hand. Having said this, for special linear spaces \mathcal{B} one can come across interesting mathematics, as we demonstrate below. We also remark in passing that the linear problem in Corollary 6.2.1 is the Haar condition of [51] in our setting. Quite generally, the sufficient Haar condition of Crain for the existence of MLE is in the uniqueness of a linear problem while our necessary and sufficient condition is in the uniqueness of a linear-programming problem. The latter is still computationally manageable but more subtle (and optimal); see also the last sentence in Example 6.2.1 for a difference between these two conditions in a very simple setting.

6.3 Applications

Maximization of likelihood is fundamental in estimation, model selection and testing. In many procedures it is important to know if MLE actually exists for given data x_1, \dots, x_n and the linear space of exponents \mathcal{B} ; see [75] (Introduction) for a list of such problems. [75] interpret the existence of MLE by using the geometry of the polyhedral cone spanned by the rows of a specific design matrix. This result is connected with the criterion of [11]. They also inquire which parameters are estimable when MLE is missing.

Below we show that the notion of the set of uniqueness is useful in characterizing the existence of MLE in discrete exponential families for specific spaces \mathcal{B} . There are two types of results we propose:

- (a) conditions for the existence of MLE for a given sample,
- (b) probability bounds for the existence of MLE for independent identically distributed samples.

To this end let \mathcal{X} and \mathcal{B} be as in Section 6.1.1. Let X_1, X_2, \dots be *i.i.d.* random variables with values in \mathcal{X} . We define the random (stopping) time

$$\nu_{\text{uniq}} = \inf\{n \geq 1 : \{X_1, \dots, X_n\} \text{ is a set of uniqueness for } \mathcal{B}_+\}.$$

We will estimate tails of the distribution of ν_{uniq} in terms of \mathcal{X} , \mathcal{B} and n . Typically we are interested in uniformly distributed X_i 's: $\mathbb{P}(X_i = x) = 1/K$, $x \in \mathcal{X}$, $i = 1, 2, \dots$, where $K = |\mathcal{X}|$. In the setting of Theorem 6.2.2 we consider $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$. We fix arbitrary $\mu > 0$ on \mathcal{X} , see Remark 6.2.3. Here is a trivial observation.

Lemma 6.3.1. *MLE for $e(\mathbb{R}^{\mathcal{X}})$ and x_1, \dots, x_n exists if and only if $\{x_1, \dots, x_n\} = \mathcal{X}$.*

Proof. By Theorem 6.2.2 it is enough to verify that \mathcal{X} is the only set of uniqueness for $\mathbb{R}_+^{\mathcal{X}}$. Obviously, \mathcal{X} is a set of uniqueness for $\mathbb{R}_+^{\mathcal{X}}$ (in fact for $\mathbb{R}^{\mathcal{X}}$). On the other hand, if $U \subset \mathcal{X}$ and $x_0 \in \mathcal{X} \setminus U$, then $\mathbb{1}_{x_0}$ vanishes on U but not on \mathcal{X} , hence U is not of uniqueness for $\mathbb{R}_+^{\mathcal{X}}$ (neither it is for $\mathbb{R}^{\mathcal{X}}$). \square

Example 6.3.1. Using notation of Section 6.2.1, we have $U_{\mathcal{B}_+} = U$, for every $U \subset \mathcal{X}$. Clearly, $U \subset U_{\mathcal{B}_+}$. On the other hand, using Equation (6.2.3), one may observe that for every $x \notin U$ the function $\phi(x) = \mathbb{1}_{\{x\}} \in \mathcal{B}_+$ and $\phi = \{0\}$ on U , but $x \notin \phi^{-1}(\{0\})$, so $U_{\mathcal{B}_+} \subset U$. In particular, $\{x_1, \dots, x_n\}_{\mathcal{B}_+} = \{x_1, \dots, x_n\}$ is the new state space $\tilde{\mathcal{X}}$.

Later on we give examples which use the full strength of Theorem 6.2.2 and the nonnegativity of functions in \mathcal{B}_+ therein. For now we propose a probabilistic consequence of Lemma 6.3.1.

Corollary 6.3.1. *Let $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$ and $K = |\mathcal{X}|$. Let X_1, X_2, \dots be independent random variables, each with uniform distribution on \mathcal{X} . Then, for every $c \in \mathbb{R}$,*

$$\lim_{K \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < K \log K + Kc) = e^{-e^{-c}}.$$

Proof. Let $\nu_{\mathcal{X}} = \inf\{n \geq 1 : \{X_1, \dots, X_n\} = \mathcal{X}\}$. The random variable $\nu_{\mathcal{X}}$ yields a connection to the classical Coupon Collector Problem, see [68], and [143]. Namely, by [68],

$$\lim_{K \rightarrow \infty} \mathbb{P}(\nu_{\mathcal{X}} < K \log K + Kc) = e^{-e^{-c}}.$$

By Lemma 6.3.1, $\nu_{\mathcal{X}} = \nu_{\text{uniq}}$, and the proof is complete. \square

We aim to cover with large probability the whole of \mathcal{X} by a sample of suitable size depending on K .

Corollary 6.3.2. *Let $\varepsilon \in (0, 1)$, $K = |\mathcal{X}|$ and $\mathcal{B} = \mathbb{R}^{\mathcal{X}}$. Let X_1, X_2, \dots be independent random variables, each with uniform distribution on \mathcal{X} . If $K \rightarrow \infty$, then*

$$\mathbb{P}(\nu_{\text{uniq}} < (1 - \varepsilon) K \log K) \rightarrow 0 \quad \text{and} \quad \mathbb{P}(\nu_{\text{uniq}} < (1 + \varepsilon) K \log K) \rightarrow 1. \quad (6.3.1)$$

Proof. By Lemma 6.3.1 and Corollary 6.3.1, for every $c \in \mathbb{R}$ we get

$$\begin{aligned} \limsup_{K \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < (1 - \varepsilon) K \log K) &\leq \limsup_{K \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < K \log K + Kc) \\ &= e^{-e^{-c}}. \end{aligned}$$

Thus $\lim_{K \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < (1 - \varepsilon) K \log K) = 0$. The second part of (6.3.1) is obtained analogously. \square

Remark 6.3.2. We summarize (6.3.1) by saying that $K \log K$ is a *sharp threshold* of the sample size for the existence of MLE for $e(\mathbb{R}^{\mathcal{X}})$ and uniform *i.i.d.* samples. Sharp thresholds are widely used in the theory of random graphs, see [67] (Equation 3). It is also convenient to use them here to indicate the minimal size of *i.i.d.* samples that guarantees the existence of MLE with high probability.

6.3.1 Rademacher functions

For $k \in \mathbb{N}$, let us consider $\mathcal{X} = Q_k := \{-1, 1\}^k$, the k -dimensional discrete cube with, say, the uniform weight $\mu(\chi) = 2^{-k}$, $\chi \in Q_k$ (but see Remark 6.2.3). Thus, $K = |\mathcal{X}| = 2^k$. For $j = 1, \dots, k$ and $\chi = (\chi_1, \dots, \chi_k) \in Q_k$ we define the Rademacher functions:

$$r_j(\chi) = \chi_j,$$

and we denote $r_0(\chi) = 1$. Let

$$\mathcal{B}^k = \text{Lin}\{r_0, r_1, \dots, r_k\}.$$

We define, as usual, the exponential family

$$e(\mathcal{B}^k) = \{e(r) : r \in \mathcal{B}^k\}.$$

Theorem 6.3.3. *MLE for $e(\mathcal{B}^k)$ and $x_1, \dots, x_n \in Q_k$ exists if and only if for all $j = 1, \dots, k$ we have $\{r_j(x_1), \dots, r_j(x_n)\} = \{-1, 1\}$.*

Proof. By Theorem 6.2.2 we only need to prove that the above condition characterizes the sets of uniqueness for \mathcal{B}_+^k . If $j \in \{1, \dots, k\}$ is such that $r_j(x_1) = \dots = r_j(x_n) = 1$, then we let $r = r_0 - r_j$. Clearly, $r \in \mathcal{B}_+^k$ and r is not identically zero, but $r(x_i) = 0$ for all $i = 1, \dots, n$. Thus, $\{x_1, \dots, x_n\}$ is not a set of uniqueness for \mathcal{B}_+^k . Similarly, if $r_j(x_1) = \dots = r_j(x_n) = -1$, then we consider the function $r = r_0 + r_j \in \mathcal{B}_+^k$. For the converse implication we consider arbitrary

$$r = \sum_{j=0}^k a_j r_j \in \mathcal{B}_+^k.$$

Let $\chi = -(\text{sign}(a_1), \dots, \text{sign}(a_k))$, where, say, $\text{sign}(0) = 1$. Obviously, $\chi \in Q_k$, and since $r(\chi) \geq 0$, we get

$$a_0 \geq \sum_{j=1}^k |a_j|. \quad (6.3.2)$$

Assume that $r = 0$ on $\{x_1, \dots, x_n\}$. Let $j \in \{1, \dots, k\}$. There are $x, x' \in \{x_1, \dots, x_n\}$ such that $r_j(x) = 1$ and $r_j(x') = -1$. We have

$$0 = r(x) + r(x') = 2a_0 + \sum_{i \neq j} a_i [r_i(x) + r_i(x')].$$

It follows that

$$a_0 \leq \sum_{i \neq j} |a_i|.$$

By (6.3.2), $a_j = 0$, for every $j \geq 1$. Thereby $a_0 = 0$ and $r \equiv 0$. We see that $\{x_1, \dots, x_n\}$ is a set of uniqueness for \mathcal{B}_+^k . \square

Example 6.3.4. Let $x \in Q_k$ be arbitrary. By Theorem 6.3.3, MLE for $e(\mathcal{B}^k)$ and $\{x, -x\}$ exists.

We define the *positive* and *negative half-cubes*, respectively:

$$H_j^+ = \{\chi \in Q_k : r_j(\chi) = 1\}, \quad H_j^- = \{\chi \in Q_k : r_j(\chi) = -1\}, \quad j = 1, \dots, k. \quad (6.3.3)$$

We note that \mathcal{B}^k is also spanned by the indicator functions of half-cubes, namely $\mathbb{1}_j^+ = (r_0 + r_j)/2$ and $\mathbb{1}_j^- = (r_0 - r_j)/2$, $j = 1, \dots, k$.

Corollary 6.3.3. *MLE for $e(\mathcal{B}^k)$ and $x_1, \dots, x_n \in Q_k$ exists if and only if $\{x_1, \dots, x_n\}$ has a nonempty intersection with each half-cube.*

The proof of Corollary 6.3.3 is immediate from Theorem 6.3.3 and the discussion above.

Example 6.3.5. If MLE fails to exist for $e(\mathcal{B}^k)$ and $x_1, \dots, x_n \in Q_k$, then the following analysis may shed some light on Theorem 6.2.5. Let

$$J = \{j \in \{1, \dots, k\} : \{r_j(x_1), \dots, r_j(x_n)\} = \{-1, 1\}\}, \quad J' = \{1, \dots, k\} \setminus J.$$

Since we consider the case when MLE does not exist, by Theorem 6.3.3, $J' \neq \emptyset$. For $j \in J'$ we let

$$H_j = \{\chi \in Q_k : r_j(\chi) = r_j(x_1) = \dots = r_j(x_n)\}.$$

Clearly, this is a half-cube, see (6.3.3). We will show that

$$\{x_1, \dots, x_n\}_{\mathcal{B}_+^k} = \bigcap_{j \in J'} H_j. \quad (6.3.4)$$

We note that for $j \in J'$, r_j is constant on the right-hand side of (6.3.4). Accordingly, the right-hand side of (6.3.4) is isomorphic to $\{-1, 1\}^{|J'|}$ or to $Q_{|J'|}$.

Now if $r = \sum_{j=0}^k a_j r_j \in \mathcal{B}_+^k$ and $r(x_1) = \dots = r(x_n) = 0$, then $r = \sum_{j \in J} a_j r_j + c \geq 0$ on $\{-1, 1\}^{|J'|}$, where $c = a_0 + \sum_{j \in J'} a_j r_j(x_1)$ is the sum of terms which are constant on $\bigcap_{j \in J'} H_j$. In the case when $J = \emptyset$, it is obvious that $\{x_1, \dots, x_n\}_{\mathcal{B}_+^k} = \bigcap_{j \in J'} H_j = \{x_1\}$, since $x_1 = \dots = x_n$. However, if $J \neq \emptyset$, then by definition of J and Theorem 6.3.3 with $k = |J|$, $r = 0$ on $\bigcap_{j \in J'} H_j$. Thus $\bigcap_{j \in J'} H_j \subset \{x_1, \dots, x_n\}_{\mathcal{B}_+^k}$. On the other hand, we observe that for each $j \in J'$, $\mathbb{1}_{H_j^c} = 0$ on the sample and $\mathbb{1}_{H_j^c} > 0$ on H_j^c , hence $H_j^c \cap \{x_1, \dots, x_n\}_{\mathcal{B}_+^k} = \emptyset$ and $\{x_1, \dots, x_n\}_{\mathcal{B}_+^k} \subset \bigcap_{j \in J'} H_j$. By Theorem 6.2.5, MLE exists for $e(\tilde{\mathcal{B}}^k)$ and x_1, \dots, x_n with the measure $\tilde{\mu} := \mu|_{\tilde{\mathcal{X}}}$ on $\tilde{\mathcal{X}} := \bigcap_{j \in J'} H_j$. Of course, $\tilde{\mathcal{X}}$ is isomorphic with $Q_{|J'|}$, if we ignore the J' coordinates of the points in $\tilde{\mathcal{X}}$. In this way we may also think that $\tilde{\mu}$ and x_1, \dots, x_n are on $Q_{|J'|}$. Thus, one may calculate the supremum of the log-likelihood function for $e(\mathcal{B}^k)$, x_1, \dots, x_n and μ as the maximum of a log-likelihood function on $Q_{|J'|}$. Of course, the total mass of $\tilde{\mu}$ is a fraction of that of μ . For instance, if μ is the uniform probability weight on Q_k then $\tilde{\mu}$ is uniform with the total mass $2^{-|J'|}$, which adds $n|J'| \log 2$ to the log-likelihood that would be obtained for $Q_{|J'|}$ with the uniform probability weight, see, e.g., (6.1.2).

Here is a probabilistic application of Theorem 6.3.3.

Corollary 6.3.4. *Let $k \in \mathbb{N}$ and X_1, X_2, \dots, X_n be independent random variables, each with uniform distribution on Q_k . Then,*

$$\begin{aligned} \mathbb{P} \left(\text{MLE exists for } e(\mathcal{B}^k) \text{ and } X_1, \dots, X_n \right) &= \left(1 - \frac{1}{2^{n-1}} \right)^k \\ &\geq 1 - \frac{k}{2^{n-1}} \rightarrow 1, \text{ as } n \rightarrow \infty. \end{aligned}$$

Proof. We have $\mathbb{P}(X_i = x) = 2^{-k}$ for all $x \in Q_k$ and $i = 1, \dots, n$. We let $R_{ij} = r_j(X_i)$ for $i = 1, \dots, n$ and $j = 1, \dots, k$. Thus, $\mathbb{P}(R_{ij} = 1) = \mathbb{P}(R_{ij} = -1) = \frac{1}{2}$ and $\{R_{ij}\}_{i,j}$ are independent. By Theorem 6.3.3,

$$\begin{aligned} &\mathbb{P} \left(\text{MLE exists for } e(\mathcal{B}^k) \text{ and } X_1, \dots, X_n \right) \\ &= \mathbb{P} \left(\{R_{ij} : i = 1, \dots, n\} = \{-1, 1\} \text{ for } j = 1, \dots, k \right) = \left(1 - \frac{2}{2^n} \right)^k. \end{aligned}$$

Applying the Bernoulli inequality finishes the proof. \square

Corollary 6.3.5. *For $k \in \mathbb{N}$ let $X_1, \dots, X_{n(k)}$ be independent random variables, each with uniform distribution on Q_k . If $n(k) = \log_2 k + b + o(1)$ for some $b \in \mathbb{R}$ as $k \rightarrow \infty$, then*

$$\lim_{k \rightarrow \infty} \mathbb{P} \left(\text{MLE exists for } e(\mathcal{B}^k) \text{ and } X_1, \dots, X_{n(k)} \right) = e^{-2^{1-b}}.$$

Proof. By Corollary 6.3.4,

$$\mathbb{P} \left(\text{MLE exists for } e(\mathcal{B}^k) \text{ and } X_1, \dots, X_{n(k)} \right) = \left(1 - \frac{1}{k 2^{b-1+o(1)}} \right)^k$$

$$\rightarrow e^{-2^{1-b}}, \text{ as } k \rightarrow \infty. \quad (6.3.5)$$

□

Corollary 6.3.6. $\log_2 k$ is a sharp threshold of the sample size for the existence of MLE for $e(\mathcal{B}^k)$ and i.i.d. uniform samples on Q_k .

Proof. Let $\varepsilon \in (0, 1)$ and (the sample size) $n = n(k) < (1 - \varepsilon) \log_2 k$. Then,

$$\mathbb{P}(\nu_{\text{uniq}} < n) \leq \mathbb{P}(\nu_{\text{uniq}} < (1 - \varepsilon) \log_2 k).$$

For every $b \in \mathbb{R}$ by the equation in (6.3.5) we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < (1 - \varepsilon) \log_2 k) &\leq \limsup_{k \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < \log_2 k + b) \\ &= e^{-2^{1-b}}. \end{aligned}$$

Since b is arbitrary, we conclude that $\limsup_{k \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} < n(k)) = 0$. Analogously, for the sample size $n = n(k) > (1 + \varepsilon) \log_2 k$ we get

$$\liminf_{k \rightarrow \infty} \mathbb{P}(\nu_{\text{uniq}} > n(k)) = 1,$$

which ends the proof. □

The above is in stark contrast to Corollary 6.3.2, as summarized in Remark 6.3.2. Indeed, in the present setting we have $K = |Q_k| = 2^k$, so the sharp threshold for the sample size needed for the existence of MLE is $\log_2 \log_2 K$. The following result on the expectation of ν_{uniq} agrees well with the sharp threshold.

Lemma 6.3.2. Let ν_{uniq} be as in Corollary 6.3.5. Let $H_k = \sum_{i=1}^k \frac{1}{i}$ be the k -th harmonic number. Then,

$$\frac{H_k}{\log 2} + 1 \leq \mathbb{E}(\nu_{\text{uniq}}) < \frac{H_k}{\log 2} + 2, \quad k = 1, 2, \dots$$

Proof. Observe that $\nu_{\text{uniq}} = \max\{\tau_1, \dots, \tau_k\}$, where

$$\tau_j = \min\{n \geq 1 : \{r_j(X_1), \dots, r_j(X_n)\} = \{-1, 1\}\}, \quad j = 1, \dots, k.$$

From the fact that X_1, X_2, \dots are independent and uniformly distributed, we deduce that

$$\mathbb{1}_{r_j(X_i) \neq r_j(X_1)}, \quad i = 2, 3, \dots, \quad j = 1, 2, \dots,$$

are independent with symmetric Bernoulli distribution. Then τ_1, \dots, τ_k are independent, and

$$\tau_j + 1 \sim \text{Geom}(1/2)$$

for $j = 1, \dots, k$. The result follows from [64]. □

In Section 6.5 we return to Rademacher functions, but for now we turn to exponential families of random graphs, a major motivation for this work.

6.4 Random graphs

In this section we focus on random graphs. Their various applications can be found in [150], [159] and [134]. What is important for us, many such models are indeed discrete exponential families. As usual, maximum likelihood can be used to select a suitable graph model within the exponential family, see, e.g., [142] (Chapter 1 and 8) and [18]. In this section we characterize the existence of MLE in such context. The theory of random graphs started with probabilistic proofs of the existence or nonexistence of specific graphs by Erdős, see, e.g., [28]. Asymptotic properties of random graphs were developed in the seminal papers of [66, 67] and [87]. [150] discuss geometric interpretations of the existence of MLE for discrete exponential families with applications to random graphs and social networks. [45] give normalizing constants that are crucial for the computation of MLE for exponential random graph models. Furthermore, they include examples when MLE fails to exist. The same authors together with Sly discuss in [46] the asymptotic probability of the existence and uniqueness of MLE for the β -model of graphs. This allows to connect the β -model with a random uniform model of graphs with a given degree sequence, which is then explored using graphons (graph limits, see [127]). They also present an algorithm for the computation of MLE in the β -model.

[141] put nonasymptotic conditions for the existence of MLE in various random graph models parameterized by vertex-specific parameters. [151] characterize the existence of MLE for β -models. They interpret the Barndorff-Nielsen's criterion using the geometry of multidimensional polytopes of vertex-degree sequences, see also [75]. [185] transfer the criterion into discrete hierarchical models, using the notion of simplicial complices. These models include, e.g., graphical models and Ising models. Wang, Rauh and Massam also improve the approximation of the set of estimable parameters in the case of the nonexistence of MLE, which is discussed in the setting of marginal polytopes.

Let us start with the notation. Graph is a pair $G = (V, E)$, where $V = \{1, \dots, N\}$, $N \in \mathbb{N}$, is the set of nodes and E is the set of edges, i.e.,

$$E \subset \binom{V}{2} := \{(r, s) : 1 \leq r < s \leq N\}.$$

We only consider simple undirected graphs (containing no loops or multiple edges). Let $m = |E|$. If $m = \binom{N}{2}$, then the graph is called complete and is denoted as K_N . On the other hand, the empty graph (with $m = 0$) is denoted as $\overline{K_N}$. For graphs $G = (V, E_1)$ and $H = (V, E_2)$ we let, as usual,

$$G \cup H := (V, E_1 \cup E_2), \quad G \cap H := (V, E_1 \cap E_2).$$

Furthermore, $G \subset H$ means that $E_1 \subset E_2$. Let \mathcal{G}_N be the family of all the graphs with N nodes, i.e., with $V = \{1, \dots, N\}$. By a random graph we understand a random variable \mathbb{G} with values in \mathcal{G}_N . The families of distributions of such random variables are called random graph models. We focus on the exponential model of random graphs $\mathcal{G}_{N,c}$ defined as follows.

For $1 \leq r < s \leq N$ and $G \in \mathcal{G}_N$, we let

$$\mathbb{1}_G(r, s) = \begin{cases} 1, & \text{if } (r, s) \in E, \\ 0, & \text{otherwise.} \end{cases}$$

We define $\chi_{r,s} : \mathcal{G}_N \rightarrow \{-1, 1\}$ by $\chi_{r,s}(G) = 1 - 2\mathbb{1}_G(r, s)$. We consider the linear space

$$\mathcal{B}^{\mathcal{G}_N} = \text{Lin} \left\{ 1, \chi_{r,s}(G) : 1 \leq r < s \leq N \right\}.$$

Let $c \in \mathbb{R}^{\binom{V}{2}}$ be a corresponding vector of coefficients. Following the setting of Section 6.1.1 we let $\mu(G) = 1$ for each $G \in \mathcal{G}_N$ (but see Remark 6.2.3) and consider the exponential family

$$\mathcal{G}_{N,c} := e(\mathcal{B}^{\mathcal{G}_N}) = \left\{ p_c := e^{\phi_c - \psi(\phi_c)} : c \in \mathbb{R}^{\binom{V}{2}} \right\}, \quad (6.4.1)$$

where

$$\phi_c(G) = \sum_{(r,s) \in \binom{V}{2}} c_{r,s} \chi_{r,s}(G), \quad \psi(\phi_c) = \log \sum_{G \in \mathcal{G}_N} e^{\phi_c(G)},$$

for $G \in \mathcal{G}_N$, see also (6.1.3). As usual, for $p_c \in \mathcal{G}_{N,c}$ we let $L_{p_c}(G_1, \dots, G_n) = \prod_{i=1}^n p_c(G_i)$, etc.

Lemma 6.4.1. *Let $c \in \mathbb{R}^{\binom{V}{2}}$ and let \mathbb{G} be a random graph with distribution $\mathcal{G}_{N,c}$. Let $1 \leq r < s \leq N$. Then the probability of the appearance of the edge (r, s) in \mathbb{G} equals*

$$p_{r,s} = \frac{e^{c_{r,s}}}{1 + e^{c_{r,s}}}. \quad (6.4.2)$$

The result is well known but for convenience a proof is given in Appendix 6.6.3.

Lemma 6.4.2. *Let $c \in \mathbb{R}^{\binom{V}{2}}$ and let \mathbb{G} be a random graph with distribution $\mathcal{G}_{N,c}$. Let $1 \leq r_1, s_1, r_2, s_2 \leq N$, $r_1 < s_1, r_2 < s_2$, and $(r_1, s_1) \neq (r_2, s_2)$. Then the appearances of edges (r_1, s_1) and (r_2, s_2) in \mathbb{G} are independent events.*

The proof of the result is similar to that of Lemma 6.4.1, and can be found in Appendix 6.6.4. For instance, if $p_{r,s} = p \in (0, 1)$ for every edge (r, s) , then the exponential random graph with distribution $\mathcal{G}_{N,c}$ is the Erdős-Rényi random graph $\mathcal{G}_{N,p}$ in [66, 67]. The latter means that $\mathbb{P}(e \in E(\mathbb{G})) = p$ for every edge $e \in \binom{V}{2}$, and the events $e \in E(\mathbb{G})$ and $f \in E(\mathbb{G})$ are independent for different edges e, f .

Theorem 6.4.1. *MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ and $G_1, \dots, G_n \in \mathcal{G}_N$ exists if and only if*

$$\bigcup_{i=1}^n G_i = K_N \quad \text{and} \quad \bigcap_{i=1}^n G_i = \overline{K_N}.$$

Proof. By Theorem 6.2.2, MLE exists if and only if $\{G_1, \dots, G_n\}$ is of uniqueness for $\mathcal{B}_+^{\mathcal{G}_N}$.

We first prove the ‘‘only if’’ part of Theorem 6.4.1. Let us assume that there exists an edge $(r_0, s_0) \notin \bigcup_{i=1}^n G_i$. Then the function $\chi_{r_0, s_0} \in \mathcal{B}_+^{\mathcal{G}_N}$ equals zero on G_1, \dots, G_n , but not on the whole \mathcal{G}_N . In addition, if there is an edge $(r_0, s_0) \in \bigcap_{i=1}^n G_i$, then the function $(1 + \chi_{r_0, s_0}) \in \mathcal{B}_+^{\mathcal{G}_N}$ vanishes for G_1, \dots, G_n , but it is not equal to zero, e.g., for the graph $\overline{K_N}$.

We next prove the ‘‘if’’ part of the theorem. Let $\phi = k_0 + \sum_{r < s} k_{r,s} \chi_{r,s} \in \mathcal{B}_+^{\mathcal{G}_N}$, where $k_0, k_{r,s} \in \mathbb{R}$ for all $1 \leq r < s \leq N$. Since $\phi(G) \geq 0$ for every $G \in \mathcal{G}_N$,

$$k_0 \geq \sum_{r < s} |k_{r,s}|. \quad (6.4.3)$$

Let $(r_0, s_0) \in \binom{V}{2}$. Let $\phi(G_1) = \dots = \phi(G_n) = 0$. Since $\bigcup_{i=1}^n G_i = K_N$ and $\bigcap_{i=1}^n G_i = \overline{K_N}$, there exists a pair of graphs $G', G'' \in \{G_1, \dots, G_n\}$ such that $\chi_{r_0, s_0}(G') = 1$, $\chi_{r_0, s_0}(G'') = -1$. Therefore,

$$0 = \phi(G') + \phi(G'') = 2k_0 + \sum_{r < s} k_{r,s} (\chi_{r,s}(G') + \chi_{r,s}(G''))$$

$$= 2k_0 + \sum_{\substack{r < s \\ (r,s) \neq (r_0,s_0)}} k_{r,s} (\chi_{r,s}(G') + \chi_{r,s}(G'')).$$

It follows that $k_0 \leq \sum_{(r,s) \neq (r_0,s_0)} |k_{r,s}|$ and eventually we get $k_{r_0,s_0} = 0$, thanks to (6.4.3). Since (r_0, s_0) is arbitrary, $k_{r,s} = 0$ for every $1 \leq r < s \leq N$. Then also $c_0 = 0$, and thus $\phi \equiv 0$. \square

In the above random graph model it is possible to compute explicitly the probability of the existence of MLE for *i.i.d.* samples of graphs in \mathcal{G}_N . To this end, for $1 \leq r < s \leq N$ we fix $c_{r,s} \in \mathbb{R}$. By Lemma 6.4.1 the probability of the appearance of the edge (r, s) in random graph \mathbb{G} with distribution $\mathcal{G}_{N,c}$ is

$$p_{r,s} = \frac{e^{c_{r,s}}}{1 + e^{c_{r,s}}}.$$

Lemma 6.4.3. *Let $\{\mathbb{G}_1, \dots, \mathbb{G}_n\}$ be *i.i.d.* with distribution $\mathcal{G}_{N,c}$. Then the probability of the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ equals*

$$\prod_{1 \leq r < s \leq N} \left(1 - p_{r,s}^n - (1 - p_{r,s})^n\right). \quad (6.4.4)$$

Proof. By Theorem 6.4.1, MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ exists if and only if among the random graphs $\mathbb{G}_1, \dots, \mathbb{G}_n$ every edge (r, s) , $1 \leq r < s \leq N$, appears at least once, but not n times. For every edge (r, s) the above condition is satisfied with probability $1 - (1 - p_{r,s})^n - p_{r,s}^n$. The independence of the occurrences of different edges in $\mathcal{G}_{N,c}$ yields the product (6.4.4). \square

In particular, if $c = 0$, then the probability of the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ equals

$$\left(1 - 2^{1-n}\right)^{\binom{N}{2}},$$

which is an analogue of Corollary 6.3.5. From the above results we can deduce asymptotic bounds for the *i.i.d.* sample size for which MLE exists with high probability. To this end we recall the classical result on $p = p(N) \in (0, 1)$ such that \mathbb{G} from $\mathcal{G}_{N,p}$ has at least one edge with high probability.

Remark 6.4.2. [80] (Lemma 1.10) Let $\mathbb{G}_{N,p(N)}$ be a random graph with distribution $\mathcal{G}_{N,p(N)}$. Then

$$\lim_{N \rightarrow \infty} \mathbb{P} \left(\mathbb{G}_{N,p(N)} \text{ has at least one edge} \right) = \begin{cases} 0 & \text{if } p(N) = o(N^{-2}), \\ 1 & \text{if } N^{-2} = o(p(N)). \end{cases}$$

The above may be summarized by saying that N^{-2} is a *threshold* for the probability p such that \mathbb{G} with distribution $\mathcal{G}_{N,p}$ has at least one edge. For more information on threshold functions in the theory of random graphs see, e.g., [80]. In particular, a sharp threshold (mentioned previously) is a threshold but the converse is not true in general.

Lemma 6.4.4. *Let $\mathbb{G}_1, \dots, \mathbb{G}_n$ be *i.i.d.* random variables with distribution $\mathcal{G}_{N,c}$. Then $\log N$ is a threshold of the sample size n for the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$.*

Proof. According to Lemma 6.4.3, the probability of the existence of MLE for $e(\mathcal{B}^{\mathcal{G}_N})$ and $\mathbb{G}_1, \dots, \mathbb{G}_n$ equals

$$P_{\text{MLE}} = \prod_{1 \leq r < s \leq N} \left(1 - p_{r,s}^n - (1 - p_{r,s})^n\right).$$

We define the function

$$f(x) = 1 - x^w - (1 - x)^w, \quad x \in (0, 1), \quad w \geq 2. \quad (6.4.5)$$

Clearly, $f(x) = f(1 - x)$ and for $w \geq 2$ we have f increasing when $0 < x < \frac{1}{2}$ and decreasing when $\frac{1}{2} < x < 1$. Using (6.4.5) we can bound P_{MLE} from above by

$$P_{\text{BIG}} := \left(1 - 2^{1-n}\right)^{\binom{N}{2}}.$$

Applying Corollary 6.3.4 and the equality in (6.3.5) for $k = \binom{N}{2}$, we observe that for every $b \in \mathbb{R}$ and for $n = n(N) = \log_2 \binom{N}{2} + b + o(1)$ we have $P_{\text{BIG}} \rightarrow e^{-2^{1-b}}$, as $N \rightarrow \infty$. Therefore, for $n(N) = o(\log N)$ we obtain $P_{\text{MLE}} \leq P_{\text{BIG}} \rightarrow 0$, as $N \rightarrow \infty$.

We consider the sample size $n = n(N)$ (depending on N). We will prove that if $\log N/n \rightarrow 0$ as $N \rightarrow \infty$, then $P_{\text{MLE}} \rightarrow 1$. To this end we bound P_{MLE} from below by

$$P_{\text{SMALL}} := \left(1 - p_{\text{max}}^n - (1 - p_{\text{max}})^n\right)^{\binom{N}{2}},$$

where $c_{\text{max}} = \max_{1 \leq r < s \leq N} |c_{r,s}|$ and $p_{\text{max}} = e^{c_{\text{max}}} / (1 + e^{c_{\text{max}}})$.

Take n independent Erdős-Rényi random graphs $\mathbb{H}_1, \dots, \mathbb{H}_n$ with distribution $\mathcal{G}_{N, p_{\text{max}}}$. Then the probability of the existence of MLE for $e(\mathcal{B}^{\mathcal{G}^N})$ and for $\mathbb{H}_1, \dots, \mathbb{H}_n$ equals exactly P_{SMALL} . Note that intersection and union of the graphs are also Erdős-Rényi random graphs, namely

$$\bigcap_{i=1}^n \mathbb{H}_i \sim \mathcal{G}_{N, p_{\text{max}}^n}, \quad \bigcup_{i=1}^n \mathbb{H}_i = \overline{\bigcap_{i=1}^n \overline{\mathbb{H}_i}} \sim \mathcal{G}_{N, 1 - q_{\text{max}}^n},$$

where

$$q_{\text{max}} := 1 - p_{\text{max}} = \frac{e^{-c_{\text{max}}}}{1 + e^{-c_{\text{max}}}}.$$

From Remark 6.4.2, with high probability we have

$$\bigcap_{i=1}^n \mathbb{H}_i = \overline{K_N} \quad \text{and} \quad \bigcup_{i=1}^n \mathbb{H}_i = \overline{K_N},$$

provided

$$p_{\text{max}}^n = o(N^{-2}) \quad \text{and} \quad q_{\text{max}}^n = o(N^{-2}).$$

By definition, $c_{\text{max}} > 0$, so $p_{\text{max}} > q_{\text{max}}$. In order to get $P_{\text{SMALL}} \rightarrow 1$ as $n \rightarrow \infty$, it suffices to have $p_{\text{max}}^n = o(N^{-2})$. If $n(N)/\log N \rightarrow \infty$ as $N \rightarrow \infty$, then the above condition is satisfied. Therefore $\log N$ is a threshold of the sample size for existence of MLE for $e(\mathcal{B}^{\mathcal{G}^N})$ and independent $\mathbb{G}_1, \dots, \mathbb{G}_n$ from $\mathcal{G}_{N,c}$. \square

6.5 Applications to Walsh functions

We return to Rademacher functions to discuss the spaces spanned by their products. Let $k \in \mathbb{N}$, $1 \leq q \leq k$, and

$$\mathcal{B}_q^k = \text{Lin} \{w_S : S \subset \{1, \dots, k\} \text{ and } |S| \leq q\},$$

where

$$w_S(x) = \prod_{i \in S} r_i(x), \quad x \in Q_k, \quad S \subset \{1, \dots, k\},$$

are the Walsh functions, see, e.g., [108].

The case $\mathcal{B}_1^k = \mathcal{B}^k$ was discussed in Section 6.3.1 and the case $q = 2$ is related to the Ising model of ferromagnetism in statistical mechanics, see [184] (Example 3.1).

Lemma 6.5.1. *The dimension of the linear space \mathcal{B}_q^k is $\sum_{j=0}^q \binom{k}{j}$.*

The proof of Lemma 6.5.1 is given in Appendix 6.6.5.

Corollary 6.5.1. *For $q \leq \frac{k}{2}$ we have*

$$\dim(\mathcal{B}_q^k) \leq 2^{kH_2(\frac{q}{k})} \leq \left(\frac{ek}{q}\right)^q,$$

where $H_2(p) = -p \log_2 p - (1-p) \log_2(1-p)$ is the binary entropy function.

The proof follows from Lemma 6.5.1 and the entropy bound for the sum of binomial coefficients, see, e.g., [84] (Theorem 3.1).

Characterization of the existence of MLE for $e(\mathcal{B}_q^k)$ and the related sharp thresholds seem to be hard for general q , even for $q = 2$, see Remark 6.5.2. In the next section we discuss the products of $k - q$ Rademacher functions for fixed $q \in \mathbb{N}$ ($q \leq k$). We especially focus on the products of $k - 1$ and k Rademacher functions. Below we characterize the existence of MLE for $e(\mathcal{B}_{k-1}^k)$. As we will see, we get a qualitatively different result than that in Section 6.3.1. Let \mathcal{E} and \mathcal{O} be the sets of all those points in Q_k that have an even and odd number of positive coordinates, respectively.

Theorem 6.5.1. *MLE exists for $e(\mathcal{B}_{k-1}^k)$ and $x_1, \dots, x_n \in Q_k$ if and only if \mathcal{E} or $\mathcal{O} \subset \{x_1, \dots, x_n\}$.*

Proof. Thanks to Theorem 6.2.2, we only need to characterize the sets of uniqueness for $(\mathcal{B}_{k-1}^k)_+$. To this end, we consider the hypercube G_{Q_k} , defined as the graph with vertices in Q_k and edges between all pairs of points which differ by exactly one coordinate. Thus,

$$V(G_{Q_k}) = Q_k \text{ and } E(G_{Q_k}) = \{\{x, y\} \in Q_k \times Q_k : |\{j : r_j(x) \neq r_j(y)\}| = 1\}.$$

Let $U = \{x_1, \dots, x_n\}$. Assume that U is a set of uniqueness. Let $e \in \mathcal{E}$ and $o \in \mathcal{O}$. The hypercube graph G_{Q_k} is connected, so there exists a path $(e, v_1, v_2, \dots, v_{2p}, o)$ in G_{Q_k} . Then

$$\left(\mathbb{1}_{\{e, v_1\}} + \mathbb{1}_{\{v_2, v_3\}} + \dots + \mathbb{1}_{\{v_{2p}, o\}}\right) - \left(\mathbb{1}_{\{v_1, v_2\}} + \mathbb{1}_{\{v_3, v_4\}} + \dots + \mathbb{1}_{\{v_{2p-1}, v_{2p}\}}\right) \quad (6.5.1)$$

$$= \mathbb{1}_{\{e\}} + \mathbb{1}_{\{o\}} \quad (6.5.2)$$

is a nontrivial nonnegative function on Q_k . Therefore, we must have $\{e, o\} \cap U \neq \emptyset$. Then we easily conclude that $\mathcal{E} \subset U$ or $\mathcal{O} \subset U$.

For the converse implication, we consider $q \in \{0, \dots, k\}$ and $(k - q)$ -subcubes defined as follows,

$$\bigcap_{i=1}^q H_{j_i}, \quad (6.5.3)$$

where $1 \leq j_1 < j_2 < \dots < j_q \leq k$ and $H_{j_i} = H_{j_i}^+$ or $H_{j_i}^-$, see (6.3.3). When $q = k - 1$, the intersection, or a 1-cube, is a pair of points in Q_k which differ by exactly one coordinate, so they have a different parity. Moreover, each such pair can be obtained in this way. Using (6.5.3), as in the proof of Lemma 6.5.1 we see that $\mathbb{1}_{\{e, o\}} \in \mathcal{B}_{k-1}^k$ for each $e \in \mathcal{E}$ and $o \in \mathcal{O}$. Furthermore,

each q -subcube of Q_k with $q \geq 1$ can be covered by disjoint pairs $\{e, o\}$ as above. Therefore, the functions $\mathbb{1}_{\{e, o\}} \in \mathcal{B}_{k-1}^k$ with $e \in \mathcal{E}$ and $o \in \mathcal{O}$ span the linear space \mathcal{B}_{k-1}^k .

We next claim that for every $f \in \mathcal{B}_{k-1}^k$,

$$\sum_{x \in \mathcal{O}} f(x) = \sum_{x \in \mathcal{E}} f(x). \quad (6.5.4)$$

Indeed, if $f = \mathbb{1}_{\{e, o\}}$ with $e \in \mathcal{E}$ and $o \in \mathcal{O}$, then the equality is true because both sides of (6.5.4) are equal to 1. Since such functions span \mathcal{B}_{k-1}^k it follows that (6.5.4) is true for every $f \in \mathcal{B}_{k-1}^k$. Finally, if nonnegative $f \in \mathcal{B}_{k-1}^k$ vanishes on \mathcal{E} , then the sum over \mathcal{O} also equals zero, hence $f \equiv 0$, and the same conclusion holds if we assume that $f = 0$ on \mathcal{O} . Thus U is the set of uniqueness if $\mathcal{O} \subset U$ or $\mathcal{E} \subset U$. \square

Remark 6.5.2. A naïve extension of Corollary 6.3.3 fails for $e(\mathcal{B}_2^k)$, if we try to replace the half-cubes with $(k-2)$ -subcubes, that is, quarter-cubes. This is seen from Theorem 6.5.1 for $k=3$. Indeed, the set

$$\{(1, 1, -1), (1, -1, 1), (-1, 1, 1), (1, -1, -1), (-1, 1, -1), (-1, -1, 1)\}$$

is not of uniqueness for $(\mathcal{B}_2^3)_+$, as follows from (6.5.1) with $e = (-1, -1, -1)$ and $o = (1, 1, 1)$, even though the set has nonempty intersection with each quarter-cube.

We will briefly treat the case of $e(\mathcal{B}_k^k)$, as follows.

Corollary 6.5.2. $k2^k \log 2$ is a sharp threshold of the sample size for the existence of MLE for $e(\mathcal{B}_k^k)$ and i.i.d. samples uniform on Q_k .

Proof. Observe that $e(\mathcal{B}_k^k)$ is isomorphic to $e(\mathbb{R}^{\mathcal{X}})$ for $|\mathcal{X}| = 2^k$. The existence of MLE for $e(\mathcal{B}_k^k)$ is characterized in (more general) Lemma 6.3.1, and the sharp threshold is given after Corollary 6.3.2. \square

Corollary 6.5.2 is in stark contrast with the result for the (smaller) space $e(\mathcal{B}_1^k)$ because for $e(\mathcal{B}_1^k)$ the sharp threshold, and so the threshold, equal $\log_2 k$, by Corollary 6.3.6.

Remark 6.5.3. Let $1 \leq q_1 \leq q_2 \leq k$. Then every set U of uniqueness for $(\mathcal{B}_{q_2}^k)_+$ is of uniqueness for $(\mathcal{B}_{q_1}^k)_+$, because $(\mathcal{B}_{q_1}^k)_+ \subset (\mathcal{B}_{q_2}^k)_+$.

A characterization of the existence of MLE for $e(\mathcal{B}_q^k)$ for arbitrary q , even for $q=2$, turned out to be difficult. Accordingly, we do not give a sharp threshold for the size of the uniform i.i.d. sample needed for the existence of MLE for $e(\mathcal{B}_q^k)$. However, the case of $e(\mathcal{B}_{k-q}^k)$ seems a little easier in the sense that we are able to give the less precise threshold for the existence of MLE for $e(\mathcal{B}_{k-q}^k)$. Moreover, for each fixed q the threshold for $e(\mathcal{B}_{k-q}^k)$ is the same as for $e(\mathcal{B}_k^k)$, namely $k2^k$ as $k \rightarrow \infty$.

Lemma 6.5.2. Fix $q \in \mathbb{N}$. Then $k2^k$ is a threshold of the sample size for the existence of MLE for $e(\mathcal{B}_{k-q}^k)$ and i.i.d. sample uniform on Q_k .

Proof. If $\lim_{k \rightarrow \infty} n(k)/(k2^k) = \infty$, then by Remark 6.5.3 and Corollary 6.5.2, for $k \rightarrow \infty$ we get

$$\mathbb{P} \left(\left\{ X_1, \dots, X_{n(k)} \right\} \text{ is of uniqueness for } \left(\mathcal{B}_{k-q}^k \right)_+ \right)$$

$$\geq \mathbb{P} \left(\left\{ X_1, \dots, X_{n(k)} \right\} \text{ is of uniqueness for } \mathcal{B}_k^k \right) \rightarrow 1,$$

as needed. On the other hand, every set U of uniqueness for $(\mathcal{B}_{k-q}^k)_+$ must intersect with every subcube defined by fixing last $k - q$ coordinates, because each q -subcube is the support of a function in $(\mathcal{B}_{k-q}^k)_+$, to wit, of its indicator. There are 2^{k-q} such q -subcubes, each of which we can suggestively denote by $(*, \dots, *, \varepsilon_{q+1}, \dots, \varepsilon_k)$, where $\varepsilon_{q+1}, \dots, \varepsilon_k = \pm 1$. Observe that the family of the above subcubes is a partition of \mathcal{Q}_k . We consider each q -subcube as a coupon in the Coupon Collector Problem. If a sample point falls into the q -subcube, we consider the coupon as collected. The probability of collecting a given coupon is 2^{q-k} . Therefore, if $n(k) = o(2^k k)$, hence $n(k) = o(2^{k-q}(k - q))$, then

$$\mathbb{P} \left(\left\{ X_1, \dots, X_{n(k)} \right\} \text{ is of uniqueness for } (\mathcal{B}_{k-q}^k)_+ \right) \rightarrow 0, \quad \text{as } k \rightarrow \infty,$$

as needed. □

6.6 Appendix

6.6.1 Proof of Lemma 6.1.1

Let $\hat{p} = e(\phi_0), \tilde{p} = e(\phi_1) \in e(\mathcal{B})$ and $\hat{p} \neq \tilde{p}$, so that $\phi_1 - \phi_0 \neq \text{const}$. Let $\phi_t = \phi_0 + t(\phi_1 - \phi_0)$, $p_t = e(\phi_t)$ for $t \in \mathbb{R}$ and $l(t) = l_{p_t}(x_1, \dots, x_n)$. We claim that l is strictly concave, that is $l'' < 0$. Indeed, since $\overline{\phi}_t = \overline{\phi}_0 + t\overline{\phi}_1$ is a linear function, by (6.1.4) we get

$$l''(t) = -n \frac{d^2}{dt^2} \log Z(\phi_t).$$

Let X be a random variable with values in \mathcal{X} such that $\mathbb{P}(X = x) = p(x)\mu(x)$. As usual, for every $f : \mathcal{X} \rightarrow \mathbb{R}$ we have

$$\mathbb{E}f(X) = \sum_{x \in \mathcal{X}} f(x)p(x)\mu(x).$$

Clearly, $(\log Z(\phi_t))' = \frac{Z(\phi_t)'}{Z(\phi_t)}$ and $(\log Z(\phi_t))'' = \frac{Z(\phi_t)''}{Z(\phi_t)} - \left(\frac{Z(\phi_t)'}{Z(\phi_t)} \right)^2$. Hence, thanks to (6.1.1),

$$\begin{aligned} Z(\phi_t)' &= \sum_{x \in \mathcal{X}} e^{\phi_t(x)} \mu(x) (\phi_1(x) - \phi_0(x)) \\ Z(\phi_t)'' &= \sum_{x \in \mathcal{X}} e^{\phi_t(x)} \mu(x) (\phi_1(x) - \phi_0(x))^2. \end{aligned}$$

Thus,

$$\frac{Z(\phi_t)'}{Z(\phi_t)} = \mathbb{E}[\phi_1(X) - \phi_0(X)] \qquad \frac{Z(\phi_t)''}{Z(\phi_t)} = \mathbb{E}[\phi_1(X) - \phi_0(X)]^2$$

and so

$$\frac{d^2}{dt^2} \log Z(\phi_t) = \mathbb{E}[\phi_1(X) - \phi_0(X) - \mathbb{E}(\phi_1(X) - \phi_0(X))]^2 > 0,$$

since $\phi_1 - \phi_0$ is not constant. Hence, l is strictly concave, in particular $l(1/2) > (l(0) + l(1))/2$. If $\sup_{p \in e(\mathcal{B})} L_p(x_1, \dots, x_n) = L_{\tilde{p}}(x_1, \dots, x_n) = L_{\hat{p}}(x_1, \dots, x_n)$, then $l(1/2) > \sup_{p \in e(\mathcal{B})} l_p(x_1, \dots, x_n)$, which is absurd; thus at most one of \tilde{p} and \hat{p} can be the MLE.

6.6.2 Control by oscillations

λ_U defined in Section 6.2 may be thought of as a specific measure of oscillation of ϕ . Of course, $\lambda_U \geq 0$. Furthermore, for every $c \in \mathbb{R}$,

$$\lambda_U(\phi + c) = \lambda_U(\phi), \quad \phi \in \mathcal{B}, \quad (6.6.1)$$

and for every (positive number) $k > 0$ we have (homogeneity),

$$\lambda_U(k\phi) = k\lambda_U(\phi), \quad \phi \in \mathcal{B}, k \geq 0. \quad (6.6.2)$$

If $U = \mathcal{X}$, then $\lambda_{\mathcal{X}}(-\phi) = \lambda_{\mathcal{X}}(\phi)$ for $\phi \in \mathcal{B}$, and so $\lambda_{\mathcal{X}}$ is a seminorm. Clearly, $\lambda_U \leq \lambda_{\mathcal{X}}$. However, if there is a nontrivial $\phi \in \mathcal{B}_+$ such that $\phi = 0$ on U , then $\lambda_U(\phi) = \sup_{\mathcal{X}} \phi > 0$ but $\lambda_U(-\phi) = 0$. The following result is the engine of Theorem 6.2.2.

Lemma 6.6.1. *$U \subset \mathcal{X}$ is the set of uniqueness for \mathcal{B}_+ if and only if λ_U is comparable with $\lambda_{\mathcal{X}}$ on \mathcal{B} , i.e., there exist constants $c_1, c_2 > 0$ such that $c_1\lambda_{\mathcal{X}}(\phi) \leq \lambda_U(\phi) \leq \lambda_{\mathcal{X}}(\phi)$ for all $\phi \in \mathcal{B}$.*

Proof. We first prove the “if” part. Assume U is not a set of uniqueness for \mathcal{B}_+ . Then there exists a nonzero function $\phi \in \mathcal{B}_+$ such that $\phi = 0$ on U . We have $\lambda_U(-\phi) = 0$ and $\lambda_{\mathcal{X}}(-\phi) > 0$, hence λ_U and $\lambda_{\mathcal{X}}$ are not comparable on \mathcal{B} .

We now prove the “only if” part, which is delicate. For all $\vartheta, \phi \in \mathcal{B}$ we have

$$\begin{aligned} \lambda_U(\vartheta + \phi) &\leq \max_{\mathcal{X}} \vartheta + \max_{\mathcal{X}} \phi - \min_U \vartheta - \min_U \phi \\ &= \lambda_U(\vartheta) + \lambda_U(\phi) \leq \lambda_U(\vartheta) + \lambda_{\mathcal{X}}(\phi). \end{aligned}$$

It follows that $\lambda_U(\vartheta) \geq \lambda_U(\vartheta - \phi) - \lambda_{\mathcal{X}}(\phi)$, hence

$$\lambda_U(\vartheta + \phi) \geq \lambda_U(\vartheta) - \lambda_{\mathcal{X}}(\phi).$$

Therefore, $\text{vert}\lambda_U(\vartheta + \phi) - \lambda_U(\vartheta)\text{vert} \leq \lambda_{\mathcal{X}}(\phi)$. As a consequence, λ_U is continuous on \mathcal{B} .

We will prove that there is a number $h > 0$ such that $\lambda_U(\phi) \geq h\lambda_{\mathcal{X}}(\phi)$ for every $\phi \in \mathcal{B}$. Let $\mathcal{S} = \{\phi \in \mathcal{B} : \min_{\mathcal{X}} \phi = 0 \text{ and } \max_{\mathcal{X}} \phi = 1\}$. Let $\phi \in \mathcal{S}$. If $\lambda_U(\phi) = 0$, then $\phi = 1$ on U . Consider $\varphi = 1 - \phi$. Clearly, $\varphi \geq 0$ and $\varphi = 0$ on U . It follows that $\varphi = 0$ on \mathcal{X} , because U is of uniqueness. Then $\phi \equiv 1$, which contradicts the assumption $\phi \in \mathcal{S}$. Therefore, $\lambda_U(\phi) > 0$. Since \mathcal{S} is compact and λ_U is continuous, $h := \min_{\mathcal{S}} \lambda_U > 0$. By (6.6.2) and (6.6.1) we obtain $\lambda_U(\phi) \geq h\lambda_{\mathcal{X}}(\phi)$ for all $\phi \in \mathcal{B}$. \square

6.6.3 Proof of Lemma 6.4.1

By (6.4.1), each $G \in \mathcal{G}_N$ appears in $\mathcal{G}_{N,c}$ with probability $p_c(G) = e^{\phi_c(G) - \psi(\phi_c)}$. Then,

$$\begin{aligned} p_{r,s} &= \mathbb{P}((r, s) \in E(G)) = \sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} \frac{e^{\phi_c(G)}}{\sum_{G \in \mathcal{G}_N} e^{\phi_c(G)}} \\ &= \frac{\sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} e^{\phi_c(G)}}{\sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} e^{\phi_c(G)} + \sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \notin E(G)}} e^{\phi_c(G)}} \end{aligned}$$

$$= \frac{\sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} e^{\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G)}}{\sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} e^{\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G)} + \sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \notin E(G)}} e^{\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G)}}. \quad (6.6.3)$$

Note that

$$\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G) = c_{r,s} \chi_{r,s}(G) + C(G),$$

where

$$C(G) = \sum_{\substack{(k,l) \in \binom{V}{2} \\ (k,l) \neq (r,s)}} c_{k,l} \chi_{k,l}(G).$$

Therefore

$$e^{\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G)} = e^{c_{r,s} \chi_{r,s}(G)} e^{C(G)}.$$

Clearly, $c_{r,s} \chi_{r,s}(G)$ is $c_{r,s}$ if $(r, s) \in E(G)$ and it is 0 if $(r, s) \notin E(G)$. Thus, (6.6.3) equals

$$\frac{e^{c_{r,s}} \sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} C(G)}{\sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \in E(G)}} e^{C(G)} + e^{c_{r,s}} \sum_{\substack{G \in \mathcal{G}_N \\ (r,s) \notin E(G)}} e^{C(G)}}.$$

Let S be the graph with only one edge (r, s) . The map $G \mapsto G \setminus S$ is a bijection between the graphs with the edge (r, s) and graphs without (r, s) . In addition, $C(G) = C(G \setminus S)$, and so we get (6.4.2).

6.6.4 Proof of Lemma 6.4.2

By (6.4.1), each $G \in \mathcal{G}_N$ appears in $\mathcal{G}_{N,c}$ with probability $p_c(G) = e^{\phi_c(G) - \psi(\phi_c)}$. Then,

$$\mathbb{P}((r_1, s_1), (r_2, s_2) \in E(G)) = \sum_{\substack{G \in \mathcal{G}_N \\ (r_1, s_1), (r_2, s_2) \in E(G)}} \frac{e^{\phi_c(G)}}{\sum_{G \in \mathcal{G}_N} e^{\phi_c(G)}}.$$

As in the proof of Lemma 6.4.1, we observe that

$$\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G) = c_{r_1, s_1} \chi_{r_1, s_1}(G) + c_{r_2, s_2} \chi_{r_2, s_2}(G) + \tilde{C}(G),$$

where

$$\tilde{C}(G) = \sum_{\substack{(k,l) \in \binom{V}{2} \\ (k,l) \neq (r_1, s_1) \\ (k,l) \neq (r_2, s_2)}} c_{k,l} \chi_{k,l}(G).$$

Thus,

$$e^{\sum_{(k,l) \in \binom{V}{2}} c_{k,l} \chi_{k,l}(G)} = e^{c_{r_1, s_1} \chi_{r_1, s_1}(G)} e^{c_{r_2, s_2} \chi_{r_2, s_2}(G)} e^{\tilde{C}(G)}.$$

Let S_1 and S_2 be the graphs with only one edge, (r_1, s_1) and (r_2, s_2) , respectively. Let

$$\mathcal{G}_{N_{12}} = \{G \in \mathcal{G}_N : S_1 \subset G, S_2 \subset G\},$$

$$\begin{aligned}\mathcal{G}_{N_{10}} &= \{G \in \mathcal{G}_N : S_1 \subset G, S_2 \not\subset G\}, \\ \mathcal{G}_{N_{02}} &= \{G \in \mathcal{G}_N : S_1 \not\subset G, S_2 \subset G\}, \\ \mathcal{G}_{N_{00}} &= \{G \in \mathcal{G}_N : S_1 \not\subset G, S_2 \not\subset G\}.\end{aligned}$$

a partition of \mathcal{G}_N . We observe that the maps

$$G \mapsto G \setminus S_1, \quad G \mapsto G \setminus S_2, \quad G \mapsto G \setminus (S_1 \cup S_2)$$

are bijections between $\mathcal{G}_{N_{10}}$, $\mathcal{G}_{N_{02}}$, $\mathcal{G}_{N_{12}}$, respectively, and $\mathcal{G}_{N_{00}}$. Also, for every $G \in \mathcal{G}_N$,

$$\tilde{C}(G) = \tilde{C}(G \setminus S_1) = \tilde{C}(G \setminus S_2) = \tilde{C}(G \setminus (S_1 \cup S_2)).$$

Put differently, $\tilde{C}(G)$ does not depend on the edges (r_1, s_1) and (r_2, s_2) . As in the proof of Lemma 6.4.1, we obtain

$$\begin{aligned}\mathbb{P}((r_1, s_1), (r_2, s_2) \in E(\mathbb{G})) \\ = \frac{e^{cr_1, s_1} e^{cr_2, s_2}}{1 + e^{cr_1, s_1} + e^{cr_2, s_2} + e^{cr_1, s_1} e^{cr_2, s_2}} = pr_{1, s_1} pr_{2, s_2}.\end{aligned}$$

6.6.5 Proof of Lemma 6.5.1

Proof. Consider the positive half-cubes H_1^+, \dots, H_k^+ . Let

$$\mathcal{B} = \text{Lin} \left\{ \prod_{i \in I_q} \mathbb{1}_{H_i^+} : I_q \subset \{0, \dots, k\} \text{ and } |I_q| \leq q \right\}.$$

We have $\mathcal{B} = \mathcal{B}_q^k$, because $r_0 = \mathbb{1}_{Q_k}$, $r_i = 2\mathbb{1}_{H_i^+} - \mathbb{1}_{Q_k}$ and by induction it is easy to see that for every $S \subset \{1, \dots, k\}$ and $|S| < q$, if Walsh function $w_S \in \mathcal{B}$ then their product with Rademacher function $w_S r_i \in \mathcal{B}$, for any $i = 0, \dots, n$. Note that for any permutation σ of $\{1, 2, \dots, q\}$,

$$\mathbb{1}_{H_{i_1}^+} \mathbb{1}_{H_{i_2}^+} \cdots \mathbb{1}_{H_{i_q}^+} = \mathbb{1}_{H_{i_{\sigma(1)}}^+} \mathbb{1}_{H_{i_{\sigma(2)}}^+} \cdots \mathbb{1}_{H_{i_{\sigma(q)}}^+}.$$

The functions $\mathbb{1}_{Q_k}$ and $\mathbb{1}_{H_{i_1}^+} \cdots \mathbb{1}_{H_{i_q}^+}$, $1 \leq i_1 \leq \dots \leq i_q \leq k$, are linearly independent. Indeed, assume that

$$r := \alpha_0 \mathbb{1}_{Q_k} + \sum_{i_1, \dots, i_q \in \{1, \dots, k\}} \alpha_{i_1 \dots i_q} \mathbb{1}_{H_{i_1}^+} \cdots \mathbb{1}_{H_{i_q}^+} = 0.$$

There are points $x_0 \in \bigcap_{i=1}^k H_i^-$, $x_{i_1 \dots i_q} \in \bigcap_{l \in \{i_1, \dots, i_q\}} H_l^- \cap \bigcap_{l \neq i_1, \dots, i_q} H_l^-$ for each $1 \leq i_1 \leq i_2 \leq \dots \leq i_q \leq k$. We obtain $\alpha_0 = r(x_0) = 0$ and $\alpha_{i_1 \dots i_q} = r(x_{i_1 \dots i_q}) = 0$ as needed. \square

6.6.6 Propagation of extrema, relative interior and the criterion of Barndorff-Nielsen

In this section we give auxiliary results, but also explain connections to the criterion of Barndorff-Nielsen. Let \mathcal{B} be an arbitrary linear subspace of $\mathbb{R}^{\mathcal{X}}$. In Corollary 6.6.1 below we adapt the criterion in Theorem 6.2.2 to such $e(\mathcal{B})$. Let \mathcal{B}' be the linear space spanned by \mathcal{B} and $\mathbb{1}$.

Lemma 6.6.2. *If $U \subset \mathcal{X}$, then $\phi = \min_{\mathcal{X}} \phi$ on U implies $\phi = \min_{\mathcal{X}} \phi$ on \mathcal{X} for every $\phi \in \mathcal{B}$ if and only if $\phi = \max_{\mathcal{X}} \phi$ on U implies $\phi = \max_{\mathcal{X}} \phi$ on \mathcal{X} for every $\phi \in \mathcal{B}$.*

Proof. The property with the minima is equivalent to the one with the maxima because \mathcal{B} is closed upon multiplication by -1 and because $\max(-\phi) = -\min \phi$. \square

Definition 6.6.1. We say that $U \subset \mathcal{X}$ propagates extrema for \mathcal{B} if $\phi = \inf_{\mathcal{X}} \phi$ on U implies that $\phi = \inf_{\mathcal{X}} \phi$ on \mathcal{X} for every $\phi \in \mathcal{B}$.

Due to Lemma 6.6.2, the property could be equivalently stated using maxima.

Lemma 6.6.3. A nonempty $U \subset \mathcal{X}$ propagates extrema for \mathcal{B} if and only if U is of uniqueness for \mathcal{B}'_+ .

Proof. Assume that U is of uniqueness for \mathcal{B}'_+ . Let $\phi \in \mathcal{B}$ and $\phi = \min_{\mathcal{X}} \phi$ on U . Then $\varphi = \phi - \min_{\mathcal{X}} \phi \in \mathcal{B}'_+$ and $\varphi = 0$ on U , so $\varphi = 0$ on \mathcal{X} and $\phi = \min_{\mathcal{X}} \phi$ on \mathcal{X} . It follows that U propagates extrema for \mathcal{B} . Conversely, assume that U propagates extrema for \mathcal{B} . Let $\phi \in \mathcal{B}$. Then $\phi = \varphi + c$ for some $\varphi \in \mathcal{B}$ and $c \in \mathbb{R}$. If $\phi \geq 0$ and $\phi = 0$ on U , then $\varphi = \min_{\mathcal{X}} \varphi = -c$ on U , hence $\varphi = -c$ on \mathcal{X} , and so $\phi = 0$ on \mathcal{X} . Thus, U is of uniqueness for \mathcal{B}'_+ . \square

Theorem 6.2.2 yields the following.

Corollary 6.6.1. MLE for $e(\mathcal{B})$ and $x_1, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ propagates extrema for \mathcal{B} .

Proof. The MLE for $e(\mathcal{B})$ and $e(\mathcal{B}')$ must be the same. Indeed, we have $e(\mathcal{B}) = e(\mathcal{B}')$ so the suprema of the likelihood functions are the same, see Section 6.1.1. Of course, if $\phi \in \mathcal{B}$ and $e(\phi)$ is the MLE for $e(\mathcal{B})$ then it is also the MLE for $e(\mathcal{B}')$. Conversely, if $\phi \in \mathcal{B}'$, then $\phi = \varphi + c$ for some $\varphi \in \mathcal{B}$ and $c \in \mathbb{R}$. If $e(\phi)$ is the MLE for $e(\mathcal{B}')$, then $e(\varphi)$ is the MLE for $e(\mathcal{B})$. Considering \mathcal{B}' , by Theorem 6.2.2 we see that MLE for $e(\mathcal{B}')$ and $x_1, \dots, x_n \in \mathcal{X}$ exists if and only if $\{x_1, \dots, x_n\}$ is of uniqueness for \mathcal{B}'_+ , and – by Lemma 6.6.3 – if and only if $\{x_1, \dots, x_n\}$ propagates extrema for \mathcal{B} . \square

The next lemma hinges on the trivial observation that if the sample mean equals the minimum, then the sample is constant.

Lemma 6.6.4. $\{x_1, \dots, x_n\}$ propagates extrema for \mathcal{B} if and only if for every $\phi \in \mathcal{B}$, $\min_{\mathcal{X}} \phi < \max_{\mathcal{X}} \phi$ implies $\min_{\mathcal{X}} \phi < \bar{\phi} < \max_{\mathcal{X}} \phi$.

Proof. Let $\{x_1, \dots, x_n\}$ propagate extrema for \mathcal{B} . If $\min_{\mathcal{X}} \phi = \bar{\phi}$, then $\phi = \min_{\mathcal{X}} \phi$ on $\{x_1, \dots, x_n\}$, hence $\phi = \min_{\mathcal{X}} \phi$ on \mathcal{X} and so $\min_{\mathcal{X}} \phi = \max_{\mathcal{X}} \phi$. A similar argument works if $\bar{\phi} = \max_{\mathcal{X}} \phi$; see also Lemma 6.6.2. Conversely, if $\{x_1, \dots, x_n\}$ does not propagate extrema for \mathcal{B} then there is $\phi \in \mathcal{B}$ such that $\phi = \min_{\mathcal{X}} \phi$ on $\{x_1, \dots, x_n\}$, but $\max_{\mathcal{X}} \phi > \min_{\mathcal{X}} \phi$. Then $\min_{\mathcal{X}} \phi = \bar{\phi} < \max_{\mathcal{X}} \phi$. \square

Recall the setting and notation of Section 6.1.2. The following theorem was essentially proved in [11] (Theorem 9.13), except that it was stated for the minimal representation of exponential families. The formulation presented in Theorem 6.6.1 below was given in [110] (Theorem 3.5), which covers the arbitrary canonical representation and does so with a more direct proof. Notably, [110] uses the notion of relative interior of a convex set. Let C be the convex hull of S . We say that $t \in \mathbb{R}^d$ is in the relative interior of C if for every $\theta \in \mathbb{R}^d$, $\min_{y \in C} \theta \cdot y < \max_{y \in C} \theta \cdot y$ implies $\min_{y \in C} \theta \cdot y < \theta \cdot t < \max_{y \in C} \theta \cdot y$.

Theorem 6.6.1. [110] (Theorem 3.5.) MLE for $e(\mathcal{B})$ and $x_1, \dots, x_n \in \mathcal{X}$ exists, if and only if \bar{T} is in the relative interior of C .

To close the circle of ideas, we give a self-contained proof of Theorem 6.6.1, which may also be used to obtain Theorem 6.2.2 from Theorem 6.6.1.

Proof of Theorem 6.6.1. By the discussion in this section we know very well that MLE for x_1, \dots, x_n and $e(\mathcal{B})$ exists if and only if for every $\phi \in \mathcal{B}$, $\min_{\mathcal{X}} \phi < \max_{\mathcal{X}} \phi$ implies $\min_{\mathcal{X}} \phi < \bar{\phi} < \max_{\mathcal{X}} \phi$. Recall that $\phi \in \mathcal{B}$ if and only if there is $\theta \in \mathbb{R}^d$ such that $\phi = \theta \cdot T$. Then $\min_{x \in \mathcal{X}} \phi(x) = \min_{y \in \mathcal{S}} \theta \cdot y = \min_{y \in C} \theta \cdot y$, $\max_{x \in \mathcal{X}} \phi(x) = \max_{y \in C} \theta \cdot y$, and, of course, $\bar{\phi} = \theta \cdot \bar{T}$. Therefore the existence of MLE for x_1, \dots, x_n and $e(\mathcal{B})$ is equivalent to \bar{T} being in the relative interior of C . \square

For clarity, we recall that we agreed in Example 6.1.2 that the existence of MLE for $x_1, \dots, x_n \in \mathcal{X}$ and $e(\mathcal{B})$ is the same as the existence of MLE for x_1, \dots, x_n and the exponential family given by the canonical statistics T and (6.1.10), and that it is equivalent to the existence of MLE for the sample $y_1 := T(x_1), \dots, y_n = T(x_n) \in \mathbb{R}^d$ and the standard exponential family in (6.1.11). From the above discussion we also see that the convex hull C and the notion of relative interior are merely auxiliary objects to express the property in Lemma 6.6.4, or the propagation of extrema property.

Chapter 7

On Laplacian of Graphical Models in Various Graphs

7.1 Introduction

The content of this chapter comes from the published article [164] of the author of the dissertation from the conference “Geometric Science of Information 2021” in Sorbonne University, Paris.

Let $G = (V, E, C)$ be a simple undirected graph, where $V = \{1, 2, \dots, n\}$ is a set of vertices, $E \subset \binom{V}{2}$ is a set of edges and $C \subset V$ is a set of source vertices, which will be called later as a “root set” or a “root”. As a degree $\deg(v)$ of a vertex $v \in V$ we treat a number of its neighbours. In our convention the graph Laplacian is defined as $\mathbf{L}(G) = \{l_{i,j}\}_{1 \leq i,j \leq n}$ with $l_{i,i} = \deg(i)$, $l_{i,j} = -1$ if $\{i, j\}$ is an edge in G and zeros otherwise, cf. e.g. [35]. Note that the graph Laplacian is a singular matrix, since its entries in each row (and each column) sum up to 0. Therefore we introduce an augmented graph Laplacian by adding 1 to every entry $l_{c,c}$ corresponding to $c \in C$. In other words,

$$\mathbf{L}_C^*(G) := \mathbf{L}(G) + \mathbf{E}_C$$

with \mathbf{E}_C being a square matrix with 1 in (c, c) for $c \in C$ and zeros everywhere else. We define $\mathbf{L}_c^*(G) = \{l_{i,j}^*\}_{1 \leq i,j \leq n}$ with $l_{i,j}^* = l_{i,j} + \mathbb{1}_{\{i=j \in C\}}$.

7.2 Trees

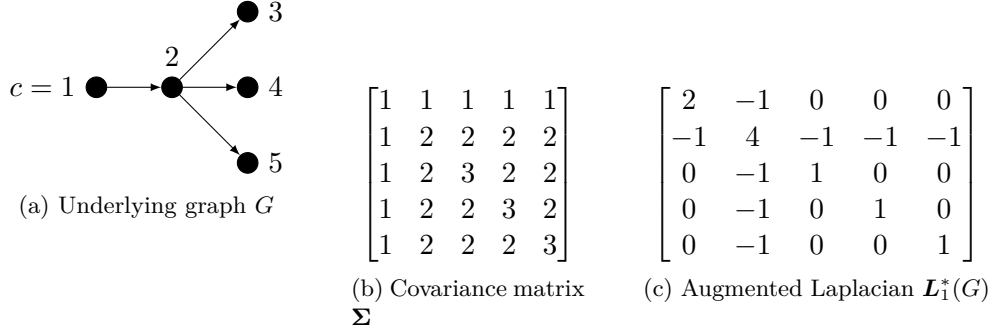
Let $T = (V, E)$ be an undirected tree. We may orient it in a following way: Choose one root vertex $C = \{c\}$. Then we orient every edge in a direction from c . Following this method we may induce a partial order \leq on the set of vertices such that $v \leq v'$ if and only if there exists a directed path from v' to v .

For every vertex $v \in V$ we define its ancestry $\text{AN}(v) := \{w \in V : v \leq w\}$ as a set of vertices in a unique path from v to c . Note that both the partial order and AN depend strictly on the choice of c .

Now consider an n -dimensional Gaussian random variable (X_1, \dots, X_n) with a covariance matrix $\Sigma = \{\sigma_{i,j}\}_{1 \leq i,j \leq n}$ such that

$$\sigma_{i,j} = |\text{AN}(i) \cap \text{AN}(j)|$$

for every $1 \leq i, j \leq n$.



Theorem 7.2.1. Let (X_1, \dots, X_n) be a Gaussian graphical model with an underlying graph G being a tree rooted in $C = \{c\}$. Assume that its covariance matrix $\Sigma = (\sigma_{i,j})$ satisfies

$$\sigma_{i,j} = |\text{AN}(i) \cap \text{AN}(j)|.$$

Then the precision matrix $\mathbf{K} = \Sigma^{-1}$ of (X_1, \dots, X_n) is equal to $L_c^* = L_c^*(G)$.

Proof. It suffices to prove that $L_c^*(G) \cdot \Sigma = \mathbf{I}_n$. At first consider the case $i = j = c$. Observe that for any $1 \leq k \leq n$ we have $\sigma_{c,k} = 1$. Therefore:

$$(L_c^* \Sigma)_{c,c} = \sum_{k=1}^n l_{c,k}^* \sigma_{k,c} = l_{c,c}^* \sigma_{c,c} + \sum_{k \sim c} l_{c,k}^* \sigma_{k,c} = (\deg(c) + 1) + \deg(c) \cdot (-1) \cdot 1 = 1.$$

Now consider $i = j \neq c$:

$$\begin{aligned} (L_c^* \Sigma)_{i,i} &= \sum_{k=1}^n l_{i,k}^* \sigma_{k,i} = l_{i,i}^* \sigma_{i,i} + \sum_{\substack{k \sim i \\ k \in \text{AN}(i)}} l_{i,k}^* \sigma_{k,i} + \sum_{\substack{k \sim i \\ i \in \text{AN}(k)}} l_{i,k}^* \sigma_{k,i} = \\ &= \deg(i) \sigma_{i,i} - (\sigma_{i,i} - 1) - (\deg(i) - 1) \sigma_{i,i} = 1. \end{aligned}$$

Now we will prove that the outside of the diagonal of $(L_c^* \Sigma)$ consists only of zeros. Observe that for $i = c, j \neq c$ we have:

$$\begin{aligned} (L_c^* \Sigma)_{i,j} &= \sum_{k=1}^n l_{i,k}^* \sigma_{k,j} = l_{i,i}^* \sigma_{i,j} + \sum_{\substack{k \sim i \\ k \notin \text{AN}(j)}} l_{i,k}^* \sigma_{k,j} + \sum_{\substack{k \sim i \\ k \in \text{AN}(j)}} l_{i,k}^* \sigma_{k,j} = \\ &= (\deg(i) + 1) \cdot 1 - (\deg(i) - 1) \cdot 1 - 1 \cdot 2 = 0. \end{aligned}$$

On the other hand, if $j = c, i \neq c$, then:

$$(L_c^* \Sigma)_{i,j} = \sum_{k=1}^n l_{i,k}^* \sigma_{k,j} = \sum_{k=1}^n l_{i,k} \cdot 1 = 0,$$

because each row of the graph Laplacian sums up to 0. Now we let $i, j \neq c, i \neq j$. If $i \in \text{AN}(j)$, then:

$$\begin{aligned} (L_c^* \Sigma)_{i,j} &= \sum_{k=1}^n l_{i,k}^* \sigma_{k,j} = \sum_{k=1}^n l_{i,k} \cdot \sigma_{k,j} = \\ &= \deg(i) \cdot \sigma_{i,j} - (\deg(i) - 2) \sigma_{i,j} - 1 \cdot (\sigma_{i,j} - 1) - 1 \cdot (\sigma_{i,j} + 1) = 0. \end{aligned}$$

Analogously, if $j \in \text{AN}(i)$, then $\sigma_{k,j} = \sigma_{i,j}$ for every $k \sim i$. Therefore:

$$(\mathbf{L}_c^* \boldsymbol{\Sigma})_{i,j} = \sum_{k=1}^n l_{i,k}^* \sigma_{k,j} = \sum_{k=1}^n l_{i,k} \cdot \sigma_{k,j} = \deg(i) \cdot \sigma_{i,j} - \deg(i) \cdot \sigma_{i,j} = 0.$$

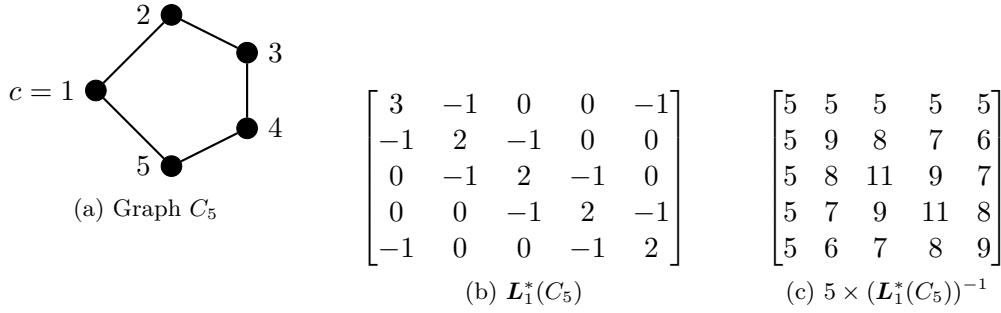
The same argument may be applied in the only case left, when $i \neq j$; $i, j \neq c$, $i \notin \text{AN}(j)$ and $j \notin \text{AN}(i)$. Therefore here also $(\mathbf{L}_c^* \boldsymbol{\Sigma})_{i,j} = 0$. \square

7.3 Discussion and non-tree graphs

7.3.1 Non-tree graphs

Cycles and complete graphs

The description of the inverse of $\mathbf{L}^*(G)$ for general G is much harder for G not being a tree. So far we are not able to present a general formula, thus we consider examples of such inverse for specific classes of graphs G . Below we show some examples of $n \times \boldsymbol{\Sigma} = [\mathbf{L}_1^*(G)]^{-1}$ for cycles C_n and complete graphs K_n with the root $C = \{1\}$:



Theorem 7.3.1. *The inverse matrix $\boldsymbol{\Sigma} = (\sigma_{i,j})$ of $\mathbf{L}^*(C_n)$ is a symmetric matrix satisfying*

$$\sigma_{i,j} = 1 + \frac{(i-1)(n-j+1)}{n}, \quad \text{for } i \leq j.$$

Proof. Again, we show that $\mathbf{L}_1^* \boldsymbol{\Sigma} = \mathbf{I}_n$, we assume that the cycle is $1 \rightarrow 2 \rightarrow 3 \rightarrow \dots \rightarrow n \rightarrow 1$ and take $C = \{1\}$. Then

$$(\mathbf{L}_1^* \boldsymbol{\Sigma})_{c,c} = \sum_{k=1}^n l_{c,k} \sigma_{k,c} = \sum_{k=1}^n l_{c,k} = 1.$$

If $j \neq c$, then we have

$$\begin{aligned} (\mathbf{L}_1^* \boldsymbol{\Sigma})_{c,j} &= \sum_{k=1}^n l_{c,k} \sigma_{k,j} = l_{c,c} \sigma_{c,j} + l_{c,2} \sigma_{2,j} + l_{c,n} \sigma_{n,j} = \\ &= 3 \cdot 1 + (-1) \cdot \left(1 + \frac{1 \cdot (n-j+1)}{n}\right) + (-1) \cdot \left(1 + \frac{(j-1) \cdot 1}{n}\right) \\ &= 3 - 1 - \frac{n-j+1}{n} - 1 - \frac{j-1}{n} = 0 \end{aligned}$$

and, analogously,

$$(\mathbf{L}_1^* \boldsymbol{\Sigma})_{i,c} = \sum_{k=1}^n l_{i,k} \sigma_{k,c} = \sum_{k=1}^n l_{i,k} = 0.$$

If $i, j \neq c$, then on the main diagonal we have

$$\begin{aligned}
(\mathbf{L}_1^* \boldsymbol{\Sigma})_{i,i} &= \sum_{k=1}^n l_{i,k} \sigma_{k,i} = l_{i,i-1} \sigma_{i-1,i} + l_{i,i} \sigma_{i,i} + l_{i,i+1} \sigma_{i+1,i} = \\
&= - \left(1 + \frac{(i-2)(n-i+1)}{n} \right) + 2 \left(1 + \frac{(i-1)(n-i+1)}{n} \right) - \left(1 + \frac{(i-1)(n-i)}{n} \right) = \\
&= \frac{-(i-2)(n-i+1) + 2(i-1)(n-i+1) - (i-1)(n-i)}{n} = \\
&= \frac{(n-i+2) + (i-1)}{n} = 1.
\end{aligned}$$

Note that the above calculations are also true for $i = n$ and the $(n+1)^{st}$ row/column of L and $\boldsymbol{\Sigma}$ being treated as c^{th} row/column. Now we only need to consider the outside of the main diagonal. Note that $\sigma_{i,n} = 1 + \frac{i-1}{n}$. Thus

$$(\mathbf{L}^* \boldsymbol{\Sigma})_{i,n} = \sum_{k=1}^n l_{i,k} \sigma_{k,n} = -\sigma_{i-1,n} + 2\sigma_{i,n} - \sigma_{i+1,n} = 0.$$

Also,

$$\begin{aligned}
(\mathbf{L}^* \boldsymbol{\Sigma})_{n,j} &= \sum_{k=1}^n l_{n,k} \sigma_{k,j} = -\sigma_{n-1,j} + 2\sigma_{n,j} - \sigma_{c,j} = \\
&= -1 - \frac{(j-1) \cdot 2}{n} + 2 + \frac{(j-1) \cdot 1}{n} \cdot 2 - 1 = 0.
\end{aligned}$$

Finally, for $i, j \neq 1, n, i \neq j$ we have

$$(\mathbf{L}^* \boldsymbol{\Sigma})_{i,j} = \sum_{k=1}^n l_{i,k} \sigma_{k,j} = -\sigma_{i-1,j} + 2\sigma_{i,j} - \sigma_{i+1,j} = (*)$$

Observe that as i, j, n are pairwise distinct, either $i+1 \leq j$ or $i-1 \geq j$. Therefore

$$\begin{cases} i+1 \leq j \Rightarrow (*) = -1 - \frac{(i-2)(n-j+1)}{n} + 2 + 2 \cdot \frac{(i-1)(n-j+1)}{n} - 1 - \frac{(i)(n-j+1)}{n} = 0, \\ i-1 \geq j \Rightarrow (*) = -1 - \frac{(j-1)(n-i+2)}{n} + 2 + 2 \cdot \frac{(j-1)(n-i+1)}{n} - 1 - \frac{(j-1)(n-i)}{n} = 0, \end{cases}$$

which ends the proof. □

Theorem 7.3.2.

$$[\mathbf{L}_1^*(K_n)]_{i,j}^{-1} = \begin{cases} 1 & \text{if } i = 1 \text{ or } j = 1, \\ 1 + \frac{2}{n} & \text{if } 1 < i = j, \\ 1 + \frac{1}{n} & \text{else.} \end{cases}$$

Proof. For proving that claim we observe that $\mathbf{L}_1^*(K_n) = \mathbf{L}(K_n) + \mathbf{E}_1 = n\mathbf{I}_n - \mathbf{J}_n + \mathbf{E}_1$. Thus $\mathbf{L}_1^*(K_n)$ belongs to an associative algebra being the 5-dimensional matrix space spanned by $\mathbf{I}_n, \mathbf{E}_1, \mathbf{J}_n, \mathbf{E}_1\mathbf{J}_n$ and $\mathbf{J}_n\mathbf{E}_1$. Therefore $(\mathbf{L}_1^*(K_n))^{-1}$ has to be looked for under the form

$$\frac{1}{n}\mathbf{I}_n + a\mathbf{E}_1 + b\mathbf{J}_n + c\mathbf{E}_1\mathbf{J}_n + d\mathbf{J}_n\mathbf{E}_1.$$

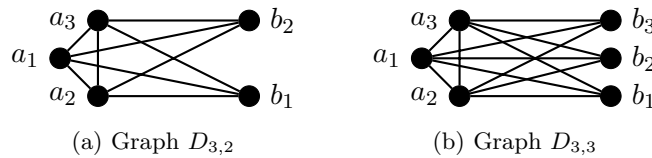
Clearly $(\mathbf{L}^*(K_n))^{-1}$ is a symmetric matrix, thus $d = c$. Solving the linear system in a, b, c gives $a = 0, b = (n+1)/n$ and $c = -1/n$. □

Daisy graphs

The daisy graphs (cf. e.g. Nakashima and Graczyk [135]) may be understood as a notion between a complete bipartite graph and a complete graph. To be more specific, the daisy graph $D_{a,b}$ is built as a sum of the complete bipartite graph $K_{a,b}$ and a complete subgraph K_a , i.e.

$$D_{a,b} := (V = V_A \sqcup V_B, E), \quad V = V_A \sqcup V_B, \quad (x, y) \in E \iff \{x, y\} \cap V_A \neq \emptyset.$$

We may interpret V_A and V_B as the internal and the external part of a daisy, respectively. To give an intuition, graphical models based on daisy graphs may be useful for analysis of a internal features of data (without knowledge of any independence among them) and b mutually conditionally independent external factors, which can influence the internal environment. Note that for $b = 1$ we have $D_{a,1}$ being a complete graph K_{a+1} and for $a = 1$ we have $D_{1,a}$ being a star graph, which is a tree (cf. 7.2).



The augmented Laplacian (and its inverse) depends on the choice of a root set $C \subset V$. Below we consider four choices of a root in $D_{a,b}$ and their augmented Laplacians:

$$\begin{aligned} L_{in}^* &:= \mathbf{L}(D_{a,b}) + \mathbf{E}_c && \text{for } C = \{c\}, \quad c \in V_A, \\ L_{in}^{**} &:= \mathbf{L}(D_{a,b}) + \sum_{c \in A} \mathbf{E}_c && \text{for } C = V_A, \\ L_{ex}^* &:= \mathbf{L}(D_{a,b}) + \mathbf{E}_c && \text{for } C = \{c\}, \quad c \in V_B, \\ L_{ex}^{**} &:= \mathbf{L}(D_{a,b}) + \sum_{c \in B} \mathbf{E}_c && \text{for } C = V_B. \end{aligned}$$

Without loss of generality we imply such ordering on vertices that the internal vertices precede the external ones. Moreover, concerning cases of one rooted internal (external) vertex we label it as the first (last) vertex.

The exact formulas for $(\mathbf{L}^*)^{-1}(D_{a,b})$ are presented below:

Theorem 4: Inverses of augmented Laplacians of daisy graphs

$[L_{in}^*(D_{a,b})]_{i,j}^{-1} = \begin{cases} 1 & \text{if } i = 1 \text{ or } j = 1, \\ 1 + \frac{1}{n} & \text{if } 2 \leq i, j \leq a, i \neq j, \\ 1 + \frac{2}{n} & \text{if } 2 \leq i = j \leq a, \\ 1 + \frac{1}{n} & \text{if } 2 \leq i \leq a < j, \\ 1 + \frac{a-1}{an} & \text{if } a < i, j, i \neq j, \\ 1 + \frac{a+n-1}{an} & \text{if } a < i = j. \end{cases}$	$[L_{in}^{**}(D_{a,b})]_{i,j}^{-1} = \begin{cases} \frac{n}{a(n+1)} & \text{if } i, j \leq a, i \neq j, \\ \frac{n+a}{a(n+1)} & \text{if } i = j \leq a, \\ \frac{1}{a} & \text{if } i \leq a < j, \\ \frac{1}{a} & \text{if } a < i, j, i \neq j, \\ \frac{2}{a} & \text{if } a < i = j. \end{cases}$
$[L_{ex}^*(D_{a,b})]_{i,j}^{-1} = \begin{cases} 1 & \text{if } i = n \text{ or } j = n, \\ 1 + \frac{1}{a} & \text{if } a < i, j < n, i \neq j, \\ 1 + \frac{2}{a} & \text{if } a < i = j < n, \\ 1 + \frac{1}{a} & \text{if } i \leq a < j < n, \\ 1 + \frac{n-1}{an} & \text{if } i, j \leq a, i \neq j, \\ 1 + \frac{a+n-1}{an} & \text{if } i = j \leq a. \end{cases}$	$[L_{ex}^{**}(D_{a,b})]_{i,j}^{-1} = \begin{cases} \frac{n+1}{bn} & \text{if } i, j \leq a, i \neq j, \\ \frac{n+b+1}{bn} & \text{if } i = j \leq a, \\ \frac{1}{b} & \text{if } i \leq a < j, \\ \frac{a}{(a+1)b} & \text{if } a < i, j, i \neq j, \\ \frac{n}{(a+1)b} & \text{if } a < i = j. \end{cases}$

Proof. Similarly to 7.3.2, the proof relies on the associative algebras being 6-dimensional (for L_{ex}^*) or 11-dimensional (for L_{in}^*) matrix spaces. The resulting inverses of augmented Laplacians follow from solving the corresponding systems of equations. \square

7.3.2 Eigenvalues of augmented Laplacian

Below we show that if the root consists of only one vertex c , then the determinant of $L_c^*(G)$ does not depend on the choice of the root vertex. Moreover, it can be proved that it is equal to the number of spanning trees of G .

Remark 7.3.3. Let $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ be the eigenvalues of L_c^* . Then

$$\prod_{i=1}^n \lambda_i = \det(L_c^*) = \#\{\text{spanning trees } (G)\}.$$

Proof. Recall that $L_c^*(G)$ differs from $L(G)$ only at $l_{c,c}$, where c is a root vertex of G . Let M be a matrix obtained from L by replacing the c^{th} column with a column with 1 at the c^{th} row and zeros elsewhere. Then

$$\det(L_c^*(G)) = \det(L(G)) + \det(M).$$

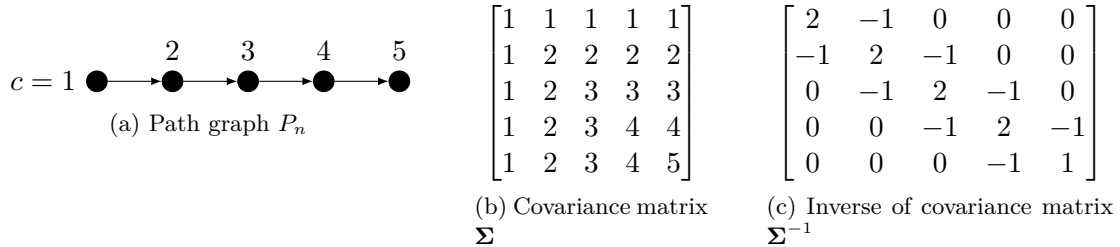
As discussed earlier, the Laplacian matrix $L(G)$ is singular. We can observe that $\det(M)$ is an $(n-1) \times (n-1)$ cofactor of $L(G)$. By the Kirchoff's matrix-tree Theorem, any $(n-1) \times (n-1)$ cofactor of $L(G)$ is equal to the number of spanning trees of G . \square

7.3.3 Discussion

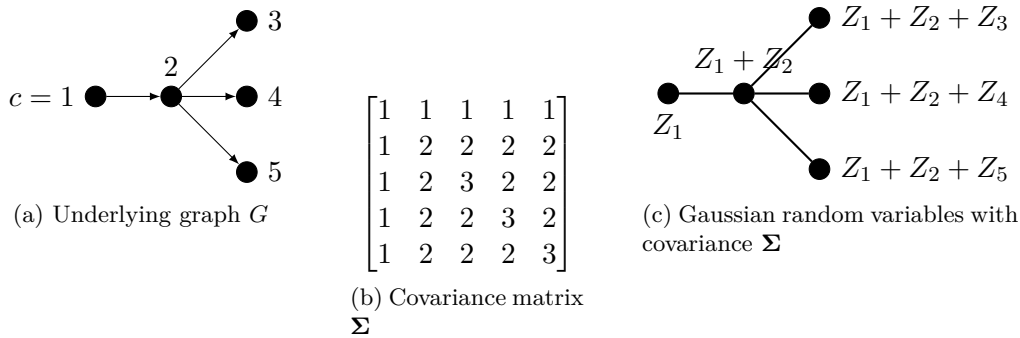
Interpretations of the inverse of the augmented Laplacian

Below we couple the obtained inverses of L^* matrices with the covariance matrices of some classical examples of random walks. At first let us remind that the covariance function of the

Wiener process $(W_t)_{t \geq 0}$ is $\text{Cov}(W_s, W_t) = \min\{s, t\}$. For example, (W_1, W_2, \dots, W_n) has its covariance matrix equal to $(\Sigma)_{i,j} = \min\{i, j\}$. At Theorem 7.2.1 we proved that this is exactly the inverse of the augmented Laplacian of a path graph with an initial vertex in one of its endpoints. This observation stays consistent with the conditional independence of W_{t_1} and



W_{t_2} under W_t for any $t_1 < t < t_2$. Similarly, the broader class of trees and their augmented Laplacians can be connected with sums of the standard Gaussian random variables 'branched' according to the underlying tree graph.



To observe the analogy of the previous examples for the cycle graph C_n , we may note that replacing the edge $(n, 1)$ with $(n, n + 1)$ (with $n + 1 \notin V$) gives a path graph

$$[c = 1] - [2] - \dots - [n + 1].$$

Therefore we may consider the model corresponding to the cycle graph as the Wiener model $(W_1, \dots, W_n, W_{n+1})$ conditioned by $W_{n+1} = W_1$. This gives a sum of a random variable $W_1 \sim \mathcal{N}(0, 1)$ and a Brownian bridge "tied down" at 1 and $(n + 1)$. Therefore the covariance matrix of (W_1, \dots, W_n) is equal to

$$\sigma_{i,j} = 1 + \frac{(i - 1)(n - j + 1)}{n}, \quad i \leq j,$$

cf. Theorem 7.3.1.

In order to find a model with a covariance matrix equal to

$$[L_1^*(K_n)]_{i,j}^{-1} = \begin{cases} 1 & \text{if } i = 1 \text{ or } j = 1, \\ 1 + \frac{2}{n} & \text{if } 1 < i = j, \\ 1 + \frac{1}{n} & \text{else,} \end{cases}$$

observe that all vertices (except of the initial one) are isomorphic and connected, therefore each of their correspondent random variables are mutually equally dependent. Let $c = 1$ and let Z_1, \dots, Z_n be i.i.d. random variables from $\mathcal{N}(0, 1)$. Therefore X_1, X_2, \dots, X_n are of the form

$$\begin{cases} X_1 = Z_1 \\ X_2 = Z_1 + \alpha Z_2 + \beta(Z_3 + Z_4 + \dots + Z_n) \\ X_3 = Z_1 + \alpha Z_3 + \beta(Z_2 + Z_4 + \dots + Z_n) \\ \dots \\ X_n = Z_1 + \alpha Z_n + \beta(Z_2 + Z_3 + \dots + Z_{n-1}). \end{cases}$$

The restriction on the covariance matrix of (X_1, \dots, X_n) induces a system of equations, which is satisfied only for $\beta = \frac{1 - \sqrt{\frac{1}{n}}}{n-1}$ and $\alpha = \beta + \sqrt{\frac{1}{n}}$.

Bibliography

- [1] A. Agresti. Foundations of linear and generalized linear models. Wiley Series in Probability and Statistics. John Wiley & Sons, Inc., Hoboken, NJ, 2015.
- [2] H. Akaike. A new look at the statistical model identification. IEEE Transactions on Automatic Control, 19(6):716–723, 1974.
- [3] A. Albert and J. A. Anderson. On the existence of maximum likelihood estimates in logistic regression models. Biometrika, 71(1):1–10, 1984.
- [4] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp. Living on the edge: Phase transitions in convex programs with random data. Information and Inference: A Journal of the IMA, 3(3):224–294, 2014.
- [5] E. B. Andersen. Sufficiency and exponential families for discrete sample spaces. J. Amer. Statist. Assoc., 65:1248–1255, 1970.
- [6] T. W. Anderson and J. B. Taylor. Strong consistency of least squares estimates in normal linear regression. The Annals of Statistics, 4(4):788–790, 1976.
- [7] J.-P. Aubin. Mathematical methods of game and economic theory, volume 7 of Studies in Mathematics and its Applications. North-Holland Publishing Co., Amsterdam-New York, 1979.
- [8] R. Bao, B. Gu, and H. Huang. Fast oscar and owl with safe screening rules. In International Conference on Machine Learning, 2020.
- [9] A. Barbara, A. Jourani, and S. Vaïter. Maximal solutions of sparse analysis regularization. Journal of Optimization Theory and Applications, 180(2):374–396, 2019.
- [10] É. Barbin and J. P. Lamarche. Histoires de probabilités et de statistiques. Ellipses Paris, 2004.
- [11] O. Barndorff-Nielsen. Information and exponential families in statistical theory. John Wiley & Sons Ltd., Chichester, 1978. Wiley Series in Probability and Mathematical Statistics.
- [12] E. M. L. Beale, M. G. Kendall, and D. W. Mann. The discarding of variables in multivariate analysis. Biometrika, 54:357–366, 1967.
- [13] P. Bellec and A. Tsybakov. Bounds on the prediction error of penalized least squares estimators with convex penalty. In Modern problems of stochastic analysis and statistics, volume 208 of Springer Proc. Math. Stat., pages 315–333. Springer, Cham, 2017.

- [14] P. C. Bellec, G. Lecué, and A. B. Tsybakov. Slope meets Lasso: improved oracle bounds and optimality. *Ann. Statist.*, 46(6B):3603–3642, 2018.
- [15] A. Belloni, V. Chernozhukov, and L. Wang. Pivotal estimation via square-root Lasso in nonparametric regression. *Ann. Statist.*, 42(2):757–788, 2014.
- [16] A. Ben-Israel and T. N. E. Greville. *Generalized inverses*, volume 15 of *CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC*. Springer-Verlag, New York, second edition, 2003. Theory and applications.
- [17] D. P. Bertsekas. *Nonlinear programming*. Athena Scientific Optimization and Computation Series. Athena Scientific, Belmont, MA, second edition, 1999.
- [18] I. Bezáková, A. Kalai, and R. Santhanam. Graph model selection using maximum likelihood. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pages 105–112, New York, NY, USA, 2006. ACM.
- [19] P. Billingsley. *Probability and measure*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons, Inc., New York, third edition, 1995. A Wiley-Interscience Publication.
- [20] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, New York, 2007.
- [21] K. Bogdan and M. Bogdan. On existence of maximum likelihood estimators in exponential families. *Statistics*, 34(2):137–149, 2000.
- [22] K. Bogdan, M. Bosy, and T. Skalski. Maximum likelihood estimation for discrete exponential families and random graphs. *ALEA Lat. Am. J. Probab. Math. Stat.*, 19(1):1045–1070, 2022.
- [23] M. Bogdan, X. Dupuis, P. Graczyk, B. Kołodziejek, T. Skalski, P. Tardivel, and M. Wilczyński. Pattern recovery by SLOPE. Technical Report 2203.12086, arXiv, 2022.
- [24] M. Bogdan, P. Graczyk, B. Kołodziejek, T. Skalski, P. Tardivel, and M. Wilczyński. Pattern recovery by slope. *arXiv preprint arXiv:2203.12086*, 2022.
- [25] M. Bogdan and T. Ledwina. Testing uniformity via log-spline modeling. *Statistics*, 28(2):131–157, 1996.
- [26] M. Bogdan, E. van den Berg, W. S. C. Sabatti, and E. J. Candès. SLOPE – adaptive variable selection via convex optimization. *Annals of Applied Statistics*, 9:1103–1140, 2015.
- [27] M. Bogdan, E. van den Berg, W. Su, and E. Candès. Statistical estimation and testing via the sorted l_1 norm, 2013.
- [28] B. Bollobás. *Modern graph theory*, volume 184 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1998.
- [29] H. D. Bondell and B. J. Reich. Simultaneous regression shrinkage, variable selection, and supervised clustering of predictors with OSCAR. *Biometrics*, 64:115–123, 2008.
- [30] H. D. Bondell and B. J. Reich. Simultaneous factor selection and collapsing levels in anova. *Biometrics*, 65:169–177, 2009.

- [31] J. M. Borwein and A. S. Lewis. Convex Analysis and Nonlinear Optimization, Theory and Examples. Springer, 2000.
- [32] S. Boucheron, O. Bousquet, and G. Lugosi. Theory of classification: a survey of some recent advances. ESAIM Probab. Stat., 9:323–375, 2005.
- [33] S. Boyd and L. Vandenberghe. Convex Optimization. Cambridge University Press, USA, 2004.
- [34] A. Brøndsted. An introduction to convex polytopes, volume 90 of Graduate Texts in Mathematics. Springer-Verlag, New York-Berlin, 1983.
- [35] A. E. Brouwer and W. H. Haemers. Spectra of graphs. Universitext. Springer, New York, 2012.
- [36] L. D. Brown. Fundamentals of statistical exponential families with applications in statistical decision theory, volume 9 of Institute of Mathematical Statistics Lecture Notes—Monograph Series. Institute of Mathematical Statistics, Hayward, CA, 1986.
- [37] D. Brzyski. Selecting relevant groups of explanatory variables via convex optimization methods with the false discovery rate control. PhD Thesis, 2015.
- [38] D. Brzyski, A. Gossmann, W. Su, and M. Bogdan. Group SLOPE—adaptive selection of groups of predictors. J. Amer. Statist. Assoc., 114(525):419–433, 2019.
- [39] D. Brzyski, M. Karas, B. M. Ances, M. Dziedzic, J. Goñi, T. W. Randolph, and J. Harezlak. Connectivity-informed adaptive regularization for generalized outcomes. Canad. J. Statist., 49(1):203–227, 2021.
- [40] Z. Bu, J. M. Klusowski, C. Rush, and W. J. Su. Algorithmic analysis and statistical estimation of SLOPE via approximate message passing. IEEE Trans. Inform. Theory, 67(1):506–537, 2021.
- [41] P. Bühlmann and S. van de Geer. Statistics for High-Dimensional Data: Methods, Theory and Applications. Springer, Heidelberg, 2011.
- [42] E. J. Candès and Y. Plan. Near-ideal model selection by ℓ_1 minimization. The Annals of Statistics, 37(5A):2145–2177, 2009.
- [43] cardinal (<https://math.stackexchange.com/users/7003/cardinal>). Proof of upper-tail inequality for standard normal distribution. Mathematics Stack Exchange. URL:<https://math.stackexchange.com/q/28754> (version: 2011-03-24).
- [44] A. Chatterjee and S. N. Lahiri. Strong consistency of Lasso estimators. Sankhya A, 73(1):55–78, 2011.
- [45] S. Chatterjee and P. Diaconis. Estimating and understanding exponential random graph models. Ann. Statist., 41(5):2428–2461, 2013.
- [46] S. Chatterjee, P. Diaconis, and A. Sly. Random graphs with a given degree sequence. Ann. Appl. Probab., 21(4):1400–1435, 2011.
- [47] S. Chen and D. Donoho. Basis pursuit. In Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers, volume 1, pages 41–44. IEEE, 1994.

- [48] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. SIAM J. Sci. Comput., 20(1):33–61, 1998.
- [49] G. Claeskens. Statistical model choice. In Annual Review of Statistics and Its Application, volume 3, pages 233–256, 2016.
- [50] B. R. Crain. Estimation of distributions using orthogonal expansions. Ann. Statist., 2:454–463, 1974. Collection of articles dedicated to Jerzy Neyman on his 80th birthday.
- [51] B. R. Crain. Exponential models, maximum likelihood estimation, and the Haar condition. J. Amer. Statist. Assoc., 71(355):737–740, 1976.
- [52] J. Cuzick. A strong law for weighted sums of i.i.d. random variables. J. Theoret. Probab., 8(3):625–641, 1995.
- [53] J. N. Darroch, S. L. Lauritzen, and T. P. Speed. Markov fields and log-linear interaction models for contingency tables. Ann. Statist., 8(3):522–539, 1980.
- [54] P. Descloux, C. Boyer, J. Josse, A. Sportisse, and S. Sardy. Robust lasso-zero for sparse corruption and model selection with missing covariates. Scandinavian Journal of Statistics, 49:1605–1635, 2022.
- [55] P. Descloux and S. Sardy. Model selection with lasso-zero: adding straw to the haystack to better find needles. Journal of Computational and Graphical Statistics, 30(3):530–543, 2021.
- [56] N. Dexheimer and C. Strauch. On lasso and slope drift estimators for lévy-driven ornstein–uhlenbeck processes. 2022.
- [57] P. Diaconis. Group representations in probability and statistics, volume 11 of Institute of Mathematical Statistics Lecture Notes—Monograph Series. Institute of Mathematical Statistics, Hayward, CA, 1988.
- [58] P. Diaconis and D. Freedman. Partial exchangeability and sufficiency. In Statistics: applications and new directions (Calcutta, 1981), pages 205–236. Indian Statist. Inst., Calcutta, 1984.
- [59] D. Donoho and J. Tanner. Counting faces of randomly projected polytopes when the projection radically lowers dimension. Journal of the American Mathematical Society, 22(1):1–53, 2009.
- [60] D. Donoho and J. Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 367(1906):4273–4293, 2009.
- [61] X. Dupuis and P. J. C. 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.
- [62] X. Dupuis and S. Vaiteer. The geometry of sparse analysis regularization. arXiv preprint arXiv:1907.01769, 2019.

- [63] J. Dziubański and A. Hejna. Remark on atomic decompositions for the Hardy space H^1 in the rational Dunkl setting. Studia Math., 251(1):89–110, 2020.
- [64] B. Eisenberg. On the expectation of the maximum of IID geometric random variables. Statist. Probab. Lett., 78(2):135–143, 2008.
- [65] C. Elvira and C. Herzet. Safe rules for the identification of zeros in the solutions of the slope problem. ArXiv, abs/2110.11784, 2021.
- [66] P. Erdős and A. Rényi. On random graphs. I. Publ. Math. Debrecen, 6:290–297, 1959.
- [67] P. Erdős and A. Rényi. On the evolution of random graphs. Magyar Tud. Akad. Mat. Kutató Int. Közl., 5:17–61, 1960.
- [68] P. Erdős and A. Rényi. On a classical problem of probability theory. Magyar Tud. Akad. Mat. Kutató Int. Közl., 6:215–220, 1961.
- [69] N. Eriksson, S. E. Fienberg, A. Rinaldo, and S. Sullivant. Polyhedral conditions for the nonexistence of the MLE for hierarchical log-linear models. J. Symbolic Comput., 41(2):222–233, 2006.
- [70] G. Ewald. Combinatorial convexity and algebraic geometry, volume 168. Springer Science & Business Media, 1996.
- [71] K. Ewald and U. Schneider. Uniformly valid confidence sets based on the Lasso. Electron. J. Stat., 12(1):1358–1387, 2018.
- [72] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. Journal of the American Statistical Association, 96:1348–1360, 2001.
- [73] W. Feller. An introduction to probability theory and its applications. Vol. II. John Wiley & Sons, Inc., New York-London-Sydney, second edition, 1971.
- [74] S. E. Fienberg and A. Rinaldo. Three centuries of categorical data analysis: log-linear models and maximum likelihood estimation. J. Statist. Plann. Inference, 137(11):3430–3445, 2007.
- [75] S. E. Fienberg and A. Rinaldo. Maximum likelihood estimation in log-linear models. Ann. Statist., 40(2):996–1023, 2012.
- [76] M. Figueiredo and R. Nowak. Ordered weighted l1 regularized regression with strongly correlated covariates: Theoretical aspects. In Artificial Intelligence and Statistics, pages 930–938. PMLR, 2016.
- [77] R. A. Fisher. The use of multiple measurements in taxonomic problems. Annals of Eugenics, 7(2):179–188, 1936.
- [78] L. Freijeiro-González, M. Febrero-Bande, and W. González-Manteiga. A critical review of LASSO and its derivatives for variable selection under dependence among covariates. Int. Stat. Rev., 90(1):118–145, 2022.
- [79] J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. Journal of Statistical Software, 33, 2010.

- [80] A. Frieze and M. Karoński. Introduction to random graphs. Cambridge University Press, Cambridge, 2016.
- [81] F. Frommlet, M. g. Bogdan, and D. Ramsey. Phenotypes and genotypes, volume 18 of Computational Biology. Springer-Verlag, London, 2016. The search for influential genes.
- [82] J.-J. Fuchs. On sparse representations in arbitrary redundant bases. IEEE Trans. Inform. Theory, 50(6):1341–1344, 2004.
- [83] J.-J. Fuchs. Recovery of exact sparse representations in the presence of bounded noise. IEEE Transactions on Information Theory, 51(10):3601–3608, 2005.
- [84] D. Galvin. Three tutorial lectures on entropy and counting, 2014.
- [85] J. Gertheiss and G. Tutz. Sparse modeling of categorical explanatory variables. Ann. Appl. Stat., 4(4):2150–2180, 2010.
- [86] C. J. Geyer. Likelihood and exponential families. ProQuest LLC, Ann Arbor, MI, 1990. Thesis (Ph.D.)—University of Washington.
- [87] E. N. Gilbert. Random graphs. Ann. Math. Statist., 30:1141–1144, 1959.
- [88] J. C. Gilbert. On the solution uniqueness characterization in the l_1 norm and polyhedral gauge recovery. Journal of Optimization Theory and Applications, 172(1):70–101, 2017.
- [89] C. Giraud. Introduction to high-dimensional statistics, volume 139 of Monographs on Statistics and Applied Probability. CRC Press, Boca Raton, FL, 2015.
- [90] T. Godland and Z. Kabluchko. Projections and angle sums of permutohedra, 2020.
- [91] G. H. Golub and C. F. V. Loan. Matrix computations. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, 2013.
- [92] P. Graczyk, T. Luks, and P. Sawyer. Potential kernels for radial Dunkl Laplacians. Canad. J. Math., 74(4):1005–1033, 2022.
- [93] P. Graczyk, U. Schneider, T. Skalski, and P. Tardivel. Pattern recovery in penalized and thresholded estimation and its geometry. HAL preprint hal-03262087, 2021.
- [94] P. Gruber. Convex and Discrete Geometry. Springer, Heidelberg, 2007.
- [95] B. Grünbaum. Convex polytopes, volume 221 of Graduate Texts in Mathematics. Springer-Verlag, New York, second edition, 2003. Prepared and with a preface by Volker Kaibel, Victor Klee and Günter M. Ziegler.
- [96] S. J. Haberman. The analysis of frequency data. The University of Chicago Press, Chicago, Ill.-London, 1974. Statistical Research Monographs, Vol. IV.
- [97] A. Hald. A history of parametric statistical inference from Bernoulli to Fisher, 1713–1935. Sources and Studies in the History of Mathematics and Physical Sciences. Springer, New York, 2007.
- [98] G. H. Hardy, J. E. Littlewood, and G. Pólya. Inequalities. Cambridge, at the University Press, 1952. 2d ed.

- [99] H. L. Harter. Expected values of normal order statistics. Biometrika, 48(1/2):151–165, 1961.
- [100] S. Helgason. Differential geometry, Lie groups, and symmetric spaces, volume 34 of Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2001. Corrected reprint of the 1978 original.
- [101] J.-B. Hiriart-Urruty and C. Lemaréchal. Convex Analysis and Minimization Algorithms I: Fundamentals, volume 305. Springer, Heidelberg, 1993.
- [102] R. R. Hocking and R. N. Leslie. Selection of the best subset in regression analysis. Technometrics, 9:531–540, 1967.
- [103] H. Hu and Y. M. Lu. SLOPE for sparse linear regression: asymptotics and optimal regularization. IEEE Trans. Inform. Theory, 68(11):7627–7664, 2022.
- [104] J.-C. Hütter and P. Rigollet. Optimal rates for total variation denoising. In Conference on Learning Theory, pages 1115–1146. PMLR, 2016.
- [105] M. Jacobsen. Existence and unicity of MLEs in discrete exponential family distributions. Scand. J. Statist., 16(4):335–349, 1989.
- [106] A. Jedrzejewski, J. Lago, G. Marcjasz, and R. Weron. Electricity price forecasting: The dawn of machine learning. IEEE Power and Energy Magazine, 20(3):24–31, 2022.
- [107] H. Jégou, T. Furon, and J.-J. Fuchs. Anti-sparse coding for approximate nearest neighbor search. In 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 2029–2032. IEEE, 2012.
- [108] J. Jendrej, K. Oleszkiewicz, and J. O. Wojtaszczyk. On some extensions of the FKN theorem. Theory Comput., 11:445–469, 2015.
- [109] W. Jiang, M. Bogdan, J. Josse, S. Majewski, B. Miasojedow, V. Rockova, and T. Group. Adaptive bayesian slope: Model selection with incomplete data. Journal of Computational and Graphical Statistics, 31(1):113–137, 2022.
- [110] S. Johansen. Introduction to the theory of regular exponential families, volume 3 of Lecture Notes. University of Copenhagen, Institute of Mathematical Statistics, Copenhagen, 1979.
- [111] S.-J. Kim, K. Koh, S. Boyd, and D. Gorinevsky. ℓ_1 trend filtering. SIAM review, 51(2):339–360, 2009.
- [112] K. Knight and W. Fu. Asymptotics of Lasso-type estimators. Annals of Statistics, 28:1356–1378, 2000.
- [113] M.-H. Ko, D.-H. Ryu, T.-S. Kim, and Y.-K. Choi. A central limit theorem for general weighted sums of LNQD random variables and its application. Rocky Mountain J. Math., 37(1):259–268, 2007.
- [114] K. Konis. Linear programming algorithms for detecting separated data in binary logistic regression models. PhD thesis, Worcester College, University of Oxford, 2007.
- [115] M. Kos and M. Bogdan. On the asymptotic properties of SLOPE. Sankhya A, 82(2):499–532, 2020.

- [116] P. J. Kremer, D. Brzyski, M. g. Bogdan, and S. Paterlini. Sparse index clones via the sorted ℓ_1 -norm. Quant. Finance, 22(2):349–366, 2022.
- [117] T. L. Lai and C. Z. Wei. A law of the iterated logarithm for double arrays of independent random variables with applications to regression and time series models. Ann. Probab., 10(2):320–335, 1982.
- [118] P. Laplace. Théorie analytique des probabilités. Courcier, Paris, 1812.
- [119] J. Larsson, M. Bogdan, and J. Wallin. The strong screening rule for slope. Proc. 34th Conf. NeurIPS 2020, Vancouver, Canada, pages 1–12, 2020.
- [120] J. Larsson, Q. Klopfenstein, M. Massias, and J. Wallin. Coordinate descent for slope. arXiv preprint, arXiv:2210.14780, pages 1–12, 2022.
- [121] S. L. Lauritzen. Extreme point models in statistics. Scand. J. Statist., 11(2):65–91, 1984.
- [122] S. L. Lauritzen. Graphical models, volume 17 of Oxford Statistical Science Series. The Clarendon Press, Oxford University Press, New York, 1996. Oxford Science Publications.
- [123] S. Lee, P. Sobczyk, and M. Bogdan. Structure learning of gaussian markov random fields with false discovery rate control. Symmetry, 11(10), 2019.
- [124] A. Legendre. Nouvelles méthodes pour la détermination des orbites des comètes. Nineteenth Century Collections Online (NCCO): Science, Technology, and Medicine: 1780–1925. F. Didot, 1805.
- [125] E. L. Lehmann and G. Casella. Theory of point estimation. Springer Texts in Statistics. Springer-Verlag, New York, second edition, 1998.
- [126] G. Letac. Lectures on natural exponential families and their variance functions, volume 50 of Monografías de Matemática [Mathematical Monographs]. Instituto de Matemática Pura e Aplicada (IMPA), Rio de Janeiro, 1992.
- [127] L. Lovász and B. Szegedy. Limits of dense graph sequences. J. Combin. Theory Ser. B, 96(6):933–957, 2006.
- [128] A. Maj-Kańska, P. Pokarowski, and A. Prochenka. Delete or merge regressors for linear model selection. Electron. J. Stat., 9(2):1749–1778, 2015.
- [129] C. Mazza-Anthony, B. Mazoure, and M. Coates. Learning Gaussian graphical models with ordered weighted ℓ_1 regularization. IEEE Trans. Signal Process., 69:489–499, 2021.
- [130] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. Annals of Statistics, 34:1436–1462, 2006.
- [131] N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. The annals of statistics, 37(1):246–270, 2009.
- [132] K. Minami. Degrees of freedom in submodular regularization: A computational perspective of Stein’s unbiased risk estimate. Journal of Multivariate Analysis, 175:104546, 2020.
- [133] S. Mousavi and J. 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.

- [134] R. Mukherjee, S. Mukherjee, and S. Sen. Detection thresholds for the β -model on sparse graphs. Ann. Statist., 46(3):1288–1317, 2018.
- [135] H. Nakashima and P. Graczyk. Wigner and wishart ensembles for graphical models, 2020.
- [136] R. Negrinho and A. Martins. Orbit regularization. Advances in neural information processing systems, 27, 2014.
- [137] S. Nowakowski, P. Pokarowski, and W. Rejchel. Group lasso merger for sparse prediction with high-dimensional categorical data, 2021.
- [138] OEIS Foundation Inc. The On-Line Encyclopedia of Integer Sequences, 2023. Published electronically at <http://oeis.org>.
- [139] M.-R. Oelker, J. Gertheiss, and G. Tutz. Regularization and model selection with categorical predictors and effect modifiers in generalized linear models. Stat. Model., 14(2):157–177, 2014.
- [140] A. Owrang, M. Malek-Mohammadi, A. Proutiere, and M. Jansson. Consistent change point detection for piecewise constant signals with normalized fused lasso. IEEE signal processing letters, 24(6):799–803, 2017.
- [141] P. O. Perry and P. J. Wolfe. Null models for network data, 2012.
- [142] E. J. G. Pitman. Some basic theory for statistical inference. Chapman and Hall, London; A Halsted Press Book, John Wiley & Sons, New York, 1979. Monographs on Applied Probability and Statistics.
- [143] A. Pósfai. Approximation theorems related to the coupon collector’s problem. PhD Thesis, 6, 2010.
- [144] A. Postnikov. Permutohedra, associahedra, and beyond. Int. Math. Res. Not. IMRN, (6):1026–1106, 2009.
- [145] J. Qian and J. Jia. On stepwise pattern recovery of the fused lasso. Computational Statistics & Data Analysis, 94:221–237, 2016.
- [146] K. R. Rao, N. Ahmed, and M. A. Narasimhan. Orthogonal transforms for digital signal processing. In Proceedings of the Eighteenth Midwest Symposium on Circuits and Systems (Concordia Univ., Montreal, Que., 1975), pages 1–6. Western Periodicals, North Hollywood, Calif., 1975.
- [147] A. Rencher. Methods of Multivariate Analysis. Wiley Series in Probability and Statistics. Wiley, 2003.
- [148] R. Riccobello, M. Bogdan, G. Bonaccolto, P. Kremer, S. Paterlini, and P. Sobczyk. Graphical modelling via the sorted l1-norm. arXiv preprint, arXiv:2204.10403, 2022.
- [149] A. Rinaldo. Properties and refinements of the fused lasso. Ann. Statist., 37(5B):2922–2952, 2009.
- [150] A. Rinaldo, S. E. Fienberg, and Y. Zhou. On the geometry of discrete exponential families with application to exponential random graph models. Electron. J. Stat., 3:446–484, 2009.

- [151] A. Rinaldo, S. Petrović, and S. E. Fienberg. Maximum likelihood estimation in the β -model. Ann. Statist., 41(3):1085–1110, 2013.
- [152] R. Rockafellar. Convex Analysis. Princeton University Press, 1997.
- [153] V. Saligrama and M. Zhao. Thresholded basis pursuit: Lp algorithm for order-wise optimal support recovery for sparse and approximately sparse signals from noisy random measurements. IEEE Transactions on Information Theory, 57(3):1567–1586, 2011.
- [154] R. Sankaran, F. Bach, and C. Bhattacharya. Identifying Groups of Strongly Correlated Variables through Smoothed Ordered Weighted L_1 -norms. In A. Singh and J. Zhu, editors, Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, volume 54 of Proceedings of Machine Learning Research, pages 1123–1131. PMLR, 20–22 Apr 2017.
- [155] R. Schneider. Convex bodies: the Brunn-Minkowski theory, volume 151 of Encyclopedia of Mathematics and its Applications. Cambridge University Press, Cambridge, expanded edition, 2014.
- [156] U. Schneider and P. Tardivel. The geometry of uniqueness, sparsity and clustering in penalized estimation. Journal of Machine Learning Research, 23(331):1–36, 2022.
- [157] A. Schrijver. Theory of linear and integer programming. Wiley-Interscience Series in Discrete Mathematics. John Wiley & Sons, Ltd., Chichester, 1986. A Wiley-Interscience Publication.
- [158] G. Schwarz. Estimating the dimension of a model. The Annals of Statistics, 6(2):461–464, 1978.
- [159] M. Schweinberger, P. N. Krivitsky, C. T. Butts, and J. Stewart. Exponential-family models of random graphs: Inference in finite-, super-, and infinite population scenarios, 2017.
- [160] A. Sepeshri and N. Harris. The accessible lasso models. Statistics, 51(4):711–721, 2017.
- [161] A. A. Sepeshri. The bayesian slope. arXiv: Methodology, 2016.
- [162] D. B. Sharma, H. D. Bondell, and H. H. Zhang. Consistent group identification and variable selection in regression with correlated predictors. Journal of Computational and Graphical Statistics, 22(2):319–340, 2013.
- [163] Y. She. Sparse regression with exact clustering. Electronic Journal of Statistics, 4:1055–1096, 2010.
- [164] T. Skalski. Remarks on Laplacian of graphical models in various graphs. In Geometric science of information, volume 12829 of Lecture Notes in Comput. Sci., pages 685–692. Springer, Cham, [2021] ©2021.
- [165] T. Skalski, P. Graczyk, B. Kołodziejek, and M. Wilczyński. Pattern recovery and signal denoising by slope when the design matrix is orthogonal. Probability and Mathematical Statistics, 42(2):283–302, 2022.
- [166] P. Sobczyk. Identifying low-dimensional structures through model selection in high-dimensional data. PhD Thesis, 2019.

- [167] U. Stadtmüller. A note on the law of iterated logarithm for weighted sums of random variables. Ann. Probab., 12(1):35–44, 1984.
- [168] B. G. Stokell, R. D. Shah, and R. J. Tibshirani. Modelling high-dimensional categorical data using nonconvex fusion penalties. J. R. Stat. Soc. Ser. B. Stat. Methodol., 83(3):579–611, 2021.
- [169] C. J. Stone. Large-sample inference for log-spline models. Ann. Statist., 18(2):717–741, 1990.
- [170] B. Stucky and S. van de Geer. Sharp oracle inequalities for square root regularization. J. Mach. Learn. Res., 18:Paper No. 67, 29, 2017.
- [171] W. Su and E. Candès. Slope is adaptive to unknown sparsity and asymptotically minimax. Annals of Statistics, 44:1038–1068, 2016.
- [172] S. Sullivant. Algebraic statistics, volume 194 of Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2018.
- [173] A. Takahashi and S. Nomura. Efficient path algorithms for clustered lasso and oscar. arXiv preprint arXiv:2006.08965, 2020.
- [174] P. J. Tardivel and M. 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, 49:1636–1668, 2022.
- [175] P. J. Tardivel, R. Servien, and D. Concordet. Simple expressions of the lasso and slope estimators in low-dimension. Statistics, 54(2):340–352, 2020.
- [176] R. Tibshirani. Regression shrinkage and selection via the Lasso. Journal of the Royal Statistical Society Series B, 58:267–288, 1996.
- [177] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 67(1):91–108, 2005.
- [178] R. J. Tibshirani. The lasso problem and uniqueness. Electronic Journal of statistics, 7:1456–1490, 2013.
- [179] R. J. Tibshirani and J. Taylor. Degrees of freedom in lasso problems. Annals of Statistics, 40:1198–1232, 2012.
- [180] B. Uniejewski, J. Nowotarski, and R. Weron. Automated variable selection and shrinkage for day-ahead electricity price forecasting. Energies, 9(8), 2016.
- [181] S. Vaïter, C. Deledalle, J. Fadili, G. Peyré, and C. Dossal. The degrees of freedom of partly smooth regularizers. Annals of the Institute of Statistical Mathematics, 69(4):791–832, 2017.
- [182] S. Vaïter, M. Golbabaee, J. Fadili, and G. Peyré. Model selection with low complexity priors. Information and Inference: A Journal of the IMA, 4(3):230–287, 2015.
- [183] M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). IEEE Transactions on Information Theory, 55:2183–2202, 2009.

- [184] M. J. Wainwright and M. I. Jordan. Graphical models, exponential families, and variational inference. Foundations and Trends in Machine Learning, 1(1-2):1–305, 2008.
- [185] N. Wang, J. Rauh, and H. Massam. Approximating faces of marginal polytopes in discrete hierarchical models. Ann. Statist., 47(3):1203–1233, 2019.
- [186] A. Weinstein, W. J. Su, M. Bogdan, R. F. Barber, and E. J. Candès. A power analysis for knockoffs with the lasso coefficient-difference statistic. arXiv preprint arXiv:2007.15346, 2020.
- [187] J. Whittaker. Graphical Models in Applied Multivariate Statistics. Wiley, 2009.
- [188] M. Yuan and Y. Lin. Model selection and estimation in the Gaussian graphical model. Biometrika, 94:19–35, 2007.
- [189] X. Zeng and M. Figueiredo. Decreasing weighted sorted ℓ_1 regularization. IEEE Signal Processing Letters, 21:1240–1244, 2014.
- [190] X. Zeng and M. A. Figueiredo. Solving oscar regularization problems by fast approximate proximal splitting algorithms. Digital Signal Processing, 31:124–135, 2014.
- [191] Y. Zhang and Z. Bu. Efficient designs of slope penalty sequences in finite dimension. In A. Banerjee and K. Fukumizu, editors, Proceedings of The 24th International Conference on Artificial Intelligence and Statistics, volume 130 of Proceedings of Machine Learning Research, pages 3277–3285. PMLR, 13–15 Apr 2021.
- [192] P. Zhao and B. Yu. On model selection consistency of Lasso. Journal of Machine Learning Research, 7:2541–2563, 2006.
- [193] G. Ziegler. Lectures on Polytopes, volume 152. Springer, New York, 2012.
- [194] H. Zou. The adaptive Lasso and its oracle properties. Journal of the American Statistical Association, 101:1418–1429, 2006.