

Multivariate statistics for dietary risk analysis Emilie Chautru

▶ To cite this version:

Emilie Chautru. Multivariate statistics for dietary risk analysis. Statistics [math.ST]. Télécom Paris-Tech, 2013. English. NNT: 2013ENST0045 . tel-01306948

HAL Id: tel-01306948 https://pastel.hal.science/tel-01306948

Submitted on 25 Apr 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.







Doctorat ParisTech

ΤΗÈSE

pour obtenir le grade de docteur délivré par

TELECOM ParisTech

Spécialité « Signal et Image »

présentée et soutenue publiquement par

Emilie CHAUTRU

le 6 septembre 2013

Statistiques multivariées

pour l'analyse du risque alimentaire

Directeur de thèse : **Stéphan CLÉMENÇON** Co-encadrement de la thèse : **Jean-Luc VOLATIER**

Jury Mme Anne RUIZ-GAZEN, Professeur, GREMAQ, Toulouse School of Economics M. Johan SEGERS, Professeur, ISBA, Université catholique de Louvain Mme Liliane BEL, Professeur, MIA, AgroParisTech Mme Anne-Laure FOUGERES, Professeur, CNRS UMR 5208, Université Lyon 1 M. Valentin PATILEA, Professeur, LSM Crest, Ensai M. Paul DOUKHAN, Professeur, AGM, Université de Cergy-Pontoise

Rapporteur Rapporteur Examinateur Examinateur Examinateur Invité

TELECOM ParisTech

école de l'Institut Mines-Télécom - membre de ParisTech 46 rue Barrault 75013 Paris - (+33) 1 45 81 77 77 - www.telecom-paristech.fr

FINANCEMENTS IN RA-ANSES-CREST









Inra - Unité Met@risk Anses - Direction de l'évaluation des risques Crest - Ensai - Laboratoire de statistiques d'enquête

ABSTRACT - RÉSUMÉ

ABSTRACT

Dietary risk analysis is a multidisciplinary field par excellence, requiring in particular a mixture of biological, sociological, chemical, cultural, economic, statistical and sanitary expertise to answer practical issues. In a global system where international exchanges are encouraged, where mass production favors cheap, profitable production strategies (GMO, pesticides, enriched animal food, addition of colorants, preservatives or artificial flavors, etc.), it is necessary to quantify the risks that result from such economic behaviors. The focus here is on the chronic (year-long) exposure to a set of food contaminants, the long-term toxicity of which is already well-known. Since food consumption is also the privileged way of supplying the body with necessary nutrients, nutritional benefits (or deficiencies) are taken into account as well. This thesis is dedicated to a collection of mathematical problems arising from both the analysis of such dietary risks and the nature of the data. First of all, the adequacy of the classical univariate long-term models that are recommended by Efsa (European Food Safety Authority) is discussed at length. They usually assume that observations have a log-normal distribution, thereby neglecting the possibility that they are in fact heavy-tailed. In terms of food risks, this would mean that a very high exposure to a single chemical element is a rare enough event to be neglected. When the component of interest clearly violates this strong assumption, extreme value theory is proposed as a relevant alternative, as revealed by a set of illustrations based on real data. Then, this sub-field of theoretical statistics is adapted to the analysis of the simultaneous exposure to many nutrients and contaminants. Following in the footsteps of classical machine learning techniques, using in particular the recent "Principal Nested Spheres" algorithm of S. Jung, I.L. Dryden and J.S. Marron, we construct a new model that identifies extreme dependencies in high dimension. This allows us to define some cocktails of chemicals that are jointly ingested in very high quantities. Remaining in a multivariate framework, we then move to the issue of dietary recommendations. In the line of the "minimum volume set" approach of C. Scott and R. Nowak, we propose an algorithm that selects food baskets realizing a compromise between toxicological risk and nutritional benefit. Finally, as consumption databases often result from complex survey schemes, the estimators constructed under the hypothesis that observations are independent and identically distributed can produce severely biased outcomes. In an attempt to take into account this preliminary sampling phase, we mimic the approaches of J Hajek and Y. Berger and focus on the specific family of Poisson-like survey plans. Under this framework, the asymptotic properties of Horvitz-Thompson empirical processes are inspected, before concluding this thesis on the introduction of a weighted version of the widely celebrated Hill estimator for the heavy-tail analysis of sampled observations.

RÉSUM É

Véritable carrefour de problématiques économiques, biologiques, sociologiques, culturelles et sanitaires, l'alimentation suscite de nombreuses polémiques. Dans un contexte où les échanges mondiaux facilitent le transport de denrées alimentaires produites dans des conditions environnementales diverses, où la consommation de masse encourage les stratégies visant à réduire les coûts et maximiser le volume de production (OGM, pesticides, nourriture enrichie donnée aux animaux en élevage, ajout de substances chimiques tels les colorants et les arômes artificiels, etc.) il devient nécessaire de quantifier les risques sanitaires que de tels procédés engendrent. Notre intérêt se place ici sur l'étude de l'exposition chronique, de l'ordre de l'année, à un ensemble quelconque de contaminants dont la nocivité à long terme est d'ores et déjà établie. Les dangers et bénéfices de l'alimentation ne se restreignant pas à l'ingestion ou non de substances toxiques, nous ajoutons à nos objectifs l'étude de certains apports nutritionnels. Les travaux réalisés au cours de cette thèse répondent à plusieurs problématiques mathématiques engendrées par l'analyse de ce type de risque alimentaire et par la nature des données utilisées à ces fins. Dans un premier temps, nous discutons la portée des modèles classiques dédiés à l'estimation de l'exposition de long terme à un unique composant alimentaire. Ces derniers, dont l'utilisation est recommandée par l'Efsa (European Food Safety Authority), requièrent généralement la log-normalité des observations et interdisent par là-même les distributions à queue plus épaisse, où l'élément chimique étudié aurait une probabilité non négligeable d'être ingéré en de très grandes quantités. Nous montrons alors, exemples à l'appui, comment la théorie des valeurs extrêmes peut être utilisée dans de telles situations. Ce pan de théorie statistique est ensuite adapté à l'étude des fortes expositions à un nombre quelconque (potentiellement très grand) de nutriments et contaminants. En nous inspirant de techniques traditionnelles de l'apprentissage statistique, notamment le récent algorithme des "Principal Nested Spheres" développé par S. Jung, I.L. Dryden et JS. Marron, nous construisons un modèle ouvrant l'étude des dépendances extrêmes à la grande dimension, qui nous permet en particulier de définir des groupes d'éléments auxquels les consommateurs sont simultanément sur-exposés. Toujours dans une optique multivariée, nous nous éloignons ensuite des expositions extrêmes pour nous placer du côté des consommations alimentaires. En nous basant sur une approche de type "ensemble de volume minimum" comme introduite par C. Scott et R. Nowak, nous proposons un algorithme pour identifier des paniers de produits réalisant un compromis entre risque toxicologique et bénéfice nutritionnel. Enfin, les données alimentaires disponibles résultant souvent de plans de sondage non triviaux, les estimateurs construits sous l'hypothèse que les données sont indépendantes et identiquement distribuées peuvent produire des résultats biaisés. Tentant de prendre en compte cette étape préliminaire d'échantillonnage, nous nous concentrons dans la lignée des travaux de J. Hàjek et Y. Berger sur la famille des plans de sondage de type Poisson et étudions le comportement asymptotique des processus empiriques pondérés à la manière Horvitz-Thompson. Suivant la même approche, nous proposons finalement une variante de l'estimateur de Hill pour l'analyse des distributions à queue épaisse adaptée au cadre des données de sondage.

PUBLICATIONS

submitted for publication

E. Chautru, S. Clémençon, Dimension reduction in multivariate heavy-tail analysis, 2013

P. Bertail, E. Chautru, S. Clémençon, Empirical Processes in Survey Sampling, 2013

P. Bertail, E. Chautru, S. Clémençon, On Tail Index Estimation Based on Survey Data, 2013

to be submitted

E. Chautru, S. Clémençon, J. Tressou, Determination of an optimal diet by minimum volume sets, 2013

E. Chautru, J.L. Volatier, Heavy-tail modeling of Usual Intakes, 2013

REMERCIEMENTS

En premier lieu, je souhaite vivement remercier les membres de mon jury d'avoir accepté d'évaluer mon travail et de me faire profiter par là-même de leurs conseils avisés. Je suis tout particulièrement reconnaissante à Johan Segers et Anne Ruiz-Gazen d'avoir sacrifié une partie de leur mois d'août pour éplucher mon mémoire en détail. Merci de l'intérêt que vous avez porté à ma thèse et merci pour les critiques constructives qui m'ont permis d'en améliorer la substance. Je suis de même très honorée par la présence dans mon jury de thèse de Liliane Bel, Anne-Laure Fougères, Valentin Patilea et Paul Doukhan, que j'admire personnellement tant sur le plan scientifique que sur le plan humain.

Il va sans dire que les lauriers devraient revenir à Jean-Luc et Stéphan, mes deux directeurs de thèse. Merci à tous deux de m'avoir introduite au risque alimentaire, l'un des champs d'étude les plus interdisciplinaires que j'aie eu l'occasion de découvrir, ainsi que de m'avoir suivie sans relâche durant ces quatre longues années. Votre collaboration aura généré un équilibre entre thématiques théoriques et appliquées dans ma recherche, point qui me tenait à cœur lorsque je me suis décidée à démarrer un doctorat et que j'ai pu réaliser grâce à vous.

Je me dois aussi de reconnaître la partie prenante qu'ont pris dans mon parcours ceux avec qui j'ai écrit la plus grande partie des papiers présentés dans cette thèse : Patrice et Jessica. Merci à tous deux pour votre gentillesse et votre patience, pour tout le temps et l'énergie que vous m'avez consacrés et qui m'ont permis d'avancer.

Bien entendu, cette thèse n'aurait pu avoir lieu sans le soutien financier de l'Inra, de l'Anses et du Crest-Ensai. Etre ainsi rattachée à un nombre non négligeable d'unités de recherche s'est avéré des plus stimulant et enrichissant d'un point de vue intellectuel. A commencer par Met@risk, où j'ai pu profiter des conseils pratiques de l'équipe informatique, découvrir l'aspect sociologique de la nutrition au contact des sociologues, partager avec les statisticiens, etc. . Un grand merci des plus enthousiastes à tous ceux que j'ai rencontrés au sein de cette unité formidable, Laura, Lydie, Damien, Stéphane, Juliette, Sandrine, Eve, Isabelle, Jessica, Mélanie, Clémence, Max et les autres.

Du côté de l'Anses, je souhaite en particulier remercier Ariane, Carine, Véronique, Amélie et Camille, non seulement pour leur accueil chaleureux dans un milieu aux codes quelque peu différents de ceux de la recherche universitaire, mais aussi pour tous leurs conseils et leur expertise sur les bases de données que j'ai utilisées dans ma thèse et les enjeux mathématiques inhérents au risque alimentaire.

Plus loin en Bretagne, dans la campagne Bruzoise, où les Ensaiens ont fait preuve d'une hospitalité remarquable, je tiens à dire un grand merci à Eric et Valentin pour m'avoir soutenue lors des auditions et permis de rejoindre les Rennais en cette dernière année de thèse. Merci aussi à François, Guillaume, Brigitte, David, Nicolas, Nicolas, Gaspar, Laurent, Stéphane, Lionel, Magalie, Myriam et tous les permanents du Crest auprès de qui j'ai passé une excellente année. N'oublions pas les Insee, tout particulièrement Philippe et Jean-Michel, grâce à qui j'ai pu payer mon loyer chaque mois. Je finirai cette liste par mes deux compagnons de route, Cyril et Samuel, sans qui Rennes aurait été bien morne. Merci à vous deux pour votre amitié et votre soutien inconditionnel durant cette période éprouvante.

Enfin, à Télécom, je souhaite remercier une fois de plus Eric pour son aide lors de la recherche de financements supplémentaires pour terminer mon doctorat. Merci à Jérémie, cobureau des débuts qui m'aura fait prendre un grand recul sur la place que je peux occuper dans la recherche et aux permanents qui m'ont aiguillée; à Sarah, Steffen, Olaf, Malika, Marine, Anne-Laure, Alexandre, Jeoffrey, Tabea, Onur, les vieux de la vieille qui sont partis beaucoup trop tôt à mon goût; à toute la bande des thésards, Emilie, Amandine, Olivier, Cristina, Sylvain, Yao, Nicolas, Benjamin, Aymeric, Sébastien et des post-doc, Charanpal, Julien, Romaric, j'en oublie et des meilleurs. Une pensée toute particulière pour mes frères de thèse, Sylvain, Ruocong, Antoine, Eric, ainsi que pour Billy et Billy; Dareau manque sérieusement de chevelus déjantés parfumés à la bière belge depuis votre départ.

Beaucoup de personnes extérieures à mes nombreux lieux de travail ont aussi eu un impact non négligeable sur ma recherche. Eric, qui m'a écoutée des heures durant parler de mathématiques un kir à la main, qui m'a épaulée, conseillée, guidée tout au long de ce parcours, je ne te dirai jamais assez merci ; Olivier, qui m'a introduite auprès de nombreux chercheurs et encouragée dans ma recherche; Xavier, sans qui je me serais sentie bien seule lors de ma toute première conférence; Julyan, Ronan, Emmanuelle, avec qui j'ai eu le plaisir de donner des TD et qui m'ont aidée à prendre confiance en moi ; M. Mano, mon professeur de mathématiques en hypokhâgne et khâgne, qui m'a donné le goût de la matière. Merci aussi à Céline et Meriem et à tous les anciens de l'Ensae, Miki, Adélaïde, Marion, Arnaud, Max, Nathanaël, Maha, Sophie, Alexandra, Simon, la JE et le forum, avec qui j'ai découvert et pris goût aux statistiques. Tout particulièrement merci à Stéphanie, une des meilleures partenaires de mémoire (et de sorties), que j'encourage vivement à quitter l'Insee pour monter une agence matrimoniale. Sans oublier les anciens d'Abu Dhabi, Nassime, Antoine, Pohe, Pierre, Cédric, toujours fidèles au poste plus de dix ans après nos aventures lycéennes.

Je remercie également toute ma famille, et plus encore mes parents, mon frère et ma sœur (à qui je dois quelques uns des dessins de mes slides), pour m'avoir toujours soutenue et encouragée, que ce soit dans mes études, ma vie professionnelle ou ma vie privée. Ce travail n'existe que grâce à vous (notamment parce que vous avez pensé à regarder les résultats d'admissibilité à l'Ensae alors que j'avais oublié que je m'étais inscrite au concours). Isabelle et Mélanie, je vous dis merci d'être mes amies depuis plus de 20 ans et de continuer encore et toujours d'occuper une place importante dans ma vie.

Enfin, à Guillaume, je dis merci et re-merci. Merci de m'avoir redonné le goût des mathématiques à un moment où je n'y croyais plus, merci d'avoir relu l'intégralité de mon mémoire, merci de m'avoir supportée, clope au bec, durant les moments les plus stressants, la liste est longue... But most of all, thank you for being yourself.

CONTENTS

Pr	efac	е		1	
р	régimes alimentaires, exposition individuelle et statistique			3	
	p.1	Qu'es	t-ce que le risque alimentaire?	4	
		p.1.1	Evaluation du risque	4	
		p.1.2	Gestion du risque	5	
		p.1.3	Communication autour du risque	5	
	p.2	Des so	purces multiples d'information	6	
		p.2.1	Données de consommation	6	
		p.2.2	Données de composition	9	
	р.З	Estimation des risques alimentaires :			
		petit a	aperçu de l'état de l'art	11	
		p.3.1	Exposition extrême de long terme	12	
		p.3.2	Exposition simultanée à un ensemble de nutriments et		
			contaminants	12	
		p.3.3	Données de sondage	13	
	p.4	Object	tifs et contributions principales	13	
		p.4.1	Estimation de l'exposition extrême de long terme à un unique		
			élément chimique	14	
		p.4.2	Evaluation de la sur-exposition simultanée à un cocktail de		
			composants alimentaires	15	
		p.4.3	Ensembles de volume minimum et approche risque-bénéfice	16	
		p.4.4	Processus empiriques dans le cadre des sondages	16	
		p.4.5	Estimation de l'indice de valeurs extrêmes à partir de données		
			de sondage	17	
Th	esis			19	
1	intı	roduct	tion: of dietary habits and exposure	21	
	1.1	Gettin	g a taste of dietary risk analysis	22	
		1.1.1	Risk evaluation	22	
		1.1.2	Risk management	22	
		1.1.3	Risk communication	23	
	1.2	Collecting information: multiple			
		sources, multiple difficulties			
		1.2.1	Consumption data	24	
		1.2.2	Composition data	26	
	1.3	Statist	ical assessment of dietary risks: a review of the literature	28	

		1.3.1	Extreme long-term exposure	29
		1.3.2	Simultaneous exposure to multiple food components	29
		1.3.3	Survey data	30
	1.4	Object	ives and main contributions	30
		1.4.1	Extreme chronic exposure to one chemical	31
		1.4.2	Simultaneous over-exposure to many food chemicals	31
		1.4.3	A minimum volume set approach to dietary risk-benefit analysis	32
		1.4.4	Empirical processes in survey sampling	32
		1.4.5	Tail index estimation based on survey data	33
2	ext r	· eme c	hronic exposureto one chemical	35
	2.1	Introd	uction	35
	2.2	Mixed	-effects models to assess usual	
		intake	\$	36
		2.2.1	The BetaBinomial-Normal model	37
		2.2.2	The LogisticNormal-Normal model	38
	2.3	Heavy	r-tail modeling of usual intakes	39
		2.3.1	Extreme value theory and heavy-tailed distributions	39
		2.3.2	Testing the heavy-tail assumption	41
		2.3.3	Assessing extreme quantities	43
	2.4	Case s	study: tails of iron, zinc, calcium and retinol	46
		2.4.1	Description of the data	46
		2.4.2	Results	46
	2.5	Discus	ssion	51
	2.6	Suppl	ements: on the second order	
		param		51
3	sim	nultaneous over -exposur e to many food chemical s 5		
3.1 Sta		Statist	ical challenges and objectives	55
	3.2	Hypot	heses and notations	57
		3.2.1	Notations	58
		3.2.2	General setting and main hypotheses	58
	3.3	A mix	ture model under multivariate	
		regula	r variation	60
		3.3.1	Exponent and spectral (probability) measures	60
		3.3.2	Mixture model of the spectral probability measure	61
		3.3.3	Latent variables representation	63
	3.4	Statist	ical inference	65
		3.4.1	Preliminaries	66
		3.4.2	Dimension reduction and clustering	66
		3.4.3	Identifying groups of asymptotically dependent variables	68
	3.5	Nume	rical experiments	70
		3.5.1	Settings	70

		3.5.2	Results	71
	3.6	Applic	cation to dietary risk assessment	73
		3.6.1	Data and required assumptions	74
		3.6.2	Analysis of extreme dependencies	74
	3.7	Discus	ssion	77
	3.8	Proofs and supplements		
		3.8.1	Intra-class regular variation	79
		3.8.2	Face-characterizing functional	83
		3.8.3	About $p_{j;}$ (k)	86
		3.8.4	A little more on PNS and spherical k-means	90
4	a mi	inimur	n volume set approach to dietary risk-benefit analysis	99
	4.1	Theore	etical analysis and methods	100
		4.1.1	Assessment of dietary exposure to chemicals and nutrients:	
			a MV-set formulation	100
		4.1.2	Uniform approximation of generalized U-statistics by their	
			incomplete versions	102
		4.1.3	Empirical MV-set estimation based on hypercubes	107
		4.1.4	Optimal regions in the consumption space	108
	4.2	Proofs	and supplements	112
		4.2.1	Maximal deviation	112
		4.2.2	Maximal deviation in dietary risk analysis	115
		4.2.3	Optimal dietary habits	115
5	emp	irical	processes in survey sampling	117
	5.1	Backgi	round and Preliminaries	118
		5.1.1	Survey sampling: some basics	118
		5.1.2	Empirical process indexed by classes of functions	124
	5.2	Empirical process in survey sampling		
		5.2.1	The Horvitz-Thompson empirical process	127
		5.2.2	Alternative estimate in the Poisson sampling case	128
	5.3	Asym	ptotic results	129
		5.3.1	Limit of the empirical process for the Poisson survey scheme	130
		5.3.2	The case of rejective sampling	133
	5.4	Applic	cation to non-parametric statistics	135
		5.4.1	Hadamard differentiable functionals	136
		5.4.2	Fréchet differentiable functionals	137
		5.4.3	Simulation-based Gaussian asymptotic confidence regions	138
	5.5	Proofs	and supplements	143
		5.5.1	Limit of the covariance operator	143
		5.5.2	FCLI in the Poisson survey case	144
		5.5.3		146
		5.5.4	FCLT in the rejective survey case	147

		5.5.5	CLT for Hadamard differentiable functionals	148
6	tail	index	estimation based on survey data	149
	6.1	Backgr	ound and Preliminaries	150
		6.1.1	Survey sampling	150
		6.1.2	Tail index inference - the Hill estimator	152
	6.2	The Hill estimator in survey sampling		
		6.2.1	The Horvitz-Thompson variant of the Hill estimator	154
		6.2.2	Consistency of $H^{\square}_{K;N}$	155
	6.3	Asymp	Distotic normality of $H^{\square}_{K;N}$	157
		6.3.1	The case of the Poisson survey scheme	159
		6.3.2	Extension to rejective sampling schemes	161
	6.4	Practical Issues and Illustrative		
		Experi	ments	163
		6.4.1	On the choice of an optimal k	163
		6.4.2	Numerical experiments	164
	6.5	Discus	sion	168
	6.6	Proofs	and supplements	169
		6.6.1	Consistency of the Horvitz-Thompson variant of the Hill	
			estimator	169
		6.6.2	Limit distribution of $H^p_{K;N}$ in the Poisson survey case	172
		6.6.3	Limit distribution of $H^{\square}_{K;N}$ in the rejective survey case $\ .\ .\ .$.	181
bibl iogr aph y				184

ACRONYMS

institutes

- Anses French agency for food, environmental and occupational health safety
- Efsa European food safety authority
- FAO Food and Agriculture Organization
- Inra French national institute for agricultural research
- Insee French national institute of statistics and economic studies
- Unep United Nations Environment Program
- WHO World Health Organization

dietary risk analysis abbr eviations

- DIL Dietary intake limit
- CIQUAL Centre d'information sur la qualité des aliments
- FFQ Food frequency questionnaire
- FPQ Food propensity questionnaire
- GMO Genetically Modified Organism
- INCA Etude nationale des consommations alimentaires
- TDS Total dietary survey
- 24H 24-hour recall

chemical symbols

Ca	Calcium
Cd	Cadmium
Fe	Iron
MeHg	Methylmercury
Na	Sodium
PCB-DL	Dioxin-like polychlorinated biphenyls
Re	Retinol
Zn	Zinc

mathematical abbreviations

cdf	Cumulative distribution function
EVT	Extreme value theory
iid	Independent and identically distributed
rv	Random variable
svf	Slowly varying function

PREFACE

RÉSUMÉ EN FRANÇAIS DES TRAVAUX DE THÈSE

Ρ

RÉGIMES ALIMENTAIRES, EXPOSITION INDIVIDUELLE ET STATISTIQUE

Variées et nombreuses sont les études scientifiques pouvant être caractérisées comme analyses du risque alimentaire. Cette désignation fort générale issue du jargon de la santé publique fait référence à toute entreprise visant à la détection, à la compréhension, et au traitement des dangers liés à l'alimentation. Parmi ces derniers, citons par exemple l'une des problématiques classiques des producteurs : le développement de bactéries au cours des différentes étapes de production, de transport, de distribution et de stockage qui précèdent la consommation (Rigaux et al., 2012). A une toute autre échelle, l'industrie agroalimentaire doit évaluer les impacts à la fois économiques et sanitaires de techniques de production de masse tels l'utilisation de pesticides (http://www.efsa.europa.eu/en/topics/topic/pesticides.htm) ou l'introduction d'organismes génétiquement modifiés dans les plantations (http:// www.efsa.europa.eu/en/topics/topic/gmo.htm). Selon leur composition chimique et leur mode d'utilisation, les ustensiles de cuisine peuvent aussi représenter une source de danger, comme en témoigne la récente polémique à propos de la présence de bisphénol A dans les biberons (cf. l'avis de l'Anses y ayant fait suite, Anses, 2013). D'autres problèmes de santé peuvent être engendrés par des troubles du comportement alimentaire (e.g. les personnes souffrant d'anorexie aiguë sont particulièrement sujettes aux maladies provoquées par des carences nutritionnelles) ou par des facteurs biologiques particuliers (e.g. le diabète). A la lumière de ces quelques exemples, il est clair que selon la question examinée, des connaissances spécifiques en chimie, biologie, médecine, sociologie, économie ou même psychologie peuvent être requises, sans oublier la modélisation probabiliste et la gestion informatique des données. L'analyse du risque alimentaire est ainsi un domaine d'étude multidisciplinaire par excellence. Le lecteur curieux d'en apprendre plus est invité à consulter Feinberg et al. (2006), manuel expliquant de manière exhaustive les tenants et aboutissants de l'analyse du risque alimentaire selon un point de vue interdisciplinaire.

Le présent travail est dédié à l'étude d'un type spécifique de risque alimentaire : indépendamment de tout processus de production, de stockage ou de cuisson, ignorant les prédispositions biologiques extraordinaires, nous nous intéressons exclusivement à la très forte (ou très faible) exposition sur le long terme à certains composants alimentaires de la population française dans son ensemble. Après avoir introduit quelques concepts élémentaires du risque alimentaire en section P.1, décrit les données disponibles en section P.2 et brièvement exposé l'état de l'art sur la modélisation probabiliste des risques alimentaires en section P.3, nous présentons en détail en section P.4 les diverses problématiques traitées au cours de cette thèse avant d'en annoncer les contributions scientifiques principales.

P.1 qu'est -ce que le risque alimentaire?

Une fois qu'un type spécifique de risque alimentaire a été porté à la connaissance de tous, son analyse consiste en trois étapes distinctes, respectivement qualifiées d'évaluation du risque, de gestion du risque et de communication autour du risque. Les divers problèmes soulevés dans chacune de ces phases de recherche sont résumés dans les paragraphes suivants. Afin d'en faciliter la compréhension, nous en illustrons la substance à l'aide de l'exemple concret des régimes riches en sel, dont l'abus récurrent peut favoriser les problèmes cardiaques (se référer par exemple à la page Internet http://www.anses.fr/en/content/salt et aux références qui y sont mentionnées).

P.1.1 Evaluation du risque

L'évaluation (ou appréciation) du risque peut elle-même être décomposée en quatre sous-étapes, qui consistent à successivement identifier puis caractériser les dangers potentiels et leur probabilité d'apparition. Dans le cas des régimes riches en sel, le danger n'est autre que celui d'ingérer de manière journalière ou hebdomadaire une trop grande quantité de sodium (Na dans le tableau périodique des éléments), ce qui sur le long terme pourrait engendrer une dégradation du système cardiovasculaire. La connaissance de ces effets nocifs est due à tout un ensemble d'études chimiques et biologiques concernant l'assimilation, l'action et l'élimination du sodium dans le corps humain. Ce processus, appelé identification du danger (Barlow et al., 2002), est suivi de recherches supplémentaires permettant de définir à partir de quel niveau de consommation le composant devient toxique, c'est-à-dire de caractériser le danger. Pour ce faire, des tests dits dose-réponse sont en général réalisés in vitro ou in vivo sur des animaux avant d'en étendre les résultats à l'espèce humaine par le biais de modèles dédiés (Dybing et al., 2002). Une fois le processus de contamination compris dans son intégralité, il devient possible d'en estimer la fréquence d'apparition au sein d'une population choisie. Cette dernière phase requiert en premier lieu la description détaillée des distributions de la consommation et de l'exposition dans la population concernée (Kroes et al., 2002), avant de les comparer à des doses maximales recommandées, déterminées par des experts à l'issue des tests dose-réponse précédemment évoqués (Renwick et al., 2003) et possiblement raffinés par des procédés mathématiques (Edler et al., 2002). En fonction de la nature des données disponibles, les modèles statistiques utilisés à ces fins peuvent inclure une dimension temporelle (Bertail

et al., 2010, 2008; Allais and Tressou, 2009), prendre en compte des caractéristiques comportementales ou biologiques individuelles, ou encore tenter de compenser le manque d'information sur les habitudes alimentaires de long terme (Dodd et al., 2006). Ce dernier problème est discuté en détail en section P.3 et au chapitre 2.

P.1.2 Gestion du risque

Après que le risque a été évalué, il est nécessaire d'en identifier les déterminants, d'en évaluer les impacts relatifs et de définir la stratégie la plus rapide et la plus efficace permettant de le réduire. Ce n'est autre que l'objectif de la gestion du risque. Par exemple, la sur-exposition au sodium peut être attribuée à la consommation répétée de produits à haute teneur en sel comme les plats tout prêts ou les biscuits. Un plan simple d'action dans ce cas serait de prévenir les consommateurs des dangers potentiels et d'imposer en parallèle aux compagnies agroalimentaires de limiter les doses de sel ajoutées à leurs préparations. Dans des situations plus extrêmes, après une évaluation rigoureuse des conséquences économiques d'une telle opération, des produits considérés trop dangereux peuvent même être retirés du marché.

P.1.3 Communication autour du risque

La dernière étape de l'analyse du risque alimentaire est appelée communication autour du risque. Elle peut être mise en place à n'importe quel moment de l'analyse et peut s'adresser aussi bien aux scientifiques et aux gestionnaires du risque qu'à l'industrie agroalimentaire ou à la population. En tant que telle, elle n'est pas cantonnée aux campagnes publicitaires de santé publique du type "évitez de manger trop salé". S'y rapportent aussi tous les rapports scientifiques concernant les procédures considérées appropriées pour détecter et quantifier les risques, la publications d'études cas-témoin, les discussions internationales sur les politiques agroalimentaires, etc. Dans cette thèse nous sommes tout particulièrement intéressés par la production de résultats permettant d'aiguiller les recherches futures et, lorsque possible, de définir des lignes directrices de consommation simples et générales à l'intention de la population.

La grande majorité de nos travaux correspond ainsi à la phase d'évaluation du risque. En quelques mots, nous proposons des méthodes statistiques visant à estimer certaines caractéristiques (comme des quantiles) de la très forte exposition à un ou plusieurs composants alimentaires sur une longue période de temps. A l'exception près du chapitre 4, dans lequel nous élaborons une procédure pour déterminer des paniers de consommation réalisant un compromis entre risque toxicologique et bénéfice nutritionnel, les problématiques relevant de la gestion du risque ne sont pas abordées. Les éléments chimiques que nous prenons en considération sont les nutri-

ments et contaminants dont les effets sanitaires liés à une sur- ou une sous-exposition chronique ont d'ores et déjà été établis. En particulier, les risque aigus sont ignorés, tels ceux impliquant la contamination bactérienne de la nourriture, qui peuvent affecter l'organisme en quelques jours ou même quelques heures seulement. L'analyse statistique est réalisée dans ce contexte à partir de formats standards de bases de données, dont une description rapide est proposée ci-après.

P.2 des sour ces multiples d'information

A fin de permettre l'analyse statistique de l'exposition chronique d'une population donnée à une collection de composants alimentaires, il est désirable d'observer de tels types d'exposition sur un large échantillon d'individus pendant une longue période de temps. Malheureusement, les quantités de nutriments et de contaminants ingérés durant un repas ne sont pas directement mesurables; ils ne peuvent qu'être estimés à l'aide de méthodes variées. L'une des procédures classiques consiste à détecter puis quantifier certains marqueurs biologiques, dont la présence est intimement liée au niveau d'exposition (voir par exemple la page Internet de l'Anses http://www. anses.fr/sites/default/files/documents/RSC1205-DossierParticipants.pdf, Où sont listées les présentations d'un colloque de mai 2012 dédié à ce sujet). D'un point de vue pratique, cela nécessite de collecter puis d'analyser la composition chimique d'échantillons de fluides corporels, de cheveux ou de peau d'un nombre important de personnes. Au niveau national, le coût de telles enquêtes peut contraindre la quantité de sondés au détriment de l'efficacité statistique. En outre, de nombreux éléments chimiques peuvent être assimilés par l'intermédiaire d'autres substances que les aliments (e.g. l'air), or les marqueurs biologiques ne permettent pas de distinguer les diverses sources d'exposition (Sirot et al., 2009). Pour pouvoir étudier exclusivement l'impact de l'alimentation sur l'exposition individuelle, il est possible d'utiliser une méthode alternative où deux bases de données sont combinées, la première listant les habitudes alimentaires d'un ensemble de consommateurs et la seconde indiquant les teneurs en composants chimiques d'une nomenclature fine de produits. Les principales caractéristiques de ces données, respectivement dites de consommation et de composition, sont décrites dans les paragraphes suivants.

P.2.1 Données de consommation

De nombreux types de bases de données peuvent être utilisées pour évaluer l'exposition chronique de la population française à des nutriments et des contaminants alimentaires, allant de cohortes de ménages dont les dépenses en nourriture sont strictement suivies (Secodip, Nichèle et al., 2008) à des sondages en ligne remplis sur la base du volontariat (Nutrinet, Hercberg et al., 2010) en passant par des enquêtes

institutionnelles dédiées à un niveau national (INCA2, Afssa, 2009). Nous nous intéressons ici à deux grandes familles de données, appelées en anglais 24-hour recalls (24H) et food frequency ou propensity questionnaires (FFQ/ FPQ). La première de ces catégories correspond à des enquêtes de grande envergure (par exemple nationale) où un échantillon de sondés note en détail les quantités d'aliments qu'ils consomment durant 2 à 7 jours, parfois consécutifs, mais la plupart du temps sélectionnés aléatoirement dans l'année. Les prises alimentaires sont indiquées soit de manière exacte, soit relativement à des photographies d'assiettes plus ou moins remplies, qui sont mises à disposition par le sondeur. Bien entendu, la période d'observation étant particulièrement courte, les 24H peuvent sembler inappropriés pour l'étude de l'exposition de long terme (Counil et al., 2006). C'est pourquoi il est souvent demandé aux sondés de remplir en parallèle lesdits FFQ/ FPQ (questionnaires de fréquence ou de propension des prises alimentaires si l'on traduit mot-à-mot). Ils permettent en particulier de distinguer les personnes ne consommant jamais certains produits de celles qui en mangent occasionnellement. Le lecteur curieux d'en apprendre davantage au sujet de ces deux grandes sources d'information est invité à se référer à van Klaveren et al. (2012); EFSA (2006); Dodd et al. (2006). Tout au long de la présente thèse nous utilisons la base INCA2 qui, ainsi que décrit ci-après, peut être assimilée à un mélange de 24H et de FFQ.

P.2.1.1 Les habitudes alimentaires en France: INCA2

INCA2 (second opus de l'enquête nationale sur les habitudes alimentaires individuelles en France) est une enquête d'envergure nationale mise en place par l'Anses (agence française de sécurité sanitaire de l'alimentation, de l'environnement et du travail) en collaboration avec l'Insee (institut national de la statistique et des études économiques) entre décembre 2005 et avril 2007. Elle collecte les informations concernant les comportements alimentaires de 2624 adultes et 1455 enfants sélectionnés aléatoirement dans la population française.

Ú Plan de sondage Les individus composant INCA2 ont été sélectionnés selon un plan de sondage complexe à plusieurs degrés, construit afin de produire un échantillon représentatif de la population française selon des critères géographiques, sociologiques et économiques. Ces mêmes variables auxiliaires ont ensuite été utilisées lors d'une étape de redressement, réalisée dans un second temps afin prendre en compte la possible non-réponse de certains sondés ainsi que les fluctuations d'échantillonnage. Une description détaillée du plan de sondage d'INCA2 est disponible dans le rapport Afssa (2009, Chapitre 2 et Appendice 2).

Ú Données alimentaires Tous les participants ont indiqué la nature et la quantité des aliments qu'ils ont consommés durant les 7 jours consécutifs de l'enquête. Ces produits ont été classifiés selon une nomenclature exhaustive de 1342 aliments, regroupés en 123 sous-classes et 45 catégories plus larges (Afssa, 2009, Section 2.2.6.3 et Appendice 1). Afin de les aider à évaluer les quantités mangées, un carnet de photographies représentant des assiettes et des verres progressivement remplis a été mis à la disposition des sondés, servant d'étalon lorsque des mesures exactes ne pouvaient être réalisées. Des informations complémentaires sur les conditions des repas ont aussi été collectées, concernant par exemple le lieu, l'heure et la durée de ces événements. Six catégories de repas sont indiquées dans la base finale, à savoir le petit déjeuner, la collation du matin, le déjeuner, la collation de l'après-midi, le dîner et la collation du soir. Pour prendre en compte les variations saisonnières des habitudes alimentaires, les participants ont été contactés à des moments aléatoires de l'année, durant l'une des 3 vagues successives de collecte des données. Un ensemble de questions concernant la prise de compléments alimentaires a aussi été introduit (Afssa, 2009, Section 2.2.6.4 et Chapitre 8).

Ú Information auxiliaire En complément des régimes alimentaires, un nombre important d'informations au sujet des sondés a été noté, allant des caractéristiques sociologiques (e.g. diplôme, profession, revenu du ménage, nationalité) aux préférences alimentaires, en passant par les activités physiques (e.g. type, fréquence, durée), l'historique médical (e.g. troubles du comportement alimentaire) et d'autres indications générales (e.g. age, poids, sexe).

Ú Consommateurs occasionnels Lorsque l'intérêt est porté sur les habitudes alimentaires et l'exposition de long terme, 7 jours peuvent sembler être une période arbitrairement courte d'observation. En effet, en une semaine seulement, les enquêtés ne peuvent couvrir l'ensemble de leur répertoire alimentaire et de nombreux produits de la nomenclature de référence ne sont pas consommés. Un tel phénomène rend difficile la distinction entre les consommateurs occasionnels et ceux qui ne mangent jamais de certains produits. Afin de remédier à ce problème, il a été demandé aux individus constituant la base INCA2 de décrire leurs régimes usuels et de déclarer clairement à quelle catégorie de consommateurs ils appartiennent. Ces questions additionnelles correspondent aux fameux FFQ/ FPQ précédemment mentionnés qui, combinés aux bases de type 24H, facilitent grandement la modélisation statistique (Dodd et al., 2006 et van Klaveren et al., 2012, Sections 3.6 et 3.8).

Ú Information incomplète Comme nous venons de l'évoquer, n'observer les enquêtés que durant 7 jours peut s'avérer être un inconvénient majeur. Cependant, du point de vue des sondés, une semaine complète peut paraître particulièrement long, et la qualité des données peut en pâtir. En effet, comme indiqué dans le rapport Afssa (2009, Section 2.4.2), plusieurs participants ont ponctuellement omis d'indiquer leurs consommations, tandis que d'autres ont largement sous-estimé les quantités ingérées

durant leurs repas. Dans nos calculs nous avons pris le parti d'ignorer ces lacunes en rapportant les informations disponibles à l'échelle de la semaine.

Ú Dépendance temporelle La seule réelle différence entre INCA2 et les bases de type 24H tient à la nature consécutive des 7 jours d'observation. Ainsi, bien qu'une semaine complète semble clairement plus appropriée que 2 jours d'enquête pour une analyse de long terme, il est alors plus difficile d'ignorer la dépendance temporelle.

Ú Traitement préliminaire d'INCA2 Les travaux réalisés dans cette thèse sont dédiés à l'analyse globale des habitudes alimentaires de long terme des adultes en France. En raison de leurs besoins nutritionnels spécifiques et du caractère temporaire de leur état, nous avons décidé de ne prendre en compte ni les femmes enceintes ni les femmes allaitant dans nos calculs. D'autres individus ont été exclus de nos analyses, notamment ceux ayant omis de renseigner des variables essentielles comme le poids corporel, amenant l'échantillon initial à 2488 unités. Bien entendu, nous aurions pu tenter d'appliquer des méthodes classiques pour remédier au problème des valeurs manquantes. Néanmoins, de telles considérations dépassent le cadre de nos travaux. En particulier, ces techniques statistiques sont en général mises en place pour éviter de dégrader l'estimation de phénomènes moyens. Or notre intérêt est porté sur les événements extrêmes (minimum et maximum). Par ailleurs, la proportion de valeurs manquantes dans l'échantillon étant infime (seuls 91 individus sont concernés, soit 3; 7% des enquêtés), nous avons préféré les ignorer.

P.2.2 Données de composition

Les bases de données dites de composition regroupent en général un ensemble de mesures chimiques réalisées sur des groupes plus ou moins raffinés d'aliments. Elles peuvent provenir d'enquêtes de natures variées. Par exemple, des plans de surveillance sont mis en place pour vérifier la sûreté de produits suspects et induisent la collecte de données à leur sujet. Bien que de nombreux éléments du répertoire alimentaire soient ainsi négligés, le nombre de mesures réalisées dans ce cadre est particulièrement conséquent, avantage non-négligeable pour toute analyse statistique. Néanmoins, en raison de la nature suspecte des produits étudiés, les bases de données résultantes ne peuvent être utilisées pour l'évaluation des risques alimentaires à grande échelle sans introduire un biais non-négligeable. A l'inverse, les études dites de l'alimentation totale tentent de couvrir la quasi-totalité des aliments consommés dans une population d'intérêt. En contrepartie, le nombre de mesures réalisées par produit est bien plus modeste (de 2 à 8 en général). Au vu de nos objectifs, nous préférons ces dernières aux plans de surveillance. Les données de contamination peuvent être combinées aux données de consommation pour ensuite approximer l'exposition à certains composants alimentaires. Malheureusement, les nomenclatures de produits utilisées dans chacune de ces bases peuvent différer substantiellement, ce qui pose un problème supplémentaire au statisticien (voir par exemple le rapport de l'Efsa concernant la mise en commun de données produites par divers pays européens à l'adresse http://www.efsa.europa.eu/en/search/doc/415e.pdf). Comme expliqué ci-après, les bases de données utilisées dans cette thèse ont été construites spécifiquement pour permettre le croisement avec INCA2, contournant par là-même le problème sus-mentionné.

P.2.2.1 Apports nutritionnels : CIQUAL 2008

La version 2008 de la base de données CIQUAL (Centre d'Information sur la QUalité des ALiments) liste les concentrations moyennes en un large panel de nutriments des 1342 produits alimentaires de la nomenclature d'INCA2. Ces teneurs sont déterminées à partir de multiples sources d'information, allant d'analyses chimiques spécifiquement commanditées pour l'étude à des rapports publiés par des instituts de recherche variés. A chaque source est attribué un poids indiquant sa fiabilité, qui est ensuite pris en compte dans un calcul produisant le résultat final. Plus de détails concernant la construction de cette base de données sont disponibles dans le rapport Anses (2008).

Ú De l'invariance nutritionnelle Bien que relativement exhaustive, CIQUAL 2008 ne donne aucune information concernant la variabilité des concentrations en nutriments à l'intérieur d'un même type d'aliment. Cela suggère par exemple que deux oranges cultivées dans des régions différentes, l'une potentiellement plus ensoleillée que l'autre, ont la même teneur en vitamine A. Si cette hypothèse semble déraisonnable de prime abord, les variations de concentration au sein d'une même famille de produits sont souvent si faibles qu'il est d'usage de considérer les teneurs fixes, à condition que la nomenclature utilisée soit assez détaillée (de Boer et al., 2009, p.1433), contrainte qui se trouve être respectée dans le cas de CIQUAL 2008.

P.2.2.2 Contamination des aliments : EAT2

Contrairement aux plans de surveillance, le second opus de l'étude de l'alimentation totale française (EAT2) a été construite dans le but d'estimer l'exposition de la population française à un large ensemble de contaminants alimentaires. Ainsi, plus de 80% des aliments consommés en France y sont pris en compte (Anses, 2011), ce qui en fait une base de donnée de choix au vu de nos objectifs. Cependant, elle possède malgré tout un certain nombre de défauts inévitables.

Ú Mélange des aliments L'une des principales lacunes d'EAT2 provient de la technique utilisée pour réaliser les mesures de concentration. En effet, les analyses chimiques ont été opérées sur des mélanges d'aliments plutôt que sur des produits bruts : différentes espèces d'une même famille d'aliments ont été mixés avant inspection, et les plats cuisinés ont été étudiés dans leur ensemble, sans en séparer au préalable les ingrédients. Afin de contourner les problèmes induits par cette procédure, des tables de recettes ont été construites en parallèle, permettant la décomposition des résultats en une nomenclature plus raffinée (Anses, 2011).

- 1. hypothèse médiane : $t_D = LDD=2$ et $t_Q = LDQ=2$,
- 2. hypothèse basse : $t_D = 0$ et $t_Q = LDD$,
- 3. hypothèse haute : $t_D = LDD$ et $t_Q = LDQ$.

L'hypothèse médiane est typiquement choisie lorsque le taux de mesures censurées dans la base de données utilisée ne dépasse pas les 60%. Nous décidons ici d'adopter ce scénario, les considérations à propos de la censure des données de contamination dépassant le cadre de nos travaux.

P.3 estimation des risques alimentaires : petit aperçu de l'état de l'art

Ces 20 dernières années, pléthore de modèles statistiques a été développé afin d'évaluer de multiples types de risques alimentaires. Construits pour répondre à des questions pratiques liées à la nature des bases de données disponibles telles que la censure mentionnée au paragraphe précédent (Tressou, 2006), ils continuent d'évoluer avec les nouvelles méthodes de collecte des données. Dans le contexte des bases de type 24H et FFQ/ FPQ, auxquelles INCA2 peut être comparée, trois grandes problématiques classiques, décrites dans les paragraphes suivants, sont abordées dans cette thèse.

P.3.1 Exposition extrême de long terme

Afin de compenser la courte durée des enquêtes de type 24H, les modèles à correction d'erreur ont récemment gagné en popularité (Tooze et al., 2006; Dodd et al., 2006; van Klaveren et al., 2012; de Boer et al., 2009; Boon et al., 2011). Ils permettent en particulier de prendre en compte deux sources de variabilité de l'exposition observée, à savoir les fluctuations individuelles autour des prises habituelles et l'hétérogénéité au sein de la population. Toutes deux sont estimées à l'aide de modèles paramétriques et la distribution finale de long terme est approchée par simulations de Monte-Carlo en ignorant la variance intra-individuelle. Les hypothèses sur lesquelles reposent les modèles paramétriques évoqués précédemment impliquent en général que la distribution de l'exposition individuelle est de type log-normale. Cependant, pour de nombreux contaminants et nutriments, la queue de distribution de la loi log-normale est souvent trop fine pour rendre compte convenablement de la probabilité d'occurrence des événements extrêmes; il y a fort à parier que la sur-exposition à de tels composants alimentaires est alors sous-estimée. Lorsque les très fortes expositions à des nutriments et contaminants alimentaires sont au cœur des préoccupations, la théorie des valeurs extrêmes (TVE) a déjà fait ses preuves dans le cadre de l'estimation des risques aigus (Tressou et al., 2004b,a; Bertail et al., 2010; Kennedy et al., 2011; Paulo et al., 2006), nous encourageant ainsi à étendre ces procédures à l'évaluation des risques de long terme.

P.3.2 Exposition simultanée à un ensemble de nutriments et contaminants

Au delà de l'opposition entre risques aigus et risques chroniques, un sujet brûlant d'actualité en évaluation du risque alimentaire est celui de l'exposition simultanée à un cocktail de composants chimiques présents dans la nourriture. En effet, si l'ingestion excessive d'un élément toxique peut avoir des effets dévastateurs sur l'organisme, les connaissances actuelles au sujet des potentiels effets synergétiques, combinés, de plusieurs contaminants sont encore faibles (Carpenter et al., 2002). En France, l'Anses a récemment initié le programme Pericles (Crépet et al., 2013), dans le but d'identifier les cocktails de composants alimentaires effectivement consommés dans la population. Les résultats de ces analyses permettront ensuite aux chimistes et biologistes d'établir un ordre de priorités quant aux recherches à effectuer sur les effets sanitaires d'une telle consommation. Pour le moment, ces travaux sont centrés autour des phénomènes moyens (Béchaux et al., 2013), au détriment des extrêmes. Faisant à nouveau appel à la théorie des valeurs extrêmes, la partie multivariée de cette branche de la statistique devrait permettre l'analyse des fortes expositions simultanées. Lorsque seul un petit nombre de composants est considéré, quelques travaux dans cet esprit

ont d'ores et déjà été réalisés (Paulo et al., 2006). Cependant, en particulier dans le cas de certains types de contaminants comme les pesticides ou les polychlorobiphényles, qui possèdent des centaines de congénères, il serait désirable d'être à même de gérer les grandes dimensions. Malheureusement, d'un point de vue théorique, la TVE ne permet pas encore de traiter des problèmes de dimension plus grande que 5 or 6. Ce problème concret issu de l'analyse du risque alimentaire fait ainsi naître une problématique tout aussi théorique qu'appliquée.

P.3.3 Données de sondage

Dans toutes les méthodes sus-mentionnées, il est de coutume de supposer que les données disponibles sont indépendantes et identiquement distribuées (iid) selon une certaine mesure de probabilité. Or les bases de données de type 24H comme INCA2 sont en général construites à l'aide d'un plan de sondage élaboré, dont l'objectif est de produire un échantillon représentatif de la population d'intérêt. L'hypothèse que les données résultantes sont iid n'est alors pas respectée, et les individus se voient chacun attribué un poids de sondage correspondant à l'inverse de leur probabilité d'être sélectionné dans l'échantillon (Droesbeke et al., 1987; Tillé, 1999). Ignorer cette étape de construction de la base de données et négliger les poids de sondage peut produire des estimateurs biaisés (Bonnery, 2011). Dans le cadre de la présente étude, la probabilité d'occurrence de la très forte exposition à un ou plusieurs composants alimentaires a de fortes chances d'être sur- ou sous-estimée. Bien que la littérature sur les sondages soit déjà très riche (Gourieroux, 1981; Droesbeke et al., 1987; Deville, 1987; Cochran, 1977; Tillé, 2006), ce n'est que très récemment que des résultats fonctionnels, qui permettent par exemple l'estimation de l'intégralité d'une fonction de répartition, ont commencé de se développer (Breslow and Wellner, 2008, 2007; Saegusa and Wellner, 2011). Quant à l'analyse des extrêmes, elle reste à notre connaissance encore inexplorée.

Essayant d'apporter des réponses à ces problématiques, nous proposons dans cette thèse un ensemble de méthodes heuristiques et de résultats théoriques développés spécifiquement pour les bases de données de type 24H. Comme expliqué plus tard en détail, plusieurs de ces procédures mathématiques s'avèrent applicables à bien d'autres domaines que l'analyse du risque alimentaire.

P.4 objectifs et contributions principal es

Au cours du présent travail, nous essayons d'apporter quelques réponses aux problématiques statistiques évoquées à la section précédente. Le mémoire est structuré comme suit : au chapitre 2, nous montrons comment la théorie des valeurs extrêmes peut être utilisée pour la modélisation de l'exposition chronique à un composant alimentaire lorsque la distribution sous-jacente est à queue épaisse. Nous tournant ensuite vers le cadre multivarié, nous nous intéressons à la forte exposition simultanée à plusieurs nutriments et contaminants. Nous proposons ainsi au chapitre 3 une nouvelle méthode, mélangeant algorithmes issus de l'apprentissage statistique et analyse de la mesure spectrale, qui permet d'identifier des groupes de variables dépendantes dans les extrêmes, réduisant par là-même la dimension initiale du problème. Toujours dans une optique multivariée, nous quittons au chapitre 4 les extrêmes en faveur des phénomènes moyens. Nous y présentons une approche en terme d'ensembles de volume minimum pour l'estimation non-paramétrique de la distribution de l'exposition à un cocktail d'éléments chimiques présents dans la nourriture. Ces résultats sont ensuite adaptés à l'identification de paniers de consommation réalisant un compromis entre risque toxicologique et bénéfice nutritionnel. Enfin, nous nous attelons dans les deux derniers chapitres au traitement des données issues d'un plan de sondage. Après avoir étendu quelques résultats fonctionnels usuels en analyse des processus empiriques au cadre des plans de sondages de type Poisson dans le chapitre 5, nous introduisons au chapitre 6 un nouvel estimateur de l'indice de valeurs extrêmes adapté aux échantillons issus de tels plans. Bien qu'encore trop restrictifs pour être directement appliqués à l'analyse du risque alimentaire, ces résultats préliminaires constituent un premier pas dans la direction de futurs développements qui, nous l'espérons, permettront bientôt la construction de modèles plus généraux, adaptés aux bases de données comme INCA2.

P.4.1 Estimation de l'exposition extrême de long terme à un unique élément chimique

Nous commençons dans ce chapitre par une réflexion sur les modèles à adopter lors du calcul de l'exposition chronique à un unique nutriment ou contaminant alimentaire. Portant un intérêt tout particulier aux phénomènes extrêmes, dans la lignée des travaux de Tressou et al. (2004a), nous remettons au goût du jour une méthode non-paramétrique généralement ignorée en faveur des modèles à correction d'erreur brièvement décrits à la section précédente. Lorsque les données sont issues d'une distribution à queue épaisse, nous montrons qu'elle comporte plusieurs avantages. La première étape de cette technique consiste à moyenner les données temporelles pour les ramener à l'échelle de la journée. Nous évitons ainsi la modélisation statistique des variances intra- et inter-individuelles qui affectent la distribution de l'exposition de long terme, encouragée par Tooze et al. (2010); van Klaveren et al. (2012); de Boer et al. (2009). Les modèles classiques de la théorie des valeurs extrêmes pour l'étude des lois à queue épaisse sont ensuite appliqués pour estimer quantiles extrêmes et faibles probabilités de dépasser un très haut seuil de recommandation. Si l'exposition est en effet à queue plus lourde qu'une loi log-normale, nous montrons qu'en procédant de la sorte, à l'inverse des modèles traditionnels, nous évitons de sous-estimer ces quantités. Comme en pratique le statisticien doit choisir le modèle le plus adapté aux données, nous recommandons l'utilisation préliminaire de quelques tests statistiques connus détectant le cas échéant la présence d'une queue lourde. L'intégralité de notre méthodologie est enfin appliquée à des données réelles pour en montrer les avantages et les inconvénients.

P.4.2 Evaluation de la sur-exposition simultanée à un cocktail de composants alimentaires

Nous étendons ensuite notre approche au cadre multivarié et considérons l'analyse de la forte exposition chronique à plusieurs nutriments et contaminants. Dans l'esprit du programme Pericles (Crépet et al., 2013), nous développons une nouvelle méthode permettant d'identifier des cocktails de composants chimiques consommés simultanément en très grandes quantités dans une population d'intérêt. Nous nous inspirons pour cela à la fois d'algorithmes d'apprentissage statistique et de concepts de la théorie des valeurs extrêmes multivariée. Nous étudions en particulier la mesure spectrale, objet mathématique qui caractérise la dépendance extrême. Elle peut être définie sur l'orthant positif de la sphère unité, i.e. le simplexe, lui-même décomposable en faces ouvertes (sommets et arêtes), chacune indiquant un groupe de variables dépendantes dans les extrêmes. L'objectif est alors d'identifier les faces du simplexe sur lesquelles la mesure spectrale est définie. Pour cela, nous commençons par mettre en œuvre le récent algorithme de Jung et al. (2012), appelé Principal Nested Spheres, qui réduit la dimension des données en les projetant successivement sur des sphères de dimension de plus en plus faible. A la manière d'une analyse en composantes principales, cette première étape permet de faciliter la suite de l'analyse en réduisant le bruit des observations. Ensuite, nous exhibons un modèle de mélange de la mesure spectrale sur chacune des faces du simplexe et définissons une variable latente caractérisant les composantes du mélange. Cette dernière est estimée de manière non-paramétrique en utilisant d'abord un algorithme de classification appelé spherical k-means (Dhillon et al., 2002), puis en utilisant un critère heuristique construit pour déterminer les faces auxquelles les classes obtenues font référence. Ce critère permet de même de choisir le nombre d'observations extrêmes qui peuvent être considérées comme représentatives de la queue de la distribution multivariée. Nous justifions notre approche à l'aide d'une étude par simulations, dont les résultats sont fort encourageants, puis l'appliquons enfin à nos bases de données réelles, INCA2, CIQUAL 2008 et EAT2. Les groupes résultants de contaminants et nutriments supposés consommés simultanément et en de très grandes quantités dans la population française font parfaitement sens et sont cohérents avec une analyse paire par paire, ce qui vient renforcer notre confiance en notre approche.

P.4.3 Ensembles de volume minimum et approche risque bénéfice

Allant plus loin encore dans l'utilisation de l'apprentissage statistique pour l'analyse du risque alimentaire, nous nous concentrons dans ce chapitre sur les phénomènes moyens, non plus extrêmes. Nous y introduisons pour la première fois la variabilité de la contamination des aliments, les teneurs ayant jusque lors été considérées fixes au sein d'une même famille de produits. Dans la lignée des travaux de Bertail and Tressou (2006), nous proposons une extension de l'estimation d'ensembles de volume minimum de Scott and Nowak (2006) au cas où le volume est inconnu et approché à l'aide d'une U-statistique. Les résultats théoriques sont ensuite appliqués à la construction non-paramétrique des ensembles de niveau de la distribution multivariée de l'exposition à un ensemble de composants alimentaires. Dans un second temps, nous montrons comment cette procédure peut être généralisée à l'identification de paniers de consommation qui réalisent un compromis entre risque toxicologique et bénéfice nutritionnel. Encore en cours de programmation, nous espérons pouvoir appliquer cette dernière technique aux bases de données INCA2, CIQUAL 2008 et EAT2 afin de définir des recommandations simples dans le même esprit que "mangez cinq fruits et légumes par jour".

P.4.4 Processus empiriques dans le cadre des sondages

Nous nous attelons dans les deux derniers chapitres de cette thèse au traitement des données issues d'un plan de sondage, comme le sont souvent les données de consommation. Il est d'un intérêt majeur pour les instituts de santé publique comme l'Anses d'être à même d'estimer convenablement la distribution de l'exposition de long terme à un ou plusieurs éléments chimiques. Or, omettre le plan de sondage dans le processus d'estimation peut induire un biais non-négligeable (Bonnery, 2011). Tentant de contribuer à l'élaboration d'une théorie générale sur les sondages garantissant la normalité asymptotique d'une large classe d'estimateurs, comprenant notamment des estimateurs fonctionnels de la fonction de répartition de l'exposition, nous commençons par étudier les plans de type Poisson. Des extensions à des plans plus complexes comme celui utilisé pour la formation d'INCA2 seront envisagées dans un futur proche. Notre approche est directement inspirée des travaux séminaux de Hàjek (1964) et Berger (1998). Nous définissons en premier lieu le processus empirique de type Horvitz-Thompson, où les observations sont pondérées par l'inverse de leur probabilité d'inclusion, dans le cadre du plan de Poisson. Sous un ensemble d'hypothèses classiques portant à la fois sur les probabilités d'inclusion et le modèle de surpopulation considéré, que nous espérons relâcher dans des travaux à venir, nous en montrons la convergence vers un processus Gaussien dont nous exhibons la covariance. Dans un second temps, nous généralisons ces résultats aux

processus empiriques à la Horvitz-Thompson impliquant des échantillons sélectionnées selon un plan de sondage à forte entropie comme le plan réjectif. Pour ce faire, nous utilisons des résultats connus exhibant la proximité entre les plans de sondage à forte entropie et le plan simple de Poisson. Les théorèmes fonctionnels de la limite centrale ainsi obtenus sont ensuite utilisés pour montrer la normalité asymptotique d'estimateurs pouvant s'écrire comme certaines fonctionnelles (Hadamard ou Fréchet différentiables) du processus empirique. Enfin, nous étudions le cas spécifique de l'estimateur Horvitz-Thompson de la fonction de répartition et, illustrations à l'appui, montrons comment utiliser nos résultats pour en construire des bandes de confiance uniformes. Au delà de l'analyse du risque alimentaire, ces résultats semblent tout à fait appropriés à la gestion de bases de données de taille gigantesque (les fameuses "big data"), dont la taille augmente sans cesse, comme les données financières, et ne peuvent être exploitées sur un unique ordinateur. Dans un tel contexte, l'échantillonnage semble une solution toute naturelle aux problèmes de mémoire informatique. Les plans de sondage que nous avons étudiés se révèlent alors d'un intérêt tout particulier, permettant d'obtenir simplement des estimateurs non biaisés et d'une efficacité optimale (il suffit pour cela de calibrer correctement les probabilités d'inclusion) sur des échantillons de taille raisonnable.

P.4.5 Estimation de l'indice de valeurs extrêmes à partir de données de sondage

Nous concluons cette thèse en revenant au problème initial de l'estimation des phénomènes extrêmes, prenant cette fois-ci en compte le plan de sondage selon lequel les données ont été collectées. A notre connaissance, la théorie des valeurs extrêmes n'a pour le moment pas été étendue au cadre des sondages. Nous commençons donc modestement par l'adaptation dans ce contexte de l'un des estimateurs les plus classiques de la TVE, à savoir l'estimateur de Hill de l'indice de valeurs extrêmes (Hill, 1975). Dans le même esprit que les travaux du chapitre précédent, nous en construisons une version Horvitz-Thompson dont nous montrons la consistance et la normalité asymptotique pour les plans de type Poisson. La vitesse de convergence de ce nouvel estimateur s'avère être la même que si la population entière avait été accessible, et la variance n'est dégradée que d'un paramètre multiplicatif dépendant du choix des probabilités d'inclusion, qu'il est théoriquement possible de minimiser. A l'aide d'expériences numériques, nous montrons que les hypothèses restrictives que nous exigeons pour établir ces résultats mathématiques sont loin d'être nécessaires et pourraient être relâchées dans de futurs travaux. Rappelant enfin les problématiques inhérentes aux "big data", nous encourageons l'extension de ces travaux préliminaires à d'autres familles d'estimateurs issus de la TVE.

THESIS

MULTIVA RIATE STATISTICS FOR DIETARY RISK ANALYSIS
1

INTRODUCTION: OF DIETARY HABITS AND EXPOSURE

Dietary risk analysis is a generic public health term that embraces as many scientific problems as there are ways for consumers to get sick by eating (or not eating). For instance, food industries are concerned with the development of bacteria during the consecutive fabrication, transportation, distribution and stocking processes that precede consumption (Rigaux et al., 2012). On another level, agribusinesses may want to evaluate the impact of economically profitable mass production methods (e.g. genetically modified organisms, pesticides) on the human organism (d. http://www.efsa.europa.eu/en/topics/topic/qmo.htm on GMO and http: //www.efsa.europa.eu/en/topics/topic/pesticides.htm on pesticides). Cooking utensils can also represent a source of danger, depending on their chemical composition and the way they are used, as was recalled by the recent polemic about the presence of bisphenol A in baby bottles (see the ensuing Anses avis, Anses, 2013). Other health issues may be caused by behavioral phenomena such as eating disorders (e.g. people suffering from severe anorexia are likely to develop diseases due to nutritional deficiency), or by biological determinants (e.g. diabetes). As suggested by this non-exhaustive list of examples, depending on the examined question, specific knowledge in chemistry, medicine, sociology, economy, biology or even psychology may be required, not to mention probabilistic modeling and computer data management, which makes dietary risk analysis a multidisciplinary field par excellence. We invite the interested reader to consult Feinberg et al. (2006) and the references therein for a comprehensive introduction to the ins and outs of dietary risk analysis from an interdisciplinary point of view.

The present work is dedicated to a very specific type of dietary risk: independently from any production, stocking, cooking, or biological predisposition phenomena, we are concerned with the very high (or low) long-term exposure of the French population as a whole to some food components. After having introduced some basic notions about dietary risk analysis in Section 1.1, described the available data in Section 1.2 and quickly reviewed the literature on probabilistic modeling of dietary risks in Section 1.3, we thoroughly present the various problems tackled in this thesis and succinctly report its main scientific contributions in Section 1.4.

1.1 getting a taste of dietary risk analysis

Once a specific type of dietary risk has been brought into focus, its analysis involves three distinct stages, usually referred to as risk evaluation, risk management and risk communication. In the next paragraphs, we give a concise overview of the various issues addressed in these successive and complementary steps. To help better understand the challenges at stake, we illustrate each introduced methodological concept with the specific example of salty diets, the recurrent abuse of which can favor cardiovascular issues (see for instance http://www.anses.fr/en/content/salt and the references therein).

1.1.1 Risk evaluation

Risk evaluation (or assessment) can be schematically decomposed into four substages, which consist in successively identifying and characterizing both potential dangers and their probability of occurrence. In the case of salty diets, the danger would be to ingest too much sodium (Na in the periodic table) on a daily or weekly basis, since in the long-run it would be likely to damage the cardiovascular system. These noxious effects were spotted by means of chemical and biological studies about the assimilation, action and elimination of Na in the human organism. This process is called hazard identification (Barlow et al., 2002). It is followed by further research designed to understand at what point the component becomes noxious, i.e. characterize the danger. In general, danger characterization involves testing dose-response effects in vitro or in vivo on animals before extending the subsequent results to humans by means of dedicated models (Dybing et al., 2002). Once the contamination process is fully comprehended, it becomes possible to assess its probability of occurrence in a given population. This requires first a thorough description of the distributions of consumption and exposure in the population (Kroes et al., 2002), then a comparison with some maximal intake limit determined by experts from the aforementioned dose-response trials (Renwick et al., 2003), possibly enhanced by mathematical designs (Edler et al., 2002). Depending especially on the data at hand, statistical models may include a temporal dimension (Bertail et al., 2010, 2008; Allais and Tressou, 2009), take into account individual biological or behavioral characteristics, or try to make up for the limited amount of available information (Dodd et al., 2006). Considerations on this matter are discussed in detail in Section 1.3 and Chapter 2.

1.1.2 Risk management

Once dietary risks have been evaluated, it is necessary to identify their determinants, evaluate their relative impact and define the best strategy to rapidly and efficiently reduce risks. This is the exact purpose of risk management. For instance, over-exposure to sodium can sometimes be traced back to the repeated consumption of some salt-saturated products such as precooked dishes or biscuits. A simple line of action there would be to simultaneously warn consumers of the potential dangers and impose food companies to limit the amount of added salt in their preparations. In more extreme cases, after a thorough evaluation of the economic impact of such measures, some foodstuffs may even be taken off the market.

1.1.3 Risk communication

The last stage of dietary risk analysis is called risk communication. It can take place at any moment of the analysis and be intended for scientists and risk managers involved in the process as well as food industries or consumers. As such, it is not restricted to public health campaigns such as "avoid eating products that contain too much salt" in our example. It also encompasses scientific reports about the methods considered appropriate to detect and quantify risks, publications of case study results, international discussions about food policies, etc. In the sequel, we are mainly interested in producing results that would help orient further research and, when possible, set general, easily understandable and applicable dietary guidelines to the population.

With this objective in mind, most of our work falls into the evaluation phase of dietary risk analysis: in a few words, we design statistical methods to assess some characteristics (e.g. quantiles) of the very high exposure to some food chemicals over a long period of time. Except in Chapter 4 where we elaborate a statistical methodology that ascertains balanced food baskets with regard to toxicological risk and nutritional benefit, risk management is not the main concern here. The components of interest are nutrients and contaminants to which the chronic over- or under-exposure has known detrimental sanitary effects. In particular, we disregard acute risks such as those involving bacterial contamination of the food, which can impact the organism in only a few days or even a few hours. Statistical analysis is based in this context on some standard types of databases, which are described in detail in the next subsection.

1.2 collecting information: multiple sources, multiple difficulties

To statistically analyze the chronic dietary exposure of a given population to a collection of food components, we need to observe such types of exposure on a large sample of individuals over a long period of time. Unfortunately, the amounts of nutrients and contaminants ingested during a meal are not directly measurable and can only be assessed by means of various methods. This can be achieved for instance by detecting and quantifying some specific biomarkers, the presence of which is ultimately linked to the level of exposure (refer for instance to the Anses website http:// www.anses.fr/sites/default/files/documents/RSC1205-DossierParticipants.pdf giving the list of presentations in a dedicated workshop held in May 2012). From a practical point of view, this necessitates collecting then chemically analyzing samples of body fluids, hair or skin of a relatively large array of people. At a national level, the cost of such procedures can constrain the sample size at the expense of statistical efficiency. Moreover, many chemicals can be assimilated via other elements than food (eg. air) and biomarkers cannot distinguish between the various sources of exposure (Sirot et al., 2009). So as to study the sole impact of nourishment on individual exposure, one may use in an alternative manner a combination of two types of databases, one listing the dietary intakes of a sample of consumers and the other the levels of components in a fine nomenclature of products. The next paragraphs provide an overview of the main characteristics of these so-called consumption and composition data.

1.2.1 Consumption data

There are many types of data that can be used to assess chronic exposure to nutrients and contaminants in France, ranging from panel cohorts on alimentary expenditures of households (Secodip, Nichèle et al., 2008) to Internet repositories filled up on a voluntary basis (Nutrinet, Hercberg et al., 2010) and dedicated nationwide institutional surveys (INCA2, Afssa, 2009). Focus is here on two major categories, namely 24-hour recalls (abbreviated 24H) and food frequency or propensity questionnaires (FFQ/ FPQ). The first category corresponds to large surveys (possibly national) where a sample of selected individuals report in detail their dietary habits during 2 to 7 days, sometimes consecutive but most of the time randomly picked within the year. Food intakes are given exactly or relative to some pictures of more or less filled plates provided by the pollster. Obviously, the short duration of these surveys can impede the estimation of long-term dietary habits (Counil et al., 2006). To help deal with this issue, food frequency (or propensity) questionnaires are proposed to the consumers, in which they declare what type of food they most commonly eat and which products they avoid or never consume. We refer the interested reader to van Klaveren et al. (2012); EFSA (2006); Dodd et al. (2006) for an account of the assets and liabilities of these complementary approaches. In this thesis, we used the INCA2 database described herein-after, which can be assimilated to a mixture of 24H and FFQ.

1.2.1.1 Dietary habits in the French population : INCA2

INCA2 (second opus of the national survey on individual dietary habits in France) is a nationwide survey conducted by Anses (French agency for food, environmental and occupational health safety) in collaboration with Insee (French national institute of statistics and economic studies) between December 2005 and April 2007. It collects information about the dietary habits of 2624 adults and 1455 children taken at random in the French population.

Ú Survey scheme Individuals in INCA2 were selected according to a complex multistage survey scheme that was designed to produce a representative sample relative to geographical, sociological and economic criteria. Post-calibration methods were applied in a second phase with respect to the same set of auxiliary variables to provide corrected survey weights relative to both non-response and sampling fluctuation. We refer to Afssa (2009, Chapter 2 and Appendix 2) for a thorough description of the survey plan.

Ú Dietary information All participants reported both the nature and the amount of food that they ate during 7 consecutive days. These products were classified according to an exhaustive nomenclature of 1342 foods grouped into 123 sub-classes and 45 wider categories (Afssa, 2009, Section 2.2.6.3 and Appendix 1). To help them assess the quantities they ate, inquired people were given a notebook displaying pictures of progressively filled plates and glasses. It served as a referential in case exact measurement was not possible. A precise description of each meal was also provided, indicating for instance when, where and with whom they occurred. In the final database, they were classified into six types of meals, namely breakfast, morning snack, lunch, afternoon snack, dinner and evening snack. So that the seasonal variation of dietary habits may be controlled, individuals were contacted randomly at different periods of the year during 3 distinct phases of data collection. A specific line of questions was also established to assess food supplement consumption in the French population (Afssa, 2009, Section 2.2.6.4 and Chapter 8).

Ú Episodic or non-consumers When interested in long-term dietary habits and food chemical exposure, 7 days of observation appear to be an arbitrarily short period of time. In particular, in only one week, individuals cannot cover their entire food repertoire and many items in the reference nomenclature are not consumed. Thus, it be-

comes difficult to discriminate between real non-consumers and episodic consumers. So as to make up for this drawback, individuals in INCA2 were asked to depict their dietary habits, and clearly declare to which category they belonged. These additional queries corresponds to the so-called food frequency/ propensity questionnaires previously mentioned, which, when combined with 24-hour recall data, facilitate the statistical modeling of long-term food or nutritional intakes (Dodd et al., 2006 and van Klaveren et al., 2012, Sections 3.6 and 3.8).

Ú Incomplete information Regarding episodic versus non-consumers, collecting data during only 7 days is clearly a liability, and a wider period of observation might be preferred. Conversely, when concerned with the quality of answers, an entire week can be considered too long to provide reliable information. Indeed, as mentioned in Afssa (2009, Section 2.4.2), some participants did not report their consumption every day and other under-estimated their intakes. In our calculations, we decided to ignore these faults and simply scaled the observed dietary habits to the entire week.

Ú Time dependence The only difference between INCA2 and classical 24-hour recall databases is that the 7 days of observation were consecutive. Consequently, although an entire week is clearly better than the usual 2 days for long-term analysis, it makes temporal dependence harder to ignore.

¹ Preliminary processing of INCA2 The present work is dedicated to the global analysis of the long-term dietary habits of French adults. Because of their very specific nutritional needs and the temporary character of their condition, we decided not to take into account pregnant or lactating women. Other individuals were excluded, namely those for whom important variables were missing (eg. body weight or consumed food amount), thereby restricting the initial sample to 2488 units. Obviously, we could have tried to apply standard techniques to infer on missing values, but this went beyond the scope of our work. In particular, such methods are usually designed to avoid degrading the estimation of average phenomena. Since we are more interested in extreme (maximum or minimum) events and incomplete data only concerned a very small proportion of the sample, we chose to ignore them instead. In the final considered database, 91 individuals (3.7%) ceased filling the questionnaire after only a few days.

1.2.2 Composition data

Composition databases usually collect an array of chemical measurements realized on more or less refined groups of foodstuffs. They can originate from surveys of very different natures. For instance, surveillance plans are designed to punctually check the safety of some suspect products. Hence, by nature, they do not cover the entire food repertoire, but are particularly thorough in the sense that the collected samples are of important size. In addition, since the inspected foodstuffs are suspected to be abnormally contaminated, they cannot be used to estimate the toxicological exposure of a large population without introducing a non-negligible bias. Alternative types of data are the total dietary surveys. Based on samples of food of very small sizes (2 to 8 items in general), they however encompass a large array of products. As such, they are more suitable for our purpose. The main difficulty when crossing them with consumption data is that the corresponding nomenclatures may substantially differ, thereby necessitating an additional step, potentially difficult, of association of the various sources. This issue is particularly pregnant with meaning for international institutes such as Efsa, which have to combine data from various countries (see the recent report on the matter at http://www.efsa.europa.eu/en/search/doc/415e.pdf). Here, we work with databases that were constructed specially to fit INCA2. They are briefly introduced below.

1.2.2.1 Nutrient supply : CIQUAL 2008

The 2008 version of the CIQUAL database (Centre d'Information sur la QUalité des ALiments in French) lists the average concentrations with regards to a large set of nutrients for each of the 1342 dietary products of the INCA2 nomenclature. These levels were determined using multiple sources of information, ranging from specifically ordered laboratory analyses to the published reports of various research institutes. Each source was attributed a weight representing its reliability before a final synthesizing calculus. We refer to Anses (2008) for more information concerning the construction of this database.

Ú About nutritional invariability Although already quite exhaustive, CIQUAL 2008 does not provide any information about the variability of concentration within a type of food. This suggests for instance that two oranges grown in different regions, one possibly sunnier than the other, contain the same level of vitamin A. Even if this assumption seems questionable at first sight, nutrient contents usually have such a small variance that it is customary to consider them fixed, provided the associated food nomenclature is detailed enough (de Boer et al., 2009, p.1433). The list of products considered in CIQUAL 2008 was designed to respect this constraint.

1.2.2.2 Food contamination : TDS2

Contrary to surveillance plans, the second opus of the French Total Dietary Survey (TDS2) is designed to assess the global exposure of the French population to a whole set of food contaminants and thereto covers more than 80% of the current food repertoire in France (Anses, 2011). Though particularly adapted to our needs, it still possesses some unavoidable drawbacks.

Ú Pooling foodstuffs One of the major limitations of TDS2 is that the analyzed products are not raw but pooled, i.e. different species of a same food item were mixed together before chemical inspection. Dishes were also treated as such, meaning that instead of separating the elements of a recipe they were considered as a whole. To overcome this issue, recipe tables were constructed in parallel, to enable the decomposition of the results of INCA2 into a more refined nomenclature (Anses, 2011).

- 1. median hypothesis: $t_D = LOD=2$ and $t_Q = LOQ=2$,
- 2. lower hypothesis: $t_D = 0$ and $t_Q = LOD$,
- 3. upper hypothesis: t_D = LOD and t_Q = LOQ.

The median hypothesis is typically chosen when the censored measures do not represent more than 60% of the entire database. In this thesis, we chose to adopt this specific approach, considerations about censorship going beyond the scope of our work.

1.3 statistical assessment of dietary risks: a review of the literature

In the last 20 years, a plethora of dedicated statistical methods have been developed to assess dietary risks of various natures. Designed to answer practical issues linked to the type of available data such as the censorship in the measurement of levels of contents (Tressou, 2006), they keep on evolving with the new methods of data collection. In the context of 24-hour recalls and food frequency questionnaires, to which INCA2 is very similar, we can distinguish three major issues, detailed in the next paragraphs, that are to be tackled in the present thesis.

1.3.1 Extremelong-term exposure

To cope with the limited duration of 24-hour recalls, statistical models of the measurement error type have recently gained popularity (Tooze et al., 2006; Dodd et al., 2006; van Klaveren et al., 2012; de Boer et al., 2009; Boon et al., 2011). They account for the presence of two sources of variance in the observed exposure, namely individual fluctuations around usual habits and populational heterogeneity. Both are estimated by means of parametric modeling and the final long-term distribution is approached with Monte-Carlo simulations that disregard the intra-individual variations. The aforementioned parametric assumptions usually implicitly require that the usual intakes have log-normal distribution. However, for many nutrients and contaminants, the tails of such probability laws are sometimes too thin to accurately account for the occurrence of extreme events; under-estimation is thus a non-negligible risk. When interested in the very high exposure to nutrients and contaminants, methods issued from extreme value theory (EVT) have already proven useful in acute risk estimation (Tressou et al., 2004b,a; Bertail et al., 2010; Kennedy et al., 2011; Paulo et al., 2006), suggesting that an extension to the long-term setting would be worth considering.

1.3.2 Simultaneous exposure to multiple food components

Beyond the issue of chronic versus acute risks, a particularly hot topic in dietary risk assessment is the analysis of simultaneous exposure to multiple chemicals. Indeed, if the excessive ingestion of toxicants can have a detrimental impact on health, knowledge about synergistic, combined effects is still poor (Carpenter et al., 2002). In France, the Anses institute has most recently launched the Pericles program (Crépet et al., 2013) in order to identify cocktails of components that are observed to be consumed in the population. The ensuing results should provide guidelines to chemists and biologists who will then be in charge of inspecting the corresponding sanitary effects on the human organism. For now, most of this work is focused on average phenomena (Béchaux et al., 2013), at the detriment of extremes. Going back to extreme value theory, using the multivariate branch of this field should provide ways of dealing with high joint exposure. When only a few dimensions are involved, attempts in that direction have already been published (Paulo et al., 2006). However, especially for some types of contaminants such as pesticides or Polychlorinated biphenyls, which possess hundreds of congeners, managing high dimensions would be desirable. Unfortunately, from a theoretical point of view, EVT does not handle dimensions higher than 5 or 6 yet, making this branch of dietary risk analysis both a practical and a methodological challenge.

1.3.3 Survey data

In all the aforementioned methods, it is usually assumed that the data at hand is independent, identically distributed (iid) according to some probability measure. However, 24-hour recalls and databases like INCA2 are typically constructed from some elaborate survey scheme, designed to produce samples that are representative of the population of interest. There, the iid hypothesis is no longer satisfied, and individuals are attributed survey weights based on their probability of being included in the sample (Droesbeke et al., 1987; Tillé, 1999). Thus ignoring the underlying design is known to produce biased estimates (Bonnery, 2011), in our case the probability of occurrence of very high exposure is likely to be either over- or under-estimated. Though the literature on survey sampling is quite rich (Gourieroux, 1981; Droesbeke et al., 1987; Toleville, 1987; Cochran, 1977; Tillé, 2006), functional results that would enable the estimation of the entire distribution function are just starting to flourish (Breslow and Wellner, 2008, 2007; Saegusa and Wellner, 2011) and extreme analysis is, to our knowledge, still unexplored.

In an attempt to answer these issues, we propose in the present work a collection of heuristic methods and theoretical results that are specially built for 24-hour recall types of databases. As shall be seen later on, some of these mathematical findings have possible applications that go beyond dietary risk analysis.

1.4 objectives and main contributions

In this thesis, we tried to bring some answers to the statistical problems mentioned in the previous section. It is structured as follows: in Chapter 2, we tackle the issue of heavy-tailed long-term exposure to one specific chemical and show how extreme value theory can be of help. Moving then to the multivariate setting, the simultaneous exposure to multiple food components is considered. A new method mixing machine learning algorithms and spectral measure estimation to reduce the dimension in the analysis of multivariate extreme values is introduced in Chapter 3. In Chapter 4, while still remaining in a multidimensional optic, we change focus and go back to average phenomena. There, statistical learning methods borrowed from the minimum volume set literature are adapted to the non-parametric estimation of the distribution of exposure to a collection of nutrients and contaminants. A natural extension to the construction of dietary habits that realize a compromise between toxicological risk and nutritional benefit is also proposed. Finally, Chapter 5 and Chapter 6 are dedicated to the treatment of survey data. After extending classical functional results on empirical processes to the analysis of observations issued from a Poisson-like sampling scheme in the former, we introduce in the latter a novel estimator of the extreme value index for survey data. Although they are not general

enough to enable direct application to dietary risk analysis, these preliminary results constitute the basis of future developments that, we hope, will soon lead to more comprehensive models that can manage databases like INCA2.

1.4.1 Extreme chronic exposure to one chemical

Before even considering multivariate types of exposure to food chemicals, we start by discussing the calculation of the usual intakes of a unique nutrient or contaminant. Particularly interested in extreme phenomena, following in the footsteps of Tressou et al. (2004a), we bring back into fashion a non-parametric method that is usually disregarded in favor of mixed-effects models. It simply relies on the preliminary averaging of temporal observations on a daily scale, thereby avoiding the statistical modeling of the between and within variances that play a role in the distribution of the long-term exposure advocated by Tooze et al. (2010); van Klaveren et al. (2012); de Boer et al. (2009). When dealing with heavy-tailed distributions, thus proceeding is shown to avoid the under-estimation of tail characteristics that is bound to occur with the classical log-normal parametric modeling of usual intakes. Supporting our arguments with a real-data analysis, we propose a systematic procedure that consists in testing first the presence of a fat tail before choosing a specific statistical procedure. For this purpose, techniques directly borrowed from the extreme value theory literature are depicted and the suitability of this field for dietary risk analysis is brought into focus.

1.4.2 Simultaneous over-exposure to many food chemicals

Extending next our approach to the multivariate level, we consider the analysis of extreme types of exposure to many nutrients and contaminants. In the same spirit as the Pericles program (Crépet et al., 2013), we propose a new method that identifies cocktails of chemicals to which individuals are simultaneously highly exposed. Inspired by techniques of both the statistical learning and the extreme value analysis fields, it consists in assessing in a non-parametric manner the elements of the support of the spectral measure, a mathematical object that characterizes extreme dependencies. It is defined on the positive orthant of the unit sphere, which is quite naturally decomposable into open faces (edges and vertices) that happen to point out the variables that are linked together. In order to detect the faces on which the spectral measure is positive (thus the groups of variables that exhibit extreme dependence), we use a novel algorithm called Principal Nested Spheres (Jung et al., 2012) that achieves a sort of Principal Components Analysis on the unit sphere. It reduces the dimension of the data by systematically projecting the cloud of points on sub-spheres of lower dimension, thereby facilitating further analyzes. Tackling this

issue from a latent variable point of view justified by a mixture model of the spectral measure, the projected data is then clustered into groups that supposedly represent the different faces forming its support. Their identification is finally handled with a new heuristic that seems to perform well on simulations. When applied to the multivariate exposure obtained with the INCA2, CIQUAL 2008 and TDS2 databases, it produces comprehensible outcomes that support this promising approach.

1.4.3 A minimum volume set approach to dietary risk-benefit analysis

Exploiting further the assets of learning procedures for dietary risk analysis, we focus this time on average phenomena. In this chapter, we introduce for the first time the variability in the contamination process of the food. Following in the footsteps of Bertail and Tressou (2006), we extend the minimum volume set approach of Scott and Nowak (2006) to the case where the volume is unknown and estimated by a U-statistic. This enables to construct in a non-parametric manner the level sets of the multivariate distribution of types of exposure to multiple chemicals. We then demonstrate how this procedure can be modified to recover in the consumption space the dietary habits that balance toxicological risk and nutritional benefit. In the near future, these theoretical results are destined to be applied to the INCA2, CIQUAL 2008 and TDS2 databases in order to provide general dietary guidelines in the spirit of "eat five fruits and vegetables a day".

1.4.4 Empirical processes in survey sampling

Our next challenge concerns the nature of the consumption data, which often results from a complex survey scheme. Since the distribution of the long-term exposure is of particular interest for public health institutes like Anses, providing functional results about such weighted data would help avoid the bias induced when ignoring the survey scheme. Though unable to develop results for the specific design employed in INCA2 yet, we make our contribution to the elaboration of a more comprehensive theory by studying the asymptotic properties of empirical processes in the context of Poisson-like survey plans. Our approach is directly inspired by the seminal papers of Hajek (1964) and Berger (1998) and exploits the proximity of large entropy designs to the simple Poisson scheme. Under some assumptions on the inclusion probabilities and a superpopulation framework, which are bound to be relaxed in the near future, we establish a functional central-limit theorem and show its implications for the asymptotic analysis of a large array of estimators. With illustrations based on simulations we present in particular how it can be applied to the construction of uniform confidence bands of the distribution function. Beyond dietary risk analysis, such results appear to be of particular interest for the management of huge databases, the sizes of which increase in permanence like financial data, and therefore cannot be fully accessed. With complete control over the sampling procedure, the Poisson and rejective plans are revealed as especially convenient for the unbiased statistical analysis of such databases.

1.4.5 Tail index estimation based on survey data

Going back to extreme phenomenons, we propose to adapt the widely celebrated Hill estimator of the extreme value index (Hill, 1975) to the survey sampling framework. Hoping again to extend this results to the complex survey scheme of INCA2, we start by establishing its consistency and asymptotic normality for sampling plans of the Poisson type. The rate of convergence of this novel estimator is found to be the same as if the entire population was available, and the variance is only depreciated by a multiplicative term that depends on the way the inclusion probabilities were chosen. With numerical experiments we show that the restrictive hypotheses that were required in our theorems could actually be relaxed. Recalling the issue of big data, we finally encourage the extension of those preliminary results to many other branches of the extreme value field.

2

EXTREME CHRONIC EXPOSURE TO ONE CHEMICAL

2.1 introduction

When combined with food composition databases, dietary surveys prove useful to assess statistical distributions within a given population of intakes of nutrients as well as different types of chemical substances such as environmental contaminants, food additives or pesticide residues present in the food. Typically, they are built out of short-term follow-ups of representative sub-populations, commonly called 24h-recalls, chosen according to some appropriate survey scheme, and therefore do not enable direct estimation of chronic (long-term) risks. Indeed, participants in cross-sectional representative dietary surveys are usually not solicited more than one week, and the current trend is even to shorten the survey duration to only two or three days (see for instance the reports of the EFCOVAL project on http://www.efcoval.eu/ and Crispim et al., 2011; De Boer et al., 2011). As a result, chronic risk evaluation may not be achieved from such data without any statistical modeling.

In answer to these practical limitations, statistical models have been developed in the last 20 years to estimate usual long-term intakes from short-term measurements (van Klaveren et al., 2012; Dodd et al., 2006). Originally constructed to assess the unbiased prevalence in some population of interest of inadequate and insufficient intakes of nutrients, they are also used now to estimate high percentiles of intake in a food safety perspective. In nutrition, there is for instance a need to verify that upper levels of intakes for vitamins and minerals are not exceeded, considering food fortification or supplements intake. Most of these models use an analysis of variance to separate between- and within-individuals variabilities of usual intakes. An assumption of normality of intakes after a Box-Cox transformation is done (Tooze et al., 2010). The present study tries to propose an adaptation of these methods to better consider high nutrient intakes, which are often not normally distributed, in order not to underestimate the risk of exceeding upper levels. This method could also be applied to chronic risk assessment to chemicals in food (Boon et al., 2011). This work corresponds to a paper currently being written in collaboration with JL. Volatier (Anses, France).

The chapter is structured as follows. We start off in Section 2.2 by setting notations and presenting the statistical background on which classical models are based. After

reviewing in further detail two popular methods and underlining their limits relative to the estimation of high quantiles and rare events, we thoroughly depict our alternative procedure in Section 2.3.1. As means of illustration, a case-study is conducted in Section 2.4. The different models previously introduced are subsequently compared on real data, which finally leads to the discussion of their respective assets and liabilities in Section 2.5. Supplementary details regarding tail estimation are provided in Section 2.6.

2.2 mixed-effects models to assess usual intakes

The setting under which usual intakes are analyzed is always the same. We observe food consumption of n ° 1 individuals, indexed by i, taken randomly within a given population, during JPt2;:::; 7u days. Let $X_{i,j}$ be the true nutrient (or contaminant) intake of individual i on day j, and $X_i = E X_{i;i} | i$ their usual intake. Here, E (.|i) represents the expectation conditional on being individual i. Because it is not possible to observe $X_{i;i}$ directly in a chosen population, we use a proxy, denoted by $\Re_{i;i}$, obtained by crossing consumption with composition data. Specifically, let C_{i+i}^h be the amount of food h ingested by individual i on day j, and Q^h the average level of nutrient (or contaminant) contained in h. We observe $\mathbf{C}_{i;i}^h$, an approximation of $C^h_{i \cdot i}$ reported by individual i during a survey, and Q^h an estimated version of Q^h obtained via multiple measurements on various samples of food h. It is assumed that E $\mathbb{Q}_{i:i}^h \square_{i:i}^h$; j = $C_{i:i}^h$ and E \mathbb{Q}^h = \mathbb{Q}^h , i.e. data collection methods produce unbiased estimates of the quantities of interest. Moreover, individuals are supposed to choose and report their food independently from their contents, which implies in particular that E $\mathbf{Q}_{a;j}^h \square \mathbf{Q}^h = \mathbf{E} \mathbf{Q}_{i;j}^h \square \mathbf{E} \mathbf{Q}^h$. Then, for a given nomenclature of H foods, $\mathbf{R}_{i;j} := \prod_{h=1}^{\infty} \mathbf{C}_{i;j}^{h} \square \mathbf{Q}^{h}$ is an unbiased estimate of $X_{i;j}$, and one can write the following measurement error model, for all i Pt1;:::;nu, j Pt1;:::;Ju:

$$\Re_{i;j} = X_i + "_{i;j};$$
(2.1)

where $"_{i;j}$ is a within-individual error with null expectation. A few additional constraints are assumed to hold:

- (C₁) X₁;:::;X_n are independent, identically distributed (iid) with cumulative distribution function (cdf) F, independent from noise terms "_{i;i},
- (C₂) within-person errors are individually (over i = 1; ...; n) and temporally (over j = 1; ...; J) independent from each other, with common cdf F_n .

Notice that in practice, (C_1) requires in particular that two people of the same household are not both included in the sample and that (C_2) excludes surveys like INCA2, where the days of observation are consecutive. Many different techniques were developed from this general background (van Klaveren et al., 2012; Dodd et al., 2006). Though all based on the same master model, they can substantially differ on many points. Depending on which family of distributions both F and F[•] are assumed to belong to, confidence intervals may or may not be computable, episodic consumption taken into account, auxiliary information exploited, etc. The final objective also plays a role in the choice of method, conservatism being sometimes preferred to realism (we would rather have a systematic over-estimation of the quantities of interest than taking the risk of under-estimating them). These techniques are regularly evaluated and compared by scientific committees mandated by public health institutes, in order to provide methodological guidelines to practitioners (see for instance Sections 2, 3 and 4 of the latest EFSA report van Klaveren et al., 2012, or the older recommendations of U.S. Environmental Protection Agency, 1999; WHO, 2000; EFSA, 2006). Here we present two major approaches advocated by EFSA, namely the BetaBinomial-Normal (BBN) and LogisticNormal-Normal (LNN) models, respectively introduced by de Boer et al. (2009) and Tooze et al. (2006, 2010).

2.2.1 The BetaBinomial-Normal model

Referring to de Boer et al. (2009), the BetaBinomial-Normal model, abbreviated BBN, is particularly appreciated because it enables to deal with episodic consumption. It relies on the simple decomposition:

$$F(x) = P(X_i = 0) + P X_i \S x X_i \circ 0 P(X_i \circ 0).$$

Specifically, it considers that the probability of consuming (or intake frequency), namely $p_0 := P(X_i \circ 0)$, has Betabinomial distribution, while positive intake amounts, denoted by X_i^+ , are normally distributed. Since it can happen that the distribution of X_i^+ is asymmetric, a preliminary logarithmic or power transformation of the data is usually recommended (Box and Cox, 1964). The BBN model further requires that p_0 and X_i^+ are independent. This prohibits for instance situations where the occasional intakes of episodic consumers are systematically small. In practice, estimation is achieved in two separate steps, one dealing with p_0 and the other with X_i^+ , as detailed herein-after.

- Step 1 Using both null and non-null observations, the parameters of the Betabinomial distribution of p₀ are assessed by means of maximum likelihood estimation.
- Step 2 Relying solely on the positive observations, denoted by $\Re_{i;j}^+$, a preliminary power transformation is achieved:

To choose an appropriate \Box , a grid of candidate values is explored and the retained value \Box^{\Box} is that minimizing the sum of squared residuals ensuing from

a regression of normal Blom scores on the transformed intakes (de Boer et al., 2009, p.1436). Then, the parameters of the following model are fitted with the maximum likelihood procedure:

$$g_{\Box} \stackrel{\square}{\overset{\square}{\mathbf{R}^+_{i;j}}} = \Box + c_i + u_{i;j};$$

where c_i is the between-person effect with Normal distribution $N(0; \square_B^2)$ and $u_{i;i}$ the within-person effect with Normal distribution $N(0; \square_W^2)$.

Once all model parameters assessed, the distribution F of the usual intakes is obtained by Monte-Carlo approximation. First, realizations of p_0 are simulated according to its estimated distribution. When they are equal to 1 (i.e. consumption is supposed to occur), the amount is drawn from a normal distribution $N(p; p_B^2)$, where p and p_B^2 are the estimated versions of \square and \square_B^2 respectfully. Concretely, this comes to annihilating the intra-variation effect encapsulated by $u_{i;j}$, which represents the individual fluctuations around the corresponding long-term dietary habits. Finally, a back-transformation is performed to go back to the initial scale (before applying the function g_{\square} to the raw data). Notice that it is also possible to add auxiliary information to the procedure by introducing covariates in the models of p_0 and $\mathcal{R}_{i;j}^+$ (de Boer et al., 2009, p.1436).

2.2.2 The LogisticNormal-Normal model

The LogisticNormal-Normal model of Tooze et al. (2010), first introduced in Tooze et al. (2006), is very similar to the BBN method depicted above. The main differences are the following.

- The transformation function is of the Box-Cox form (Box and Cox, 1964)

and the choice of an optimal \Box is directly handled in the maximum likelihood estimation of the parameters of the distribution of $\Re_{i,i}^+$.

- Correlation between p_0 and $\boldsymbol{\Re}^+_{i;j}$ can be introduced.
- The parametric model of p₀ is no longer Betabinomial, but Logistic-Normal.

In both methods (BBN and LNN), statistics of interest such as percentiles can be estimated using the simulated usual intakes. Unfortunately, none of these approaches provides associated confidence intervals, since it would require controlling at the same time the variance of the estimates of the various model parameters, that of the Monte-Carlo approach and the error due to the back-transformation.

Obviously, many more methods than these two seminal contributions have been proposed in the literature. We refer to van Klaveren et al. (2012, Sections 2, 3 and 4)

and Dodd et al. (2006) for a more detailed account of the available models for usual intakes.

2.3 heavy-tail modeling of usual intakes

All the methods that were just introduced share a common characteristic: up to a transformation, both usual intakes and the noise terms are supposed normally distributed. If fitting a Gaussian distribution can convey acceptable approximations of average phenomena, it rarely is the case when focus is on extreme objects such as the 95th and 99th percentiles, or the probability of getting over some maximum intake limit. This issue is clearly mentioned in de Boer et al. (2009, p.1438): "departures from normality may give biased estimation of the model parameters and, hence, may give wrong inference about the usual intake distribution". In the next paragraph, we review the basics of extreme value theory that are of help to understand the validity of this statement. Given this background, we then propose an alternative model dedicated to the assessment of the tail of the distribution of positive usual intakes (we are not interested in episodic consumption).

2.3.1 Extreme value theory and heavy-tailed distributions

In a univariate context, extreme value theory is dedicated to the analysis of the tail of some distribution of interest, corresponding here to $\overline{F}(x) := 1 \square F(x)$ for all large intakes $x P R_+$. Notice that it is intricately linked to the tail quantile function, defined for all $x \bullet 1$ as

$$U(x) := \inf y PR_+ : \overline{F}(y) \S \frac{1}{x}$$

by setting p := 1=x for some p P(0; 1), it is easy to see that U(1=p) is no other than the quantile of order $1 \square p$ of the studied usual intakes, also denoted by

$$Q(p) := infty PR_+ : F(y) \bullet pu.$$

$$\lim_{n \to 1} P \frac{X_{n;n} \square b_n}{a_n} \S x = G(x); \qquad (2.2)$$

for any x P R at which the limiting distribution function G is continuous. If such normalizing constants exist, then G is called an extreme value distribution and F is said

to be in the maximum domain of attraction of G (abbreviated FPMDA(G)). Going back to the tail of the distribution and referring to De Haan and Ferreira (2006, Theorem 1.1.2), Equation (2.2) is equivalent to

$$\lim_{t \to 1} t \overline{F} a_{ttu} x + b_{ttu} = \Box \log G(x); \qquad (2.3)$$

with ttuthe integer part of t. Fortunately, (univariate) extreme value distributions can be written in a parametric form, introduced first by Fisher and Tippett (Fisher and Tippett, 1928) then by Gnedenko (Gnedenko, 1943):

$$G(x) = G_{\Box}(x) := \exp^{\left[(1 + \Box x)^{\Box = \Box} \right]}; \text{ for } 1 + \Box x \circ 0;$$
 (2.4)

with, by convention, $G_0(x) = \exp t \Box e^{\Box x} u$ and $\Box P R$. Called the extreme value index (EVI), defines three types of maximum domains of attraction. Specifically, if FPMDA(G_{\Box}) with \Box ° 0, then F is said to be in the Fréchet maximum domain of attraction, which characterizes heavy-tailed distributions like Pareto laws. Conversely, thin-tailed distributions such as the Normal or Lognormal laws belong to the Gumbel maximum domain of attraction, with $\Box = 0$. The last maximum domain of attraction, called Weibull, verifies \Box † 0 and encompasses distributions with bounded upper tails. The most classical example of the latter would be the Uniform distribution on some closed interval. Of course, it can so happen that some distributions belong to neither of these three maximum domains of attraction. Log-Pareto distributions with cdf $F(x) = 1 \square (1 + \square(\log x \square u) = \square)^{\square 1 = \square}$; $\square \circ 0$, u P R and $\square \circ 0$, marked as super-heavy-tailed, are part of these exceptions (Cormann and Reiss, 2009). Further notice that only some types of thin-tailed (resp. bounded) distributions belong to the Gumbel (resp. Weibull) maximum domain of attraction (Embrechts et al., 2011, Sections 3.3.2 and 3.3.3). In fact, F can have a bounded upper tail and still verify $F P M DA(G_0)$. Conversely, light-tailed distributions with bounded upper tail can belong to the Weibull maximum domain of attraction. The interested reader may refer to Embrechts et al. (2011, Chapter 3) for a very thorough introduction to the fluctuations of maxima, enlivened by many detailed examples. In particular, Tables 3.4.2, 3.4.3 and 3.4.4 from p.153 to p.157 in this manual display a list of characteristics of the most common distributions in each maximum domain of attraction, including some possible normalizing constants.

In the present study, we are interested in heavy-tailed distributions that respect F P MDA(G_{\Box}), \Box ° 0. They are characterized as follows (Embrechts et al., 2011, Section 3.3.1):

$$1 \square F(x) = x^{\square 1 = \square} L(x); \qquad (2.5)$$

where $\square \circ 0$ is the EVI and L(x) a slowly varying function (svf), i.e. for all t $\circ 0$, L(t x)=L(x) — 1 as x — +1. Distributions of that form are said to be regularly varying with index $\square 1=\square$ (Resnick, 2007); the set of such functions is denoted by R_{$\square 1=\square$}. They enjoy quite a few useful properties, e.g. their tail quantile function is

also regularly varying with index \Box . When concerned with the estimation of the EVI, Equation (2.5) is a typical semi-parametric model, where the svf L plays the role of the perturbation. Actually, to guarantee its accuracy, statistical inference often requires this disruptive function to fulfill some extra regularity criteria. Called Von Mises or second order conditions, they are usually written as follows (cf. Goldie and Smith, 1987 and De Haan and Ferreira, 2006, Chapter 2).

Assumption 2.1 The regularly varying tail quantile function U P R_{\Box} with \Box ° 0 is such that there is a real parameter \Box † 0, referred to as the second order parameter, and a positive or negative function A with $\lim_{x \to -1} A(x) = 0$ such that for any t ° 0,

$$\frac{1}{A(x)} \frac{U(t x)}{U(x)} \Box t^{\Box} \sum_{x=1}^{\Box} t^{\Box} \frac{t^{\Box} \Box 1}{\Box}$$

or equivalently

$$\frac{1}{A \frac{1}{\overline{F}(x)}} \stackrel{\Box}{\overline{F}(t x)} \underbrace{t^{\Box 1=\Box}}_{x-1} \underbrace{t^{\Box 1=\Box}}_{x-1} \underbrace{t^{\Box = \Box} 1}_{\Box}.$$

This ensures that L is almost of the form of a Hall slowly varying function $1 \square cx^{\square}$, which makes its influence over the Pareto form $x^{\square 1=\square}$ of \overline{F} controllable. There, the parameter \square controls the speed of convergence in extremes and, as we shall see later on, plays a crucial role in the asymptotic analysis of most estimators.

In terms of nutrient or contaminant intakes, assuming X_i has cdf F as in Equation (2.5) means that the larger \Box , the less the probability that X_i reaches high levels is negligible. On the contrary, constraining F to be in the Gumbel maximum of attraction comes to considering that extreme events can be rare enough to be disregarded. Actually, even if normality is not attributed to the raw data but to its Box-Cox-transformed counterpart, usual intakes will still be implicitly assumed thintailed (Teugels and Vanroelen, 2004; Wadsworth et al., 2010). Depending on which of these two hypotheses is privileged, corresponding large percentiles and probabilities of exceeding some maximum intake limit may differ significantly. In particular, dietary risks are bound to be drastically under-estimated if heavy tails are ignored. So as to decide which model to apply to our data, one may use statistical testing procedures borrowed from the literature in EVT such as those depicted in the next section.

2.3.2 Testing the heavy-tail assumption

Because constraining F to be thin-tailed when it is in fact heavy-tailed may induce serious bias in tail estimation, we would like to test the null hypothesis

$$H_0$$
: FPMDA(G_); \Box ° 0 versus H_1 : FPMDA(G_); \Box § 0.

Since we do not observe X_i directly but rather a collection of $(\Re_{i;j})_{1 \le j \le J}$ as in Equation (2.1), we cannot apply testing procedures directly. An extra assumption has to be made to enable further analysis and we will consider the following conditions.

Assumption 2.2 The cumulative distribution functions F and F^a are such that

 $D^{`}_{0} P[0; +1]: \frac{1 \Box F_{"}(x)}{1 \Box F(x)} \underset{x \to 1}{\Box} ^{`}_{0} \text{ and } @t P(0; 1); \limsup_{x \to +1} \frac{1 \Box F(t x)}{1 \Box F(x)} \uparrow 1.$

Assumption 2.3 The cumulative distribution functions F and F⁺ are such that

 $\frac{1 \Box F^{_{\!\!\!\!\!\!\!\!\!}}\left(x\right)}{1 \Box F(x)} \underset{x \to 1}{\sqsubseteq} 0 \quad \text{and} \quad Dx_0 \mbox{ } \mbox{ } 1 \quad : @x \ ^{\circ} \ x_0 \ ; \ F^{_{\!\!\!\!\!\!\!\!\!\!\!\!\!\!}}\left(x\right) = 1.$

If either Assumption 2.2 or Assumption 2.3 holds, then under the setting introduced at the beginning of Section 2.2 the cdf of $\Re_{i;j}$, denoted by $F_{i;j}$, is in the Fréchet maximum domain of attraction with EVI $\square^{\square} \circ 0$ if and only if F P M DA(G $_{\square^{\square}}$), as shown in Maddipatla et al. (2011, Theorem 3.3 and Remark 3.4). The first mathematical constraint in Assumption 2.2 simply means that the tail of Fr should either be equivalently thick or thinner than that of F. Without this assumption, we cannot tell anything about F based solely on $\Re_{i:i}$. Though it may seem quite restrictive and cannot be tested, such a requirement encompasses a large class of distributions for both F and F. In particular, F. can still be heavy-tailed, as long as its EVI does not exceed that of F. Notice that if the former is indeed in the Fréchet maximum domain of attraction with EVI D, then to allow the existence of its expectation we need to have $\square_{i} \ddagger 1$. Similarly, for "_{i;j} to have a finite variance, it is required that $\square_{i} \ddagger 1=2$. We refer to Embrechts et al. (2011, Section A3), De Haan and Ferreira (2006, Chapter 1) or Beirlant et al. (2004, Chapter 2) for more details on these properties. Hence, desirable regularity conditions on individual errors limit the thickness of the tail of their distribution, thereby advocating the reasonableness of our required assumption. At this point, we still cannot use classical testing techniques, because our sample of observations t $\Re_{i;i}$; 1 § i § n; 1 § j § Juis not iid. However, it suffices to average the data over time to get into an appropriate setting. According to the model depicted in Section 2.2, we have

$$\overline{X}_{i} := \frac{1}{J} \prod_{j=1}^{P} \Re_{i;j} = X_{i} + \frac{1}{J} \prod_{j=1}^{P} "_{i;j} = :X_{i} + "_{i}; \qquad (2.6)$$

where all \overline{X}_i , 1 § i § n, are iid with cdf $F_{\overline{X}}$. Further denote by F_{F} the cdf of $J^{\Box 1} \overset{\circ}{}_{j=1}^{J}$, then if F_{F} fulfills Assumption 2.2 or Assumption 2.3 instead of F_{F} , we still have the equivalence (F P M DA(G_{\Box}); $\Box^{\Box} \circ 0$) Ù ($F_{\overline{X}}$ P M DA(G_{\Box})). This requirement is met for a large set of distributions. In particular, if F_{F} has a right upper bound, then so does F_{F} and Assumption 2.3 is verified. Other examples of accurate

distributions can be found, for instance, in Embrechts et al. (2011, Sections 1.3 and A.3.2).

From now on, we will consider the model in Equation (2.6) and assume that Fand F fulfill either Assumption 2.2 or Assumption 2.3. To test whether F (or equivalently $F_{\overline{X}}$) is in the Fréchet maximum domain of attraction, we first have to check beforehand that it does indeed belong to a maximum domain of attraction. For this purpose, one can use the statistical test introduced by Dietrich et al. (2002) and studied at length in Hüsler and Li (2006). It is designed to check H₀: FPMDA(G₀) for some \Box PR, and relies on the statistic

$$\mathsf{E}_{n\,;k} \coloneqq k \overset{a}{\underset{0}{\overset{1}{\longrightarrow}}} \frac{\log \overline{X}_{n\,\square\,tkt\,\mathfrak{u}\,n}\,\square\log \overline{X}_{n\,\square\,k;n}}{p_{+}} \square \frac{t^{\square\,p_{\square}}\,\square\,1}{p_{\square}} (1 \square p_{\square}) \overset{!}{\xrightarrow{2}} t^{\square}dt;$$

where $\overline{X}_{1;n} \$ $\vdots \$ $\vdots \$ $\overline{X}_{n;n}$ denotes the order statistics relative to the sample of average intakes, \Box ° 0 is a parameter to be chosen, \overline{p}_{+} and \overline{p}_{\Box} are the moment estimators of Dekkers et al. (1989), and k is a fixed number of upper values of $(\overline{X}_{1}; \ldots; \overline{X}_{n})$ that are considered to be representative of the tail of their distribution. Under H₀ and for some appropriate choice of k, the quantiles of the asymptotic distribution of $E_{n;k}$ can be computed and compared to its actual value. For more details on this procedure, we refer to Hüsler and Li (2006), who also provide a R package implementing this test. In practice, we compute $E_{n;k}$ for various choices of k (from 10 to 0.3 \Box n for instance), and accept H₀ if the resulting function almost always remains below the corresponding asymptotic quantiles.

Once it has been established that F is in some maximum domain of attraction, we may focus on testing H_0 : F P M DA(G_0), $\square \circ 0$ versus H_1 : F P M DA(G_0), $\square \S 0$. Many procedures have been developed in the literature, as reviewed in Neves and Alves (2008). We propose to use the statistic introduced in Beirlant et al. (2006); it is based on a modified version of the Jackson statistic (Jackson, 1967) that tests the "exponentiality" in the tail of the log-transformed data. Under the second order conditions in Assumption 2.1 and provided that the parameter \square is consistently assessed (eg. with an estimator of the class presented in Alves et al. (2003)), it converges as n - +1 to a Gaussian distribution. Because it has a complicated form, for clarity purposes we do not provide the explicit formula of this statistic, but rather refer to the original paper (Beirlant et al., 2006). However, we precise that just like the aforementioned $E_{n;k}$, it depends on the choice of some number k of largest values. Again, we compute this statistic for k ranging from 10 to 0.3 \square n, and accept H_0 if the resulting function generally stays in the acceptance interval.

2.3.3 Assessing extreme quantities

In this section we deal with the case where both tests depicted herein-before were accepted, i.e. where $F_{\overline{X}}$, thus F, was found to be heavy-tailed, in the Fréchet maximum

domain of attraction. We provide examples of widely-used estimators to assess the corresponding EVI , extreme percentiles and small probabilities of exceeding some maximum intake limit. Of course, there is a plethora of such estimators, that are eluded here in favor of the general advantage of our probabilistic model for dietary risk assessment. In particular, it is well-known that the statistical tools introduced herein-after can perform quite poorly in a number of situations (Embrechts et al., 2011, Remarks p.337-338), and alternative estimators may be preferred. However, the latter are often quite complicated and based on intricate mathematical results that go beyond the scope of the present chapter. For the sake of clarity, and because we are only interested here in the methodological aspect of the modeling of usual intakes, we deliberately choose to introduce simple, well-known, techniques as means of illustration. The interested reader may refer to comprehensive textbooks such as Embrechts et al. (2011); De Haan and Ferreira (2006); Beirlant et al. (2004); Resnick (2007); Reiss and Thomas (1997) or to Caeiro and Gomes (2008); Beirlant et al. (2012); Scarrott and MacDonald (2012) to get a wider overview of the state of the art on this topic.

Ú The EVI One of the most famous estimators of the EVI for heavy-tailed distributions is the Hill estimator (Hill, 1975), computed on a number k of largest values in the sample, supposedly representative of the tail of the underlying distribution:

$$H_{k;n} := \frac{1}{k} \prod_{i=1}^{[n]} \log \overline{X}_{n \ \square \ i+1;n} \ \square \ \log \overline{X}_{n \ \square \ k;n} \,.$$
(2.7)

It corresponds to the maximum likelihood estimator of \Box when the tail of the distribution is assumed to be exactly Pareto, i.e. $\overline{F}(x) = cx^{\Box 1=\Box}$ for some c° 0 and $x^{\circ} x_{\Box}^{\circ}$ 0 such that $\sum_{x_{\Box}}^{z_{\Box}} c = x^{\Box 1=\Box 1} dx = 1$. Under Assumption 2.1 and assuming k = k(n) - 1, k=n - 0, $\overline{k} A(n=k) - \Box + 1$ as n - 1, then

$$\frac{?}{k}(H_{k;n} \square \square) \underset{n-1}{\square} N \frac{\square}{1 \square \square}; \square^{2} \text{ in distribution,}$$

see for instance Embrechts et al. (2011, Example 4.1.12 and Section 6.4.2, Method 2), Beirlant et al. (2004, Section 4.4) or De Haan and Ferreira (2006, Section 3.2). Obviously, one needs to select some acceptable value k on which to calculate $H_{k;n}$ while ensuring its asymptotic normality. This is a typical bias-variance dilemma: if k is small, the selected observations are likely to belong to the tail of the distribution, producing slightly biased estimates. However, they should be in small number, thereby implying a wide variance. On the contrary, as k increases more observations are used and the variance reduces, but the bias may grow, because observations that are not located in the tail of the distribution, and do not correspond to the stipulated model, may contaminate the estimation. To overcome this issue in practice, we propose to use the double-bootstrap algorithm of Danielsson et al. (2001), recommended in Gomes and Oliveira (2001). Assessing \Box usually represents a preliminary step to the estimation of some extreme quantities, such as high quantiles or small probabilities of exceeding a large threshold. We now review some classical estimators of both objects, based on H_{k:n}.

 \acute{U} Extreme quantiles and the probability of exceeding maximum intake limits From the estimation of \Box , by referring to Equation (2.4) one may easily build estimators to assess extreme quantiles or the probability of getting over some high threshold. Indeed, combining it with Equation (2.3) yields the following approximation, for some large enough x PR₊ and well-chosen t ° 0:

$$P \stackrel{\Box}{X_{i}} \circ x \stackrel{\Box}{=} P (X_{i} \circ x) \stackrel{\Box}{=} \frac{1}{t} \stackrel{\Box}{=} 1 + \stackrel{\Box}{=} \frac{x \stackrel{\Box}{=} b_{ttu}}{a_{ttu}} \stackrel{\Box}{=} 1 = 0$$
(2.8)

Notice that by setting $a_{ttu} = \Box U(t)$, $b_{ttu} = U(t)$, \Box° 0 and replacing x by $(x \Box 1) = \Box$ in Equation (2.3), we recover Equation (2.5). Moreover, if we set t = n = k for some number k of largest observations, then the empirical counterpart of U(n=k) is simply

$$U_{n}(n=k) := \inf^{\pi} y PR_{+} : \frac{1}{n} \prod_{i=1}^{n} ItX_{i} \circ yu \S \frac{k}{n} = X_{n \square k;n}$$

Coupling these remarks to the foregoing approximation in Equation (2.8) conveys quite natural estimators of both $Q(p) := U(1=(1 \Box p))$, $p \bullet 0.95$, the quantile of order p, and P (X_i ° `□) for some maximum intake limit `□ (Embrechts et al., 2011, Section 6.4.2, Method 2):

$$\mathbf{P}(X_{i} \circ \mathbf{\hat{G}}) \coloneqq \frac{\mathbf{k}^{\Box}}{n} \frac{\mathbf{\hat{G}}_{\Box \mathbf{k}^{\Box};n}}{\overline{X}_{n \Box \mathbf{k}^{\Box};n}}; \qquad (2.9)$$

$$\mathbf{Q}(\mathbf{p}) \coloneqq \inf^{!} \mathbf{y} \mathbf{P} \mathbf{R}_{+} : \mathbf{P}(\mathbf{X}_{i} \circ \mathbf{y}) \S \mathbf{1} \Box \mathbf{p}^{'} = \overline{\mathbf{X}}_{n \ \Box \ \mathbf{k}^{\Box}; n} \stackrel{\Box}{=} \frac{\mathbf{n}}{\mathbf{k}^{\Box}} (\mathbf{1} \Box \mathbf{p}) \stackrel{\Box \ \mathbf{H}_{k \ \Box; n}}{=}; \quad (2.10)$$

where k^{\Box} is the optimal number of largest observations obtained with the doublebootstrap algorithm of Danielsson et al. (2001). Confidence intervals may be computed via the asymptotic result below (Caeiro and Gomes, 2008, Proposition 1.2), assuming again that k = k(n) - 1, k=n - 0 and k = k(n-1).

$$\frac{\stackrel{?}{\overline{k}}}{\log \frac{k}{n(1-p)}} \quad \frac{\mathbb{Q}(p)}{\mathbb{Q}(p)} \square 1 \quad \square \\ \stackrel{n-1}{\longrightarrow} N \quad \frac{\square}{1 \square \square}; \square^{2}.$$
(2.11)

Equipped with these statistical tools, we now illustrate our approach on a set of nutrients, namely iron, zinc, calcium and retinol. Results are subsequently compared to those obtained with the BBN and LNN methods recommended by EFSA, which were briefly presented in Section 2.2.

2.4 case study: estimating the tail of usual intakes of iron, zinc, calcium and retinol in the french population

2.4.1 Description of the data

To assess tail characteristics of the distribution of usual intakes for iron (Fe), zinc (Zn), calcium (Ca) and retinol (arbitrarily abbreviated Re), we crossed the INCA2 consumption database with the levels of nutrients within foods of the CIQUAL base. They are both described at length in Chapter 1. So as to evaluate the impact of the number of reporting days in our model, we built a 24h-recall like data from the weekly information of INCA2. Specifically, among all 7 days considered in the latter, we selected one day of the week-end and one day of the week, with at least 3 days of interval, to mimic the sampling scheme of 24h-recall surveys. These declared amounts of ingested food were then combined with informations on their nutritional composition to produce estimates of individual intakes. We precise that no null intake was observed on the considered data, hence there is no need to model intake frequencies, denoted by p_0 in the presentation of the BBN and LLN methods (Section 2.2).

2.4.2 Results

We start this statistical analysis by computing mean intakes over all dates J = 2; 7, as in Equation (2.6). Before applying a heavy-tail model to all 4 types of nutritional intakes, we use the testing procedure of Dietrich et al. (2002) mentioned in Section 2.3.2 to check if their distributions belong to one of the three possible maximum domains of attraction. Results are displayed in Figure 2.1. Given that the corresponding test statistics mostly remain in the acceptance area when k varies, the distributions of intakes of Fe, Zn and Ca may well respect Equation (2.2). However, it is obviously not the case with Re; given the shape of its histograms in Figure 2.2, we suspect it is in fact super-heavy-tailed. Because our model does not encompass such distributions, we remove retinol from further analyses.



Figur e 2.2 - Histograms of intakes of retind

Nutrients that passed the former test are then subjected to the second statistical test of Beirlant et al. (2006), designed to detect distributions belonging to the Fréchet maximum domain of attraction (heavy-tailed). Results are presented in Figure 2.3. Again, statistics remain in the acceptance region whatever k, especially with Fe, thereby suggesting that the distributions of iron, zinc and calcium are all heavy-tailed.



Figure 2.3 – Fluctuations of the test statistic of Beirlant et al. (2006) with k on the 30% largest values of the sample, and the corresponding acceptance regions (light pink area for $\Box = \overline{p}$ and dark pink area for $\Box = \Box 2$) for iron, zinc and calcium intakes over 2 days (left hand plots) and 7 days (right hand plots) of observation

It is now possible to estimate the EVI \Box for all three distributions, extreme quantiles and probabilities of getting over some maximum intake limit, with the techniques introduced in Section 2.3.3. For comparison purposes, we also assess the last two mentioned quantities with the LLN methods recommended by EFSA, using the source SAS code available at http://riskfactor.cancer.gov/diet/usualintakes/macros. html. Results are displayed in Table 2.1 and Table 2.2.

Table 2.1 – Estimates of th	he EVI based on the Hill	estimator ((standard	errorsin	parentheses)	and	
associated optimal numbers of largest values							

	ł	<□	H _k	H _{k[□];n}		
Nutrient	2 days	7 days	2 days	7 days		
Iron	693	554	0.268 (0.010)	0.218 (0.009)		
Zinc	19	6	0.158 (0.036)	0.134 (0.055)		
Calcium	336	82	0.192 (0.010)	0.143 (0.016)		

Table 2.2 – Estimates of the 95-th and 99-th percentiles using the Hill estimator (standard errors in parentheses) and the LLN method (in mg/day)

	Quantiles						
_	Hill (7	(7 days) Hill (2 days)		2 days)	LLN (2 days)		
Nutrient	95%	99%	95%	99%	95%	99%	
Iron	20.51	29.13	22.36	34.41	19.40	23.53	
	(0.014)	(0.029)	(0.017)	(0.034)			
Zinc	14.59	18.09	17.52	22.60	15; 25	17; 97	
	(0.165)	(0.078)	(0.068)	(0.010)			
Calcium	1427	1798	1479	2014	1371	1645	
	(0.007)	(0.019)	(0.010)	(0.027)			

In view of Table 2.1, if the various tests performed on the data revealed that they were heavy-tailed, their EVI is relatively small, never exceeding 1=3. Provided these results are correct, this suggests that moments of order at least 3 exist for the distributions of iron, zinc and calcium. Moreover, it seems that observing only 2 days of consumption results in a significant increase in the heaviness of the tail. This phenomenon is naturally passed on to the quantile estimates, as can be observed in Table 2.2. This is relatively natural, since the smaller the period of observation, the further away we are from long-term habits by averaging the daily intakes. Even with 7 days, using the Hill estimator makes quantile estimates significantly higher than with the LLN model. Given that these usual intakes were found to be heavy-tailed, it is probable here that the LLN approach under-estimates extreme quantiles. From

a risk assessment point of view, it is thus safer to rely on extreme value theory than on classical methods to assess tails of distributions.

2.5 discussion

We have shown in this chapter how classical methods like the BetaBinomial-Normal and LogisticNormal-Normal presented in Section 2.2 can severely under-estimate the characteristics of the tail of the distribution of heavy-tailed usual intakes. Far from rejecting these approaches, we underlined in Section 2.3.2 that preliminary testing of the presence of a fat tail can be of substantial help when trying to decide which method to apply to the data. The nutrients that were taken as examples in Section 2.4 happened to pass these tests, but there are many other substances for which normality (up to a Box-Cox transformation) is a reasonable assumption. When it is not, we proposed in Section 2.3.3 simple methods to estimate extreme quantities such as high percentiles. Of course, EVT overflows with more refined statistics, which help avoid the classical problems of the Hill estimator (asymptotic bias). In particular, for distributions in the Fréchet domain of attraction with a very small index , Cai et al. (2011) have recently developed a family of dedicated estimators. Regression models for extreme values also exist and would permit the insertion of covariates such as the age or the sex of individuals (Beirlant et al., 2004, Chapter 7). Such extensions are of major interest, since nutritional recommendations usually depend on auxiliary variables (women or children do not have the same nutritional needs as adult men). Beyond improvement of the estimation, more work remains to be done, in particular to include the modeling of intake frequencies in our approach. This is left for future research.

2.6 supplements: on the second or der parameter □

Let us dwell for a moment on the estimation of the second order parameter \Box . To approximate its value, we used the family of estimators proposed by Alves et al. (2003) and the tuning parameters they recommend:

$$\mathbf{p}(\mathbf{k};\Box) = \mathbf{p}_{\mathsf{n}\,|\mathsf{T}}^{(1;2;3;\Box)}(\mathbf{k}) = \frac{3 \mathsf{T}_{\mathsf{n}}^{(1;2;3;\Box)}(\mathbf{k}) \Box 1}{\mathsf{T}_{\mathsf{n}}^{(1;2;3;\Box)}(\mathbf{k}) \Box 3};$$

with □ Pt 0; 0.5; 1; 2u,

$$T_{n}^{(1;2;3;)}(k) = \begin{cases} M_{n}^{(1)}(k) & M_{n}^{(2)}(k) = 2 \\ M_{n}^{(2)}(k) = 2 & M_{n}^{(3)}(k) = 6 \end{cases} \quad \bigcirc \quad 0;$$

$$\int_{0}^{0} \frac{\log M_{n}^{(1)}(k)}{\log M_{n}^{(2)}(k) = 2} \frac{\log M_{n}^{(2)}(k) = 2}{2 \log M_{n}^{(3)}(k) = 6} \quad \bigcirc \quad 0;$$

and

$$M_{n}^{(m)}(k) = \frac{1}{k} \prod_{i=1}^{m} \log X_{(n \square i+1)} \square \log X_{(n \square k)} \prod_{i=1}^{m} ; m Pt1; 2; 3u$$

For the results can be very volatile with k, they recommend to choose a final p where $p(k; \Box)$ is the most stable, usually for large k. Therefore, we start by selecting the value of \Box P t 0; 0.5; 1; 2u that yields the most stable sample path tk; $p(k; \Box)u$; we consider that a path is stable if the difference between two successive occurrences is small for a relatively long period. Thus, for a fixed \Box , we calculate

$$D_k^{\Box} := p(n \Box k + 1; \Box) \Box p(n \Box k; \Box)$$

for k Pt1; ...; n 1 u and get the cumulative empirical variance:

$$V_{k}^{\Box} = \frac{1}{k} \prod_{i=1}^{[n]} (D_{i}^{\Box})^{2} \Box = \frac{1}{k} \prod_{i=1}^{[n]} D_{i}^{\Box}; \text{ k Pt 1; ...; n } \Box 1u.$$

The \Box that minimizes V_k^{\Box} for a maximum of occurrences of k is then selected and denoted by \Box_0 . In a second step, we identify the values of k for which $V_k^{\Box_0}$ is small, eg. below the median of $V_{k-1\S k\S n \Box 1}^{\Box_0}$, and calculate our final estimate p as the empirical mean of $p(k; \Box_0)$ on these k. Results on exposures to our 4 nutrients are displayed in Figure 2.4.



Figur e 2.4 – Fluctuations of the estimated second order parameter with k on and the corresponding optimal value (red line) for iron, zinc and calcium intakes over 7 days of observation

3

SIMULTANEOUS OVER-EXPOSURE TO MANY FOOD CHEMICALS

In Chapter 2, we focused for a while on univariate extreme types of exposure to food nutrients or contaminants. The objective was to analyze the tail of the distribution of the chronic exposure to some chemical. It allowed in particular to estimate the small proportion of individuals in a given population who exceed some tolerable doses (the so-called dietary intake limits) and are thus likely to develop serious health problems in the long run. However, by analyzing only one component at a time, we ignored further noxious effects that may be caused by possible interactions between elements that are ingested simultaneously (Carpenter et al., 2002). Evaluating the sanitary impact of cocktails of nutrients and contaminants is a hot and complex topic public health institutes are currently mobilizing efforts on. For instance, Anses recently launched the Pericles research program (PEsticide Residue In vitro Combined Level of Exposure Study), which is dedicated to the identification and quantification of the risk due to the exposure to mixtures of pesticides (Crépet et al., 2013; Crépet et al., 2012; Crépet and Tressou, 2011; Béchaux et al., 2013). The present chapter corresponds to a paper written in collaboration with S. Clémencon (Télécom ParisTech, France) and recently submitted for publication, in which we develop a statistical methodology to find groups of any number of chemicals that are jointly absorbed in high quantity in the population of interest.

3.1 statistical challenges and objectives

High dimension raises important issues in applied multivariate statistics; while sample sizes are finite, the set on which probability measures are defined can be so large that extrapolation is intricate. Referred to as the curse of dimensionality (Donoho, 2000), this phenomenon makes the variance of classical estimators explode, thereby impeding inference. In extreme value analysis, the quality of estimation is all the more degraded as it is not carried out on the entire sample, but on some relatively small number of largest observations that are considered representative of the tail of the distribution. Whereas a plethora of techniques has been developed in the field of statistical learning to overcome this issue (Friedman et al., 2001), multivariate extremes in dimension larger that 2 are still handled with difficulty. It is the main purpose of the present chapter to address this issue, by developing a
non-parametric technique for identifying groups of variables exhibiting asymptotic dependence. Beyond a possible overall description of the tail dependence structure, when these classes are of small dimension, our method would enable further and more efficient assessment of multivariate tails. It combines novel statistical learning algorithms with multivariate extreme value theory (MEVT). From a practical perspective, it should be also pointed out that it includes a heuristic criterion to help select the sub-sample of extreme observations on which inference should be performed.

From a theoretical perspective, non-parametric assessment of multivariate extreme dependencies is already well documented. It relies on the necessary but quite mild assumption that there exists a tail dependence distribution, or equivalently that once marginal distributions have been transformed into standard Pareto, the cumulative distribution function of the resulting random vector is multivariate regularly varying. Then, its limit measure characterizes the extreme behavior of the original variables and possesses useful properties that facilitate its investigation. In particular, when switching to a pseudo-polar representation of the data, it can be expressed as a tensor product of two measures, one related to the radius, the other to the angles. The limit measure of the angles, termed spectral or angular measure, exhaustively embodies extreme dependencies. In the bivariate setting, the classical estimators introduced in the literature of this angular measure may vary depending on how marginals are standardized, which radius norm is picked and how measures are assessed. For instance, Einmahl et al. (2001); Resnick (2007); Beirlant et al. (2004) use the rank transform for standardization, then alternatively use the L_1 , L_2 and L_1 norms and ground estimation on the basic empirical measure. However, the spectral measure is required to be Lebesgue-dominated, thereby failing to encompass situations where it is degenerate on some points. Breaking the barrier, Einmahl and Segers (2009) introduced a maximum empirical likelihood statistic, while extending theoretical results to the full set of L_p-norms, p • 1. Bayesian models have also flourished (Boldi and Davison, 2007; Guillotte et al., 2011). In the same vein, Sabourin and Naveau (2012) recently proposed a novel algorithm that handles moderate dimensions. Though it can be viewed as a subsequent improvement in multivariate extremes analysis, their technique is only efficient when all variables considered are asymptotically dependent; higher-complexity spectral measures may unfortunately not be studied by their method. Hence the need to first identify groups of dependent variables in regard to their extreme behavior: once this preliminary analysis carried out, the aforementioned estimators would enable more precise estimation up to dimension 5.

Lately, Haug et al. (2010) have ingeniously adapted one of the most celebrated dimension reduction method, namely Principal Components Analysis (PCA), to multivariate extremes analysis. Under an elliptical copula assumption, they recover the set of straight lines summarizing best the extreme covariance function, thereby leading to a clustering of variables based on extreme dependence. Following in their footsteps, we propose to borrow concepts from statistical learning to achieve dimen-

sion reduction, without making any parametric assumption in contrast. We rather base our analysis on a mixture model of the spectral measure exploiting its specific geometry, tackled from a latent variable point of view, under which useful properties arise and enable identification and interpretation of hopefully small groups of asymptotically dependent variables. Inference mimics classical non-parametric spectral measure estimation and focuses on the cloud of observation angles related to the L₂-norm. Since they belong to the positive orthant of the unit hypersphere, also called simplex, their structure is explored through a recent algorithm fitting PCA to Riemannian manifolds (Jung et al., 2012). Not only does this procedure respect the intrinsic combinatorial geometry of the simplex, but it also enriches the set of eligible summarizing sub-manifolds compared to standard PCA. Identification of the variables exhibiting asymptotic dependence is subsequently achieved using an appropriate clustering technique (Dhillon et al., 2002) on the obtained sub-space. We also provide a heuristic to help select the number of upper values most representative of extremes, thereby circumventing a traditionally intricate issue in MEVT.

To illustrate the assets and liabilities of our method, we perform numerical experiments and conduct a real case study for long-term dietary risk assessment. Extreme value theory (EVT) has already proven useful in studying high exposures to a single toxicant (Tressou et al., 2004a; Paulo et al., 2006), but to our knowledge the question of simultaneous extreme exposures to multiple chemical elements has never been addressed from a statistical point of view.

The chapter is organized as follows: we start off in Section 3.2 with the introduction of a few notations and hypotheses, subsequently used throughout the methodological part of our work in Section 3.3. There, after recalling a few basic notions in spectral measure analysis, we introduce a mixture model for the spectral probability measure and emphasize the ensuing fruitful properties it enjoys, when viewed as a latent variable model. Then we turn to the practical aspects of the approach we promote, and thoroughly depict our strategy for statistical inference under the assumed model, based on dimension reduction algorithms, in Section 3.4. It is supported by numerical experiments carried through in Section 3.5, and subsequently applied for illustration purposes to dietary risk assessment in Section 3.6. In view of both simulation and case study results, assets, liabilities, natural extensions and required improvements of our method are finally listed and discussed in Section 3.7.

3.2 hypotheses and notations

We start by introducing a few essential notations used throughout the chapter, followed by a short listing of the main hypotheses involved in the subsequent analysis.

3.2.1 Notations

$$x = \bigcup_{x_d}^{x_1} \bigcup_{x_d}^{x_1} X_a$$
 and its transpose is denoted by $x^1 = (x_1; \dots; x_d)$.

In particular, we write $0 = (0; ...; 0)^1$ to mean the null vector in \mathbb{R}^d , and $\mathbf{e}_1; ...; \mathbf{e}_d$ the vectors of the canonical basis of \mathbb{R}^d . Operations between vectors should be interpreted matricially, eg. $\mathbf{x}^1 \mathbf{x} = \sum_{i=1}^{d} \mathbf{x}_i^2$.

 $\acute{\upsilon}$ Norms When working on \mathbb{R}^d , recall that all norms are equivalent. For any collection of norms $\}$, $_{(1)};$, $_{(2)};$::: in \mathbb{R}^d , we denote their corresponding unit spheres by $S_{(1)}^{d \square 1};$, $S_{(2)}^{d \square 1};$::: respectively, thereby emphasizing the topological dimension of these objects.

 \acute{v} Random variables For any sample $Z_1; \ldots; Z_n$ of n ° 1 independent and identically distributed (iid) random vectors on a space product of d ordered vector spaces $E_1 \square \blacksquare E_d$ with multivariate cumulative distribution function (cdf) F, dimensions are indexed by j Pt1; \ldots; du and observations by i Pt1; \ldots; nu. Order statistics are denoted $Z_{(1;j)}$ § \blacksquare § $Z_{(n;j)}$, for all j Pt1; \ldots; du. This notion is intrinsically linked to that of ranks; we define the rank function

Rank:
$$\begin{array}{cccc} E_{j} & \Box & N & \Box \\ Z_{i;j} & fi & \underset{m=1}{\overset{\infty}{\longrightarrow}} I & Z_{m;j} & Z_{i;j} \end{array}$$

where for any condition A, I tAu = 1 if A is true and 0 otherwise. Then we have $Rank(Z_{(i;j)}) = i$.

3.2.2 General setting and main hypotheses

Throughout this article, we consider a d-dimensional random vector

$$X := (X_1; \ldots; X_d)^1;$$

d • 2, with Lebesgue-dominated probability distribution P on the positive orthant $C^d := [0; 1]^d$ and cumulative distribution function (cdf) F, the tail structure of which we wish to assess. For 1 § j § d, we denote by P_j the j-th 1-dimensional marginal distribution of P, i.e the probability distribution of X_j, with corresponding continuous cdf F_j(x) := P_j([0; x]), x • 0. Statistical inference on the extreme behavior of F will be based on the observation of a sample X₁;...;X_n, n ° 1 (we shall write X_i = (X_{i;1};...;X_{i;d})¹ for 1 § i § n), supposedly drawn independently from P. We do not assume that F is characterized by its marginals, as would be the case in a situation where the X_j's are independent, or when considering copulas for modeling the dependence structure. Additionally, neither F nor any F_j, 1 § j § d, are supposed to belong to the maximum domain of attraction of an extreme value distribution. The only imposed regularity constraint, apart from continuity of marginal cdfs, is the existence of a Radon measure \Box , not identically zero and not degenerate at a point, concentrated on the blunt convex cone C^d₁ := [0; 1]^d zt 0u such that

$$t \mathsf{P} \stackrel{\square}{=} \frac{1}{t} \stackrel{\square}{=} \frac{1}{(1 \square \mathsf{F}_1(\mathsf{X}_1))}; \dots; \frac{1}{(1 \square \mathsf{F}_d(\mathsf{X}_d))} \stackrel{\square}{=} \mathsf{P} \stackrel{\square}{=} \frac{\mathsf{v}}{t-1} \square(.).$$
(3.1)

$$t \to \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ \hline 1 & 1 & F_1(X_1) \end{bmatrix}; \dots; \frac{1}{(1 \to F_d(X_d))} \to \begin{bmatrix} a & a \\ \Box & \Box & a \\ t \to 1 & C_{\Box}^d \end{bmatrix} f \to C_{\Box}^d$$

In words, Equation (3.1) simply states that there is, indeed, an extreme dependence structure between the random variables $X_1; \ldots; X_d$, exhaustively described by measure \Box , see for instance Section 8.2.3 in Beirlant et al. (2004) or Section 6.5.6 in Resnick (2007). Alternatively, consider the random vector $Z := (Z_1; \ldots; Z_d)^1$ of standardized components

$$Z_j \coloneqq \frac{1}{(1 \Box F_j(X_j))}; j Pt1; \dots; du;$$

$$(3.2)$$

and the corresponding transformed sample Z_1 ;:::; Z_n . Written this way, all Z_j with j P t1;:::;du, are standard Pareto distributed, i.e. @x • 1, P $Z_j \circ x = x^{-1}$. Let R refer to the set of regularly varying functions with index \Box , then by definition P $Z_j \circ x PR_{-1}$, with tail quantile function inftx $PR_+ : P Z_j \S x • 1 \Box 1$ =tu = t. For more details on univariate regular variation, we refer to Chapter 2 in Resnick (2007). Then, Equation (3.1) defines the distribution of Z as regularly varying in the multivariate sense. A larger overview of multivariate variation can be found for instance in Resnick (2007, Chapters 3 and 6).

Because we are interested in summarizing the dependence structure of the random variables X_1 ;:::; X_d , focus is here on the analysis of measure \Box in Equation (3.1). Theoretical properties of this mathematical object are detailed in Section 3.3, as a preliminary to the proposed inference technique, subsequently introduced in Section 3.4.

3.3 a mixture model under multivariate regular variation

In the present section, after reviewing a few known properties of the limit measure in Equation (3.1) and defining spectral (probability) measures, we introduce a useful mixture model of the latter on which inference is next based. For simplicity, we shall work exclusively with the vector Z of standardized random variables defined in the previous section.

3.3.1 Exponent and spectral (probability) measures

Under the regular variation hypotheses listed herein-before, the limit measure \Box in Equation (3.1), called exponent measure, exhibits some convenient properties that we shall exploit later on. Specifically, it is homogeneous, i.e.

for all 0 † s † 1 and Borel subset B of
$$C_{\Box}^{d}$$
;

$$\Box(sB) = s^{\Box 1} \Box(B); \qquad (3.3)$$

and fulfills d marginal constraints expressing the nature of the marginal survival functions, namely

for all j = 1;:::; d and 0 † z † 1 ;

$$\Box x P C_{\Box}^{d} : x_{j} \circ z^{(\Box)} = z^{\Box 1}; \qquad (3.4)$$

see for instance Section 8.2.2 in Beirlant et al. (2004) and Section 6.1.4 in Resnick (2007). Consequently, \Box can be expressed as a tensor product of two measures when switching to pseudo-polar coordinates. Indeed, choose two norms $\}$. $\}_{(1)}$ and $\}$. $\}_{(2)}$ on R^d and define the following mapping:

$$T: \begin{array}{cccc} C^{d}_{\square} & \square & (0; 1] \square S^{d \square 1}_{(2)} \\ x & fi \longrightarrow (\square; !) = {x}_{(1)}; x = {x}_{(2)} \end{array}$$

with $T^{\Box 1}(\Box; !) = \Box! = ! _{1)} = x$. Typical choices of norms include the L_p -norm or the sup-norm L_1 . Then, the homogeneity property stated in Equation (3.3) implies that

$$\Box \Box T^{\Box 1} := \Box_{\Box 1} b S; \tag{3.5}$$

н

where the radius measure \Box_{\Box_1} , defined on (0; 1], is such that for all x ° 0, we have $\Box_{\Box_1}((x; 1]) = x^{\Box_1}$, and the angle measure S, referred to as the spectral measure, has support on $\Box^{d \Box_1} := S_{(2)}^{d \Box_1} X C_{\Box}^d$ and satisfies

$$S(B) = \Box \overset{\sqcup}{\mathrm{tx}} PC_{\Box}^{\mathsf{d}} : \{x\}_{(1)} \bullet 1; x = \{x\}_{(2)} PBu^{\Box}$$
(3.6)

for all Borel subsets B of $\Box^{d \Box 1}$. A simple normalization of S yields the so-termed spectral probability measure Q on $\Box^{d \Box 1}$,

$$Q := S = S(\Box^{d \Box 1}). \tag{3.7}$$

Let us set $! = Z = Z_{(2)}$ and $\Box = Z_{(1)}$, then Equations (3.1), (3.6) and (3.7) imply

$$t \mathsf{P} (! \mathsf{P}_{\cdot}; \Box \bullet t) \underset{t=1}{\Box^{\vee}} \mathsf{S}(.); \tag{3.8}$$

$$P \stackrel{\square}{!} P \stackrel{\square}{!} \bullet t \stackrel{\square}{!} \frac{P}{t-1} Q(.).$$
(3.9)

In words, the latter expression stipulates that Q is the limit distribution of the angles when the radius gets infinitely large. It thereby encapsulates the extreme (or asymptotic) dependence structure between the d variables in dimension d \Box 1. Observe that Equation (3.4) can be expressed in terms of moment constraints for S and Q respectively. Namely, for all j Pt1;:::;du, we have:

$$a^{d d 1} \frac{\frac{j}{j}}{j!} S(d!) = 1;$$
(3.10)

$$\frac{! j}{1 + j} Q(d!) = 1 = S(\Box^{d \Box 1}).$$
(3.11)

3.3.2 Mixture model of the spectral probability measure

The extreme dependence structure between the variables $X_1; \ldots; X_d$ can be expressed in terms of the geometry of the support of Q (or S), which we denote by supp(Q). Indeed, recall that supp(Q) is included in $\Box^{d}\Box^1$, the positive orthant of the unit hypersphere $S_{(2)}^{d}\Box^1$, or the simplex associated with $\}.\}_{(2)}$. The latter can be partitioned into $2^d \Box 1$ non-empty and disjoint open faces with dimensions ranging from 0 up to d $\Box 1$. They are identified by the collections of indexes $I_1; \ldots; I_{2^d} \Box^1$ forming $P^{\Box}(t1; \ldots; du)$: for any h Pt1; $\ldots; 2^d \Box^1$ u such that $I_h PP^{\Box}(t1; \ldots; du)$, the open face generated by $t_{e_i}; j PI_h$ u is

$$\Box_{h}^{d \Box 1} \coloneqq \Box^{d \Box 1}(I_{h}) \coloneqq \Box^{d \Box 1} X \operatorname{Vect}(te_{j}; j PI_{h}u);$$

with dimension $m_h \square 1$, where $m_h \coloneqq \#I_h$. The star of vertex $te_j u$, i.e. the reunion of all open faces the closure of which contains $te_j u$, is denoted by $star(te_j u)$ with corresponding set of indexes $S(j) \coloneqq th Pt1; \ldots; 2^d \square 1u : te_j u \tilde{N} \square_h^{d \square 1} u$. Notice that we also have $S(j) = th Pt1; \ldots; 2^d \square 1u : j PI_h u$. See Figure 3.1 for an illustration in dimension 3. By extension, we denote by $I_0 \coloneqq to the empty face$ $\square_0^{d \square 1} \coloneqq H$, with $m_0 = 0$.



Figur e 3.1 – The 7 nonempty open faces in the L₂-norm simplex \Box^2 : the 3 vertices (left), the 3 edges (right), and the interior (bottom). The set star(te₁u) corresponds to the reunion of the 4 open faces $\Box^2(t1u)$, $\Box^2(t1;2u)$, $\Box^2(t1;3u)$ and $\Box^2(t1;2;3u)$.

Given this decomposition, for any h Pt1;:::; $2^d \square 1u$, supp(Q) X $\square_h^d \square 1 \square$ H means that all X_j such that j P I_h exhibit asymptotic dependence (see Beirlant et al., 2004, Section 8.2.3). Consequently, recovering the set of faces intersecting the support of Q suffices to identify the sets of variables which are dependent in the extremes and those which are not. This motivates the following mixture model:

$$Q(.) = \sum_{h=1}^{2^{h}} \Box_{h} Q_{h}(.); \qquad (3.12)$$

where for all h Pt1;::: $2^d \Box 1u$,

$$\Box_{h} := Q(\Box_{h}^{d \Box 1}).$$

In addition, denote by H the set made of all open faces intersecting supp(Q), i.e.

$$H := th Pt1; \dots; 2^d \square 1u : supp(Q) X \square_h^{d \square 1} \square Hu;$$

Now that the theoretical framework has been set out, following in the footsteps of standard mixture model analysis (see for instance McLachlan and Peel, 2000), we shall exploit the properties deriving from Equation (3.12) when reasoning in terms of latent variables and intrinsic clustering, in order to identify all \Box_{h}^{d} with h PH.

3.3.3 Latent variables representation

We consider now the iid copies Z_1 ; :::; Z_n of the standardized random vector Z defined in Section 3.2. Going back to the model in Equation (3.12) and setting H := #H, one would expect the empirical distribution of these observations to reflect the decomposition of Q as a mixture of H probability measures on $I_{hPH} \square_{h}^{d\square1}$. Hence, there should be an intrinsic (unknown) clustering of the data into H classes leading to an identification of these open faces. Formally, define n unobserved random vectors $\square := (\square_{;0}; \square_{;1}; :::; \square_{;2^d \square 1})^1$ with standard multinomial distribution such that $P(\square_{;h} = 1) =: p_h P[0; 1]$, where $\sum_{h=0}^{2^d \square1} p_h = \sum_{h=0}^{2^d \square1} \square_{;h} = 1$. Consider each $\square_{;h}$ as an indicator of whether observation Z_i is drawn from a distribution with spectral probability measure Q_h or not, i.e when $\square_{;h} = 1$, individual i may reach extreme values on the m_h variables identified by I_h alone. Conversely, for any other h RH, the event $\square_{;h} = 1$ indicates that each of the d coordinates of Z_i should generally have small to moderate values. Mathematically, for all h P t0; :::; $2^d \square 1$ u, i P t1; :::; nu, we define the latent vector \square_i as satisfying

$$P \stackrel{\square}{!}_{i} P \stackrel{\square}{\cdot}_{i} \bullet t; \square_{i;h} = 1 \stackrel{\square}{\overset{D}{\overset{D}{\cdot}}_{t-1}} \stackrel{\$}{\overset{Q}{}_{h}(.) \text{ if h PH;}}_{0 \text{ otherwise;}} (3.13)$$

where ! i and \Box are respectively the angle and radius of individual i. In fact, for the sake of interpretation, as soon as all $\Box_{;h}$, h PH, are null, we set $\Box_i = (1;0;:::;0)$. In other words, we impose the following equivalence:

@h Pt1;:::;2^d □ 1ư; (
$$\Box_h = 0$$
) Ù (p_h = 0). (3.14)

Some useful properties can be established in such a setting. We display here two results which are subsequently exploited for inference, as shall be seen in the next section. Proofs and technical details are deferred to Section 3.8. The proposition below exhibits the asymptotic behavior of conditional marginals under the latent variable model depicted herein-before.

Proposition 3.1 - Intra-class regular variation. We place ourselves in the framework of Section 3.3 and denote by H(j) the set th PH : j Pl_h u, i.e. the intersection between H and S(j).

Then, for all j Pt1;:::;du, h Pt0;:::;2^d □ 1u, x • 1,

$$t P Z_{i;j} \circ x t = 1 Z_{i;h} = 1 Z_{i;h} c_{j;h} x^{-1}; \qquad (3.15)$$

where $c_{j;h} P[0; 1=p_h]$ is non-null if and only if h PH(j), and $\sum_{h=0}^{\infty} p_h c_{j;h} = 1$.

In words, if $\Box_{h}^{d \Box 1}$ is an open face where Q is null or if its closure does not contain vertex te_j u, then it cannot project any mass on the j-th dimension. On the contrary, as soon as h PH(j), the marginal distribution of $Z_{i;j}$ given that it was drawn from a distribution with spectral probability measure Q_{h} is tail equivalent with the non-conditional distribution of $Z_{i;j}$. Hence, nonempty open faces intersecting with supp(Q) are identifiable by remaining only in the univariate level. In particular, the following result reveals that one can build a discriminative function of conditional marginal distributions the asymptotic behavior of which enables the characterization of tH(j); 1 § j § du, and by extension of tI_h; h PH u.

Proposition 3.2 – Face-characterizing functional. We place ourselves in the framework of Section 3.3 and consider H(j) as in Proposition 3.1. For all j in t1;:::; du, h Pt1;:::; $2^d \square$ 1u, x • 1, define the functional

$$\Box_{j;h}(t) \coloneqq \int_{1}^{a} t P(\Box \bullet t) P Z_{i;j} \circ xt = 1 dx;$$

and assume that there exist some constants $\Box^{\Box} P(0; 1)$, $c^{\Box} \bullet 0$ and $t^{\Box} \circ 1$, such that for all j Pt1;:::; du, h RH(j),

Then

As a consequence, for fixed dimension j P t1;:::;du, the set H(j) consists of all indexes h such that $\Box_{j;h}(t)$ diverges as t tends to infinity, instead of converging towards a finite constant, possibly zero. Since for all h PH, $I_h = tj : h PH(j)u$, we have

$$\Box_{h}^{d \Box 1} \coloneqq \Box^{d \Box 1}(I_{h}) = \Box^{d \Box 1} tj : \Box_{j;h}(t) \Box_{t-1} + 1 u; h PH; \qquad (3.17)$$

and H is the set of all h Pt1;:::; $2^d \square$ 1u such that there exists at least one dimension j Pt1;:::;du for which $\square_{;h}(t) - +1$ as t - 1.

Remark 3.3 The assumption in Proposition 3.2 simply requires that the extreme dependence structure is reached at a reasonably fast rate. It can be directly linked to the concept of hidden regular variation introduced in Heffernan and Resnick (2005); Resnick (2002, 2007, 2008); Das and Resnick (2011); Das et al. (2013). Roughly speaking, if the distribution of Z had hidden regular variation, there would be an angular measure on $\Box d \Box 1 z^{I}_{hPH} \Box d \Box 1$ when making the radius increase with some regularly varying function b(t) = o(t) with index 1= \Box § 1 instead of t. In that case, our assumption guarantees that 1= \Box § $\Box \dagger 1$, which is a rational condition for hidden regular variation not to be mistaken for multivariate regular variation in practice.

The proposed approach to statistical inference is based on Proposition 3.2 and Equation (3.17), as explained at length in the next section. Numerical experiments illustrating the relevance of the method we promote here are subsequently presented in Section 3.5.

3.4 statistical inference

Relying on the probabilistic framework detailed in Section 3.3, we now review the various steps of the proposed methodology to assess the dependence structure governing the extreme values of X_1 ;...; X_d . In short, it combines techniques borrowed from multivariate extreme value theory with clustering algorithms. Its declared purpose is to try to circumvent the classical curse of dimensionality that gravely deteriorates estimator variances (Massart, 1989). In particular, under some sparsity-like hypothesis, there is real hope of improvement in the estimation of the spectral measure: if supp(Q) is condensed on small manifolds of $\Box^{d \Box 1}$, recovering its geometrical structure should be manageable and would ultimately enable inference in lower, well-identified dimensions. Though high dimension is a classical issue in multivariate statistics (Friedman et al., 2001; Donoho, 2000), it is even more pregnant with meaning in extreme value analysis, where statistical inference only relies on a small sub-sample of most extreme observations. This justifies our approach, detailed step by step in the next three subsections: after a brief review of standard preliminaries, we depict our clustering algorithm given that the number H of open faces intersecting with supp(Q) is known. Finally, we propose some heuristic tools to choose both the aforementioned H and the number of upper values on which to base statistical analysis appropriately.

3.4.1 Preliminaries

Just as in classical spectral measure assessment, we consider that for some high enough threshold t, asymptotic relations such as in Equations (3.8), (3.9), (3.15) and 3.8.2 are sufficiently well approached to enable estimation. We use t = n = k in practice, where k represents a fixed number of upper radii, that has to be chosen carefully as shall be seen later on. Other essential elements of the aforementioned asymptotics are not known a priori and need to be estimated beforehand, namely the marginals $F_1; \ldots; F_d$, used to standardize the original variables in Equation (3.2). To avoid restrictive hypotheses, we privilege here a non-parametric procedure, usually referred to as the rank transform: for all i Pt1; \ldots; nu and j Pt1; \ldots; du, set

$$\mathbf{P}_{j}(X_{i;j}) = \frac{1}{n} \operatorname{Rank}(X_{i;j}) \square 1^{\square}; \qquad (3.18)$$

and pursue the analysis with $\mathbf{P}_{i;j} = 1 = (1 \square \mathbf{P}_j(X_{i;j})), 1 \S i \S n, 1 \S j \S d (Einmahl et al., 2001; Beirlant et al., 2004; Resnick, 2007; Einmahl and Segers, 2009). The generic random vector associated with this sample is written <math>\mathbf{P} = (\mathbf{P}_1; \ldots; \mathbf{P}_d)^1$. Angles and radii are subsequently denoted by $!_i^a$ and $\hat{\Box}_i^a$ respectively. For geometrical reasons explained in the next subsection, we set $\}_{(2)}$ as the L₂-norm. In addition, we use the L₁ -norm for $\}_{(1)}^{(1)}$ because of its natural adequacy with marginal analysis. Observe that whereas it is unimportant regarding the angle, selecting a specific norm for the radius can have major implications. This is due to the selection process of "tail observations", defined as those with radius larger than n=k; clearly, different norms are bound to produce different sub-samples (Einmahl and Segers, 2009). However, such issues go beyond the scope of our analysis and are not discussed further here.

3.4.2 Dimension reduction and dustering

Equipped with the objects and notations introduced herein-before, we now proceed with the analysis of the spectral measure. Assume for the moment that H is known, k is appropriate, and that our only task is to identify the set of open faces $\Box_h^{d\Box 1}$, h PH, or equivalently the corresponding collections of indexes I_h . For this, we propose to mimic a classical approach in statistical learning, namely Principal Components Analysis (PCA, Friedman et al., 2001). Actually, Haug et al. (2010) have already extended PCA to extreme dependence analysis, but they assume an elliptical copula to describe the dependence structure in extremes. Because we would like to avoid any parametric restriction, we propose to work on the angles instead of the raw data and carry out PCA directly on the simplex. Then, choosing the L₂-norm enables the use of algorithms that respect the intrinsic distance of the unit hypersphere $S_{(2)}^{d\Box 1}$, now identified with the usual unit hypersphere $S_{(2)}^{d\Box 1}$ of R^d. We shall refer to the geodesic distance introduced hereafter.

Definition 3.4 The geodesic distance between two points x and y of R^d located on $S^{d_{1}}$ is written

 $d_G(x; y) \coloneqq \arccos x^1 y^{\Box}.$

Among all algorithms that were proposed in the literature (Jung et al., 2011; Huckemann and Ziezold, 2006; Fletcher et al., 2004), we consider the most general one, namely the Principal Nested Spheres (PNS) technique developed by Jung et al. (2012). In short, it is an iterative procedure that projects the data on smaller and smaller hyperspheres ($S^{d \square 2}$, $S^{d \square 3}$;:::), until the unit circle S^1 is reached. These so-called PNS enrich the set of summarizing sub-manifolds compared to classical PCA. In Figure 3.2, we provide an illustration of the better adequacy of PNS compared to PCA when studying spectral measures. Among all d \square 1 resulting PNS, we may pick one of reasonably small dimension, which can be considered as representative of the geometry of supp(Q). In practice, this is achieved with a simple rule of thumb argument: we iteratively calculate the marginal level of geodesic variance encapsulated in S^1 ; S^2 ;::: as in Jung et al. (2012), and select the biggest PNS still providing a gain in variance that exceeds a given threshold (eg. 10% of the total variance). As was pointed out by Jung et al. (2012), S^1 and S^2 usually capture a high enough level of variance to ignore higher dimensions.



Figure 3.2 – A PNS on the L₂-norm simplex \Box^2 : assume data points are uniformly concentrated around the greyed areas, then there is no straight line that can achieve a better separation than the represented dotted circle.

Once a sub-sphere has been selected, we still need to analyze the structure of the projected points to identify $t \Box_h^{d} \Box^1$; h P Hu. For this, we propose to use an accurate clustering procedure such as spherical k-means (Dhillon et al., 2002; Maitra and Ramler, 2010), based on the geodesic distance again. Though it may not seem necessary, the first stage of PNS makes the subsequent clustering more robust by getting rid of any misleading noise while best preserving the geometry of the original cloud of points. In fine, because of the structure of $\Box^{d} \Box^1$, we can expect an adequacy between the obtained set of clusters and the underlying mixture model stated in Equation (3.12). Specifically, individuals reaching extremes on all tZ_j ; j PI_h u should

be concentrated near $\Box_{h}^{d\Box_{1}}$ and projected onto a different zone of the PNS than observations taking high values on $tZ_{j^{1}}$; $j^{1}PI_{h^{1}}$; $h^{1}\Box$ hu, and be affected to different classes. Hence, identifying the set of open faces intersecting supp(Q) comes down to finding out on which dimensions individuals of each class reach extreme values. Technicalities about this last process are detailed in the next subsection. Obviously, many other natural techniques in mixture model analysis could have been adopted here (McLachlan and Peel, 2000); our preference for geometrical methods is based on a strong belief that Riemannian geometry is a key concept for understanding the structure of the spectral (probability) measure, as suggested by the encouraging results of the numerical experiments conducted in Section 3.5.

3.4.3 Identifying groups of asymptotically dependent variables

Let us begin by still assuming that H is known and k is well chosen. In view of the procedure explained above, we have at our disposal a clustering into H groups of the set of most extreme observations $@_k := \text{ti } P \text{t1}; \ldots; \text{nu} : \square \bullet \text{ n=ku}$. These correspond to estimates of H coordinates of the unobserved vectors \square_i introduced in Section 3.3.3, for all observations i P $@_k$. We denote them by $P_{i,:}$, ` P t1; …; Hu, i P $@_k$, where $P_{i,:} = 1$ if i is in group `, and a fixed ` corresponds to some unknown h P t1; …; 2^d \square 1u. Observe that the case h = 0 is neglected, since Proposition 3.1 suggests that P ($\square \bullet \text{ n=k}$) \square 0 for large enough n=k. Unfortunately, we are not able to comprehend to which open faces the events $P_{i,:} = 1$ are referring yet. In order to recover them, we propose to take advantage of the marginal properties stated in Proposition 3.2, and start by assessing H(j) with the empirical counterpart of $\square_{j:h}(t)$ in Equation (3.17). Formally, define

where $n_k := \#@_k$ is the number of observations the radius of which exceeds n=k, and $n' := \bigcap_{i \in @_k} It P_{i;} = 1$ uthe size of class `. A more explicit version of this statistic and a short discussion regarding its accuracy are available in Section 3.8.3. Then, we set

$$\mathbf{H}(j) := \mathbf{P}(1; \dots; Hu : \mathbf{p}_{j;}(k) = \mathbf{0}^{(j)}$$

In practice, to decide which `Pt1;:::; Hu provide a large enough $\bar{p}_{j;}$ (k) to be selected, we perform a scree test-like analysis (Cattell, 1966). Specifically, we compute $\bar{p}_{j;}$ (k) on all `Pt1;:::; Hu, and j Pt1;:::; du. The V := H \Box d resulting values, denoted for instance by \bar{p}^1 ;:::; \bar{p}^V , are subsequently ordered so that $\bar{p}^{(1)}$ § $\Box\Box$ § $\bar{p}^{(V)}$. In fine, we say \bar{p}^v " 0, 1 § v § V, if

$$V \Box v + 1 \S \operatorname{argmax}_{1 \S w \S V} \overline{p}^{(V \Box w + 1)} \Box \overline{p}^{(V \Box w)}.$$

Because to each result index v corresponds a couple (j;`), we obtain in this way a collection of such couples from which all $\mathbb{P}(j)$, 1 § j § d, are determined. Then, identification of the corresponding open faces is straightforward: for all `Pt1;:::;Hu, a natural estimator of I ` is

subsequently characterizing the desired spaces $\Box \Box \Box (P)$, $1 \S \ S H$.

Unfortunately, in practice H is usually unknown and k has to be picked at hand. To overcome this issue, we develop a heuristic criterion to measure the quality of a clustering, given a couple $(k; \mathbf{\hat{H}}(k))$, where $\mathbf{\hat{H}}(k)$ is a number of clusters fixed a priori. Consider $\mathbf{\hat{H}}(j)$ as before, but replace H with $\mathbf{\hat{H}}(k)$ in its original definition, and set

$$\Box(\mathbf{k};\mathbf{P}(\mathbf{k})) = \prod_{j=1}^{\mathbf{p}_{j} \cdot \mathbf{P}_{j}(\mathbf{k})} (\Box \mathbf{1})^{\mathsf{It} \cdot \mathsf{R}^{\mathsf{H}}(j) \mathsf{u}} \mathbf{p}_{j; \cdot}(\mathbf{k}).$$

This statistic is simply built from Section 3.8.2: after having computed the empirical counterpart of \Box_{j} , (k) on all `Pt1;:::; $\mathbf{\hat{P}}(k)u$, we add up all quantities corresponding to `P $\mathbf{\hat{H}}(j)$ (which should be large) and substract the others (supposedly close to zero). When (k; $\mathbf{\hat{P}}(k)$) provides an accurate clustering of our data, $\Box(k; \mathbf{\hat{P}}(k))$ should reach high values. To avoid possible practical errors, we further refine this criterion with some additional constraints. Specifically, classes should contain more that 1 individual, groups should each identify a different open face and no set $\mathbf{\hat{H}}(j)$, 1 § j § d, should be empty. Observe that while the first two conditions are just common sense, the last one is necessary to respect the theoretical properties of Q: were there any empty $\mathbf{H}(j)$, marginal distributions could not be standard Pareto, and finding extra meaningless classes would come down to stating that the chosen threshold was not high enough to get rid of the empty face. Finally, we retain the partition inherited from (k; $\mathbf{\hat{P}}(k)$)^{\Box}, defined below.

$$(\mathbf{k}; \mathbf{\hat{P}}(\mathbf{k}))^{\Box} = \underset{(\mathbf{k}; \mathbf{\hat{H}}(\mathbf{k}))}{\operatorname{argmax}} \Box (\mathbf{k}; \mathbf{\hat{H}}(\mathbf{k})) \Box \qquad \overset{\mathbf{\hat{H}}_{ff}(\mathbf{k})}{= 1} \qquad 1 \text{ tr}^{\circ} 1 \text{ tr}^{1$$

On account of the nice properties of the rank transform (Resnick, 2007; Heffernan and Resnick, 2005; Das and Resnick, 2011), we expect these statistical objects to converge to the true quantities they approximate as n - 1. Unfortunately, due to the lack of probabilistic results on PNS and spherical k-means, which were originally introduced as geometrical techniques, we cannot provide here a thorough asymptotic analysis of the solution output by the statistical procedure described above. Further developments are the object of an ongoing work. Nonetheless, as shall be seen in the next section, numerical experiments provide strong empirical evidence of the efficiency of the approach we propose.

3.5 numerical experiments

We tested our method through a number of numerical experiments, for various values of n, d, and H. In doing so, we tried to handle various types of extreme dependence structures, to illustrate the impact of the complexity of supp(Q) on our algorithm. In the next two subsections, we first describe the different scenarios analyzed, then present and comment on the simulation results.

3.5.1 Settings

We generated n i.i.d. copies of a d-dimensional random vector $(X_1; \ldots; X_d)^1$ with varying degrees of extreme dependence. Observations were drawn using a symmetric multivariate logistic model, with function rmvevd in R package evd (Stephenson, 2003). Asymptotic dependence was controlled via a parameter r P (0; 1], which indicates the strength thereof. In particular, r = 1 gives asymptotic independence, whereas asymptotic perfect dependence occurs when $r \square 0$. We repeated 100 trials of our algorithm under 3 scenarios, listed in Table 3.1.

		<u> </u>	<u> </u>		
Scenario 1	Scenario 2		Scenario 3		
H = 2 d = 20	H = 4 d = 20		H = 4 $d = 6$		
Open faces r	Open faces	r	Open faces	r	
□ ^d [□] ¹ (t 1; 2; 3u) 0.1	□ ^d [□] ¹ (t 1; 2u)	0.1	□ ^{d □ 1} (t 1; 2u)	0.1	
□ ^{d □ 1} (t 4; : : : ; 20u) 0.1	□ ^{d □ 1} (t 2; 3u)	0.1	□ ^d [□] ¹ (t 2; 3u)	0.1	
	□ ^{d □ 1} (t 4u)	1	□ ^{d □ 1} (t 4u)	1	
	□ ^{d □ 1} (t 5; : : : ; 20)u) 0.2	□ ^{d □ 1} (t 5; 6u)	0.2	

Table 3.1 – List of scenarios considered in our numerical experiments; open faces intersecting with supp(Q) are filled in, with the corresponding extreme dependence coefficient r

To limit computation time, we tested its performance on 5 different sample sizes, namely $n = 5 \square 10^2$; 10^3 ; $5 \square 10^3$; 10^4 , and 10 thresholds t = n = k, with corresponding $k = n \square 0.001$; $n \square 0.002$; :::; $n \square 0.01$. Not all possible number of classes were exploited either, because the constraint ltn, \circ 1u appearing in the definition of $(k; \mathbf{P}(k))^{\square}$ in Equation (3.19) restricts the maximum number of clusters to $n_k = 2$. In practice, we iteratively compute our criterion for n = 1; 2; ::: and stop as soon as the next 5 iterations cease improving it. For the same reasons, we disregarded situations where d \circ 20. However, in multivariate EVT, d = 6 and d = 20 can already be considered as high dimensions. Our code is based on the PNS algorithm provided by its authors at http://www.stat.pitt.edu/sungkyu/MiscPage.html, as well as the spherical k-means version of Dhillon et al. (2002), available in R package skmeans by setting method = "pclust", and start = "S".

3.5.2 Results

Results are displayed in Table 3.2, Table 3.3 and Table 3.4. The highlighted row reports the number of trials where we managed to exactly recover the set of open faces intersecting with the support of the spectral probability measure. To better understand the assets and liabilities of our algorithm, we also provide a detailed account of all inaccurate results, ordered relatively to their impact on the final interpretation.

Table 3.2 - Results of our numerical experiments in Scenario 1, repeated on 100 trials

			n		
	₽;`Pt1;:::;ฅ(k)u	500	1000	5000	10000
Accurate sets	t 1; 2; 3u; t 4; : : : ; 20u	90	96	100	100
Other inaccurate sets		10	4	0	0

Table 3.3 - Results of our numerical experiments in Scenario 2, repeated on 100 trials

				11	
	₽;`Pt1;:::; P (k)u	500	1000	5000	10000
Accurate sets	t 1; 2u; t 2; 3u; t 4u; t 5; : : : ; 20u	54	72	97	93
Extra sets	t 1; 2u; t 2; 3u; t 4u; t 5; : : : ; 20u; t 1u	2	3	0	0
with cardinal	t 1; 2u; t 2; 3u; t 4u; t 5; : : : ; 20u; t 3u	7	2	0	1
1	t 1; 2u; t 2; 3u; t 4u; t 5; : : : ; 20u; t 1u; t 3u	1	0	0	0
Missing oot	t 1; 2u; t 3u; t 4u; t 5; : : : ; 20u	13	12	0	1
t 1: 2u or	t 1u; t 1; 2u; t 3u; t 4u; t 5; : : : ; 20u	1	0	0	0
t 2:3u	t 1u; t 2; 3u; t 3u; t 4u; t 5; : : : ; 20u	8	5	1	1
12,00	t 1u; t 2; 3u; t 4u; t 5; : : : ; 20u	0	0	0	0
	Other inaccurate sets	14	6	2	4

Table 3.4 - Results of our numerical experiments in Scenario 3, repeated on 100 trials

				n	
	₽;`Pt1;:::;ฅ(k)u	500	1000	5000	10000
Accurate sets	t 1; 2u; t 2; 3u; t 4u; t 5; 6u	39	65	85	88
Fortage and a with	t 1; 2u; t 2; 3u; t 4u; t 5; 6u; t 1u	0	1	0	0
cardinal 1	t 1; 2u; t 2; 3u; t 4u; t 5; 6u; t 3u	1	1	0	0
	t 1; 2u; t 2; 3u; t 4u; t 5; 6u; t 1u; t 3u	1	0	0	0
	t 1; 2u; t 3u; t 4u; t 5; 6u	16	11	6	2
Missing set	t 1u; t 1; 2u; t 3u; t 4u; t 5; 6u	0	0	0	0
t 1; 2u or t 2; 3u	t 1u; t 2; 3u; t 3u; t 4u; t 5; 6u	23	14	5	6
	t 1u; t 2; 3u; t 4u; t 5; 6u	0	0	0	1
Oth	er inaccurate sets	20	8	4	3

As expected, in all scenarios, results improve when n increases, and success rates become particularly satisfactory as soon as n • 5000, for they then exceed 85% in all 3 scenarios. The best performance is obtained in scenario 1, where d = 20 and H = 2. Indeed, even with a very small sample (n = 500), only 10 trials out of 100 fail to recover the true decomposition of supp(Q), while in scenario 3, where d = 6 and H = 4, this rate never goes below 12% whatever n. This suggests that rather than the dimension, the complexity of supp(Q) may be one of the principal determinants of the performance of our procedure. Actually, given two spectral probability measures with equivalently complex supports, increasing dimensionality can produce better outcomes. This is the case with scenarios 2 and 3, where supp(Q) is contained on small subsets of 4 open faces, but d = 6 in the former while d = 20 in the latter. These results are not surprising and illustrate a typical phenomenon called the blessing of dimensionality (Donoho, 2000); as d increases, observations occur in relatively small subsets of the original space and are therefore easier to detect and separate. This property is the basis for common techniques in statistical learning, such as the widely celebrated Support Vector Machine (Friedman et al., 2001, Chapter 12), which projects the data onto some space with higher dimension in which they are well divided. In our numerical experiments, switching from Scenario 3 to Scenario 2 significantly reduces the risk of overriding either $\Box^{d \Box 1}(t 1; 2u)$ or $\Box^{d \Box 1}(t 2; 3u)$, which are very close to one another in the unit hypersphere and may be wrongfully confused during the PNS procedure. Observe nonetheless that these simulations were performed for very small values of parameter r, i.e. all dependencies were strong. Since we used the multivariate logistic model, this means that for all h PH, subsets supp(Q) X $\square_{h}^{d \square 1}$ did not cover the entire open faces $\Box_h^{d \Box 1}$ but were concentrated around small neighborhoods of one of their points. Had we considered less obvious extreme dependencies, these results would have probably been significantly degraded. This remark can be linked to the influence of the hidden spectral measure on inference (Resnick, 2002), for it controls the rate at which extreme structure is reached and thus dangerously impacts statistical analysis if the chosen threshold n=k is too small.

In fine, these results are quite encouraging, and underline the usefulness of methods from the field of statistical learning for multivariate EVT. Our next step will be to conduct a full theoretical analysis of our approach, which would demonstrate its rate of convergence in terms of sample size, strength of extreme dependencies, and complexity of the underlying spectral probability measure. This would for instance enable the construction of some confidence intervals and may help design a more efficient algorithm in terms of computation time. Indeed, at this early stage of development, our procedure is very long to implement; this forced us to limit the number of trials to 100 and prevented us from exploring the whole range of largest values on which our criterion may be calculated. Hence, results might have improved, had we been able to reach its true maximum.

3.6 application to dietary risk assessment

While eating is the privileged way of providing the necessary nutrients for the human organism, it also conveys toxic elements that, due to various environmental causes, contaminate the food. When consumed over certain tolerable doses, called dietary intake limits (DIL), these toxic elements can have a non-negligible impact on health. Similar phenomena also occur when diets are either too rich or too poor in nutrients. More importantly, further noxious effects may be caused by possible interactions between elements that are ingested simultaneously (Carpenter et al., 2002). For international institutes concerned about public health issues such as the WHO (World Health Organization), FAO (Food and Agriculture Organization), Unep (United Nations Environment Program), Efsa (European Food Safety Authority) or for national agencies such as the Anses (the French agency for food, environmental and occupational health safety), it is then of major interest to identify cocktails of food chemicals to which populations are indeed highly exposed. EVT has already proven useful to assess the probability of getting over a single dietary intake limit, in both univariate (Tressou et al., 2004a) and bivariate settings (Paulo et al., 2006). Here, we propose to apply the method detailed in Section 3.3 and Section 3.3.3 to examine the relationships between high simultaneous long-term exposure to 6 common nutrients and contaminants, namely iron (Fe), calcium (Ca), sodium (Na), methylmercury (MeHg), cadmium (Cd) and dioxins and dioxin-like polychlorinated biphenyls (PCB-DL). Their long-term toxicity is well-known, see for instance Anses (2011) and Carpenter et al. (2002). Methylmercury, cadmium, and PCB-DL are three contaminants found mainly in seafood products. While cadmium was recognized in 2004 as a type 2 carcinogen by the European Union, methylmercury and PCB-DL can attack the nervous system. Sodium, calcium and iron are three minerals principally found in animal products such as meat or dairy products. Long-term over-exposure to these nutrients is also harmful, e.g. consuming too much calcium can provoke urinary and renal calculi and excessive ingestion of sodium favors cardiovascular issues. As for iron, some studies have underlined a probable link between its excessive ingestion and Parkinson disease (Jenner et al., 1992). The current knowledge about possible synergistic effects between these chemicals, which may increase sanitary risks, is still quite poor, due to the complexity of these phenomena. Only methylmercury and PCB-DL have been studied jointly, and their simultaneous consumption was observed to amplify health issues in a number of experimental surveys (Bemis and Seegal, 1999; Carpenter et al., 2002). Henceforth, recovering groups of nutrients or contaminants to which the population is observed to be simultaneously over-exposed can help orient future biological and chemical research, which would in turn provide a better understanding of dietary risks. In terms of statistical analysis, thus reducing the dimension would also enable a more accurate estimation of the complex relationships between these types of exposure. Indeed, even though they are clearly linked

by the type of food (fish or meat) introduced in the diet, there are differences of composition between species — like tuna or salmon — that can imply independence between types of extreme long-term exposure. In particular, exceeding the DIL of more than 3 of these elements is an event never observed in the data. Because of the variety of individual dietary habits and the complexity of the contamination process, simultaneous types of high exposure are not an obvious phenomenon, are rarely observed, and need to be analyzed in detail. In the next paragraphs, after a brief presentation of the data, we apply the procedure introduced in the previous sections to the 6 aforementioned nutrients and contaminants.

3.6.1 Data and required assumptions

Our vectors of 6 types of exposure were calculated on the n := 2488 non-pregnant, non-lactating adults of the INCA2 database for which no important variable was missing, as described in Section 1.2. Levels of nutrients within each of the 1342 food items were given in the CIQUAL database and equivalents for contaminants were found in TDS2, both described in Section 1.2.2. In keeping with Chapter 2, the vectors of exposure X_1 ;:::; X_n were obtained by multiplying amounts of food with average contents then averaging over the number of reported days. Using similar notations and hypotheses as in the previous chapter, for all components j Pt1;:::;du and consumers i Pt1;:::;nu, we assume that

$$X_{i;j} = X_{i;j}^{\Box} + U_{i;j};$$

where $X_{i;j}^{\Box}$ is the long-term individual daily exposure to contaminant or nutrient j (also called usual intake) and $U_{i;j}$ is an independent noise with lighter tail than $X_{i;j}^{\Box}$. We also suppose that the n vectors X_i are iid and that for any $i \Box i^1$ and $j \Box j^1$, $U_{i;j}$ is independent from $U_{i^1;j^1}$. Under this setting, it is clear that the extreme dependence between the d components of vector X is determined solely by that between the d components of X[□], the vector of interest. Relaxing these quite restrictive hypotheses would require further work; this is left for future research.

3.6.2 Analysis of extreme dependencies

Results of our method on the aforementioned sample are summarized in Figure 3.3. For various choices of threshold t = n = k such that $k P [10; n \square 0.3]$, we selected the partition that maximized our criterion $\square(k; \square(k))$ over all $(k; \square(k))$, and represented schematically the corresponding dependence structure. To get further confidence in this outcome, we summarize in Table 3.5 the strongest relationships that were found over all thresholds. The evolution of our criterion $\square(k; \square(k))$ with k is displayed on Figure 3.4. Its maximum is reached when k = 564, i.e. when calculations are based on the 1591 observations with largest radii.

In fact, the dependence structure represented in Figure 3.3 is found on all 16 largest values of \Box (k; \mathbf{P} (k)). The corresponding number of largest values k can be divided into two groups, one where k is in a neighborhood of 360, and another where k is around 560, as illustrated by the highlighted regions in Figure 3.4. Moreover, Table 3.5 shows that some dependencies are spotted whatever the number of largest values. In particular, methylmercury is almost always associated to PCB-DL, while cadmium and calcium get separated from all other chemicals. Concerning iron and sodium, uncertainty remains quite high, and a complementary bivariate analysis seems necessary to confirm the nature of their relationship. Figure 3.5 shows the estimated bivariate spectral probability measures of joint exposure first to MeHg and PCB-DL, then to Fe and Na. They were obtained using the maximum empirical likelihood (abbreviated MEL) approach of Einmahl and Segers (2009).



Table 3.5 - N umber of times extreme dependencies occur among all thresholdst = n=k, k P[10; n \Box 0.3] (in %)

MeHg	MeHg&	PCB-DL		
	PCB-DL			
7.60	97.15	45.32		
Cd		Ca		
95.52		78.83		
= 0 1				
re& N	ia	Na		
50.88	48.85			

Figur e 3.3 – Dependence structure between the 6 nutrients and contaminants of interest, on k = 564 that maximizes our criterion; arrows indicate extreme dependencies and the number of observations within each dass is given in parentheses

Clearly, the strong asymptotic dependence between methylmercury and PCB-DL is confirmed, on whatever value of k the estimation may be carried out. The presence of a sub-population reaching extreme exposure to PCB-DL alone is also suggested by the form of \hat{Q} , which gets close to vertical height on the extreme left part of the plot for many values of k. However, methylmercury does not exhibit such a behavior, and given that a specific class of independent exposure to MeHg only occurs for 7.60% of the largest values, we decide to disregard it. Actually, in terms of dietary habits, getting two clusters of individuals, one highly exposed to both MeHg and PCB-DL, and another solely to PCB-DL makes perfect sense. Indeed, contrary to PCB-DL, methylmercury is a contaminant found exclusively in seafood products. Hence, it is possible to get over-exposed to PCB-DL without ingesting high amounts of MeHg.

Now, let us turn to iron and sodium. According to the evolution of \hat{Q} with k shown on Figure 3.5, if these two types of exposure exhibit asymptotic dependence, the latter is clearly weak. In fact, we are more inclined to believe in the presence of a mixture of three sub-populations, one ingesting high amounts of both Fe and Na, and the other two getting over-exposed to only one of these nutrients. It is also possible that k = 564 being quite high, the relationship appearing in Figure 3.3 corresponds not to extreme but moderately high levels of exposures. This inconclusive example suggests that extending our approach to the analysis of the hidden spectral measure (Resnick, 2002, 2008) would be of major interest.



Figur e 3.4 – Evolution in log-scale of $\Box(k; \mathbf{P}^{\Box}(k))$ with the number of largest values k, where $\mathbf{P}^{\Box}(k)$ is the number of classes maximizing our criterion for some fixed k. The two dashed lines indicate the location of $\Box(k; \mathbf{P}(k))^{\Box}$, while the graved areas highlight regions where the 16 best criteria are obtained



Figur e 3.5 – MEL estimator of the bivariate spectral probability measure Q, obtained for various values of k (grey lines) up to k = 564 (black line), the optimal number of largest values selected by our criterion; the horizontal dashed line represents asymptotic independence, and the vertical one perfect asymptotic dependence

3.7 discussion

Non-parametric analysis of extreme dependencies via the spectral measure in high dimension d is still an open issue in multivariate extreme value theory. Though the bivariate setting has already been thoroughly investigated (Beirlant et al., 2004; Resnick, 2007; de Haan, 1985; Einmahl et al., 2001; Einmahl and Segers, 2009; Guillotte et al., 2011), and moderate dimensions are now accessible when all variables are asymptotically dependent (Sabourin and Naveau, 2012), the matter is still unresolved for d ° 5. Following in the footsteps of Haug et al. (2010), who adapted the most celebrated Principal Components Analysis to extreme dependence assessment, we proposed a method combining multivariate extreme value theory with statistical learning and data mining standards so as to identify sub-groups of variables exhibiting asymptotic dependence. Once these dusters are identified, if they each encompass less than 5 variables, it then becomes possible to further estimate the corresponding sub-parts of the spectral measure with any existing method, for instance those that were cited herein-before.

We started in Section 3.3 by developing the theoretical context under which our approach was constructed. First of all, contrary to Haug et al. (2010), we did not make any parametric assumption on the extreme dependence structure. This led us to focus on the spectral measure itself, or more specifically its standardized version called spectral probability measure Q. After recalling that this can be viewed as the limit distribution of observation angles given that their radius is getting infinitely large, we underlined the adequacy between the geometry of its support on the positive orthant of the unit hypersphere and the nature of extreme dependencies. Indeed, if a group of variables, say Z₁ and Z₂, are asymptotically dependent, then Q will have positive mass on the open face generated by the corresponding dimensions, here $\Box^{d \Box 1}$ (t1;2u). Therefore, we proposed a natural model of the spectral probability measure as a mixture of angular distributions with supports on each of the $2^d \square 1$ non-empty open faces of the simplex. Tackled from a latent variable point of view, this model provided particularly useful properties, formulated in Proposition 3.1 and Proposition 3.2, that reduced an initially d-dimensional problem to d univariate ones. In particular, we showed that open faces intersecting the support of Q, namely $t \square_{h}^{d \square 1}$, h PHu, could be identified by means of a simple functional $\square_{j;h}(t)$, 1 § j § d, 1 § h § $2^d \square$ 1, introduced in Proposition 3.2.

Then, we moved to the practical part of our method in Section 3.4. Because we had pinpointed the major role of geometry for analyzing spectral measures in the preceding section, we adopted geometrical techniques suited for Riemannian objects such as the unit hypersphere for statistical inference. Borrowed from the statistical learning field, they consist in first projecting the initial cloud of points on a lower-dimensional space by means of the Principal Nested Spheres algorithm of Jung et al.

(2012), then clustering the obtained data with spherical k-means (Dhillon et al., 2002; Maitra and Ramler, 2010). By first implementing PNS, we reduced potential noise and enabled more efficient classification. Resulting clusters were then considered as representative of the open faces that intersect supp(Q), and analyzed as such. To recover to which \Box_h^{d-1} , 1 § h § $2^d \Box$ 1, they were referring, we constructed estimators based on the empirical counterpart of the functional $\Box_{j;h}(t)$. The latter was also exploited to build a heuristic statistic that selects both the appropriate numbers of groups of dependent variables and of "extreme" observations. Unfortunately, due mainly to the absence of probabilistic analysis of PNS and spherical k-means in the literature, we were not able to provide asymptotic results about the aforementioned objects (this is the object of an ongoing work). Hence, assets and liabilities of our technique were discussed based solely on numerical experiments.

In Section 3.5, we tested our method on a set of simulated data bases. Three scenarios were considered, which try to encapsulate as many different situations as possible: they differed depending on d, on the number H of open faces containing mass, and on the complexity of supp(Q). In spite of a clearly improvable practical algorithm, the encouraging results we obtained enabled us to define which characteristics of Q have most influence on estimation. In particular, we saw that unlike H, d is of negligible importance to the complexity of supp(Q) and the strength of extreme dependence. The closer t $\Box_{h}^{d \Box 1}$; h P H u are to one another (e.g. both $\Box^{d \Box 1}$ (t 1; 2u) and $\Box^{d \Box 1}$ (t2; 3u) intersect supp(Q)), the harder it is to separate and correctly identify each of them. Estimation may also be impeded if H is large in comparison with n, or if asymptotic dependencies are weak. Indeed, to easily spot the desired open faces, the corresponding angular distributions denoted by Qh should concentrate most of their mass on a small neighborhood of the middle point of $\Box_{h}^{d \Box 1}$; the weaker dependencies are, the farther we are from this ideal situation. Though they were not considered in the simulations, we added some comments on rates of convergence to the asymptotic dependence structure that were sensed as a determining factor in assessment efficiency. Specifically, we insisted on the role that the hidden spectral measure may play when selecting an optimal number of largest values and suggested the interest of generalizing our approach to its analysis.

Further insight into our method was provided by a case-study illustration. Applied to real databases about exposures to 6 food contaminants, it produced stable outcomes, thereby giving confidence in the results. We were able to conclude that only two pairs of chemicals are actually linked in extremes, namely methylmercury and PCB-DL on the one hand, and iron and sodium on the other hand. These associations were confirmed by further computing the MEL estimator of Einmahl and Segers (2009) on the two pairs of variables. In addition, our method spotted a configuration usually hard to notice with traditional estimators, but quite natural given the underlying mixture model on which we based the analysis: it underlined the presence of a mixture of populations, some being jointly over-exposed to a couple of

elements, while others ingest high quantities of only one of them (PCB-DL or Na). In terms of public health implications, this means that people who are over-exposed to methylmercury tend to ingest simultaneously high amounts of dioxins and PCB. Knowing that these two toxicants have similar noxious effects on the human organism (Fischer et al., 2008; Weihe et al., 1996), and that when combined, synergistic effects can occur (Bemis and Seegal, 1999; Carpenter et al., 2002), this suggests paying particular attention to the populations that do not respect the corresponding DIL. It also justifies the need for specific research on potential combined effects of these two contaminants, which would help in assessing the sanitary risks brought upon the concerned population.

In view of these results, one advantage of our multivariate approach is that people in the data are dispatched into multiple classes that embody different types of extreme dependencies. In our case-study example, it facilitates the understanding of over-exposure categories by allowing classical discriminant analyses. An interesting alternative would be to model the various \Box_h appearing in the mixture model of the spectral probability measure in function of auxiliary covariates, eg. some sociologic or economic variables here. More than providing easily interpretable results, this would probably increase the performance of our procedure by helping discriminate between the various clusters. Such generalizations of the present work will be the subject of further investigation in the near future.

3.8 proofs and supplements

3.8.1 Intra-class regular variation

We shall start the proof of Proposition 3.1 by exhibiting two preliminary results. The first one, given in the lemma below, states that \Box_h can be viewed as the limit probability that $\Box_{;h}$ equals 1, 1 § i § n, when the radius becomes infinitely large.

Lemma 3.5 Consider the same framework as in Proposition 3.1, then for all $h Pt0; \ldots; 2^d \square 1u, i Pt1; \ldots; nu$,

$$P \square_{i;h} = 1 \square_{i} \bullet t \square_{t-1} \square_{h}.$$
(3.20)

Proof First of all, extend Q to the whole sphere by setting Q $S_{(2)}^{d_{1}} z_{2}^{d_{1}} z_{2}^{d_{1}} = 0$, then consider the following neighborhoods of each of the 2^{d} 1 open faces of the simplex: for any \circ 0, h Pt1;:::; 2^{d} 1 u and the geodesic distance d_G(.;.) on \circ d^{-1} , set

 $V_{\Box}(\Box_{h}^{d\Box 1}) := {\stackrel{!}{} ! \quad PS_{(2)}^{d\Box 1} : inf \ d_{G}(! ; x); x P \Box_{h}^{d\Box 1} ({\stackrel{)}{\$} \Box }.$

We shall prove that for all h Pt1;:::; $2^d \square 1u$,

$$\lim_{t \to 1} \mathsf{P} \stackrel{\square}{!} \mathsf{P} \mathsf{V}_{\square}(\square_{h}^{d \square 1}) \stackrel{\square}{\square} \bullet t \stackrel{\square}{=} \mathsf{Q} \stackrel{\square}{\mathsf{V}_{\square}(\square_{h}^{d \square 1}) \stackrel{\square}{:};$$
(3.21)

for an arbitrary small \Box . This result can be obtained by applying the Portmanteau theorem to Equation (3.9), provided that we find at least a decreasing sequence of positive constants \Box_1 ; \Box_2 ; ::: that tends to 0 such that for any $\mathbf{m} \cdot \mathbf{1}$ and open face $\Box_h^{d \Box 1}$, the frontier of $V_{\Box_m} (\Box_h^{d \Box 1})$ has null measure relative to Q. Since Q is a finite measure, its associated cdf admits at most countably many discontinuity sets, hence the requirement is met.

Now we shall prove that for all h Pt 1; : : : ; $2^d \square 1u$,

$$\lim_{\Box \to 0} Q V_{\Box} (\Box_{h}^{d \Box 1})^{\Box} = Q \Box_{h}^{d \Box 1}.$$
(3.22)

Observe that $Q V_{\Box}(\Box_{h}^{d \Box 1})^{\Box} = {}^{\geq}I ! PV_{\Box}(\Box_{h}^{d \Box 1})^{(Q(d!))}$; and that $V_{\Box}(\Box_{h}^{d \Box 1})$ tends to $\overline{\Box_{h}^{d \Box 1}}$ as \Box tends to 0. Therefore, Equation (3.22) can be deduced from the dominated convergence theorem, using $I ! PV_{\Box}(\Box_{h}^{d \Box 1})^{(T)} + 1$.

By combining Equation (3.21) and Equation (3.22), we obtain

$$\lim_{\omega \to 0} \lim_{t \to 1} \mathsf{P} \stackrel{\square}{!} \mathsf{P} \mathsf{V}_{\square}(\square_{h}^{d \square 1}) \stackrel{\square}{\frown} \bullet t \stackrel{\square}{=} \mathsf{Q} \stackrel{\square_{d \square 1}^{d \square 1}}{!};$$
(3.23)

for all h Pt1;:::; $2^d \square 1u$.

Now let D_j , $j \in t_0; \ldots; d \square$ 1u denote the set of indexes $h \in t_1; \ldots; 2^d \square$ 1u that identify j-dimensional open faces. We shall prove Lemma 3.5 by strong induction. First, observe that for all $h \in D_0$ we have $\square_h^{d \square 1} = \square_h^{d \square 1}$ and that the events $t \square_{i;h} = 1$; $h \in t_0; \ldots; 2^d \square$ 1uu are disjoint by construction. Hence, Equation (3.23) can be rewritten as follows:

$$Q(\square_{h}^{d\square1}) = \lim_{\substack{0 \to 0 \ t \to 1}} \lim_{\substack{i \to 0 \ t \to 1}} P_{i}^{2f\square1} P_{i}^{i} P_$$

Since Equation (3.15) ensures that $\lim_{t \to 1} P$! $PV_{\Box}(\Box_{h}^{d \Box 1}) = \bullet$ t; $\Box_{i;h} = 1 = 1$ for all \Box_{h} o and that for all \Box_{h} h, $\lim_{\Delta \to 0} \lim_{t \to 1} P$! $PV_{\Box}(\Box_{h}^{d \Box 1}) = \bullet$ t; $\Box_{i;\Sigma} = 1 = 0$, we can conclude that

$$\lim_{n \to 0} \lim_{t \to 1} P \stackrel{\square}{\hookrightarrow}_{;h} = 1 \stackrel{\square}{\longrightarrow} \bullet t \stackrel{\square}{=} \lim_{t \to 1} P \stackrel{\square}{\hookrightarrow}_{;h} = 1 \stackrel{\square}{\longrightarrow} \bullet t \stackrel{\square}{=} Q(\square_{h}^{d} \square).$$

Lemma 3.5 is thus true for all h P D₀. Now fix some J P t 1; :::; d \Box 2u and assume that it holds for all h P¹ $_{j=0}^{J}$ D_j. Set

$$F(h) \coloneqq {\overset{\circ}{\underset{\sim}{0}}} Pt1; \ldots; 2^{d} \square 1u : \square \overset{d}{\underset{\sim}{0}} P \square {\overset{d}{\underset{h}{0}}} I \overset{d}{\underset{h}{0}} J \overset{h \hbox{Heesess}}{\underset{h}{\overset{d}{\underset{h}{0}}} I \overset{d}{\underset{h}{0}} I \overset{d}{\underset{h}{} I \overset{d}{} I \overset{d}{\underset{h}{0}} I \overset{d}{\underset{h}{0}}$$

Using the same arguments as before, for all h PD_{J+1} we have:

$$Q(\boxed{a^{d}}) = \lim_{\substack{n \to 0 \\ h}} \lim_{\substack{n \to 0 \\ t \to 1}} P! PV_{a}(a^{d}) + e t; a_{i;h} = 1 P^{d}_{a;h} = 1 P^{d$$

Invoking again Equation (3.15), $\lim_{t \to 1} P \stackrel{\square}{!} PV_{\square}(\square_{h}^{d\square1}) \stackrel{\square}{=} \bullet t$; $\square_{i;h} = 1 \stackrel{\square}{=} 1$ for all $\square \circ$ 0 and for all $\square h$,

$$\lim_{n \to 0} \lim_{t \to 1} P \stackrel{\square}{!} P V_{n}(\square_{h}^{d \square 1}) \stackrel{\square}{=} t; \square_{i;} = 1 \stackrel{\square}{=} \frac{\overset{\&}{}_{1} if \ PF(h);}{\overset{0}{}_{0} if \ RF(h).}$$

Combined with the induction hypothesis, these equations entail

$$\lim_{t \to 1} P \square_{;h} = 1 \square \bullet t \square = Q(\square \square \square) \square \square \square \square \square \bullet = \square_h.$$

This concludes the proof.

The second preliminary result in the lemma below states that the distribution of vector Z_i given that $\Box_{;h} = 1$ is multivariate regularly varying when h PH.

Lemma 3.7 Consider the same framework as in Proposition 3.1, then for all h P H, there is a Radon measure \Box_h , non identically zero and not degenerate at a point, concentrated on the blunt convex cone $C_h^d := tx P C_{\Box}^d : x= x_{(2)} P \Box_h^{d \Box 1} u$, such that

$$t \mathsf{P} \stackrel{\sim}{\tilde{t}} \frac{\mathsf{Z}_{\mathsf{i}}}{\mathsf{t}} \mathsf{P} \stackrel{\sim}{\tilde{t}} \stackrel{\simeq}{\underset{t-1}{\overset{\vee}{t}}} \Box_{\mathsf{h}} (.).$$
(3.24)

Proof By Lemma 3.5, Equation (3.8) and Equation (3.13), we have that

$$t P \stackrel{\square}{!}_{i} P .; \square \bullet t \stackrel{\square}{=}_{i;h} = 1 \stackrel{\square}{=}_{t-1}^{v} S_{h} (.) \coloneqq \frac{\overset{\vee}{\&} S(. X \square_{h}^{d \square 1}) = p_{h} \text{ if } h P H;}{\overset{\otimes}{}_{0} \text{ otherwise,}} (3.25)$$

where by definition,

$$S(.X \square_{h}^{d \square 1}) = \square \stackrel{\square}{t} x PC_{\square}^{d} : \{x\}_{(1)} \bullet 1; x= \{x\}_{(2)} P.X \square_{h}^{d \square 1} u^{\square}$$

$$= \Box tx PC_{\Box}^{d} : \{x\}_{(1)} \bullet 1; x= \{x\}_{(2)} P . uX tx PC_{\Box}^{d} : x= \{x\}_{(2)} P \Box_{h}^{d} \Box_{u}^{\Box}.$$

Recall that $C_h^d := t x P C_{\square}^d : x= x_{\{2\}} P \square_h^d \square^1 u$, and set

then we can rewrite S_h in function of \Box_h as below:

$$S_h(.) = \Box_h(tx PC_{\Box}^d : x_{(1)} \bullet 1; x_{x} x_{(2)} P.u).$$

Since C_h^d is a cone, the homogeneity property of \Box stated in Equation (3.3) is passed on \Box_h , h PH. Indeed, for all $0 \dagger s \dagger 1$ and Borel subset B of C_{\Box}^d ,

$$\Box_{h} (sB) = \Box (sB) X (C_{h}^{d}) = p_{h} = \Box s(B X C_{h}^{d}) = p_{h} = s^{\Box 1} \Box B X C_{h}^{d} = p_{h} = s^{\Box 1} \Box_{h} (B).$$

According to Theorem 6.1 in Resnick (2007), it naturally follows that for all i Pt1;:::; nu,

$$t P \stackrel{\Box}{\xrightarrow{t}} P \stackrel{\Box}{\underset{i,h}{=}} 1 \stackrel{\Box}{\underset{t-1}{=}} n_h(.);$$

where, just like \Box , \Box _h can be written as the product of a measure on the radius with a measure on the angles when switching to pseudo-polar coordinates:

$$\Box_{\mathsf{h}} \Box \mathsf{T}^{\Box 1} = \Box_{\Box 1} \Box \mathsf{S}_{\mathsf{h}}.$$

We can now tackle the proof of Proposition 3.1, which is recalled below for convenience.

Proposition – Intra-class regular variation. We place ourselves in the framework of Section 3.3 and denote by H(j) the set th PH : j Pl_h u, i.e. the intersection between H and S(j).

Then, for all j Pt1;:::; du, h Pt0;:::; $2^d \square 1u, x \bullet 1$,

$$t P Z_{i;j} \circ x t = 1 Z_{i;h} = 1 C_{j;h} x^{-1};$$

where $c_{j;h}$ P[0; 1=p_h] is non-null if and only if h PH(j), and $\sum_{h=0}^{\infty} p_h c_{j;h} = 1$.

Proof Going back to the marginal level, multivariate regular variation of conditional distributions gives for all $x \cdot 1$, $1 \leq i \leq n$, $1 \leq j \leq d$,

$$\mathsf{t} \mathsf{P} \stackrel{\square}{\xrightarrow{}} \frac{\mathsf{Z}_{i;j}}{\mathsf{t}} \circ \mathsf{x} \stackrel{\square}{\xrightarrow{}}_{;h} = 1 \stackrel{\square}{\xrightarrow{}} \frac{\mathsf{v}}{\mathsf{t}-1} \square_{\mathsf{h}} (\mathsf{t} \mathsf{z} \mathsf{P} \mathsf{C}^{\mathsf{d}}_{\square} : \mathsf{z}_{j} \circ \mathsf{xu})$$

Notice that we now have a null limit for all h RH(j), i.e. the intersection between H and S(j), which identifies the star of vertex te_i u. Indeed, if te_i u is not included in the

closure of $\Box_h^{d \Box_1}$, then by definition of Q_h , S_h and \Box_h , after projection no mass is put on the j-th dimension. Furthermore, for all h PH(j), we have

$$\Box_{h} (t z P C_{\Box}^{d} : z_{j} \circ xu) = \bigcup_{\substack{d \subseteq 1 \\ \Box \ a}^{d \subseteq 1} (0; 1]} I \qquad \Box_{j! \ j(1)}^{l} \circ x \qquad \Box_{\Box 1} (d \Box) S_{h} (d!)$$
$$= x^{\Box 1} \underbrace{\frac{l \ j}{l \ j}}_{l \ j} S_{h} (d!) .$$
$$I = \bigcup_{\substack{d \subseteq 1 \\ \Box \ a}} S_{h} (d!) .$$

Hence, for all h Pt0;:::;2^d 1u, i Pt1;:::;nu, j Pt1;:::;du, x • 1, we can write

$$\mathsf{t} \mathsf{P} \overset{\Box}{\mathsf{Z}}_{i;j} \circ \mathsf{x} \mathsf{t} \overset{\Box}{\mathrel{\sqcup}}_{i;h} = \mathbf{1}^{\Box} \underset{\mathfrak{t}-1}{\overset{\Box}{\mathrel{\sqcup}}} c_{j;h} \mathsf{x}^{\Box 1};$$

where $c_{j;h} \circ 0$ when h P H(j), and $c_{j;h} = 0$ otherwise. Based on the marginal constraints on S stated in Equation (3.10) and because $t \Box_{h}^{d \Box 1} u_{0\S h\S 2^{d} \Box 1}$ forms a partition of $\Box_{h=0}^{d \Box 1}$, we have that $c_{j;h} P [0;1=p_{h}]$ for all h P t0;:::;2^d \Box 1u and $\Box_{h=0}^{2^{d} \Box 1} p_{h} c_{j;h} = \int_{hPH(j)}^{\infty} p_{h} c_{j;h} = 1$.

3.8.2 Face-characterizing functional

Before tackling its proof, Proposition 3.2 is recalled for convenience.

Proposition – Face-characterizing functional. We place ourselves in the framework of Section 3.3 and consider H (j) as in Proposition 3.1. For all j in t1;:::; du, h Pt1;:::; $2^d \square 1u, x \bullet 1$, define the functional

$$\Box_{j;h}(t) \coloneqq \int_{1}^{1} t P (\Box \bullet t) P Z_{i;j} \circ xt \Box \bullet t; \Box_{j;h} = 1^{\Box} dx;$$

and assume that there exist some constants $\Box^{\Box} P(0; 1)$, $c^{\Box} \bullet 0$ and $t^{\Box} \circ 1$, such that for all j Pt1;:::;du, h RH(j),

$$(x^{\circ} 1; (t^{\circ} t^{\circ}) \dot{O} = \frac{P Z_{i;j}^{\circ} x t}{P Z_{i;j}^{\circ} t} = 1 \\ (t^{\circ} t^{\circ}) \dot{O} = \frac{P Z_{i;j}^{\circ} x t}{P Z_{i;j}^{\circ} t} = 1 \\ (t^{\circ} t^{\circ}) \dot{O} = 1 \\ (t^$$

Then

Proof We shall handle the situations where h PH(j) and h RH(j) separately. To simplify notations, for all h Pt0;:::; $2^d \square$ 1u and x • 0 we will denote by $\overline{F}_{j;h}(x)$ the conditional probability that $Z_{i;j}$ exceeds x given $\square_{;h}$ equals 1, for any i Pt1;:::;nu:

$$\overline{F}_{j;h}(x) \coloneqq P \stackrel{\Box}{Z}_{i;j} \circ x \stackrel{\Box}{\Box}_{j;h} = 1^{\Box}.$$

From Equation (3.15) in Proposition 3.1, it is straightforward that $\overline{F}_{j;h}$ is regularly varying with index $\Box 1$, i.e. for any x • 1,

$$\frac{\overline{F}_{j;h}(x t)}{\overline{F}_{j;h}(t)} \underset{t \to 1}{\square} x^{\square 1}.$$

Hence, $\overline{F}_{i;h}$ may be written as follows:

$$\overline{F}_{j;h}(x) = x^{\Box 1} L_{j;h}(x);$$

where $L_{j;h}(x)$ is a slowly varying function $(L_{j;h} PR_0)$ that converges to $c_{j;h}$ as x - 1.

Remark 3.11 Define $x_{j;h}^{\Box} := \inf tx \cdot 1 : \overline{F}_{j;h}(x) = 0u$, the right endpoint of survival function $\overline{F}_{j;h}$ for any j Pt1;...;du and any h Pt0;...; $2^d \Box$ 1u. Then for all h PH(j), $x_{i;h}^{\Box} = +1$, that is $@ \cdot 1$, $\overline{F}_{j;h}(t) \circ 0$.

Since Bayes' formula gives

$$\begin{array}{c} \overset{1}{\overset{1}{\underset{1}{1}}} t P (\Box \bullet t) P \overset{1}{Z}_{i;j} \circ xt \overset{1}{\overset{1}{\underset{1}{1}}} \bullet t ; \Box_{;h} = 1 \overset{1}{\overset{1}{\underset{1}{0}}} dx = \\ \overset{a}{\overset{1}{\underset{1}{1}}} \underbrace{t P \overset{1}{Z}_{i;j} \circ xt ; \Box \bullet t \overset{1}{\overset{1}{\underset{1}{1}}} \underbrace{1 - 1 \overset{1}{\underset{1}{0}} \underbrace{1 - 1 \overset{1}{\underset{1}{0}}} e \underbrace{1 - 1 \overset{1}{\underset{1}{0}} e \underbrace{1 - 1 \overset{1}{\underset{1}{0}}} e \underbrace{1 - 1 \overset{1}{\underset{1}{0}} e \underbrace{1 - 1 \overset{1}{\underset{1}{0}} e \underbrace{1 - 1 & \overset{1}{\underset{1}{0}}} e \underbrace{1 - 1 & \overset{1}{\underset{1}{0}} e \underbrace{1 - 1 & \overset{1$$

Fix some $\square \circ 0$, small enough to verify $c_{j;h} \square \square \circ 0$, and some $t_\square \circ 0$ such that @t • t_\square , we have simultaneously $\square \square_{j;h} = 1 \square \square \bullet t \square \square_{h} \square \uparrow \square$ (Lemma 3.5) and $\square(t) \square c_{j;h} \square \uparrow \square$ Obviously, as soon as t • t_\square , we also have $\square(xt) \square c_{j;h} \square \uparrow \square$ for all $x \bullet 1$, and

$$0 \dagger \frac{c_{j;h} \Box \Box}{c_{j;h} + \Box} \dagger \frac{L_{j;h}(xt)}{L_{j;h}(t)}.$$

Hence, $@t \cdot t_{\Box}$,

$$\frac{a}{1} t P (\Box \bullet t) P Z_{i;j} \circ xt = 1 dx \circ$$

$$\frac{(c_{j;h} \Box)^2 p_h}{(\Box_h + \Box)(c_{j;h} + \Box)} \int_{1}^{1} x^{\Box 1} dx = +1;$$

or equivalently,

$$\int_{1}^{1} t P (\Box \bullet t) P Z_{i;j} \circ xt \Box \bullet t; \Box_{;h} = 1 dx \Box_{t-1} + 1.$$

Contrary to the case where h PH(j), we no longer have $\overline{F}_{j;h}$ PR₁. In particular, the conditional cdf can have either finite or infinite right endpoint. When its support is bounded, relying on the Bayes decomposition exhibited in the previous paragraph, the desired result is straightforward: because there exists some t₀ $^{\circ}$ 1 such that for all t $^{\circ}$ t₀, $\overline{F}_{j;h}(t) = 0$, then as t — 1, the integral also becomes null. If on the contrary, $\overline{F}_{j;h} ^{\circ}$ 0 for all t $^{\circ}$ 1, then, as previously, we can rewrite the quantity of interest in the following form:

$$\begin{array}{c} a \\ 1 \\ 1 \end{array} t P (\Box \bullet t) P \\ Z_{i;j} \circ xt \\ \hline \Box \bullet t; \\ \Box_{;h} = 1 \\ \hline dx = \\ \hline \frac{t \overline{F}_{j;h}(t) p_{h}}{P \\ \Box_{;h} = 1 \\ \hline \Box \bullet t \\ \end{array} \begin{array}{c} a \\ 1 \\ \hline \overline{F}_{j;h}(xt) \\ \hline \overline{F}_{j;h}(t) \\ \hline \overline{F}_{j;h}(t) \\ \hline \end{array} dx.$$

Since as t tends to infinity t $\overline{F}_{j;h}(t)$ tends to 0 (Proposition 3.1), P $\Box_{i;h} = t$ tends to $\Box_h \circ 0$ (Lemma 3.5) and since $p_h \circ 0$, for the integral of interest to converge to 0 it suffices to prove that there exists some $t_0 \circ 1$ such that for all t $\circ t_0$,

$$\int_{1}^{a} \frac{\bar{F}_{j;h}(xt)}{\bar{F}_{j;h}(t)} dx \dagger 1 .$$

According to the assumption in Proposition 3.2, there exists some constants \Box^{\Box} in (0; 1), $c^{\Box} \bullet 0$ and $t^{\Box} \circ 1$ such that

$$(t \circ t) \dot{O} \frac{\overline{F}_{j;h}(xt)}{\overline{F}_{j;h}(t)} \S c^{\Box} x^{\Box 1 = \Box^{\Box}}.$$

Hence, for all t $^{\circ}$ t^{\Box},

$${}^{1}_{1} \ \frac{\bar{F}_{j;h}(x\,t)}{\bar{F}_{j;h}(t)} \, dx \ \S \ c^{\Box} {}^{1}_{1} \ x^{\Box \, 1 = \Box^{\Box}} \, dx = \ \frac{\Box \, c^{\Box}}{1 \, \Box \ 1 = \Box^{\Box}} \ \dagger \ 1 \ ;$$

which produces the desired outcome.

 $\mathbf{\acute{U}}$ h PH^cztOu: p_h = $\Box_{h} = 0$

а

By definition, for all h Pt1;:::; $2^d \square$ 1u, the equivalence below holds true:

(h P H^czt 0u) Ù (
$$\Box_h = 0$$
) Ù (p_h = 0).

Consequently, when h P H ^czt 0u, we have P $Z_{i;j} \circ x = 1 = 0$ for all x • 0, and by extension

$$\int_{1}^{1} t P (\Box \bullet t) P Z_{i;j} \circ xt = 0; \quad t; \Box_{i;h} = 1 dx = 0;$$

for all t $^{\circ}$ 0. This remains true as t — 1 .

 $\mathbf{\acute{U}}$ h = 0 : p_h \Box 0 and \Box_{h} = 0

Let us start again with the following decomposition :

$$\Box_{j;0}(t) = \frac{t p_0 \bar{F}_{j;0}(t)}{P \Box_{j;0} = 1} \begin{bmatrix} a & & \\ & & \\ & & \\ \end{bmatrix}_{1}^{a} \frac{\bar{F}_{j;0}(xt)}{\bar{F}_{j;0}(t)} dx.$$

Contrary to the case where h P H zH (j), we cannot guarantee the convergence of $\Box_{;0}(t)$ to 0 as t grows to infinity, since P $\Box_{;0} = 1 \Box \bullet t$ now tends to 0 instead of a positive constant. Nonetheless, it is still possible to prove that it does not diverge to 1. Indeed, notice that

$$\frac{t p_0 \bar{F}_{j;0}(t)}{P_{i;0} = 1 \cdots t} = \frac{t P(\cdots t) \bar{F}_{j;0}(t)}{P_{i;0} = 1 \cdots t}$$

and that $\overline{F}_{j;0}(t)$ P $\square \bullet t = 1$. Hence,

$$\Box_{j;0}(t) \S t P (\Box \bullet t) \int_{1}^{a} \frac{\overline{F}_{j;0}(xt)}{\overline{F}_{j;0}(t)} dx.$$

We have already seen that according to the assumption in Proposition 3.2, there exists some constants $\Box^{\Box} P(0; 1), c^{\Box} \bullet 0$ and $t^{\Box} \circ 1$ such that for all t $\circ t^{\Box}$,

Moreover, by virtue of Equation (3.8), for all \square° 0 there exists some t \square° 0 such that for all t \degree t \square , $\stackrel{l}{\notin}$ P (\square° t) \square S($\square^{d}\square^{1}$) $\stackrel{l}{\longrightarrow}$ T is some \square° 0 and set $\square^{\circ} := \frac{\square(1 \square 1 = \square^{\circ})}{c^{\square}}$, then for all t \degree max(t \square ; t \square), we have

$$\Box_{j;0}(t) \ \ \frac{\Box S(\Box^{d \Box 1}) c^{\Box}}{1 \Box 1 = \Box^{\Box}} + \Box^{\Box} + 1 .$$

Observe that the smaller \Box^{\Box} , i.e. the faster the limit dependence structure is reached, the smaller the bound of $\Box_{j;0}(t)$. Ideally, when all $\overline{F}_{j;0}$, 1 § j § d, are rapidly varying, i.e. $c^{\Box} = 0$, we obtain the same result as in the case where h P H zH (j). This would correspond in fact to the absence of hidden regular variation, like mentioned in Section 3.5 and Section 3.7 (Resnick, 2002; Heffernan and Resnick, 2005; Resnick, 2008).

3.8.3 About \overline{p}_{i} (k)

For the sake of clarity, we give here a more explicit version of the statistic $\bar{p}_{j;}$ (k), which was defined as

$$\overline{p}_{j;}(k) \coloneqq \prod_{1}^{n} \frac{1}{k} \frac{n_k}{n^{n}} \prod_{i \in \mathbb{Q}_k} I \stackrel{!}{\stackrel{!}{\stackrel{!}{l}} \mathbf{Z}_{i;j} \circ x \frac{n}{k}; \underline{P}_{i;} = 1 \quad dx; 1 \S ` \S H; 1 \S j \S d;$$

where $n_k := #@_k$ is the number of observations the radius of which exceeds n = k, and $n := \bigcap_{i \in @_k} It P_{i;} = 1$ uthe size of class `. Recall that ` is supposed to refer to some h Pt1; :::; du, that indexed the open face $\Box_h^{d \Box 1}$.

Let us begin by considering that k is fixed, and set

$$f_{j;}(x) = \frac{1}{k} \frac{n_k}{n_k} \frac{\prod_{i \in \mathbb{Q}_k} I! \mathbf{\mathcal{P}}_{i;j} \circ x \frac{n}{k}; \mathbf{\mathcal{P}}_{i;k} = 1; 1 \ ; 1 \ S \ S \ H; 1 \ S \ J \ S \ d.$$

Our statistic of interest, $\overline{p}_{j;\hat{}}(k)$ is none other than the integral over $x \cdot 1$ of $f_{j;\hat{}}(x)$. Actually, because it relies on a finite set of n § 1 observations $f_{j;\hat{}}(x)$ is a step function with support on the interval $\min_{1 \le i \le n} \mathbf{P}_{i;j} k=n; \max_{1 \le i \le n} \mathbf{P}_{i;j} k=n$. As $f_{j;\hat{}}$ only takes into account observations verifying $\mathbf{P}_{i;\hat{}} = 1$, we denote by $(\mathbf{P}_{1;j}; \ldots; \mathbf{P}_{n;j})$ the sub-sample of $n_{\hat{}}$ observations within $(\mathbf{P}_{1;j}; \ldots; \mathbf{P}_{n;j})$, for any $j \in 1, \ldots; j$. Further consider

$$\mathbf{\mathcal{P}}_{(n \, \cdot \, \Box \, u^{\Box}; j)} = \inf \mathbf{\mathcal{P}}_{i;j}^{\tilde{}}; 1 \S i \S n^{\tilde{}} : \mathbf{\mathcal{P}}_{i;j}^{\tilde{}} \bullet \frac{n}{k} u;$$

the smallest observation $\mathbf{\mathcal{P}}_{i;j}$ that exceeds n=k, and arbitrarily set $\mathbf{\mathcal{P}}_{(n^{+} \square u^{-} \square 1;j)} = 1$, then $f_{j;}$ can be expressed as follows:

$$f_{j;}(x) = \frac{1}{k} \frac{n_k}{n} \frac{u \vec{\uparrow} \vec{\uparrow}^1}{u = 1} u \vec{\mid} x P \vec{\not{P}}_{(n \cdot u;j)} \frac{k}{n}; \vec{\not{P}}_{(n \cdot u+1;j)} \frac{k}{n}$$

In particular, when $x \cdot \mathbf{Z}_{(n_{ij},j)}^{\circ} k=n$, there is no $\mathbf{Z}_{i;j}^{\circ}$, 1 § i § n, such that $\mathbf{Z}_{i;j}^{\circ} k=n^{\circ} x$, and conversely, when $x \in 1; \mathbf{Z}_{(n \cdot \Box u^{\Box};j)}^{\circ} k=n^{\circ}$, there are exactly $u^{\Box} + 1$ observations $\mathbf{Z}_{i;j}^{\circ}$ in the sub-sample defined by $\mathbf{P}_{i;i}^{\circ} = 1$ that exceed x n=k. Therefore, the integral of $f_{j;i}^{\circ}(x)$ over all $x \cdot 1$ verifies

$$\begin{array}{c} a \\ \max_{1 \atop j \ i \ j \ n} \boldsymbol{\mathcal{B}}_{i;j}^{\circ} \ k = n \\ 1 \end{array} \quad f_{j;\gamma}(x) \ dx = \ \overline{p}_{j;\gamma}(k) = \ \frac{n_k}{n} \ \frac{1}{n^{\gamma}} \ \frac{u \ \overline{\Gamma} \ 1}{n^{\gamma}} \ u = 1 \\ u = 1 \end{array} \quad \underbrace{\boldsymbol{\mathcal{B}}_{(n \ \neg \ u + 1;j)}^{\circ} \ \Box \ \boldsymbol{\mathcal{B}}_{(n \ \neg \ u + 1;j)}^{\circ} \ \Box \ \boldsymbol{\mathcal{B}}_{(n \ \neg \ u + 1;j)}^{\circ} \end{array}$$

Let us dwell for a moment on this expression. The part

$$\frac{1}{n} \frac{u \vec{\Gamma} \mathbf{1}}{u = 1} u \mathbf{2}_{(n \cdot \Box u + 1;j)} \Box \mathbf{2}_{(n \cdot \Box u;j)}$$

represents the integral under the empirical survival function of variable $\mathbf{P}_{i;j}$ conditional on i being in cluster ` and \mathbf{P}_i ₍₁₎ • n=k. When ` P H(j), there should be a lot of extreme observations $\mathbf{P}_{i;j}$ in cluster `, and this quantity should be very large. Conversely, in all clusters `¹ R H(j), there should be very few to no extreme values on the j-th dimension, and the corresponding integral should be very small. Figure 3.6 and Figure 3.7 give an illustration of this phenomenon on exposures to the 6 nutrients and contaminants investigated in Section 3.6. Notice that dividing by n. enables comparison between classes and avoids systematically selecting poor classifications. In a similar way, the term $n_k = n$ penalizes small values of k, which would otherwise always be preferred to higher ones and provide non-explicable groups, in the sense that they would contain too few observations to be interpreted. In terms of bias-variance compromise, intuitively it would generate overly wide variances for the final estimates of $t \square_h^{d\square1}$; h PH uto be reliable.



Figure 3.6 – Log-scaled marginal distributions of exposures to the 6 chemical elements within clusters 1, 2, 3 and 4, obtained for the couple $(k; \mathbf{P}(k))^{\Box}$ as defined in Section 3.4: distributions of contaminants with extreme exposures are displayed in red with black contours, while the others are white with grey contours. The thin horizontal line indicates k^{\Box}



Figur e 3.7 – Log-scaled marginal distributions of exposures to the 6 chemical elements within clusters 5, 6 and 7, obtained for the couple $(k; P(k))^{\Box}$ as defined in Section 3.4: distributions of contaminants with extreme exposures are displayed in red with black contours, while the others are white with grey contours. The thin horizontal line indicates k^{\Box}

3.8.4 A little more on PNS and spherical k-means

As a complement to the succinct description of Principal Nested Spheres and spherical k-means in Section 3.4.2, we provide here a more detailed overview of these algorithms accompanied by illustrative figures. With these additional specifications we are then able to discuss some technical choices in the implementation that were only briefly mentioned in the core of the chapter.

3.8.4.1 Principal Nested Spheres

Recall from Definition 3.4 that the geodesic distance between two points x and y of S^{d_1} (the unit sphere in \mathbb{R}^d) is written

$$d_G(x; y) = \arccos^1 y;$$

where x¹ stands for the transpose of vector x. Now consider any (d \square 2)-dimensional sub-sphere A_{d \square 2} in S^{d \square 1}. Relative to the geodesic distance on the sphere, its center and radius are respectively a point v PS^{d \square 1} and a distance r P(0; \square =2] such that

$$A_{d \square 2} := A_{d \square 2}(v; r) := x P S^{d \square 1} : d_G(v; x) = r^{(}.$$

Given this representation, the signed distance between any point x P S^{d \square 1} and a sub-sphere A_{d \square 2}(v;r) is naturally defined as

$$d_{S}(x; A_{d \square 2}(v; r)) \coloneqq d_{G}(v; x) \square r.$$

Equipped with these tools, the main steps of the PNS algorithm can be depicted as follows.




 $A_{0} :=$ $\underset{x \in S^{1}}{\operatorname{argmin}} \int_{i=1}^{n} d_{G}^{2}(! \quad (d \square 2); x);$

the Fréchet mean of the data.



The computational algorithm designed to solve the least squares problem that defines each sub-sphere can be found in Section 3 of Jung et al. (2012) and the explicit formulas corresponding to the successive transformations of each PNS in Section 2. Many more details are provided in the supplementary materials associated with this seminal paper, concerning in particular the geometry of PNS. More importantly, a penalized version of the initial procedure is proposed to decide at each step whether small sub-spheres are more relevant candidates than those with maximal radius. Actually, both the numerical experiments in Section 3.5 and the case study in Section 3.6 are using this refined version of the PNS algorithm. We refer to the end of Section 1 (p.7) in the aforementioned supplementary materials for more details on the subject.

At the end of the procedure, we obtain a collection of unit spheres with dimensions ranging from 1 (the unit circle) to d \Box 1 (the space of the original data), which can be understood as a spherical equivalent of the principal components in PCA. In order to proceed with the rest of the analysis, we need to choose one of these d \Box 1 PNS and work with the corresponding projected angles. Recall that for all j Pt1;:::;du and i Pt1;:::;n_ku,! $i_{j}^{(j \Box 1)}$ denotes the projection of angle! i_{j} on S^{d \Box j} and define

$$\Box_{i}^{(d \square j)} := \int_{i=1}^{j_{1} + 1} \sin r \cdot d_{G} \downarrow_{i}^{(j)}; A_{d \square j}(v_{j}; r_{j});$$

the corresponding scaled residual. In short, scaling enables comparison between the deviations as if they were all measured on S_2^1 (see Jung et al., 2012, Sections 2.1 to 2.4). Then the relative variance encapsulated by PNS d \Box j, 1 § j § d, is understood as

$$V_{d \square j} := \frac{ \begin{pmatrix} \infty & n_k & & & \\ i = 1 & & & \\ \hline & & & \\ &$$

Given these notations, the selection heuristic evoked in Section 3.4.2 simply consists in picking the smallest sub-sphere $S^{d \ j}$ such that $V_{d \ j \ 1} \bullet 0.1$, $V_{d \ j} \bullet 0.1$ and $V_{d \ j + 1} \dagger 0.1$. Observe that contrary to PCA, there is no obvious link here between

 V_1 ;:::; $V_{d \ \Box 1}$ and the variance of the angles V (!), ! PR^d . Obviously, many refinements could and should be brought to our method in the near future, starting with a more adaptive way of identifying the "optimal" PNS. Moreover, though the situation was never encountered in our applications, many practical difficulties can arise and should receive appropriate attention. For instance it can so happen that there is no unique Fréchet mean: imagine a cloud of two points (0; \Box 1) and (0; 1), then there are two possible candidates for A_0 , namely (\Box 1; 0) and (1; 0). These intricate issues are left for future research.

3.8.4.2 Spherical k-means

Once a selected PNS, clustering is achieved using the spherical k-means algorithm. Before getting into detail, let us introduce a few additional notations. First, denote by $\binom{1}{1}$;...; $\binom{1}{n_k}$ the cloud of angles projected onto the optimal PNS. For some fixed number of clusters H, the objective is to estimate the n_k vectors of class indicators \Box_1 ;...; \Box_{n_k} , where @ P t1;...; n_k u, $\Box_i := (\Box_{;1};...;\Box_{;H})$. In spherical k-means, clusters are represented by their barycenter: for any set of $n \cdot 1$ points x_1 ;...; x_n on the unit sphere S^{d $\Box 1$} Ä R^d such that for all i Pt1;...; n_u , $x_i := (x_{i;1};...;x_{i;d})^1$, the barycenter function is written

$$B(x_{1}; \ldots; x_{n}) := \frac{1}{n} \prod_{i=1}^{n} x_{i;1}; \ldots; \frac{1}{n} \prod_{i=1}^{n} x_{i;d};$$

and its projection on the unit sphere

$$SB(x_1;\ldots;x_n) \coloneqq \frac{B(x_1;\ldots;x_n)}{B(x_1;\ldots;x_n)}$$

For simplicity, we denote by b_h the barycenter of all angles in class h, i.e.

and by ch its projection on the unit sphere, also called the concept vector:

$$c_h := SB \stackrel{\square}{!} \stackrel{:}{}_1 I t \square_{;h} = 1u; \dots; \stackrel{:}{!} \stackrel{:}{}_{n_k} I t \square_{;h} = 1u \stackrel{\square}{=} \frac{b_h}{\}b_h\}_2.$$

Given these notations, the spherical k-means algorithm tries to find the collection of indicators $\mathbb{P}_{1,h}$, 1 § i § n, 1 § h § H, that minimize the intra-class geodesic variance, namely

$$\begin{array}{c} \square \\ \text{GV} \quad t \square_{i;h} \underbrace{u_{1 \text{sisn}}}_{1 \text{shsH}} \stackrel{\square}{\Rightarrow} \stackrel{\text{I}}{=} \underbrace{I_{i}}_{h=1} \stackrel{\text{I}}{=} d_{G}^{2}(! \stackrel{i}{:};c_{h}) \mid t \square_{i;h} = 1u;$$

by implementing the following basic steps.





The option "S" in the R package skmeans that we used in our applications stipulates how the initial concept vectors $c_1^{(0)}$; :::; $c_H^{(0)}$ are to be chosen. Specifically, for H • 2 desired clusters, they are successively selected according to the following recurrence relation:

In words, the first initial concept vector is set as the Fréchet mean of the sample of projected angles, and the rest as the H \Box 1 observations farthest away from all already picked concept vectors. This produces an initial clustering with centers as scattered as possible. Obviously, many other initialization techniques may have been applied, e.g. picking the first concept vectors at random. The present version has been chosen in accordance with our belief that after projection, the angles corresponding to different faces are disseminated on different regions of the retained PNS while those belonging to the same face are concentrated on a common neighborhood.

The main advantage of the spherical k-means algorithm is that it is very simple to implement. However, it can often happen that it remains stuck at a local minimum of

the intra-class geodesic variance function. To counteract this undesirable effect, many refinements have been proposed in the literature (see for instance Dhillon et al., 2002 and the references therein). As a first go, we confined ourselves to the basic version of spherical k-means, leaving further considerations on algorithmic improvement for future research.

4

A MINIMUM VOLUME SET APPROACH TO DIETARY RISK-BENEFIT ANALYSIS

In the same spirit as in Chapter 3, we propose an alternative method to inspect the multivariate distribution of the exposure to multiple food chemicals, which accounts for the variability of the contamination of foodstuffs. Directly inspired from typical statistical learning techniques, this non-parametric approach is no longer focused on extreme events. One of its advantages is that it can be very naturally extended to the identification of optimal dietary habits, in the sense that they realize a compromise between toxicological risk and nutritional benefit. From a practical point of view, such results would facilitate public communication of general dietary recommendations. The present chapter corresponds to a paper currently being written in collaboration with JTressou (INRA Met@risk, France) and S. Clémençon (Télécom ParisTech, France). It is not complete yet, in particular we are still working on the practical applications of the theoretical results introduced therein.

It is the major purpose of this chapter to show how to adapt recent (unsupervised) machine-learning techniques, specifically introduced to deal with very highdimensional data in a non-parametric manner (avoiding thus the curse of dimensionality), to food risk/ benefit analysis. Precisely, a variant of the minimum-volume set methodology (MV-set in abbreviated form), originally investigated in Polonik (1997) (see also Scott and Nowak, 2006), is proposed in order to determine confidence or predictive regions for the joint dietary exposure to a variety of chemicals and nutrients present in the food. The main originality of this problem lies in the fact that a natural empirical counterpart of the dietary exposure is of the form of a generalized U-statistic, based on the combination of consumption survey data with a database gathering measures of contents of a variety of chemicals and nutrients in most food items. Mainly due to the large number of foodstuffs involved in the observed diets, this statistic cannot be calculated in general, its computation requiring to sum over a prohibitive number of terms. Following in the footsteps of Bertail and Tressou (2006), we replace the latter with an incomplete U-statistic (Blom, 1976) the computation of which is numerically feasible, and we establish a novel uniform deviation result, which shows that this approximation stage does not damage the learning rate of the MV-set procedure. Next, similar concepts and results are applied to identify regions where the multivariate distribution of dietary habits is mostly concentrated and where types of exposure simultaneously remain within toxicological values of reference (limitation of the risk) and recommended dietary allowances (preservation

of the benefit) with maximum probability. Statistical results, i.e. rate bounds guaranteeing the performance of the generic learning techniques we propose, are stated.

The chapter is structured as follows. A detailed account of the statistical issues related to food risks and benefits tackled in the chapter and the learning methods proposed to deal with them is given in Section 4.1, together with theoretical results claiming their validity. Technical proofs are postponed to Section 4.2.

4.1 theoretical analysis and methods

4.1.1 Assessment of dietary exposure to chemicals and nutrients: a MV-set formulation

We are concerned here with dietary types of exposure to d • 1 different food chemicals, contaminants or nutrients, over a certain statistical population of interest during a given period of time, say a week like in the INCA2 database. Foods are classified according to some given nomenclature that accounts for H families of products, indexed by label h Pt1;:::;Hu. The joint dietary exposure can be then described by a random vector X := $(X_1; :::; X_d)$, where

$$X_j := \prod_{h=1}^{[n]} C_h \square Q_{h;j}.$$
(4.1)

for 1 § j § d, denoting by C_h the quantity of food item h consumed per week by an individual drawn at random in the studied population and by $Q_{h;j}$ the (random) content related to food item h and component j. In the field of food safety, risk assessors are interested in building confidence/ predictive regions for the exposure X in R_+^d :

 $R_{\Box} := \operatorname{arg\,min} L(R) : P(XPR) \bullet \Box; R\ddot{A} R_{+}^{d} \operatorname{Borelian}'; \quad (4.2)$

where $\Box P(0; 1)$ and the Lebesgue measure on R^d_+ is denoted by L. For values of the level \Box dose to 1, such minimum volume sets (MV-sets in short, see Scott and Nowak, 2006) describe regions where the distribution of exposure is most concentrated, those lying in their complementary sets being possibly interpreted as "abnormal". The construction of predictive/ confidence regions for the dietary exposure is based on the observation of the dietary habits of $n \cdot 1$ individuals independently drawn from the population, yielding an iid sample t $C_i := (C_{i;1}; \ldots; C_{i;H}); 1 \S i \S nu$. They are combined with databases where $m_{h;j} \cdot 1$ iid measures of the amount of pollutant or nutrient j present in food h are listed for all h Pt1; \ldots ; Hu and j Pt1; \ldots ; du. The corresponding vectors of contents are denoted by $Q_{h;j} := (Q^1_{h;j}; \ldots; Q^{m_{h;j}}_{h;j})$. Usually, the following hypotheses are supposed to hold true.

Assumption 4.1 For all couples $(h; \tilde{h}) Pt1; :::; Hu^2$ such that $h \Box \tilde{h}$ and any j Pt1; :::; du, level $Q_{h;j}$ is independent from $Q_{\tilde{h};j}$.

Assumption 4.2 For all couples $(j; \tilde{j}) P t 1; \dots; du^2$ such that $j \square \tilde{j}$ and any h Pt1; \dots ; Hu, level $Q_{h;j}$ is independent from $Q_{h;\tilde{j}}$.

Assumption 4.3 For any (h;j) Pt1;:::;Hu \Box t1;:::;du, level $Q_{h;j}$ is independent from consumption vector C_i .

Assumption 4.1, which formally states that nutrients and contaminants are independently assimilated by foodstuffs, could easily be relaxed without any substantial impact on the present approach. Its applicability depends mostly on the level of detail provided by the database(s) at hand. On the contrary, Assumption 4.2 and Assumption 4.3 are necessary. They respectively stipulate that the contents of one foodstuff are independent from that of others and that dietary habits are not dictated by nutritional or toxicological characteristics of the food. From a more practical point of view, Assumption 4.2 is likely to be true if consumers get their supplies from various productions and Assumption 4.3 if they do not base their consumption decision on a systematic scrutiny of the composition of the food.

Based on these data, the probability involved in the constraint of the MV-set problem stipulated in Equation (4.2) is estimated by

$$\mathbf{P}(X P R) := \frac{1}{\Box} \begin{bmatrix} \mathbf{n} & \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{i=1}^{i_{1}} & \mathbf{n} \\ \vdots \\ \mathbf{n} \end{bmatrix}_{i=1}^{i_{1}} \\ \vdots \\ \mathbf{n} \end{bmatrix}_{H;d=1}^{i_{1}} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{H;d=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \begin{bmatrix} \mathbf{n} \\ \mathbf{n} \end{bmatrix}_{h=1}^{i_{1}} \\$$

with

and I t.u the indicator function. In practice, all types of sets cannot be explored; the search is restricted to a class R of Borelian sets, the complexity of which is controlled, like hypercubes or ellipses (see the next subsection). The level \Box is in turn replaced by $\Box \Box$, where \Box is some tolerance level that depends, roughly speaking, on the order of magnitude of \Box

Hence, one should ideally try to solve the constrained minimization problem:

$$\min_{\mathsf{RPR}} \mathsf{L}(\mathsf{R}) \text{ subject to } \mathbf{P}(\mathsf{X} \mathsf{P} \mathsf{R}) \bullet \Box \Box \Box.$$
(4.4)

The major difference with the formulation in Scott and Nowak (2006) lies in the fact that the estimate of the probability involved in the mass constraint is here of the form

of a (generalized) U-statistic. Thus, the study of the performance of solutions of Equation (4.4) includes the proof of concentration results for U-processes (i.e. collections of U-statistics). Unfortunately, averaging over the n $\Box \stackrel{\pm}{}_{j=1}^{d} \stackrel{\pm}{}_{h=1}^{H} m_{h;j}$ terms appearing in Equation (4.3) is generally numerically unfeasible, even for moderate sample sizes. In Bertail and Tressou (2006) for instance, where the estimation of the probability that the exposure to Ochratoxin A exceeds a critical threshold is considered, this corresponds to $4 \Box 10^{21}$ terms! As shall be seen below, the statistic in Equation (4.3) can be uniformly approximated by a "Monte-Carlo" version, the computation cost of which is drastically reduced.

Remark 4.4 – Alternative approaches. We underline that the MV-set methodology is by no means the sole possible approach for constructing predictive regions. For instance, density sub-level sets can be built by means of non-parametric density estimation techniques, see Tsybakov (1997) and the references therein. However, when trying to implement such "plug-in" alternatives, even for moderate dimensions, one faces significant computational problems inherent to the curse of dimensionality. This motivates the machine-learning approach promoted here, which avoids a preliminary density estimation stage, while focusing directly on performance optimization. One may also refer to Vert and Vert (2006) or Steinwart et al. (2005) for closely related techniques.

4.1.2 Uniform approximation of generalized U-statistics by their incomplete versions

For clarity, we recall below the definition of generalized U-statistics. Properties and asymptotic theory of U-statistics can be found in Lee (1990).

Definition 4.5 – Generalized U-statistic. Let $K \bullet 1$, $(d_1; \ldots; d_K) P N^{\Box K}$ and consider the vectors $X_1^{(k)}; \ldots; X_{n_K}^{(k)}$, $1 \S k \S K$, corresponding to K independent samples of iid random variables, taking their values in some space X_k with distribution P_k respectively. The generalized (or K-sample) U-statistic of degrees $(d_1; \ldots; d_K)$ with kernel $X_1^{d_1} \Box \Box X_K^{d_K} - R$, square integrable with respect to the probability distribution $P_1^{b d_1} b \Box b P_K^{b d_K}$, is defined as

$$U_{n}() := \frac{1}{\frac{K}{k=1}} \prod_{\substack{n_{k} \\ k = 1}} \prod_{\substack{n_{k} \\ k = 1}} \prod_{\substack{n_{k} \\ k = 1}} (X_{l_{1}}^{(1)}; X_{l_{2}}^{(2)}; \dots; X_{l_{K}}^{(K)});$$
(4.5)

where ${}^{\infty}_{l_k}$ refers to the summation over all ${}^{n_k}_{d_k}$ subsets $X^{(k)}_{l_k} := X^{(k)}_{i_1}; \ldots; X^{(k)}_{i_{d_k}}$ related to a set l_k of d_k indexes 1 § $i_1 + \cdots + i_{d_k}$ § n_k . It is said symmetric when is permutation symmetric in each set of d_k arguments $X^{(k)}_{l_k}$.

Observe that the functional in Equation (4.3) corresponds to a K-sample U-statistic of degrees (1; 1; :::; 1), with $K = d \Box H + 1$ and kernel given by:

$$\begin{array}{c} & & & \\ & & & \\ & & & \\ R(c;q) := I & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & & \\ & & & \\$$

for c := $(c_1; \ldots; c_H) P R_+^H$ and q := $(q_{h;j}; \ldots; q_{H;d})_{j=1;\ldots;d} P R_+^{H \square d}$. Beyond this example, many statistics used for pointwise estimation or hypothesis testing are actually U-statistics (e.g. the sample variance, the Gini mean difference, the Wilcoxon Mann-Whitney statistic, the Kendall tau). Their popularity mainly arise from their "reduced variance" property: the statistic $U_n()$ has minimum variance among all unbiased estimators of the parameter

$$\Box() := \mathsf{E}^{ \square} (\mathsf{X}_{1}^{(1)}; \ldots; \mathsf{X}_{d_{1}}^{(1)}; \ldots; \mathsf{X}_{1}^{(K)}; \ldots; \mathsf{X}_{d_{k}}^{(K)})^{ \square}.$$

Classically, the limit properties of these statistics (law of large numbers, central limit theorem, etc.) are investigated in an asymptotic framework stipulating that, as the full sample size

$$n \coloneqq n_1 + \square + n_K$$

tends to infinity, we have $n_k = n - \Box_k \circ 0$ for all k P t1;:::;Ku. They can be established by means of a linearization technique (see Hoeffding, 1948), permitting to write $U_n()$ as a sum of K basic sample mean statistics (of the order $O_P(1=?n)$) each, after recentering), plus possible degenerate terms called degenerate U-statistics.

As mentioned before, in practice, the number $\stackrel{\pm}{k}_{k=1}^{K} \stackrel{n_k}{d_k}^{n_k}$ of terms to be summed up to compute Equation (4.5) is generally prohibitive. As a remedy to this computational issue, in the seminal contribution of Blom (1976), the concept of incomplete generalized U-statistic has been introduced: the summation in Equation (4.5) is replaced by a summation involving much less terms, extending only over low cardinality subsets of the $\stackrel{n_k}{d_k}^{n_k}$ d_k-tuples of indexes, 1 § k § K. In the simplest formulation, the subsets of indexes are obtained by sampling with replacement, leading to the following definition.

Definition 4.6 – Incomplete Generalized U-statistic. Let B • 1. The incomplete version of the U-statistic in Equation (4.5) based on B terms is defined by:

$$\mathbb{U}_{B}(\) \coloneqq \frac{1}{B} \frac{||}{||_{1}; \dots; ||_{K}||_{PD}|_{B}} (X_{||_{1}}^{(1)}; \dots; X_{||_{K}}^{(K)});$$
(4.6)

where D_B is a set of cardinality B built by sampling with replacement in the set

$$L := \begin{pmatrix} \# \\ (i_1^{(1)}; \dots; i_{d_1}^{(1)}); \dots; (i_1^{(K)}; \dots; i_{d_K}^{(K)}) \end{pmatrix}_{\substack{1 \le i_1^{(K)} \ddagger \dots \ddagger i_{d_k}^{(K)} \le n_k; \ 1 \le k \le K}}^{+}$$

Remark 4.7 – Alternative sampling schemes. We point out that, as proposed in Janson (1984), other sampling schemes could be considered, in particular sampling without replacement or Bernoulli sampling. The results of this chapter could be extended to these situations. For the sake of simplicity, we restrict our attention here to the sampling with replacement scheme.

 $V U_{B}() = 1 \Box \frac{1}{B} V (U_{n}()) + O \frac{1}{B} as B - +1;$

refer to Lee (1990, p.193). Incidentally, we underline that the empirical variance of Equation (4.5) is not easy to compute since it involves summing approximately #L terms and bootstrap techniques should be used for this purpose, as proposed in Bertail and Tressou (2006). The asymptotic properties of incomplete U-statistics have been investigated in several articles, see Janson (1984); Brown and Kildea (1978); Enqvist (1978). The angle embraced in the present chapter is of different nature: the key idea we promote here is to use incomplete versions of collections of U-statistics in learning problems such as that described in Section 4.1.1. The following result shows that this approach solves the numerical problem, while not damaging the learning rates. It reveals that, under adequate complexity assumptions on the considered collection \Box of (symmetric) kernels (refer to Dudley, 1999), concentration results established for U-processes (i.e. collections of U-statistics) may extend to their incomplete versions.

Theor em 4.8 – Maximal deviation. Let \Box be a collection of bounded symmetric kernels on $\Box := {}^{\pm} {}_{k=1}^{K} X_{k}^{d_{k}}$. Suppose that \Box is a VC major class of functions with finite Vapnik-Chervonenkis dimension V and that $M_{\Box} := \sup_{(x) \in X} |x| + 1$.

Then, the following assertions hold true.

i) For all \square° 0, we have: $@n = (n_1; \dots; n_K) PN^{\square K}$, $@B \bullet 1$,

P sup
$$\mathbb{U}_{B}() \square U_{n}() \stackrel{!}{\square} \square \$ 2(1 + \#L)^{\vee} e^{\square B \square^{2} = M \square^{2}}$$

ii) For all $\Box P(0; 1)$, with probability at least $1 \Box \Box$, we have: $@_{1k} \bullet 1, 1 \S k \S K$,

$$\sup_{P_{-}} \underbrace{\mathbb{U}_{B}(\)}_{P_{-}} = \underbrace{\mathbb{U}$$

where

 $\Box := \min t \mathfrak{m}_1 = d_1 \mathfrak{q} : ::; \mathfrak{m}_K = d_K \mathfrak{u}$

and txudenotes the integer part of any real number x.

We refer to Section 4.2 for the proof. Observe that, with the asymptotic settings previously specified, \Box n and log(#L) \Box log(n) as n — +1. The bounds stated above show that, for a number B := B_n of terms tending to infinity as n — +1 at a rate O(n), the maximal deviation sup $_{P\Box} |U_B() \Box \Box |$ is asymptotically of the order $O_P(\stackrel{a}{\Box og(n)=n})$, just like sup $_{P\Box} |U_n() \Box \Box |$. Remarkably, except in the case K = 1 and d_K = 1 solely, using such incomplete U-statistics thus yields a significant gain in terms of computational cost and fully preserves the order of the probabilistic upper bounds for the uniform deviations. Before showing how these results apply to the analysis of the dietary exposure distribution, a few remarks are in order.

Remark 4.9 – On the complexity assumption. We point out that, as can be seen by examining their proof in Section 4.2, the results above could be extended to other complexity measures than the VC dimension, such as Rademacher averages (Boucheron et al., 2005). However, the confidence regions we shall consider in practice being of the form of the union of a limited number of hypercubes, the finite VC dimension hypothesis is sufficient to provide a validity framework to the present analysis.

Remark 4.10 – Learning tasks based on optimization of U-statistics. The uniform deviation result stated in Theorem 4.8 can be proved very useful much beyond the framework described in Section 4.1.1. Indeed, statistical learning problems where the empirical performance criterion one seeks to optimize is of the form of a (generalized) U-statistic have recently been the subject of a good deal of attention in the literature: supervised ranking (Clémençon et al., 2008), learning on graphs (Biau and Bleakley, 2006) or dissimilarity-based clustering (Clémençon, 2011) for instance. The result above permits to show that, in such problems, the empirical criterion can be replaced by an incomplete version of much simpler computation, with only a slight impact on the learning rate, provided that the parameter B is suitably chosen.

Let us now come back to the specific learning task formulated in Section 4.1.1. In order to avoid the computation of Equation (4.3), which involves summing over $\Box := n \stackrel{\pm}{}_{j=1}^{d} \stackrel{\pm}{}_{h=1}^{H} m_{h;j}$ terms based on N := n + $\stackrel{\infty}{}_{j=1}^{d} \stackrel{\infty}{}_{h=1}^{H} m_{h;j}$ observations, Theorem 4.8 suggests to draw with replacement B times in the set of indexes

$$L := t 1; :::; n u \square \prod_{j=1}^{T^{ej}} t 1; :::; m_{h;j} u; L = \square;$$

yielding an index set D_B of cardinality B. For any borelian R Ä R^d, the probability that the dietary exposure lies in the region R is estimated by the incomplete U-statistic:

$$\mathbf{P}_{B} (X P R) := \frac{1}{B} \frac{\prod_{(i; \hat{h}; j) P D_{B}} \prod_{k=1}^{\infty} \frac{1}{k} \sum_{h=1}^{\infty} \frac{1}{C_{i; h} Q_{h; j}^{\hat{h}; j}} \sum_{1 \le j \le d} \frac{1}{P R}$$

Suppose that the class R is of finite VC dimension V. Let $\Box P(0; 1)$ be the target mass and $\Box P(0; 1)$ the desired confidence level. Notice that in the present case,

 $\Box = \min n; \min m_{h;j} : 1 \S j \S d; 1 \S h \S H^{(()}$

Further define the complexity penalty by

$$\Box(\mathsf{B};\mathsf{N};\Box) \coloneqq 2^{\mathsf{C}} \frac{\overline{2 \vee \log(1+\Box)}}{\Box} + {}^{\mathsf{C}} \frac{\overline{\log(2=\Box)}}{\Box} + {}^{\mathsf{C}} \frac{\overline{\vee \log(1+\Box) + \log(4=\Box)}}{\mathsf{B}}; \quad (4.7)$$

and consider the solution \mathbf{R}_{\Box} of the constrained optimization problem:

$$\min_{\mathsf{RPR}} \mathsf{L}(\mathsf{R}) \text{ subject to } \overset{\mathsf{P}}{\mathsf{P}} (\mathsf{X} \mathsf{P} \mathsf{R}) \bullet \Box \Box \Box (\mathsf{B}; \mathsf{N}; \Box).$$
(4.8)

The result below shows that, if the number B of exposure values computed through the sampling scheme and \Box are of the order O(N), the performance of \mathbf{R}_{\Box} is then comparable to that of the region the selection of which is based on the quantities in Equation (4.3) (see Section 4.2 for a sketch of proof).

Corollary 4.11 For all
$$P(0; 1)$$
, we have with probability at least 1 \square
 $L \stackrel{\square}{\mathbf{R}}_{\square}$ $\stackrel{\square}{\$}$ $\inf_{\text{RPR}: P(XPR)^{\bullet}}$ $L(R)$ and $P \stackrel{\square}{X} P \stackrel{\square}{\mathbf{R}}_{\square}$ \bullet $\square 2 \square (B; N; \square).$

From a practical perspective, the constrained optimization program in Equation (4.8) will be here performed over a collection of sets obtained as the union of small hypercubes by the means of a simplistic sorting algorithm, detailed at length in the next subsection. Alternative techniques, whose implementation is less immediate, such as those based on Dyadic Decision Trees (DDT) could be also considered, see Scott and Nowak (2006, Subsection 6.2) for further details.

4.1.3 Empirical MV-set estimation based on hypercubes

We now turn to the issue of solving Equation (4.8) from a practical perspective. For simplicity, such regions are built by binding together hypercubes of the positive orthant R_{+}^{d} . Suppose that the exposure X takes its values in the compact set $[0; 1]^{d}$, even if it means dividing each component of the exposure random vector by a supposedly finite essential upper bound. Observe that this assumption makes the present approach quite different from that introduced in Chapter 3, where the support of the distribution of individual types of exposure could extend to infinity. Since extremes are no longer of interest here, such refinements are of little concern. Now let $k \cdot 1$ and consider the partition of the unit cube $\frac{t}{j} \frac{d}{j=1}[s_j=k; (s_j + 1)=k]$ made of sub-cubes of side length 1=k, with $k \cdot 1$ and $s_j P t 0; \ldots; k \Box 1 u$ for $1 \leq j \leq d$. Denote by $C_1; \ldots; C_M$ these sub-cubes, with $M = k^d$, and consider the collection R_k of subsets obtained as the union of such cubes. Observe that R_k is of finite cardinality: $\#R_k = 2^{k^d}$. In this situation, the theoretical results established in the previous subsection apply, with the penalty

$$\Box_{k}(\mathsf{B};\mathsf{N};\Box) \coloneqq 2^{\mathsf{C}} \frac{\overline{2k^{\mathsf{d}} \log(1+\Box)}}{\Box} + {\mathsf{C}} \frac{\overline{\log(2=\Box)}}{\Box} + {\mathsf{C}} \frac{\overline{k^{\mathsf{d}} \log(1+\Box) + \log(4=\Box)}}{\mathsf{B}};$$

and a solution of the constrained minimization problem in Equation (4.8) can be obtained in two steps, as follows. Let \Box , \Box in (0; 1) and B • 1.

1. Sort the sub-cubes C_m , 1 § m § M , so that:

$$\mathbf{P}_{B}(X \mathsf{PC}_{1;M}) \bullet \square \bullet \mathbf{P}_{B}(X \mathsf{PC}_{M;M}).$$

 Bind together the cubes sequentially, until the incomplete U-statistic estimating the mass of the resulting set exceeds □ □ _k (B; N; □), yielding the region:

$$\mathbf{\hat{R}}_{k;\Box} := \sum_{m=1}^{N_{s}} C_{m;M};$$

where

$$\overset{\#}{\mathsf{M}} = \min \overset{\#}{\mathsf{M}} \bullet 1: \overset{\mathsf{M}}{\underset{m=1}{\overset{\mathsf{P}}{\mathsf{B}}}} (\mathsf{X} \mathsf{P} \mathsf{C}_{\mathsf{m};\mathsf{M}}) \bullet \Box \Box_{\mathsf{k}} (\mathsf{B};\mathsf{N};\Box)$$

Let 1 § k_{\Box} † k^{\Box} † +1. The issue of selecting the "resolution level" k Ptk_;:::; k^{\Box}u automatically can be handled through complexity penalization, as shown in Scott and Nowak (2006, Section 4). Precisely, one should pick the value

$$\mathbf{R} \coloneqq \underset{\substack{k \subseteq S \ k \subseteq S \ k}}{\operatorname{argmin}} \stackrel{!}{\overset{!}{\underset{k \subseteq S \ k \subseteq S \ k}}} - \left[\begin{array}{c} \overset{\;}{\underset{k \subseteq S \ k \subseteq S \ k}} - \begin{array}{c} \overset{\;}{\underset{k \subseteq S \ k \subseteq S \ k}} - \begin{array}{c} \overset{\;}{\underset{k \subseteq S \ k \subseteq S \ k \subseteq S \ k}} - \begin{array}{c} \overset{\;}{\underset{k \subseteq S \ k \subseteq S \ k$$

minimizing thus a complexity-penalized version of the volume, in order to approximate the MV-set over ${}^{I}_{k} R_{k}$ as accurately as possible without overfitting the data. An oracle inequality showing that the chosen set $R_{R;\square}$ corresponds to an optimal trade-off between excess volume and missing mass can be straightforwardly derived from the analysis carried out in Section 4.1.2, just like in Scott and Nowak (2006, Theorem 7). Details are left to the reader.

4.1.4 Optimal regions in the consumption space in regard to dietary risks and benefits

In food safety, one topic of crucial interest is to determine optimal regions in the dietary consumption space R_{+}^{H} in the sense that they achieve a compromise between toxicological risk and nutritional benefit. Indeed, for each nutrient and contaminant, experts define maximum threshold levels of exposure, generally called dietary intake limits (DIL), above which health issues due to excessive supply of chemicals are likely to occur. Equivalent lower bounds are defined for nutrients, which indicate nutritional deficiency. From a public health point of view, it is then of particular interest to identify the dietary habits that offer the best chance of respecting both lower and upper DIL and draw general, easily comprehensible dietary guidelines to the concerned population, in the same spirit as "eat at least five fruits and vegetables a day". Since it is very difficult to recommend people to completely alter their habits, public health institutes are more interested in determining realistic food baskets, already consumed by a non-negligible amount of persons, which they could set as examples for the rest of the population. This problem can be addressed by means of concepts and tools very similar to those investigated in Section 4.1.1. Consider d • 1 pollutants and nutrients, present in a nomenclature of foodstuffs indexed by h Pt1;:::;Hu at (random) concentration levels $Q_{h;1}$;:::; $Q_{h;d}$. We denote by C(j)(resp. $\[(resp. \[(resp. \(resp.$ By convention, if j is a contaminant, we set (j) = 0. Notice that whereas those limits are given in terms of amounts of nutrients, regarding contaminants they usually also depend on the body weights of consumers. Hence, from now on, when we write $Q_{h;i}$ to designate the amount of chemical j in food family h, if j is a contaminant then we implicitly refer to its standardized version $Q_{h;i} = w_i$, with w_i the body weight of individual i.

Equipped with these notations, the "safe" situation in regard to nutritional benefits and dietary chemical contamination is described by the random subset of the consumption space

$$S_{Q} := c P R_{+}^{K} : @ Pt1; :::; du; c^{1}Q_{j} P[`_{(j)}; `_{(j)}]$$

$$(4.9)$$

with c¹ the transpose of vector c, i.e the set of dietary habits that yield types of exposure respecting all d considered DIL. Conversely, the unsafe consumption zone is denoted by S_Q^c , the complementary of S_Q in R_+^H . The subscript Q emphasizes their dependence on contents $Q_{h;j}$, 1 § h § H, 1 § j § d, (and body weights for contaminants), which justifies the random nature of both sets. Further set

$$S_{Q_j} := c P R_+^{K} : c^1 Q_j P[[](j); [](j)]]^{(j)}; \qquad (4.10)$$

the safe zone for chemical element j and $S_{Q_j}^c$ its complementary, then according to Equation (4.9), we have $S_Q := {\stackrel{\hat{l}}{i}}_{j=1}^d S_{Q_j}^c$ and $S_Q^c := {\stackrel{\hat{l}}{i}}_{j=1}^d S_{Q_j}^c$. Graphical examples of some possible realizations of S_Q^c are given in Figure 4.1.



Figure 4.1 – Illustration in dimension 2 of the form of S_q^c for some fixed contamination levels q_1 and q_2 of 2 nutrients. The exterior of the two plain (resp. dotted) lines defines the set $S_{q_1}^c$ (resp. $S_{q_2}^c$). Consequently, S_q^c coincides with $A_1 Y A_2$ in the left hand graph and with A_3 (the entire space) in the right hand graph.

Again, we assume in this section that we dispose of independent samples: iid consumption vectors $tC_i := (C_{i;1}; \ldots; C_{i;H}) : 1 \\$ i duare gathered together with a number $m_{h;j}$ of iid measures of content of component j occurring in food item h, namely $tQ_{h;j} := (Q_{h;j}^1; \ldots; Q_{h;j}^{m_{h;j}})$ u for 1 h h H, 1 j d. For simplicity, dietary habits are divided by some large enough constant so that they may fit into the unit cube of R_+^d . The issue mentioned above can be then formulated as follows. Denote by R^H the set of all subspaces in $[0; 1]^H$ that can be written as a finite (with reasonable cardinal, say inferior to some $r^{\Box} PN$) union of hyperrectangles, and fix $\Box P[0; 1]$, a desired level of diet frequency in the studied population. Optimal diets \Box^{\Box} are then defined as solutions of the optimization program:

$$\square := \operatorname{argmin}_{\square PR^{H}} \square (\square) \text{ subject to } P (C P \square) \bullet \square \square \square; \qquad (4.11)$$

where \Box is a dietary risk measure. Notice that exploring sub-spaces in the form of hyperrectangles has the advantage of facilitating the eventual communication of the results: it enables interpretations such as "one should eat more or less than this and that amount of specific food". The main question now is to choose an explicit expression for the volume \Box . One very natural way of modeling dietary risks would be to define $\Box(\Box)$ as the conditional probability $P \ C PS_Q^c \ C P \Box$ for any $\Box PR^H$. Going back to Equation (4.11), notice that in the case where solutions \Box are such that $P(C_P \Box \Box) = \Box \Box \bullet \Box$, it is equivalent to minimizing the joint probability $P \ C PS_Q^c X \Box$ over the set R^H of region candidates, subject to $P(C P\Box) \bullet \Box^c$. This problem is very similar to that tackled in Section 4.1.1, except that the target criterion one seeks to optimize now is not a (Lebesgue) volume anymore, but an unknown probability measure. As in Scott and Nowak (2005) and Clémençon and Vayatis (2010), the criterion must then also be replaced by an empirical estimate. Following the approach proposed in Section 4.1.2, we consider the incomplete U-statistic

where D_B is the index set of cardinality B obtained by drawing with replacement B times in the set of indexes L := $t1; :::; nu \Box \stackrel{\pm}{}_{h=1}^{H} \stackrel{\pm}{}_{j=1}^{d} t1; :::; m_{h;j}u$. This leads to the constrained optimization problem:

$$\min_{\mathbf{P}_{\mathsf{R}}} \stackrel{\Box}{\mathsf{P}}_{\mathsf{B}} \stackrel{\Box}{\mathsf{C}} \mathsf{P} \square \mathsf{X} \overset{\Box}{\mathsf{S}}_{\mathsf{Q}} \text{ subject to } \stackrel{P}{\mathsf{P}} (\mathsf{C} \mathsf{P} \square) \bullet \square \square \square; \qquad (4.12)$$

where

$$\mathbf{P}(\mathbf{C} \mathbf{P} \Box) \coloneqq \frac{1}{n} \prod_{i=1}^{m} \mathbf{I} \mathbf{t} \mathbf{C}_{i} \mathbf{P} \Box \mathbf{u}.$$
(4.13)

The following theorem describes the properties of solutions of the problem in Equation (4.12), involving statistical quantities the computation of which is feasible, the issue of finding such a solution in practice shall be tackled subsequently.

Theor em 4.12 Let $(\Box; \Box) P(0; 1)^2$. Suppose that the collection R^H is of finite VC dimension V $\dagger + 1$ and set:

$$\square^{1}(n;\Box) := 4 \frac{2 \log(8) + V \log(n+1) + \log(2=\Box)}{n}$$

Then, if ${}^{f}_{B}$ is a solution of Equation (4.12), we have with probability at least 1 \square

$$\mathsf{P} \quad \mathsf{C} \ \mathsf{P} \stackrel{\mathbb{I}}{=} \mathsf{B} \quad \bullet \quad \Box + 2\Box \stackrel{1}{=} (\mathsf{n}; \Box = 2) \text{ and } \Box (\stackrel{\mathbb{I}}{=} \mathsf{B}) \\ \$ \quad \Box (\Box \stackrel{\Box}{=}) + 2\Box (\mathsf{B}; \mathsf{N}; \Box = 2); \quad (4.14)$$

denoting by \Box the penalty given by Equation (4.7) and by \Box^{\Box} any minimizer of \Box (.) among the elements \Box of R^{H} such that $P(CP\Box) \bullet \Box$.

The (sketch of the) proof is given in Section 4.2. The statistical procedure described in Section 4.1.3 cannot be extended in a straightforward manner, since the quantity $P C P S_Q^c X C P \Box$ is far from being constant over the collection of cubes C of fixed side length 1=k, which paves the supposedly compact support of the distribution of C.

 $\acute{\textbf{v}}$ Statistically equivalent hypercubes Our proposal to solve Equation (4.12) approximately is to partition [0; 1]^H into a finite number of hypercubes C_1 ;:::; C_M (of variable side length) such that $\mathbb{P}_B \quad C \ P \ C_m \ X \ S^c_Q$, m P t1;:::;Mu, remains (approximately) constant, equal to ! := $M^{\Box 1} \ \mathbb{P}_B \quad C \ P \ S^c_Q$. This can be achieved by means of

a variety of greedy procedures, see section 21.4 in Devroye et al. (1996) and the references therein. Here we use a slightly modified version of the celebrated Gessaman's rule (Gessaman, 1970) as follows. Recall that the quantities involved are built from the sampled dataset $D_B := t(i^{(b)}; (\hat{b}); \dots; \hat{b}) : 1$ b § Bu of cardinality B. For all i Pt1;:::; Bu, assign to the observation C_i the weight

$$!_{i} := \frac{1}{B} \frac{[f]}{b=1} [i = i^{(b)}] I \prod_{h=1}^{\#} C_{i^{(b)};h} Q_{h;j}^{(b)} R[\hat{U}(1); \hat{U}(1)] \text{ or } :::$$

$$::: \text{ or } \prod_{h=1}^{H} C_{i^{(b)};h} Q_{h;j}^{(b)} R[\hat{U}(0); \hat{U}(0)] .$$

Equipped with this notation, observe that $P_B C P S_Q^c X C_m = \prod_{i=1}^{n} !_i I t C_i P C_m u$ for any m P t1;:::; M u. Starting with all the consumption data, consider the food item h and a "split value" s P [0; 1] and consider the regions of the consumption space

 $\Box_{i}(h;s) := tc P[0;1]^{H} : c_{h} \S \text{ su and } \Box_{2}(h;s) := tc P[0;1]^{H} : c_{h} \degree \text{ su.}$

Then find the threshold value s¹ that solves

min	Π	!	i		П	!	i	
sP[0;1]	└-i:C _i P□₁(h;s)		i:C _i P⊡₂(h;s)					

over $t(C_{i;h}): 1$ § i § nu. For each food item h, the minimization problem in s can be solved very rapidly by scanning through all of the data projected onto the h-th axis. Having thus determined the best "split value", one repeat the binary splitting procedure on both regions in a recursive manner. Let $d_1; \ldots; d_H \cdot 1$. Starting from the whole consumption space $[0; 1]^H$, identified as the root node, we propose to build a recursive partition, which can be represented by a binary tree of depth $D = {\stackrel{\infty}{h}}_{h=1}^H d_h$, by using first C_1 as "split variable" to build a complete binary tree of depth d_1 , then using C_2 to grow the tree until depth $d_1 + d_2$ and continuing in the same way with the remaining food items. Its terminal leaves correspond to the cells $C_1; \ldots; C_{2^D}$ of the data-dependent partition. The collection of such regions, obtained by combining hierarchically $2^D \Box 1$ splits perpendicular to the coordinate axes is of finite VC dimension $V_D \dagger + 1$ (namely, its shatter coefficient for m • 1 points is classically bounded by $(m + 1)^{2H(2^D \Box 1)}$). In addition, by construction, the latter are such that $P_B \ C P C_m X S_Q^c$ is approximately equal to $P_B \ C P S_Q^c = 2^D$ for $m = 1; \ldots; 2^D$.

Now, from a practical perspective, the target region of the consumption space can be assessed by using a strategy similar to that described in the previous subsection, as follows.

1. Sort the terminal leaves of the tree representing the partition in a way that

$$\mathbf{P}(\mathbf{C} \mathbf{P} \mathbf{C}_{1:2^{\mathsf{D}}}) \bullet ::: \bullet \mathbf{P}(\mathbf{C} \mathbf{P} \mathbf{C}_{2^{\mathsf{D}}:2^{\mathsf{D}}}).$$

2. Bind together the cells sequentially, until the empirical estimate of the mass of the resulting set exceeds □ □ □ (B; N; □), yielding the region:

Beyond its computational efficacy and simplicity, a crucial advantage of the approach described above lies in its capacity to produce regions which can be visually summarized by a binary tree, the terminal leaves of which can be described by combining elementary rules of the form " c_h ° s" or " c_h § s" in a hierarchical manner. This point is of major importance when designing dietary guidelines to improve nutrition over the population of interest.

4.2 proofs and supplements

4.2.1 Maximal deviation

where

We start by establishing the following intermediary result, which extends Corollary 3 in Clémençon et al. (2008) to the K-sample situation.

Lemma 4.13 Suppose that the hypotheses in Theorem 4.8 are fulfilled. For all $\Box P(0; 1)$, we have with probability at least $1 \Box \Box$,

$$\sup_{P_{\Box}} |U_{n}() \Box \Box()| \S M_{\Box} 2 \frac{\# c}{2 \sqrt{\log(1 + \Box)}} + \frac{c}{\Box} \frac{1}{\log(1 = \Box)} +$$

Proof Set
$$\square$$
 = mint $tn_1 = d_1 u :::; tn_K = d_K u$ and let

$$V X_1^{(1)}; ...; X_{n_1}^{(1)}; ...; X_1^{(K)}; ...; X_{n_K}^{(K)} \square$$

$$= \frac{1}{X_1^{(1)}}; ...; X_{d_1}^{(1)}; ...; X_1^{(K)}; ...; X_{d_K}^{(K)} \square$$

$$+ X_{d_1+1}^{(1)}; ...; X_{2d_1}^{(1)}; ...; X_{d_{K+1}}^{(K)}; ...; X_{2d_K}^{(K)} + ...;$$

$$+ X_{d_1 \square d_1 + 1}^{(1)}; ...; X_{\square d_1}^{(1)}; ...; X_{\square d_K \square d_{K+1}}^{(K)}; ...; X_{\square d_K \square d_{K+1}}^{(K)} \square$$

for any $P \square$. Recall that the K-sample U-statistic $U_n()$ can be expressed as

$$U_{n}() = \frac{1}{n_{1}! \prod n_{K}!} \prod_{\substack{\square_{1} PS \\ \square \\ \square_{K} PS \\ n_{K}}} V X_{\square_{1}(1)}^{(1)}; \dots; X_{\square_{1}(n_{1})}^{(1)}; \dots; X_{\square_{K}(1)}^{(K)}; \dots; X_{\square_{K}(n_{K})}^{(K)}$$

where S_m denotes the symmetric group of order m for any m • 1. This representation as an average of sums of \Box independent terms is known as the (first) Hoeffding's decomposition (Hoeffding, 1948). Then, using Jensen's inequality in particular, one may easily show that, for any nondecreasing convex function G : R₊ — R, we have:

$$E \quad G \quad \sup_{P_{\Box}} \bigcup_{n} (\bar{})^{[1]} \qquad \S$$

$$E \quad G \quad \sup_{P_{\Box}} \bigcup_{-} (X_{1}^{(1)}; \dots; X_{n_{1}}^{(1)}; \dots; X_{1}^{(K)}; \dots; X_{n_{K}}^{(K)})^{[1]}; \quad (4.15)$$

where we set $:= \Box \Box$ () for all P \Box . Now, using standard symmetrization and randomization arguments (see Giné and Zinn (1984) for instance) and Equation (4.15), we obtain that !!

$$E \quad G \quad \sup_{P_{\Box}} \overline{\bigcup}_{n} (\overline{})^{\Box} \quad \S \quad E \quad (G (2R_{\Box})); \qquad (4.16)$$

where

$$\mathsf{R}_{\Box} \coloneqq \sup_{\mathsf{P}_{\Box}} \frac{1}{\Box} \bigcap_{\mathsf{i}=1}^{\Box} \square \overset{\Box}{\underset{\mathsf{i}=1}{}} X^{(1)}_{(\mathsf{i}^{\Box}_{1})\mathsf{d}_{1}+1}; \dots; X^{(1)}_{\mathsf{i}_{d}_{1}}; \dots; X^{(K)}_{(\mathsf{i}^{\Box}_{1})\mathsf{d}_{K}+1}; \dots; X^{(K)}_{\mathsf{i}_{d}_{K}};$$

is a Rademacher average based on the Rademacher chaos $\Box_1; \ldots; \Box_1$ (independent random symmetric sign variables), independent from the $X_i^{(k)}$'s. We now apply the bounded difference inequality (McDiarmid, 1989) to the functional R_{\Box} , seen as a function of the iid random variables ($\Box; X_{(_\Box])d_1+1}^{(1)}; \ldots; X_{d_1}^{(1)}; \ldots; X_{(_\Box])d_{K+1}}^{(K)}; \ldots; X_{d_K}^{(K)}$), 1 § `§ \Box changing any of these random variables changes the value of R_{\Box} by at most $M_{\Box} = \Box$. One thus obtains from Equation (4.16) with $G(x) = expt \Box xu$, where \Box ° 0 is a parameter which shall be chosen later, that:

$$\mathsf{E} \exp \bigcup_{\mathsf{P}_{\square}}^{\#} \operatorname{exp}_{\mathsf{P}_{\square}}^{+!} \operatorname{exp}_{\mathsf{S}}^{"} \operatorname{exp}_{\mathsf{Z}}^{"} \operatorname{E}(\mathsf{R}_{\square}) + \frac{\mathsf{M}_{\square}^{2} \operatorname{c}^{2}}{4_{\square}}^{*}$$

Applying Chernoff's method, one then gets:

Using the bound (see Equation (6) in Boucheron et al. (2005) for instance)

$$E(R_{\odot}) \S M_{\odot} = \frac{2V \log(1 + \Box)}{\Box}$$

and taking $\Box = 2 \Box (\Box \Box 2E(R_{\Box})) = M_{\Box}^2$ in Equation (4.17), one finally establishes the desired result.

Now we shall prove Theorem 4.8, the statement of which is recalled below.

Theor em – Maximal deviation. Let \Box be a collection of bounded symmetric kernels on $\Box := {\overset{\pm}{\overset{K}{k=1}}} X_k^{d_k}$. Suppose that \Box is a VC major class of functions with finite Vapnik-Chervonenkis dimension V and that $M_{\Box} := \sup_{(x) \in \Box \times V} |x| + 1$. Then,

the following assertions hold true.

i) For all \square° 0, we have: $@n = (n_1; \dots; n_K) PN^{\square K}, @B \bullet 1$,

$$P \sup_{P \subseteq \square} \bigcup_{B}^{I} (D) \supseteq U_{n} (D) \bigcup_{B}^{I} \square$$
 § 2(1+ #L)^V e ^{\square} B $\square^{2} = M^{2}$.

ii) For all $\Box P(0; 1)$, with probability at least $1 \Box \Box$, we have: $@n_k \bullet 1, 1 \S k \S K$,

$$\begin{array}{c} \sup_{P^{\circ}} \mathbb{U}_{B}(\cdot) = \mathsf{E}^{\circ} \mathbb{U}_{B}(\cdot) = \mathbb{S} \\ \# \circ \mathsf{C} \\ \mathbb{M} = 2 \end{array} \xrightarrow{\mathbb{Z} \vee \log(1 + \Box)}_{=} + \mathbb{C} \frac{1 \log(2 = \Box)}{\Box} + \mathbb{C} \frac{\nabla \log(1 + \#\mathsf{L}) + \log(4 = \Box)}{\mathsf{B}}^{+}; \end{array}$$

where

 $\Box \coloneqq \min t \mathfrak{m}_1 = d_1 \mathfrak{q} :::; \mathfrak{m}_K = d_K \mathfrak{w}$

and txudenotes the integer part of any real number x.

Proof For convenience, we introduce the random sequence $\Box := ((\Box_b(I))_{I \neq L})_{1 \leq b \leq B}$, where $\Box_b(I)$ is equal to 1 if the tuple I := $(I_1; \ldots; I_K)$ has been selected at the b-th draw and to 0 otherwise: the \Box_b 's are iid random vectors and, for all (b;I) in t1;...; Bu \Box L, the random variable $\Box_b(I)$ has a Bernoulli distribution with parameter 1=#L. We also set $X_I := (X_{I_1}^{(1)}; \ldots; X_{I_K}^{(K)})$ for any I in L. Equipped with these notations, observe first that one may write: @B • 1, @h P(N $\Box)^K$,

$$\mathbb{U}_{\mathsf{B}}(\) \Box \, \mathsf{U}_{\mathsf{n}}(\) = \frac{1}{\mathsf{B}} \frac{\mathsf{P}}{\mathsf{b}_{\mathsf{b}}} Z_{\mathsf{b}}(\);$$

where $Z_b() := \bigcap_{I \neq L} (\Box_b(I) \Box 1 = \#L) (X_I)$ for any (b; I) Pt1;:::; Bu \Box L. It follows from the independence between the X_I's and the $\Box(I)$'s that, for all P \Box , conditioned upon the X_I's, the variables $Z_1()$;:::; $Z_B()$ are independent, centered and almost-surely bounded by 2M \Box (notice that $\bigcap_{I \neq L} \Box_b(I) = 1$ for all b • 1). By virtue of Sauer's lemma, since \Box is a VC major class with finite VC dimension V, we have, for fixed X_I's:

$$\#t((X_1))_{1 PL}$$
: P $\Box u \S (1 + \#L)^V$.

Hence, conditioned upon the X_I 's, using the union bound and next Hoeffding's inequality applied to the independent sequence $Z_1(); \ldots; Z_B()$, for all $\square \circ 0$, we obtain that:

$$P \sup_{P \cap D} \overline{D}_{B}() \cap U_{n}() \stackrel{P}{\to} \cap \left[(X_{1})_{1 P L} \right] \stackrel{P}{\otimes} P \sup_{P \cap D} \left[\frac{1}{B} \right]_{b=1} Z_{b}() \stackrel{P}{\to} \cap \left[(X_{1})_{1 P L} \right] \stackrel{P}{\longrightarrow} \left$$

§
$$2(1 + \#L)^{\vee} e^{\Box B \Box^2 = M^2};$$

which proves the first assertion of the theorem. Notice that this can be formulated: for any $\Box P(0; 1)$, we have with probability at least $1 \Box \Box$

$$\sup_{P \subseteq \Box} \overline{U}_{B}() \Box U_{n}() \stackrel{\square}{\Longrightarrow} M_{\Box} \stackrel{\square}{\longrightarrow} \frac{\overline{V \log(1 + \#L) + \log(2 = \Box)}}{B}.$$
 (4.18)

Turning to the second part of the theorem, it straightforwardly results from the first part combined with Lemma 4.13.

4.2.2 Maximal deviation in dietary risk analysis

We shall prove Corollary 4.11, the statement of which is recalled below.

Proof Observe first that the assumptions of Theorem 4.8 are fulfilled when taking $\Box = t_R$: R P Ru: the collection \Box of indicator functions is a VC major class of functions with finite VC dimension V and $M_{\Box}^2 = 1$. Applying thus Theorem 4.8, the proof is derived by following line by line the argument of Corollary 6 in Scott and Nowak (2006). Details are left to the reader.

4.2.3 Optimal dietary habits

We shall prove Theorem 4.12, the statement of which is recalled below.

Theorem Let $(\Box; \Box) P (0; 1)^2$. Suppose that the collection R^H is of finite VC dimension V $\dagger + 1$ and set:

$$\Box^{1}(n;\Box) := 4^{c} \frac{2\log(8) + V\log(n+1) + \log(2=\Box)}{n}.$$

Then, if $\[\[\]_B \]$ is a solution of Equation (4.12), we have with probability at least 1 \square P C P $\[\[\]_B \]$ • \square + 2 \square ¹(n; \blacksquare =2) and \square ($\[\]_B \]$) § \square (\square) + 2 \square (B;N; \blacksquare =2); (4.19)

denoting by \Box the penalty given by Equation (4.7) and by \Box^{\Box} any minimizer of \Box (.) among the elements \Box of R^{H} such that $P(CP\Box) \bullet \Box$.

Proof The result immediately follows from the argument of Theorem 10 in Clémençon and Vayatis (2010) combined with Theorem 4.8 and Vapnik-Chervonenkis inequality (see Theorem 12.5 in Devroye et al., 1996 for instance) to control the deviations of the supremum:



5

EMPIRICAL PROCESSES IN SURVEY SAMPLING

Like INCA2, most consumption databases are now constructed with some survey design to produce representative samples. Formally, this means that in the population of interest, the probability that an individual may be selected is taken into account in the form of a survey weight. For institutional data, these weights often correspond to the so-called true inclusion probabilities, but statisticians may sometimes have at their disposal calibrated or post-stratification weights (e.g. minimizing some discrepancy with the inclusion probabilities subject to some margin constraints). In most cases, the survey scheme is ignored, potentially yielding a significant sampling bias. When considering some functional of the empirical process such as the empirical distribution function, this may cause severe drawbacks and completely jeopardize the estimation, as can be revealed by simulation experiments. In the context of dietary risk analysis, the impact of such an omission would be for instance an erroneous, greatly biased estimation of the true proportion of over-exposed people. To avoid such undesirable outcomes, many estimators have been developed in the branch of survey sampling theory (Tillé, 2006; Gourieroux, 1981; Droesbeke et al., 1987), which take into account these survey weights and make up for the induced bias of the sampling phase. Unfortunately, to our knowledge, except in the specific case of stratified sampling, there is still no general functional result that would guarantee the asymptotic normality of a large family of estimators in the context of survey sampling. In particular, when estimating the distribution function of the exposure to some food chemical, the construction of confidence bands, as opposed to point confidence intervals, has not been made possible yet. We started addressing this issue in the paper presented from Section 5.1 to Section 5.4, in the specific case of Poisson-like survey plans with no post-calibration. It is the result of a collaboration with P. Bertail (Université Paris X, France) and S. Clémençon (Télécom ParisTech, France) and has been submitted for publication. Unfortunately, the results that are presented there do not apply to the complex sampling of INCA2; they have to be understood as a first step towards the elaboration of a more general theory that would encompass a wider range of survey techniques.

The main goal of this chapter is to investigate how to incorporate the survey scheme into the inference procedure dedicated to the estimation of a probability measure P on a measurable space (viewed as a linear operator acting on a certain class of functions F), in order to guarantee its asymptotic normality. This problem

has been addressed by Breslow and Wellner (2007) and Saegusa and Wellner (2011) in the particular case of a stratified survey sampling, where individuals are selected at random (without replacement) in each stratum, by means of bootstrap limit results. Our approach is different and follows that of Hàjek (1964), extended next by Berger (1998, 2011), and is applicable to more general sampling surveys, namely those with unequal first order inclusion probabilities which are of the Poisson type or sequential/rejective. The main result of the chapter is a Functional Central Limit Theorem (FCLT) describing the limit behavior of an adequate version of the empirical process (referred to as the Horvitz-Thompson empirical process throughout the article) in a superpopulation statistical framework. The key argument involved in this asymptotic analysis consists in approximating the distribution of the extended empirical process by that related to a much simpler sampling plan. In order to illustrate the reach of this result, statistical applications are considered, where the extensions of the empirical process are used to construct confidence bands around the Horvitz-Thompson estimator of the cumulative distribution function.

The chapter is organized as follows. In Section 5.1 and Section 5.2, the statistical framework is described at length, notations are set out and some basics on survey sampling theory are recalled, together with important examples of survey schemes to which the subsequent asymptotic analysis can be applied. The main result of the chapter, a FCLT for the Horvitz-Thompson empirical process, is stated in Section 5.3, while applications of the latter to non-parametric functional estimation are displayed in Section 5.4. Finally, technical details are deferred to Section 5.5.

5.1 background and preliminaries

We start off with recalling some crucial notions in survey sampling and in modern empirical process theory, which shall be extensively used in the subsequent analysis. Throughout the article, the Dirac mass at x in some vector space X is denoted by \Box_x and the indicator function of any event E by I t Eu. We also denote by #E the cardinality of any finite set E, and by P(E) its power set.

5.1.1 Survey sampling: some basics

The purpose of survey sampling is to study some characteristics of a population U_N of $N \cdot 1$ units (or individuals) identified by an arbitrary collection of labels: $U_N := t1; \ldots; Nu$. For various reasons (limited budget, geographical constraints, etc.), it is usually not possible to reach the whole population, and the features of interest have to be estimated from a finite, relatively small number of its elements, namely a sample s := $ti_1; \ldots; i_{n(s)}u \ A U_N$ of size $n(s) \ N$. So as to provide handy ways of controlling the accuracy of estimation, sample units are picked randomly among

 U_N (see for instance Tillé, 2006, Chapter 1, Tillé, 1999 or Gourieroux, 1981 for an introduction to the origins of random sampling). Equipped with this representation, a sampling scheme (design/ plan) is determined by a discrete probability measure R_N on $P(U_N)$, the set of all possible samples in U_N . Depending on the adopted point of view, like in superpopulation models, the characteristics of the population can be considered random too. In the next paragraphs, while introducing crucial concepts and notations, we shall discuss both sources of hazard and their classical modeling in survey sampling theory.

5.1.1.1 Survey schemes without replacement

Consider a sampling scheme R_N where individuals are only selected once, i.e. a design without replacement. Our analysis is restricted to this popular family of survey plans. By definition, the two conditions below are always fulfilled,

1.
$$\bigotimes_{s \in P} P(U_N), R_N(s) \bullet 0,$$

2. $\underset{s \in P}{s \in U_N} R_N(s) = 1,$

and the mean survey sample size is given by

$$E_{R_N}(n(S)) = \prod_{sPP(U_N)} n(s) R_N(s).$$

Here, the notation E_{R_N} (.) denotes the expectation taken with respect to the random sample S with distribution R_N . In a similar fashion, P_{R_N} (SPS) refers to the probability of the event tSPSu with SÄ P(U_N), when S is drawn from R_N . In particular, R_N (s) = P_{R_N} (S= s). Such distributions are entirely characterized by the concepts listed below.

$$\Box_{i}(R_{N}) := P_{R_{N}}(i PS) = \prod_{sPP(U_{N})}^{|I|} R_{N}(s) I ti Psu;$$

is the probability that the individual labeled i belongs to a random sample S under the survey scheme R_N . When there is no ambiguity on the sampling design, we will simplify notations and write \Box_i instead of $\Box_i (R_N)$. In the subsequent analysis, first order inclusion probabilities are assumed to be strictly positive: @ PU_N, $\Box_i (R_N) \circ 0$. We shall even require the stronger hypothesis that they never get either too small or too large, as formally stated below.

Assumption 5.1 There exist \Box_0 ° 0 and N₀ PN \Box such that for all N • N₀ and i PU_N,

$$\Box_{i}(\mathsf{R}_{\mathsf{N}})^{\circ} \Box_{\Box}$$

In addition,

$$\limsup_{N \to +1} \frac{1}{N} \prod_{i=1}^{N} \Box_{i}(R_{N}) + 1.$$

When the first condition holds, the rate of convergence of the estimators considered in Section 5.2 and Section 5.3 will be shown to be typically of order $1 = \frac{1}{N}$. One could possibly relax it and allow \Box_{\Box} to depend on N, with $\Box_{\Box} = \Box_{\Box}(N)$ decaying to zero as N tends to infinity at a specific rate, and still be able to establish limit results. The analysis would be however much more technical; this is left for further research.

Conditions involving the second order inclusion probabilities shall also be used in our asymptotic analysis. They are denoted by

$$\Box_{i;j}(\mathsf{R}_{\mathsf{N}}) \coloneqq \mathsf{P}_{\mathsf{R}_{\mathsf{N}}} \stackrel{\square}{(i;j)} \mathsf{PS}^{2} = \prod_{\mathsf{sPP}(\mathsf{U}_{\mathsf{N}})} \mathsf{R}_{\mathsf{N}}(\mathsf{s}) \mathsf{Itti;ju} \mathsf{\ddot{\mathsf{A}}} \mathsf{su};$$

for all (i; j) $P U_N^2$. In other words, $\Box_{i;j}(R_N)$ is the probability that two distinct individuals labeled i and j are jointly selected under design R_N . Again, we may eventually write $\Box_{i;j}$ when there is no need to emphasize the dependency on the sampling plan R_N . Notice that higher order inclusion probabilities may be defined in a similar way, up to the maximal order for which the entire population is selected.

 $\Box_{i} := I \text{ ti } P \text{ Su} = \begin{cases} \$ & 1 \\ \$ & 1 \\ \% & 0 \end{cases} \text{ with probability } \Box_{i};$

Notice indeed that the set $P(U_N)$ of all possible samples is in one-to-one correspondence with t0; $1u^N$, which provides a handy alternative representation of sampling schemes. Again, for simplicity, we will omit the subscript (N) when no ambiguity is possible. By definition, the distribution of $\Box := \Box_{(N)}$ has univariate marginals that correspond to the Bernoulli distributions $B(\Box_i)$, i PU_N , and covariance matrix given by

$$\Box_{N} := \Box_{i;j} \Box \Box_{i} \Box_{j} \overset{(}{}_{1 \leq i;j \leq N} \cdot$$

Incidentally we have $\overset{\infty}{}_{i=1}^{N} \Box_{i} = n(S)$ and thus $\overset{\infty}{}_{i=1}^{N} \Box_{i} = E_{R_{N}} (n(S))$.

Before considering the issue of extending the concept of empirical process in the context of survey sampling, we recall a few important classes of survey schemes, to which the results established in Section 5.2 and Section 5.3 can be applied. One may refer to Deville (1987) for instance for an excellent account of survey theory, including many more examples of sampling designs.

Example 5.2 – Simple Random Sampling Without Replacement. A simple random sampling without replacement (SRSWOR in abbreviated form) is a sampling design of fixed size n(S) = n, according to which all samples with cardinality n in the population U_N are equally likely to be chosen, with probability $(N \square n)!=n!$. It follows that all units of U_N have the same chance of being selected, n=N namely, and all second order probabilities are equal to $n(n \square 1)=(N(N \square 1))$.

Example 5.3 – Poisson survey sampling. The Poisson sampling plan without replacement (POISSWOR), denoted here by T_N , is one of the simplest survey schemes. In this case, the N elements of \Box are independent Bernoulli random variables with respective parameters $\Box_i (T_N) =: p_i$, i Pt1;:::;Nu so that for any sample s PP(U_N),

$$T_{N}(s) = \prod_{i Ps}^{\Pi} \prod_{i Rs}^{\Pi} (1 \Box p_{i}).$$

Notice that the size n(S) of sample S with distribution T_N is random (except in the sole situation where $p_i P t 0$; 1u for i = 1; ...; N) and that the corresponding survey plan is fully characterized by the first order inclusion probabilities. In the specific situation where they are all equal, i.e $p_1 = \Box \equiv p_N = p$, the design is called Bernoulli.

Example 5.4 – Stratified sampling. A stratified sampling design permits to draw a sample S of fixed size $n(S) = n \S N$ within a population U_N that can be partitioned into K • 1 distinct strata $U_{N_1}; \ldots; U_{N_K}$ (known a priori) of respective sizes $N_1; \ldots; N_K$ adding up to N. Let $n_1; \ldots; n_K$ be non-negative integers such that $n_1 + \ldots + n_K = n$, then the drawing procedure is implemented in K steps: within each stratum U_{N_k} , k P t1; \ldots ; Ku, perform a SRSWOR of size $n_k \S N_k$ yielding a sample S_k . The final sample is obtained by assembling these sub-samples: $S = \begin{bmatrix} 1 & K \\ k=1 & S_k \end{bmatrix}$. The probability of drawing a specific sample s by means of this survey scheme is

$$\mathsf{R}_{\mathsf{N}}^{\mathsf{str}}(\mathsf{s}) = \frac{\mathsf{r}_{\mathsf{N}}}{\mathsf{n}_{\mathsf{k}}} \frac{\mathsf{n}_{\mathsf{k}}}{\mathsf{n}_{\mathsf{k}}}$$

Naturally, first and second order inclusion probabilities depend on the stratum to which each unit belong: for all i \Box j in U_N,

$$\Box_{i}(\mathsf{R}_{\mathsf{N}}^{\mathsf{str}}) = \prod_{k=1}^{[\mathsf{N}]} \frac{\mathsf{n}_{k}}{\mathsf{N}_{k}} \mathsf{I} \mathsf{ti} \mathsf{PU}_{\mathsf{N}_{k}} \mathsf{u} \text{ and } \Box_{i;j}(\mathsf{R}_{\mathsf{N}}^{\mathsf{str}}) = \prod_{k=1}^{[\mathsf{N}]} \frac{\mathsf{n}_{k}(\mathsf{n}_{k} \Box 1)}{\mathsf{N}_{k}(\mathsf{N}_{k} \Box 1)} \mathsf{I} \quad (i;j) \mathsf{PU}_{\mathsf{N}_{k}}^{2} (\mathsf{n}_{k}^{\mathsf{str}}) = \mathsf{I}_{\mathsf{N}_{k}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}_{k}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}}^{\mathsf{n}}^{\mathsf{n}}} \mathsf{I}_{\mathsf{n}}^{\mathsf{n}}^{\mathsf{n}_{k}} \mathsf{I}_{\mathsf{$$

Example 5.5 – Canonical Rejective Sampling. Let n § N and consider a vector $\Box^R := (\Box_1^R; \ldots; \Box_N^R)$ of first order inclusion probabilities. Further define $S_n := ts P P(U_N) : #s = nu$, the set of all samples in population U_N with car-

dinality n. The rejective sampling (Hàjek, 1964; Berger, 1998), sometimes called conditional Poisson sampling (CPS), exponential design without replacement or maximum entropy design (Tillé, 2006, Section 5.6), is the sampling design R_N^R that selects samples of fixed size n(s) = n so as to maximize the entropy measure

$$H(R_{N}) = \Box \prod_{s \neq S_{n}}^{\prod} R_{N}(s) \log R_{N}(s);$$

subject to the constraint that its vector of first order inclusion probabilities coincides with \Box^R . It is easily implemented in two steps.

- Draw a sample S with a POISSWOR plan T_N = T^p_N, with properly chosen first order inclusion probabilities vector p := (p₁;...;p_N). The representation is called canonical if ^N_{i=1} p_i = n. In that case, relationships between each p_i and □^R_i, 1 § i § N, are established in Hàjek (1964).
- 2. If $n(S) \square n$, then reject sample S and go back to step one, otherwise stop.

Vector p must be chosen in a way that the resulting first order inclusion probabilities coincide with \Box^R , by means of a dedicated optimization algorithm (Tillé, 2006, Algorithms 5.5 to 5.9). The corresponding probability distribution is given for all s PP(U_N) by

$$R_{N}^{R}(s) = \frac{T_{N}^{p}(s) \mid t \# s = nu}{s^{1} P S_{n} T_{N}^{p}(s^{1})} 9 \prod_{i P S}^{\pi} p_{i} \prod_{i R S}^{\pi} (1 \Box p_{i}) \Box \mid t \# s = nu;$$

where 9 denotes the proportionality. We refer to Hàjek (1964, p.1496) for more details on the links between rejective and Poisson sampling plans.

Example 5.6 – Rao-Sampford Sampling. The Rao-Sampford sampling design generates samples s P P(U_N) of fixed size n(s) = n with respect to some given first order inclusion probabilities $\Box^{RS} := (\Box_1^{RS}; \ldots; \Box_N^{RS})$, fulfilling $\overset{\sim}{\underset{i=1}{\overset{N}{\longrightarrow}} = n$, with probability

$$\mathsf{R}^{\mathsf{RS}}_{\mathsf{N}}(\mathsf{s}) = \Box_{\mathsf{i}} \Box_{\mathsf{i}}^{\mathsf{RS}} \frac{\mathsf{\pi}}{\mathsf{j}_{\mathsf{Rs}}} \frac{\Box_{\mathsf{j}}^{\mathsf{RS}}}{\mathsf{1} \Box \Box_{\mathsf{j}}^{\mathsf{RS}}}.$$

Here, \Box ° 0 is chosen such that $\underset{sPP(U_N)}{\overset{\circ}{}} R_N^{RS}(s) = 1$. In practice, the following algorithm is often used to implement such a design (Berger, 1998):

- 1. select the first unit i with probability $\Box_i^{RS}=n$,
- 2. select the remaining n \Box 1 units j with drawing probabilities proportional to $\Box_i^{RS} = (1 \Box \Box_i^{RS}), j = 1; \dots; N,$
- 3. accept the sample if the units drawn are all distinct, otherwise reject it and go back to step one.

5.1.1.2 Superpopulation models

The characteristics of interest in population U_N are modeled as follows. We consider the probability space $(U_N; P(U_N); P)$ and a random variable/ vector X defined on the latter, taking its values in a Banach space $(X; \}.$), with probability measure P. We set

$$X: \begin{array}{cccc} U_{N} & \Box & X \\ i & fi & X(i) = X_{i} \end{array}$$

and the \Box -algebra induced by the normed vector space topology structure of X is denoted by A. For instance, X could represent the amounts of K food products consumed on a specific day. In that case, we would have $X = R_{+}^{K}$, A the associated Borel algebra, }.} the euclidean norm and $X_i = (X_{i;1}; \ldots; X_{i;K})$ would give the daily intakes of individual i P U_N. Then, the studied features correspond to some synoptic mapping $(X_1; \ldots; X_N)$ fi- $f(X_1; \ldots; X_N)$. In our example, we could consider $f(X_1; \ldots; X_N) = (N^{\Box 1} \bigcap_{i=1}^N X_{i;1}; \ldots; N^{\Box 1} \bigcap_{i=1}^N X_{i;K})$, the average consumption of the K foods in U_N.

In survey sampling, a superpopulation is basically an imaginary infinite population, U_1 say, from which U_N is supposed to be issued. In a model-based approach, it is assumed that the random vectors of interest X_1 ;:::; X_N are in fact realizations of N random vectors \tilde{X}_j : $U_1 - X$, 1 § j § N, with joint distribution Q. Then, a superpopulation model is simply a set of conditions that characterize Q (Droesbeke et al., 1987, Chapter 4). The main advantage of such a framework is that it often facilitates statistical inference; in particular, it permits the development of an asymptotic theory, when sample and population sizes grow conjointly to infinity. The superpopulation model we consider here stipulates that all N random vectors X_i , i P U_N, are independent and identically distributed (iid) with common distribution P, i.e. Q = P^{b N}, where b denotes the tensor product of measures.

Remark 5.7 The most celebrated iid superpopulation model that we adopt here establishes a setting very similar to that of weighted bootstrap (Arcones and Giné, 1992; Barbe and Bertail, 1995): the original iid N-sample there would correspond to the complete vector (X_1 ;:::; X_N), from which sub-samples are drawn according to some procedure likened to the survey scheme. Actually, both approaches are completely equivalent if the survey weights ($\Box_1 = \Box_1$;:::; $\Box_N = \Box_N$) are exchangeable (i.e. the N-variate distribution of this vector is invariant to the order of its elements). For instance, in the specific case of stratified sampling, drawing units with equal probabilities in each stratum (with a finite and given stratum-size) amounts to bootstrapping (without replacement) in some given cell. It is not surprising then that both Breslow and Wellner (2007) and Saegusa and Wellner (2011) construct a general asymptotic theory in two-phase sampling by using bootstrap type results.

5.1.1.3 Auxiliary information

In practice, sampling from a population U_N is only possible if all individuals are listed somehow, and can be identified once selected. Such documents are called survey frames; in the case of social surveys, they are collected by government institutions and often provide some minimal information about its components. For instance, in France, the geographical situation, the age, the genre and the profession and socio-professional category ("Profession et Catégorie Socioprofessionnelle", PCS, in French) of citizens are displayed in files managed by INSEE (Institut National de la Statistique et des Etudes Economiques). These auxiliary variables, supposedly known for all i P U_N , can sometimes be used to optimize in some sense the survey scheme. In a superpopulation framework, we denote by W the auxiliary random vector, valued in some measurable space W, and set $W_{(N)} := (W_1; \ldots; W_N)$. Again, the subscript (N) shall be dropped when no confusion is possible. As soon as W is correlated with X, the vector of interest, it becomes possible to boost the efficiency of estimators by defining inclusion probabilities as a function of $W_{(N)}$ (Droesbeke et al., 1987).

In the present analysis, we denote by $P_{X;W}$ the joint distribution of (X;W) and by P_W the marginal distribution of W. Like in most applications, we assume that the W_i's are independent (or exchangeable) random variables/ vectors, linked to the variable of interest X through a linear model (notice that W may be constant over the population). It is required though that W is not proportional to X (in a deterministic sense) to avoid degenerate situations; in such a case, knowing W on the whole population would mean knowing the empirical process without any error. For the sake of simplicity, the dependence of survey weights in W will only be emphasized when it is necessary, starting in Section 5.3.

5.1.2 Empirical process indexed by classes of functions

In the context of iid realizations X_1 ;:::; X_N of a probability measure P, empirical process theory (Ledoux and Talagrand, 1991) consists in the study of the fluctuations of random processes of the type

$$tG_N f$$
; $f PFu$; where $G_N := P_N \square P$

There, class F designates a certain set of P-integrable real-valued functions,

$$\mathsf{P}_{\mathsf{N}} \coloneqq \frac{1}{\mathsf{N}} \prod_{i=1}^{|\mathsf{N}|} \Box_{\mathsf{X}_{i}}$$

is the "classical" empirical measure, and for any signed measure Q on a measurable space (X; A), Qf := $\int_X f(x) Q(dx)$ when the integral is well-defined. We assume that class F admits a square integrable envelope H as defined below.

Assumption 5.8 There exists a measurable function H : X — R such that $_{x}$ H²(x) P(dx) † 1 and |f(x)| § H(x) for all x PX and any f PF.

As a consequence, F is a subset of the space

$$L_2(P) := h : X - R; h \text{ measurable and } h_{2;P}^2 := E_P h^2(X)^+ + 1$$

Notice that we may assume without loss of generality that there exists $\square \circ 0$ such that $H(x) \circ \square$ for every x PX, even if it entails replacing H by H + \square in the condition above.

5.1.2.1 Donsker classes

When viewed as a linear operator acting on F, a probability measure P satisfying Assumption 5.8 may be considered as an element of 1 (F), i.e. the space of all maps \Box : F — R such that

$$|\Box|_{F} := \sup_{f \in F} |\Box(f)| + +1 ;$$

equipped with the uniform convergence norm (or, equivalently, with Zolotarev metric), namely

$$P \square Q_F := d_F(P;Q) = \sup_{h \in F} \left[\begin{array}{c} a \\ h & dP \end{array} \right]^a h dQ;$$

for any couple of probability measures P and Q. The main purpose of empirical process theory is to find conditions on the class of functions F guaranteeing that the distribution of $\overline{N} G_N$ converges, as N - +1, to that of a Gaussian, Banach space valued process in `1 (F). Such collections of functions are called Donsker classes by analogy to the classical results on the empirical distribution function that analyze $\overline{N} (F_N \Box F)$, where

$$F_{N}(x) := \frac{1}{N} \prod_{i=1}^{|V|} I t X_{i} P(\Box 1 ; x_{1}] \Box \Box \Box \Box (\Box 1 ; x_{d}] u$$

and

$$F(x) := P(X P(\Box 1 ; x_1] \Box \Box \Box \Box (\Box 1 ; x_d])$$

for $x := (x_1; :::; x_d) P R^d$ (see Example 5.9). In particular, the study of the uniform deviations over F

$${\mathbf{N}} P_{\mathsf{N}} \square P_{\mathsf{F}}$$

is of great interest, with a variety of applications in statistics, see Shorack and Wellner (1986). A nearly exhaustive review of asymptotic results ensuring that F is a Donsker class of functions is available in van der Vaart and Wellner (1996). The purpose of this chapter is to extend typical empirical processes results obtained for iid data to the framework of survey sampling.

5.1.2.2 On measurability issues

Recall that the normed vector space $(1 (F); \}_{F}$ is (generally) a non-separable Banach space. The major problem one faces when dealing with sums of random variables taking their values in such an infinite-dimensional non-separable space congerns the measurability of events. For instance, the "classical" empirical process \overline{N} (F_N \Box F), which can be viewed as a random sequence in the Skorokhod space D ([0; 1]) of càd-làg functions endowed with the supremum norm, is not Borel-measurable. In this specific case, the topology induced by the sup-norm on D ([0; 1]) can be classically replaced by the Skorokhod metric in order to overcome this technical difficulty. Alternative approaches can be found in Pollard (1984). The ideas developed in Hoffmann-Jørgensen (1991) have led to a general solution, based on the concept of outer probability, extending the original probability measure P to nonmeasurable events by setting $P^{\Box}(A) := \inf t P(B) : A \ddot{A} B$; B measurable u. Then, the related concept of Hoffman-Jørgensen weak convergence permits somehow to forget the measurability assumptions. Hence, expectations and probabilities must now be understood as outer expectations and probabilities for non-measurable events. For simplicity, the same notations are kept to denote original and outer probabilities (resp. expectations). Here, weak convergence is metrized through the bounded Lipchitz metric on the space 1 (F): for all random functions X and Y valued in 1 (F),

$$d_{\mathsf{BL}}(X;Y) = \sup_{\mathsf{b}\,\mathsf{PBL}_1(^{1}(\mathsf{F}))} \stackrel{\sqcup}{\in} (\mathsf{b}(X)) \Box \mathsf{E} (\mathsf{b}(Y)) \stackrel{\sqcup}{;}$$

where $BL_1(^1(F))$ is the set of all 1-Lipchitz functions on $^1(F)$ bounded by 1. In the following we define the \square_P semi-metric under P as

$$\Box_{\mathsf{P}}(\mathsf{f};\mathsf{g}) := \mathsf{E}_{\mathsf{P}} \left(\mathsf{f}(\mathsf{X}) \Box \mathsf{g}(\mathsf{X}) \right)^2 =: \{\mathsf{f} \Box \mathsf{g}\}_{2;\mathsf{P}}^2.$$

We refer to van der Vaart and Wellner (1996) for technical details and general results.

5.1.2.3 Uniform covering numbers

A key concept in the study of empirical process is the covering number N(";F;|.|), which corresponds to the minimal number of balls of radius " ° 0 for a given semi metric |.| needed to cover F. Donsker classes of functions are often characterized by some integrability conditions of the form

arising from maximal inequalities. Such a condition essentially ensures that the size of class F is not too big and that one may be able to approximate any of its elements (up to ") by functions in a set of finite cardinality. In our non-iid setting, we will essentially consider the $L_2(P)$ norm for |.| and use uniform covering numbers

where D is the set of all discrete probability measures Q such that $0 \ddagger H^2 dQ \ddagger +1$. Explicit calculus of (uniform) covering numbers for general classes of functions may be found in several textbooks, see van der Vaart and Wellner (1996) or van de Geer (2000).

5.2 empirical process in survey sampling

We now introduce two different empirical processes built from survey data, the asymptotic behaviors of which shall be investigated at length in Section 5.3.

5.2.1 The Horvitz-Thompson empirical process

In the context of survey data drawn through a general survey plan R_N , the empirical measure P_N cannot be computed since the whole statistical population is not observable. Hence, a variant based on the observations must be naturally considered. For any measurable set M Ä X, the Horvitz-Thompson estimator of the empirical probability $P_N(M) = N^{\Box 1} \stackrel{\sim}{\underset{i=1}{\sim}} \Box_{X_i}(M)$ based on the survey data described above is defined as follows, see Horvitz and Thompson (1951):

$$\mathsf{P}_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{R}_{\mathsf{N}})}(\mathsf{M}) \coloneqq \frac{1}{\mathsf{N}} \frac{|\mathsf{N}|}{|\mathsf{I}|} \xrightarrow{\Box_{\mathsf{i}}} \Box_{\mathsf{i}}(\mathsf{M}) = \frac{1}{\mathsf{N}} \frac{\mathsf{\Pi}}{|\mathsf{P}\mathsf{S}|} \frac{\mathsf{I} \operatorname{ti} \mathsf{P} \operatorname{Su}}{|\mathsf{I}|} \Box_{\mathsf{i}}(\mathsf{M}).$$
(5.1)

We highlight the fact that the measure $P_{R_N}^{\Box(R_N)}$ is an unbiased estimator of P (resp. P_N , when conditioned upon (X_1 ;:::; X_N)) although it is not a probability measure. For a fixed subset M, the consistency and asymptotic normality of the estimator in Equation (5.1) are established in Robinson (1982) and Berger (1998), as N tends to infinity. When considering the estimation of measure P_N (the measure of interest in survey sampling) over a class of functions F, we are led to the asymptotic study of the collection of random processes

$$G_{R_N}^{\square(R_N)} \coloneqq G_{R_N}^{\square(R_N)} f_{fPF}$$

where

$$G_{R_{N}}^{\Box(R_{N})}f := \frac{?}{N} \frac{\Box}{P_{R_{N}}^{\Box(R_{N})}} \Box P_{N} \frac{\Box}{f} = \frac{?}{N} \frac{[N]}{\Box_{i=1}} \frac{\Box}{\Box_{i}(R_{N})} \Box 1 f(X_{i}); \quad (5.2)$$

which shall be referred to as the F-indexed Horvitz-Thompson empirical process (HTempirical process, in short). The seemingly redundant notation $G_{R_N}^{\square(R_N)}$ is motivated by the fact that extensions involving first order probabilities related to a different sampling scheme T_N will be considered in the sequel. Precisely, $G_{R_N}^{\square(T_N)}$ shall denote the process obtained when replacing all $\square_i(R_N)$ by $\square_i(T_N)$, 1 § i § N, in Equation (5.2).
The main purpose of this chapter is to establish the convergence of the re-weighted empirical process $(G_{R_N}^{\Box} f)_{fPF}$ under adequate hypotheses involving some properties of measure P, certain characteristics of the sequence of sampling plans (R_N) , and the "complexity" of class F (in the classical metric entropy sense) as well. In particular, such a result would permit to describe the asymptotic behavior of the quantity below (assumed to be almost-surely finite, see Assumption 5.8):

$$G_{R_{N}}^{(R_{N})} = \sup_{f \in F} G_{R_{N}}^{(R_{N})} f$$

By virtue of Cauchy-Schwarz inequality and Assumption 5.1 and Assumption 5.8, we almost-surely have, @N • 1,

Under Assumption 5.1 and Assumption 5.8, the F-indexed HT-empirical process in Equation (5.2) may thus be seen as a sequence of random elements of $1^{(F)}$.

Example 5.9 – Empirical cumulative distribution function. In the case where $X = R^d$ with d • 1 for instance, a situation of particular interest is that where F is the class of indicator functions of rectangles of the type

$${}^{\#}(\Box 1 ; x] := {}^{T {f^{\sharp}}} \Box 1 ; x_{j} \Box 1 ; x_{j} \Box : x_{j} : x_{j}$$

Then, the empirical process can be identified with the Horvitz-Thompson version of the empirical cumulative distribution function (cdf) $F_{R_N}^{\Box(R_N)}(x) \coloneqq P_{R_N}^{\Box(R_N)}(\Box 1 ; x]$, x P R^d, and the goal pursued boils down to investigating conditions under which uniform versions of the Law of Large Numbers (LLN) and of the Central Limit Theorem (CLT) hold for $F_{R_N}^{\Box(R_N)}(x) \Box F_N(x)$, where $F_N(x) \coloneqq P_N(\Box 1 ; x]$. As shall be seen later, the study of the asymptotic behavior of this empirical process lies at the center of the validity of the confidence band construction considered in Section 5.4.

5.2.2 Alternative estimate in the Poisson sampling case

The Poisson sampling scheme T_N (see Example 5.3) has been the subject of much attention, especially in Hàjek (1964), where asymptotic normality of (pointwise) Horvitz-Thompson estimators have been established in this specific case. Following in the footsteps of this seminal contribution, we consider the following Poisson version of the empirical process rather than the original process :

$$\mathbb{G}_{\mathsf{T}_{\mathsf{N}}}^{\mathsf{p}} \mathsf{f} \coloneqq \frac{2}{\mathsf{N}} \prod_{i=1}^{\mathsf{N}} (\Box_{i} \Box_{\mathsf{p}_{i}})^{\Box} \frac{\mathsf{f}(\mathsf{X}_{i})}{\mathsf{p}_{i}} \Box_{\mathsf{N};\mathsf{p}}(\mathsf{f})^{\Box}; \mathsf{f} \mathsf{P}\mathsf{F};$$
(5.3)

where for all f PF,

$$\Box_{N;p}(f) \coloneqq \frac{1}{d_N} \bigcap_{i=1}^{[n]} (1 \Box p_i) f(X_i) \text{ and } d_N \coloneqq \bigcap_{i=1}^{[n]} p_i (1 \Box p_i).$$

Under the assumption that $d_N - +1$ as N - +1, it has been established in Hàjek (1964, Lemma 3.2) that conditioned upon $(X_1; \ldots; X_N)$, for fixed f P F, when N tends to infinity and under a Lindeberg-Feller type condition, the weighted sum of independent random variables in Equation (5.3) can be approximated by a centered Gaussian random variable with (conditional) variance

$$V_{N}^{2}(f) = \frac{1}{N} \frac{\prod_{i=1}^{N} \frac{f(X_{i})}{p_{i}}}{\prod_{i=1}^{N} \frac{f(X_{i})}{p_{i}}} \Box \Box_{N;p}(f) \frac{u^{2}}{p_{i}} p_{i}(1 \Box p_{i}).$$

As claimed by Theorem 5.14 in the next section, this result can be extended to a functional framework under adequate hypotheses.

Although the subscript T_N in $\mathfrak{G}_{T_N}^p$ f could have been dropped since the process above only depends on vector p, we keep it in order to emphasize that the corresponding inclusion vector \Box is distributed according to the sampling scheme T_N . In this subsection, the weights $p_i := \Box_i (T_N)$, 1 § i § N, correspond to the inclusion probabilities of the Poisson sampling plan. Later on, when investigating a general sampling scheme R_N , we shall consider the Poisson-like empirical process defined by

$$\mathbb{G}_{\mathsf{R}_{\mathsf{N}}}^{\mathsf{p}} \mathsf{f} \coloneqq \frac{2}{\mathsf{N}} \frac{\mathsf{P}}{\mathsf{i}_{\mathsf{i}}} (\Box_{\mathsf{i}} \Box \mathsf{p}_{\mathsf{i}}) \frac{\mathsf{P}}{\mathsf{p}_{\mathsf{i}}} \Box \Box_{\mathsf{N};\mathsf{p}}(\mathsf{f}) ;$$

where $p := (p_1; \ldots; p_N)$ is the vector of first order inclusion probabilities of a Poisson design. In general, it will not coincide with those of R_N , namely $\Box(R_N)$, but the subscript specifies that \Box is still distributed according to R_N (in particular, $E(\Box_i) = \Box_i(R_N)$ for $i = 1; \ldots; N$). In the subsequent analysis, we start off by establishing that the process $G_{T_N}^p$ can be asymptotically approximated by a Gaussian process.

5.3 asymptotic results

The main results of the chapter are stated in the present section. As a first go, we establish a FCLT for the empirical process variant of Equation (5.3) in the Poisson survey scheme case. Combined with an approximation result, it will serve as the main tool for proving next a similar result in the context of rejective sampling.

5.3.1 Limit of the empirical process for the Poisson survey scheme

The purpose of this section is to obtain a Gaussian approximation of the empirical process $\mathfrak{G}_{T_N}^p$ related to a Poisson survey plan T_N with first order inclusion probabilities $p = (p_1; \ldots; p_N)$ depending on some auxiliary variable W (see Section 5.1.1.3). The proof relies on Theorem 2.11.1 in van der Vaart and Wellner (1996), applied to the triangular collection of independent variables defined for all f PF by

$$Z_{N;i}(f) \coloneqq Z_{N;i}(f; \Box) \coloneqq \frac{2}{N} (\Box_i \Box p_i) \stackrel{\Box}{\longrightarrow} \frac{f(X_i)}{p_i} \Box \Box_{N;p}(f) \stackrel{\Box}{\longrightarrow} \text{for i Pt1; :::; Nu.}$$

For clarity, the result is recalled below.

Theorem 5.10 - Triangular arrays (van der Vaart and Wellner, 1996).

Let $Z_{N;i}(f)$, 1 § i § N be independent F-indexed stochastic processes defined on the product probability space $\stackrel{\pm}{}_{i=1}^{N}(t_{0}; 1_{u}; P(t_{0}; 1_{u}); B(\Box_{i}(R_{N})))$ where the process $Z_{N;i}(f) := Z_{N;i}(f; \Box)$ only depends on the i-th coordinate of $\Box := (\Box_{1}; \ldots; \Box_{N})$. Assume that the maps

$$(\Box_{1}; \ldots; \Box_{N}) \stackrel{\text{fi-}}{=} \sup_{\Box_{P}(f;g) \uparrow \Box} \stackrel{\text{fi-}}{=} e_{i} (Z_{N;i}(f;\Box) \Box Z_{N;i}(g;\Box))$$

and

$$(\Box_{1}; \ldots; \Box_{N}) \stackrel{\text{fi-}}{=} \sup_{\Box_{P}(f;g) \uparrow \Box} \stackrel{P}{=} e_{i} (Z_{N;i}(f) \Box Z_{N;i}(g))^{2}$$

are measurable for every \square° 0, every $(e_1; \ldots; e_N)$ Pt \square 1; 0; 1u^N and every N PN. Further define the random semi-metric

$$d_{N}^{2}(f;g) \coloneqq \prod_{i=1}^{|N|} (Z_{N;i}(f) \Box Z_{N;i}(g))^{2};$$

and suppose that the following conditions are fulfilled.

- i) $\begin{array}{c} \overset{\bullet\bullet}{\mathsf{E}} \\ \underset{i=1}{\overset{\bullet\bullet}{\mathsf{E}}} \\ \overset{\bullet\bullet}{\mathsf{E}} \\ \overset{\bullet\bullet}{\mathsf{Z}}_{N;i}(f) \\ \overset{\bullet\bullet}{\mathsf{F}} \\ \overset{\bullet\bullet}{\mathsf{I}} \\ \overset{\bullet\bullet}{\mathsf{I} \\ \overset{\bullet\bullet}{\mathsf{I}} \\ \overset{\bullet\bullet}{\mathsf{I}} \\ \overset{\bullet\bullet}{\mathsf{I}} \\ \overset{\bullet\bullet}{\mathsf{I}} \\ \overset{$
- iii) $\sum_{0}^{\geq} a \overline{\log N("; F; d_N)} d" \supseteq_{N-1} 0$ as $\Box = 0$.
- iv) The sequence of covariance functions $cov(Z_{N;i}(f); Z_{N;i}(g))$ converges pointwise on $F \square F$ as N 1 to a non degenerate limit $\square(f;g)$.

Then the sequence $\sum_{i=1}^{\infty} (Z_{N;i}(f) \square E(Z_{N;i}(f)))$ is \square_{P} -equicontinuous and converges in i (F) to a Gaussian process with covariance function $\square(f;g)$.

5.3.1.1 Convergence of the covariance operator

The following intermediary results show that condition iv) in Theorem 5.10 is fulfilled in the particular case of Poisson survey plans. For $(f;g) P F^2$, set

$$\operatorname{cov}_{N;p}(f;g) \coloneqq \frac{1}{N} \prod_{i=1}^{|f|} \frac{f(X_i)}{p_i} \Box_{N;p}(f) \frac{g(X_i)}{p_i} \Box_{N;p}(g) p_i (1 \Box p_i).$$

Due to the independence of the \Box_i 's, it is clear that

$$\operatorname{cov}_{\mathsf{T}_{\mathsf{N}}} \overset{\square}{\mathbb{G}}_{\mathsf{T}_{\mathsf{N}}}^{\mathsf{p}}(\mathsf{f}); \overset{\square}{\mathbb{G}}_{\mathsf{T}_{\mathsf{N}}}^{\mathsf{p}}(\mathsf{g}) \overset{\square}{:=} \operatorname{cov} \overset{\square}{\mathbb{G}}_{\mathsf{T}_{\mathsf{N}}}^{\mathsf{p}}(\mathsf{f}); \overset{\square}{\mathbb{G}}_{\mathsf{T}_{\mathsf{N}}}^{\mathsf{p}}(\mathsf{g}) \overset{\square}{=} (\mathsf{X}_{\mathsf{i}}; \mathsf{W}_{\mathsf{i}})_{1 \\ {\S i } \\ {\S N}} \overset{\square}{=} \operatorname{cov}_{\mathsf{N}; \mathsf{p}}(\mathsf{f}; \mathsf{g}).$$

We thus essentially have to determine conditions ensuring that $cov_{N;p}(f;g)$ has a nondegenerate limit. The following assumptions are by no means necessary but provide a useful framework to derive such conditions. Similar types of assumptions may be found in Bonnéry et al. (2011) or Hàjek (1964) for instance.

Recall that inclusion probabilities were defined relative to some auxiliary variable W. An additional assumption on the latter is required in the subsequent result.

Assumption 5.11 The couples of random vectors $(X_1; W_1); \ldots; (X_N; W_N)$ are iid (exchangeable at least) with distribution $P_{X;W}$. Moreover, the conditional inclusion probabilities $p := (p_1; \ldots; p_N)$ are then given for all i P t1; \ldots; Nu and $W_{(N)} P W^N$ by

$$p_i := p(W_i) := E \square W_{(N)}$$

Remark 5.12 It can happen that p_i not only depends on W_i , but on the entire vector $W_{(N)}$. It is the case, for instance, when there is a unique auxiliary variable W to which weights are proportional:

$$p_i := n \underbrace{\otimes W_i}_{\substack{i=1\\j=1}} W_j$$

In such situations the iid property of the vectors $(X_i; W_i)$, 1 § i § N, can be used to bypass the part involving all $(W_1; :::; W_N)$ in the subsequent asymptotic analysis.

Under this supplementary condition, we have the following result, the proof of which can be found in Section 5.5.1.

Lemma 5.13 – Limit of the covariance operator. Suppose that Assumption 5.1, Assumption 5.8 and Assumption 5.11 are fulfilled and that

0†
$$p^{2}(w) P_{W}(dw) \dagger 1$$
.

а

Then we have

$$\frac{1}{N} d_{N} \bigoplus_{N=1}^{a} D_{p} \coloneqq \bigcup_{W} (1 \Box p(w)) p(w) P_{W}(dw) \circ 0$$

and

$$\operatorname{cov}_{N;p}(f;g) \underset{N-1}{\Box} \Box(f;g);$$

where for all $(f;g) PF^2$,

$$\Box(f;g) \coloneqq \prod_{X \square W}^{a} f(x)g(x) \square \frac{1}{p(w)} \square 1 \square P_{X;W}(dx;dw) \square \square_{p}(f) \square_{p}(g) D_{p}; \quad (5.4)$$

with

with

$$\Box_{p}(f) \coloneqq \frac{1}{D_{p}} \int_{X \square W}^{a} (1 \square p(w)) f(x) P_{X;W}(dx; dw).$$

5.3.1.2 Functional Central Limit Theorem

Applying Theorem 5.10 to the empirical process $\mathbb{G}_{T_N}^p$ f defined in Equation (5.3) thus leads to the theorem below, proved in Section 5.5.2.

Theor em 5.14 – FCLT in the Poisson survey case. Suppose that Assumption 5.1, Assumption 5.8 and Assumption 5.11 hold, as well as the following conditions.

i) Lindeberg-Feller type condition: @ 0,

$$\begin{array}{c} \square \\ \mathsf{E} \end{array}^{(Z_{N};i)}^{2} \mathsf{I} \overset{!}{Z_{N};i} \circ \overset{?}{\square} \overset{?}{\overset{}{\square}} \overset{)}{\overset{}{\square}} \underset{\overset{}{\overset{}{\square}}{\overset{}{\square}} 0; \\ \mathbb{Z}_{N;i} \coloneqq (\square \ \square \ \mathsf{p}(\mathsf{W}_{i})) \ \mathsf{sup}_{\mathsf{f}\,\mathsf{PF}} \overset{\overset{}{\overset{}{\overset{}{\vdash}}} (X_{i})}{\overset{}{\overset{}{\square}} \square \overset{}{\overset{}{\overset{}{\square}}} N_{;\mathsf{p}}(\mathsf{f}) \overset{\overset{}{\overset{}{\square}}}{\overset{}{\overset{}{\overset{}{\square}}} \end{array}$$

ii) Uniform entropy condition: let D be the set of all finitely discrete probability measures defined in Section 5.1.2.3, and assume

$$d_{1} = b_{1} = b_{1} = b_{1} = b_{1} = b_{2;Q};F; ..., b_{2;Q} = b_{1} = b_$$

Then there exists a \Box_P -equicontinuous Gaussian process G in 1 (F) with covariance operator \Box given by Equation (5.4) such that

 ${{\tt G}}^{p}_{{\sf T}_{\sf N}}$ Ò G weakly in `1 (F); as N — 1 .

Remark 5.15 – On the Lindeberg-Feller condition. Observe that, as can be proved using Hölder's inequality, condition i) in Theorem 5.14 can be replaced by the simpler condition: $D\Box^{\circ}$ 0 such that

i)
$$E_{P_{X;W}} \xrightarrow{H(X_i)} E_{T_N} \xrightarrow{(\Box_i \Box p(W_i))^{2+\Box}} (X_i;W_i)^{1/2} + 1.$$

5.3.2 The case of rejective sampling

As shall be shown herein-after, the result obtained above in the case of a Poisson sampling scheme may carry over to more general survey plans, as originally proposed in the seminal contribution of Hajek (1964).

5.3.2.1 Reduction to simpler sampling designs

The lemma stated below, following in the footsteps of Berger (1998), shows that the study of the empirical process related to a general sampling design R_N may be reduced to that related to a simpler sampling design, T_N say, which is close to R_N with respect to some metric and entirely characterized by its first order inclusion probabilities. The only "drawback" is that the estimator involved in this approximation result is not the Horvitz-Thompson estimator, since it does not involves the inclusion probabilities of the sampling plan of interest but those related to a Poisson scheme (Hàjek, 1964). However, as will be shown next, the two estimators may asymptotically coincide, as N tends to +1.

In order to formulate the approximation result needed in the sequel, we introduce, for two sampling designs R_N and T_N , the total variation metric

$$\{R_{N} \Box T_{N} \}_{1} \coloneqq \frac{\prod_{s \in P(U_{N})} |R_{N}(s) \Box T_{N}(s)| }{s \in P(U_{N})}$$

as well as the entropy

$$\mathsf{D}(\mathsf{T}_{\mathsf{N}};\mathsf{R}_{\mathsf{N}}) \coloneqq \frac{\prod_{\mathsf{s}\mathsf{PP}(\mathsf{U}_{\mathsf{N}})} \mathsf{T}_{\mathsf{N}}(\mathsf{s}) \log^{-1} \frac{\mathsf{T}_{\mathsf{N}}(\mathsf{s})}{\mathsf{R}_{\mathsf{N}}(\mathsf{s})}^{-1}.$$

In practice, T_N will typically be the Poisson sampling plan investigated in the previous subsection and $G_{T_N}^{\Box(T_N)}$ the corresponding empirical process.

Lemma 5.16 – Approximation result. Let R_N and T_N be two sampling designs and assume that T_N is entirely characterized by its first order inclusion probabilities, $\Box(T_N)$. Then, the empirical processes $\mathbb{G}_{T_N}^{\Box(T_N)}$ and $\mathbb{G}_{R_N}^{\Box(T_N)}$ valued in `1 (F) satisfy the relationships:

$$\mathsf{d}_{\mathsf{BL}} \overset{\Box}{\mathfrak{G}}_{\mathsf{T}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})}; \overset{\Box}{\mathfrak{G}}_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})} \overset{\Box}{\$} \ \mathsf{R}_{\mathsf{N}} \ \Box \ \mathsf{T}_{\mathsf{N}} \,\mathsf{I}_{\mathsf{1}} \, \, \mathsf{\$}^{\mathsf{a}} \, \, \overline{\mathsf{2D}(\mathsf{T}_{\mathsf{N}}\,; \mathsf{R}_{\mathsf{N}})}.$$

Consequently, if the sequences $(R_N)_{N \cdot 1}$ and $(T_N)_{N \cdot 1}$ are such that $R_N \square T_N_1$ tends to 0 or $D(T_N; R_N) - 0$ as N - 1 and if there exists a Gaussian process G such that

$$d_{\mathsf{BL}}(\mathbf{\check{G}}_{\mathsf{T}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})};\mathsf{G}) \bigsqcup_{\mathsf{N}=1} 0;$$

then we also have

$$\mathsf{d}_{\mathsf{BL}}(\mathbf{G}_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})};\mathsf{G}) \square_{\mathsf{N}=\mathsf{I}} 0.$$

The same result holds true when replacing $G_{T_N}^{\square(T_N)}$ and $G_{R_N}^{\square(T_N)}$ by $G_{T_N}^{\square(T_N)}$ and $G_{R_N}^{\square(T_N)}$ respectively.

This result, proved in Section 5.5.3, reveals that as soon as a possibly complicated survey design R_N can be approximated by a simpler one T_N through some coupling argument ensuring that the $\}.\}_1$ distance between them decays to zero (as in Hàjek, 1964), then an asymptotic approximation result possibly holding true for the empirical process related to T_N immediately extends to that related to R_N , when built with the inclusion probabilities $p = \Box(T_N)$. As shall be seen below, a typical situation where this result applies corresponds to the case where R_N is a rejective sampling design, while T_N is a Poisson sampling design, as in Hàjek (1964). Other natural applications may arise in the framework of post-stratification, which can be connected with empirical likelihood results.

5.3.2.2 Empirical process for the rejective sampling and its variants

The Central Limit Theorem for rejective sampling and some variants of this survey scheme has been studied at length in Hàjek (1964) and Berger (1998, 2011). Consider the rejective sampling scheme defined in Example 5.5 from a given vector \Box^R corresponding to the vector $p := (p_1; \ldots; p_N) = (p(W_1); \ldots; p(W_N)) =: p(W)$. Assume in addition that the representation is canonical, i.e. is such that $\prod_{i=1}^{N} p(W_i) = n$. The key argument for proving a CLT in the rejective sampling case consists in exhibiting a certain coupling $((\Box_1; \ldots; \Box_N); (\Box_1^\Box; \ldots; \Box_N^\Box))$ of the Poisson sampling scheme with inclusion probabilities $p(W_1); \ldots; p(W_N)$ and the rejective sampling scheme with corresponding inclusion probabilities \Box^R such that $R_N \Box T_N = 0$, see Hàjek (1964, p. 1503-1504) for further details. A straightforward application of Lemma 5.16 will then immediately yield a functional CLT in our framework. We point out that, under the rejective sampling scheme, the survey size is fixed, so that

$$\bigcap_{i=1}^{n} (\Box_i \Box p(W_i)) = n \Box n = 0.$$

Thus, we have:

$$\mathbf{G}_{\mathsf{R}_{\mathsf{N}}}^{\mathsf{p}} \mathsf{f} \coloneqq \frac{2^{\mathsf{1}}}{\mathsf{N}} \prod_{i=1}^{\mathsf{N}} (\Box_{i} \Box_{\mathsf{p}}(\mathsf{W}_{i}))^{\mathsf{D}} \frac{\mathsf{f}(\mathsf{X}_{i})}{\mathsf{p}(\mathsf{W}_{i})} \Box_{\mathsf{N};\mathsf{p}}(\mathsf{f})^{\mathsf{D}}$$

$$= \frac{2^{1}}{\overline{N}} \int_{i=1}^{|V|} \frac{\overline{Q}_{i}}{p(W_{i})} = 1 f(X_{i})$$
$$=: G_{R_{N}}^{p(W)} f.$$

Hence, the Poisson-like empirical process coincides, in that case, with the original HT-empirical process where the weights p(W) are involved instead of the true inclusion probabilities \Box^R , the latter being however asymptotically equivalent to the former, see Hàjek (1964).

The next theorem is obtained by combining Lemma 5.16 with Theorem 5.14.

Theor em 5.17 – FCLT in the rejective survey with Poisson weights case. Suppose that Assumption 5.1, Assumption 5.8, Assumption 5.11 and conditions i) and ii) of Theorem 5.14 are satisfied. Then, there exists a \Box_P -equicontinuous Gaussian process G in 1 (F) with covariance operator \Box given by Equation (5.4) such that

$$G_{R_{N}}^{p(W)}$$
 Ò G weakly in ¹ (F); as N – 1.

It has been established in Berger (1998) that for a variety of sampling plans R_N , including the Rao-Sampford scheme defined in Example 5.6, we have $D(R_N; T_N) - 0$ as N - 1. By virtue of Lemma 5.16, Theorem 5.17 naturally extends to these sampling schemes.

Going back to the original HT-empirical process in Equation (5.2) related to the plan R_N , the corollary below reveals that the asymptotic result still holds true for the latter (see the proof in Section 5.5.4). This essentially follows from the fact that the weights p(W) and the inclusion probabilities corresponding to the rejective sampling are asymptotically equivalent.

Corollary 5.18 – FCLT in the rejective survey case. Suppose that Assumption 5.1, Assumption 5.8, Assumption 5.11 and conditions i) and ii) of Theorem 5.14 are satisfied. Then, there exists a \Box_P -equicontinuous Gaussian process G in `1 (F) with covariance operator \Box given by Equation (5.4) such that

 $G_{R_{M}}^{\square(R_{N})}$ Ò G weakly in `1 (F); as N — 1.

5.4 application to non-parametric statistics

For illustration purpose, we consider now several statistical applications of the asymptotic results previously established.

5.4.1 Hadamard differentiable functionals

We first highlight that the FCLT stated above permits to establish the asymptotic normality of any statistic that can be expressed as the empirical version of some Hadamard differentiable functional, see Shorack and Wellner (1986). For the sake of clarity, we recall the definition of uniform Hadamard differentiability in Definition 5.19, adapted from Pons and de Turkheim (1991). Our results apply to many situations considered in their paper, related in particular to certain functionals of censored data. Other examples are treated in Gill (1989), van der Vaart and Wellner (1996) (see Chapter 3.9 p. 379 therein, in particular refer to the discussion about the validity of the bootstrap for uniform Hadamard differentiable functionals). Define B(F; P) as the set of measures Q in 1 (F) whose paths f P F fi- Qf := fdQ are }.}_{2;P}-uniformly continuous and bounded. This is the smallest natural space containing G. We consider the uniform Hadamard differentiability tangentially to the subspace B(F; P) because it weakens the notion of differentiability tangentially to check in practice.

Definition 5.19 A functional T:¹ (F) — R^q is said to be uniformly Hadamard differentiable at P tangentially to B(F; P); if and only if there exists a continuous linear mapping dT_P such that for any sequence P_N converging to P, any h_N converging to h PB(F; P) and every t_N converging to 0 such that P_N + t_N.h_N P¹ (F), we have:

$$\frac{\mathsf{T}(\mathsf{P}_{\mathsf{N}} + \mathsf{t}_{\mathsf{N}}.\mathsf{h}_{\mathsf{N}}) \Box \mathsf{T}(\mathsf{P}_{\mathsf{N}})}{\mathsf{t}_{\mathsf{N}}} \Box \mathsf{d}\mathsf{T}_{\mathsf{P}}.\mathsf{h} \underset{\mathsf{t}_{\mathsf{N}} = 0}{\Box = 0} \mathsf{0}.$$

Notice that T may be defined not on the entire space 1 (F) but on a subset L only. In this case, one must check that $P_N + t_N h_N PL$.

Remark 5.20 We may in addition assume that the differential dT_P admits an integral representation, i.e.

$$dT_{\rm P}.h = T^{(1)}(x; P) h(dx);$$

where $T^{(1)}(.; P)$ is the influence function defined from X to B_1 such that we have

$$\mathsf{E}_{\mathsf{P}} \stackrel{\square}{\mathsf{T}^{(1)}}(\mathsf{X};\mathsf{P}) = \mathbf{0}.$$

We recall that in the robustness terminology (Hampel et al., 1986), the influence function of the parameter T(P) may be calculated directly by computing the derivative of the functional taken at the contaminated distribution $(1 \square t)P + t \square_{x}$, i.e.

$$\mathsf{T}^{(1)}(\mathsf{x};\mathsf{P}) \coloneqq \lim_{\mathsf{t} \to 0} \frac{\mathsf{T}((1 \Box \mathsf{t})\mathsf{P} + \mathsf{t}\Box_{\mathsf{x}}) \Box \mathsf{T}(\mathsf{P})}{\mathsf{t}}$$

In this case, the limiting distribution may be calculated more easily.

Theorem 5.21 – CLT for Hadamard differentiable functionals. Suppose that the assumptions of Theorem 5.14 hold and that functional T : L Ä 1 (F) — R^q is Hadamard differentiable at P with differential dT_P and influence function T⁽¹⁾(x; P). Then, as N — +1 , we have:

$$\frac{2}{N} \prod_{R_N} (P_{R_N}^{\square(R_N)}) \square T(P_N) \stackrel{\square}{} \dot{O} dT_P.G;$$

where G is a Gaussian process with covariance operator \Box , as in Equation (5.4).

The result above, the proof of which is available in Section 5.5.5, applies in particular to the following statistics.

Example 5.22 – Expectation and variance. It is well-known that for some appropriate choice of F, the functionals $T(P) = E_P(X)$ and $T(P) = V_P(X)$ are uniformly Hadamard-differentiable. When $T(P) = E_P(X)$, Theorem 5.21 exactly reduces to the Central Limit Theorem established in Hàjek (1964).

Example 5.23 – Cumulative distribution function. In a univariate setting, the functional $T(P) = F(x) := P(XP(\Box 1; x])$ can be dealt with by simply considering the class of indicator functions u fi– I tu § xu with x PR and applying next Theorem 5.17 and Corollary 5.18. We provide illustrations of this specific example in Section 5.4.3.

5.4.2 Fréchet differentiable functionals

Hadamard differentiability is sometimes difficult to prove and it does not yield a precise control of the remainder for further approximations like Edgeworth expansions. Another approach followed by Dudley (1990) and Barbe and Bertail (1995) is to assume Fréchet differentiability with respect to a metric d_F indexed by a class of function F, for which some uniform entropy conditions hold. A functional is said to be Fréchet differentiable at P for such a metric if there exists a gradient (for instance the influence function $T^{(1)}(x; P)$, which fulfills E_P $T^{(1)}(x; P) = 0$ and a continuous function "(.), null at 0, such that for any probability Q,

$$T(Q) \ \Box \ T(P) = \quad T^{(1)}(x;P) \ (Q \ \Box \ P)(dx) + \ d_F \ (Q;P) \ " \ (d_F \ (Q;P)).$$

It is generally possible to choose the class of functions according to the functional of interest, see for instance Arcones and Giné (1992) for general classes of M-estimators. Notice that in that case, by applying Fréchet differentiability twice, we have

$$\widehat{\mathsf{N}}^{\square}\mathsf{T}(\mathsf{P}_{\mathsf{R}_{\mathsf{N}}}^{\square(\mathsf{R}_{\mathsf{N}})}) \square \mathsf{T}(\mathsf{P}_{\mathsf{N}}) \stackrel{\square}{=} \widehat{\mathsf{N}}^{\square} \mathsf{T}^{(1)}(x;\mathsf{P}) (\mathsf{P}_{\mathsf{R}_{\mathsf{N}}}^{\square(\mathsf{R}_{\mathsf{N}})} \square \mathsf{P}_{\mathsf{N}})(dx) + \mathsf{r}_{\mathsf{N}}$$

$$= \frac{21}{\overline{N}} \prod_{i=1}^{|V|} \frac{\Box_i}{\Box_i} T^{(1)}(X_i; P) + r_N;$$

with a remainder

$$r_{N} = \frac{?}{N} d_{F}(P_{R_{N}}^{\Box(R_{N})}; P) "(d_{F}(P_{R_{N}}^{\Box(R_{N})}; P)) + \frac{?}{N} d_{F}(P_{N}; P) "(d_{F}(P_{N}; P)))$$

By virtue of the results in Hàjek (1964), it is then obvious that the linear term in this approximation is asymptotically Gaussian with known variance. Controlling the remainder essentially amounts to controlling the behavior of $\overline{N} d_F(P_{R_N}^{\Box(R_N)}; P)$ or alternatively, by the triangular inequality, that of $\overline{N} d_F(P_{R_N}^{\Box(R_N)}; P_N)$, which was the purpose of Section 5.2 and Section 5.3.

5.4.3 Simulation-based Gaussian asymptotic confidence regions

A straightforward application consists in the building of Gaussian confidence regions for the (univariate) empirical cumulative distribution function in the entire population, denoted by $F_N(x)$, x PR, when the survey scheme is of the rejective type. Indeed, consider the class of functions $F := tf_x(.) := It.$ § xu; x P Ru. Provided Assumption 5.1 is fulfilled, it respects the required conditions for Corollary 5.18 to hold (see Van der Vaart, 2000, Example 19.16 for the uniform entropy condition and take H(x) = 1 and $\Box = 1$ when checking condition i^{\Box}) in Remark 5.15), which implies in particular that $\{G_{R_N}^{\Box(R_N)}\}_F$ converges in distribution to $\{G\}_F$ as N — +1 (Van der Vaart, 2000, Corrolary 19.21). This yields the following asymptotic uniform confidence band of level $\Box P(0; 1)$ for the population cdf F_N :

$$\mathsf{CB}_{\square} := \overset{\square}{\mathsf{F}_{\mathsf{R}_{\mathsf{N}}}^{\square}(\mathsf{R}_{\mathsf{N}})} \square \overset{\mathfrak{A}_{\square}}{\overset{\square}{\overset{\square}{\overset{\square}{\mathsf{N}}}}}; \, \mathsf{F}_{\mathsf{R}_{\mathsf{N}}}^{\square}(\mathsf{R}_{\mathsf{N}}) + \overset{\mathfrak{A}_{\square}}{\overset{\square}{\overset{\square}{\mathsf{N}}}};$$

where $F_{R_N}^{\Box(R_N)}$ is the Horvitz-Thompson estimator of the cdf based on the rejective sample and q_{\Box} the \Box -quantile of random variable $\{G\}_F$. Since in practice q_{\Box} is unknown, it needs to be estimated. It can be achieved by means of Monte-carlo simulations, using a simple technique based on the Cholesky decomposition of the covariance matrix (Kroese et al., 2011, Algorithm 5.1).

Algorithm 5.24 – Simulation of the limit process G and estimation of q_{\Box} .

 Choose a grid of real values tx₁;:::;x_Ku, K ° 1, and compute the Horvitz-Thompson estimator of □(f_{xk};f_{xk□}) for each couple (x_k;x_{k□}) in tx₁;:::;x_Ku², namely

$$\Box_{\mathsf{R}_{\mathsf{N}}}^{\mathsf{p}}(\mathsf{f}_{\mathsf{x}_{k}};\mathsf{f}_{\mathsf{x}_{k}^{\Box}}) \coloneqq \frac{1}{\mathsf{N}} \prod_{i=1}^{|\mathsf{N}|} \Box_{i} \frac{1 \Box p_{i}}{p_{i}^{2}} \mathsf{I} \mathsf{t} \mathsf{X}_{i} \ \S \ \mathsf{min}(\mathsf{x}_{k};\mathsf{x}_{k^{\Box}})\mathsf{u}$$

$$\Box \frac{1}{N} \frac{\prod_{i=1}^{N} \Box_{i} \frac{1 \Box p_{i}}{p_{i}} | t X_{i} \S x_{k} u \overset{T}{\underset{i=1}{N}} \Box_{i} \frac{1 \Box p_{i}}{p_{i}} | t X_{i} \S x_{k} \Box u}{\prod_{i=1}^{N} \Box_{i} (1 \Box p_{i})}.$$

- 2. Derive the Cholesky decomposition $\Box_{R_N}^p(f_{x_k}; f_{x_k \Box}) = LL^1$, where L is a lower-triangular Cholesky matrix.
- Generate B ° 1 independent copies Y₁; ...; Y_B of the Gaussian random vector Y := (Y₁; ...; Y_K)¹ with null expectation and covariance I, the identity matrix.
- Compute Z_b := LY_b, b Pt1;:::;Bu, which are considered as realizations of the limit process G.
- 5. For each b P t 1;:::; Bu, calculate $Z_b^{\Box} \coloneqq \{Z_b\}_1$ the maximum absolute distance to 0 of the path Z_b and sort the obtained sample $Z_{1:B}^{\Box}$ § $\Box \Box$ § $Z_{B:B}^{\Box}$.
- 6. Set $\hat{q}_{\square} := Z_{\mathfrak{B}_{\square u+1;B}}^{\square}$, the empirical \square -quantile of the sample of maximum absolute deviations, where tudenotes the floor function.

In the next subsections, a set of numerical experiments is performed to provide illustrative examples of this technique.

5.4.3.1 Experiment setting

Simulations were based on the following model, chosen for its simplicity in terms of both computation and interpretation:

$$X = \Box W + U; \Box Pt0; 1u$$

$$W ; TN(\Box; \Box_W^2; w_{\Box}; w^{\Box});$$

$$U ; N(0; \Box_U^2);$$

$$P(W \S w; U \S u) = P(W \S w) P(U \S u);$$

where X is the variable of interest, W the auxiliary information, U a white noise independent from W, and TN(0; \Box_W^2 ; w ; w) refers to the truncated Normal distribution over [w; w], with expectation \Box and variance \Box_W^2 . Such a representation enables a simple control of the dependence between X and W, since their correlation is then

$$\operatorname{corr}(X; W) = \Box = \frac{\Box_W}{\Box_W^2 + \Box_U^2}.$$

For a given population U_N of size N, where it is assumed that $t W_i$; i $P U_N u$ (resp. $t U_i$; i $P U_N u$) are independent (hence exchangeable) realizations of W (resp. U), inclusion probabilities of the Poisson sampling scheme are defined as

$$p_i = p(W_i) = n \frac{W_i}{\sum_{j=1}^{N} W_j};$$
 (5.5)

with n the desired expected sample size (Hàjek, 1964, Section 6, p.1512). When the inclusion probabilities are proportionate to the auxiliary variable like in Equation (5.5), the stronger the correlation between X and W, the smaller the variance of the estimator of the population mean $\frac{1}{N} \sum_{i=1}^{\infty} X_i$ (or, equivalently, of the total $\sum_{i=1}^{\infty} X_i$). Recall that under Assumption 5.1, we have n=N — c P (0; 1) as both n and N tend to infinity. Hence, p_i can be viewed as the empirical version in the population of

$$\mathsf{p}(\mathsf{W}) \coloneqq \mathsf{W} \ \frac{\mathsf{c}}{\mathsf{E}(\mathsf{W})}.$$

Observe that thus defined, $p(W) P [p_{\Box}; p^{\Box}]$, where $p_{\Box} = c w_{\Box} = 0$ and $p^{\Box} = c w^{\Box} = 0$, which offers an easy way of ensuring Assumption 5.1 is fulfilled.

Numerical experiments were conducted on a set of populations with increasing sizes N = 10², 5 \square 10², 10³, 5 \square 10³ and 10⁴. Though the latter may seem quite small to study asymptotic properties, they are in fact representative of many practical situations, where populations under the microscope have moderate sizes in comparison to nationwide surveys. Several scenarios were investigated depending on both the variance parameter \square_{U}^{2} and the coefficient \square , so as to cover situations where corr(X; W) is high, low or null. They are summarized in Table 5.1. For each scenario, two sample sizes were considered: one small with n = 0.1 \square N and one relatively large with n = 0.5 \square N. Parameters of the distribution of W were chosen to ensure that for all i P U_N, p_i P [0.01; 1]. Specifically, we set \square = 1, \square_{W}^{2} = 0.09, w $_{\square}$ = 0.1 and w \square = 2, thereby implying that (p $_{\square}$; p \square) = (0.01; 0.02) when n = 0.1 \square N and (p $_{\square}$; p \square) = (0.05; 1) when n = 0.5 \square N.

Scenario		\Box_{U}^{2}	corr(X;W)
S ₁	1	0.01	0.95
S ₂	1	35.91	0.05
S3	0	35.91	0

Table 5.1 – List of scenarios depending on \Box and \Box_{U}^{2} , and corresponding model characteristics

For each scenario, we drew 1000 samples according to a rejective sampling scheme, following Algorithm 5.9 in Tillé (2006). The true inclusion probabilities, denoted by \Box_i , 1 § i § N, could have been deduced from their Poisson equivalents defined in Equation (5.5) using Formula (5.13) in Tillé (2006). Though a very popular and natural algorithm, due to the limits of computer precision, the successive approximations it involves can lead to unexpected results like negative inclusion probabilities. Other algorithms have been developed to compute exact inclusion probabilities in a rejective sampling scheme (eg. Aires, 1999 or Tillé, 2006, Algorithms 5.8 and 5.9), but again, especially when N is large, they have a tendency to produce illogical estimates. This is why we adopted a simpler, although computationally expensive Monte-Carlo approximation technique, based on the repetition (10⁵ times) of the basic algorithm stated in Example 5.5. Notice that since rejective sampling is a Poisson sampling conditioned upon its size, we have ($p_i = 1$) \dot{O} ($\Box_i = 1$).

We constructed asymptotic uniform 95% confidence bands of the population cdf F_N using Algorithm 5.24, with B = 1000 and K = 10, 20 or 100 depending on the sample size n. More precisely, for n • 100, the grid was made of the standard empirical percentile estimators of variable X based on each artificial samples. To enable computation of the lower-triangular Cholesky matrix, which requires that the covariance matrix has full rank, we confined ourselves to deciles for n = 10 and to the quantiles of levels 0.05; 0.1; :::; 0.95; 1 for n = 50.

5.4.3.2 Experiment results

The average and maximal width of the confidence bands over the 1000 simulated samples for each scenario are given in Table 5.2. Coverage probabilities were also estimated, the results of which are displayed in Table 5.3. Finally, some graphical illustrations are provided in Figure 5.1.

Ν 10² 5 🗆 10² 10³ 5 🗆 10³ 10⁴ Scenario c Av Мx Αv Мx Αv Мx Av Мx Αv Мx 0.5 31.24 37.12 15.68 20.70 11.43 14.61 5.10 5.76 3.58 3.94 Sı 87.68 116.62 42.67 68.80 32.67 45.64 10.19 12.33 0.1 14.46 18.15 27.74 33.55 13.77 16.52 9.95 12.88 4.48 5.05 3.56 0.5 3.17 S_2 0.1 76.43 103.01 37.63 53.10 28.40 40.36 12.70 15.95 8.99 10.81 0.5 27.53 32.79 13.67 16.12 9.85 13.36 4.43 5.13 3.14 3.41 S3 0.1 75.65 105.18 37.20 53.93 27.93 38.51 12.57 15.92 8.90 10.61

Table 5.2 – Average (Av) and maximal (Mx) width of confidence bands multiplied by 100

Table 5.3 – Estimated coverage probabilities (in %)

		Ν						
Scenario	С	10 ²	5 🗆 10 ²	10 ³	5 □ 10 ³	10 ⁴		
S.	0.5	96.30	96.80	96.80	97.80	97.00		
9	0.1	94.74	97.20	97.12	98.30	98.50		
C .	0.5	92.50	97.30	95.60	96.40	96.40		
52	0.1	89.12	93.80	93.55	95.60	96.90		
<u>د</u>	0.5	91.80	96.50	95.50	96.10	95.60		
J 3	0.1	87.82	93.50	93.48	94.80	96.70		



Figure 5.1 – Example of the 95% confidence bands of the empirical distribution function in the population F_N (black line) constructed on one of the 1000 simulated samples under scenario S_1 with c = 0.1 (dark pink area) and c = 0.5 (light pink area) for $N = 5 \square 10^2$ (left hand plot) and $N = 10^4$ (right hand plot)

As expected, the larger N and c, the smaller the confidence bands. Regarding coverage probabilities, they appear to be close to 95%, the desired level, for any N and c. The most remarkable variability is that observed between scenarios: confidence bands get significantly tighter as the correlation between X and W decreases. As a consequence, estimated coverage probabilities are systematically smaller in scenarios S₂ and S₃ than in scenario S₁, especially when N = 10^2 and c = 0.1. This phenomenon is due to the formula used to construct inclusion probabilities, in Equation (5.5). Let us dwell for a moment on this expression. It ensures that the Horvitz-Thompson estimator (based on the Poisson inclusion probabilities) of the expectation of W coincides with the classical empirical mean in the entire population:

$$\frac{1}{N} \prod_{i=1}^{|\mathbf{n}|} \frac{\Box_i}{p_i} W_i = \frac{1}{N} \prod_{i=1}^{|\mathbf{n}|} \frac{\Box_i}{n \frac{\infty W_i}{\sum_{j=1}^{N} W_j}} W_i = \frac{1}{N} \prod_{i=1}^{|\mathbf{n}|} \frac{\Box_i}{n} \prod_{j=1}^{|\mathbf{n}|} W_j = \frac{1}{N} \prod_{j=1}^{|\mathbf{n}|} W_j;$$

since $\sum_{i=1}^{N} \square_i = n$ by definition. It is no surprise then that the stronger the correlation between X and W, the closer (in terms of variance) the weighted mean $\frac{1}{N} \sum_{i=1}^{N} \square_{p_i} X_i$ is to its population counterpart. However, when considering empirical distribution functions, the standard and sample estimators for W are no longer equal. Hence, not only does the model in Equation (5.5) fail to improve the variance of the HT-cdf of X, but the deviations of $F_{R_N}^p$ are expected to grow as the link between X and W tightens. To counterbalance this drawback, we could for instance choose the inclusion probabilities p_i , 1 § i § N, that minimize the uniform difference between the HT and the empirical cdf of W (see for instance Rueda et al., 2007). Such refinements are left for further research. Although not optimal, the confidence bands constructed on our numerical experiments are still satisfactory and advocate the utility of our asymptotic results whatever the available inclusion probabilities.

5.5 proofs and supplements

5.5.1 Limit of the covariance operator

We shall prove Lemma 5.13, the statement of which is recalled below.

Lemma Suppose that Assumption 5.1, Assumption 5.8 and Assumption 5.11 are fulfilled and that

Then we have

$$\frac{1}{N} d_{N} \bigoplus_{N=1}^{a} D_{p} \coloneqq \bigoplus_{W}^{a} (1 \Box p(w)) p(w) P_{W}(dw) \circ 0$$

and

$$\operatorname{cov}_{N;p}(f;g) \underset{N-1}{\Box} (f;g)$$

where for all $(f;g) PF^2$,

$$\Box(\mathsf{f};\mathsf{g}) \coloneqq \left(f(\mathsf{x}) \mathsf{g}(\mathsf{x}) \right) = \left(f(\mathsf{x}) \mathsf{g}(\mathsf{x}) \right) \left(\frac{1}{\mathsf{p}(\mathsf{w})} \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) \left(f(\mathsf{w}) \right) = \left(f(\mathsf{w}) \right$$

with

$$\Box_{p}(f) \coloneqq \frac{1}{D_{p}} (1 \Box p(w)) f(x) P_{X;W}(dx; dw).$$

Proof Notice first that under Assumption 5.11, we have: @(f;g) PF²,

а

$$cov_{N;p}(f;g) = \frac{1}{N} \frac{\prod_{i=1}^{N} \frac{f(X_i)}{p(W_i)} \frac{g(X_i)}{p(W_i)} p(W_i) (1 \Box p(W_i))}{\Box_{N;p}(f) \Box_{N;p}(g) \frac{1}{N} \frac{\prod_{i=1}^{N} p(W_i) (1 \Box p(W_i))}{i=1}}$$

with

Now it is sufficient to apply the Strong Law of Large Numbers for exchangeable vectors to obtain that

$$\frac{1}{N} d_{N} = \frac{1}{N} \prod_{i=1}^{N} (1 \Box p(W_{i})) p(W_{i})$$

а

$$\square_{N-1} (1 \square p(w)) p(w) P_{W}(dw) \text{ almost-surely.}$$

The limit above is finite, positive under Assumption 5.1 (that implies there exists $p_{\Box} \circ 0$ such that $p(w) \circ p_{\Box}$). Additionally, we have with probability one

By virtue of Assumption 5.8, the latter integral is finite. Finally, observe that we almost-surely have

$$\Box_{N;p}(f) \Box_{N-1} \Box_{p}(f);$$

and the desired result follows. In particular notice that the limiting variance $V^2(f)$ is given by

$$V^{2}(f) \coloneqq \left[\begin{array}{c} a \\ X \boxtimes W \end{array} \right]_{a} f(x)^{2} \left[\begin{array}{c} 1 \\ p(w) \end{array} \right]_{w} \left[\begin{array}{c} 1 \\ p(w) \end{array} \right]_{w} P_{X;W}(dx;dw)$$
$$\Box_{p}(f)^{2} \\ W \\ W \\ \end{array} \right]_{w} (1 \boxtimes p(w)) p(w) P_{W}(dw);$$

which is strictly positive except in the degenerate case where f(x) = p(w). Typically, this occurs when the inclusion probabilities are based directly on the variable of interest (or W = c \square X for some c P R). Positivity of the operator results from Cauchy-Schwarz inequality.

5.5.2 FCLT in the Poisson survey case

We shall prove Theorem 5.14, the statement of which is recalled below.

Theorem Suppose that Assumption 5.1, Assumption 5.8 and Assumption 5.11 hold, as well as the following conditions.

i) Lindeberg-Feller type condition: @_ ° 0,

with

$$\begin{array}{c} \square \\ \mathsf{E} \end{array}^{(Z_{N};i)^{2}} \mathsf{I} \overset{!}{Z_{N};i} \circ \quad \square \overset{?}{\mathsf{N}} \overset{)}{\underset{N \to 1}{\overset{\square}{\longrightarrow}}} 0; \\ \\ \mathsf{Z}_{N;i} \coloneqq (\square \mathsf{p}(\mathsf{W}_{i})) \operatorname{sup}_{\mathsf{f}\,\mathsf{PF}} \overset{f(X_{i})}{\underset{p(\mathsf{W}_{i})}{\overset{\square}{\longrightarrow}}} \square \overset{(\mathsf{N};p(\mathsf{f})) \overset{\square}{\overset{\square}{\longrightarrow}} \end{array}$$

- ii) Uniform entropy condition: let D be the set of all finitely discrete probability measures defined in Section 5.1.2.3, and assume
 - $\begin{array}{cccc} a & & & & \\ & & & sup & \\ & & & sup & log(N("\}H\}_{2;Q};F;\}.\}_{2;Q}) d" \ \ \ 1 \ . \end{array}$

Then there exists a \Box_P -equicontinuous Gaussian process G in 1 (F) with covariance operator \Box given by Equation (5.4) such that

$$\mathbf{G}_{T_{N}}^{p}$$
 $\hat{\mathbf{O}}$ G weakly in $\hat{\mathbf{I}}$ (F); as N — 1

Proof We essentially have to check hypotheses i) \Box iv) of Theorem 5.10.

Concerning hypothesis i), the Lindeberg-Feller condition can be written as

$$\frac{1}{N} \prod_{i=1}^{[N]} E \xrightarrow{Z_{N;i}^{2}} I \xrightarrow{I} Z_{N;i}^{2} \circ \xrightarrow{?} \overrightarrow{N} \xrightarrow{\square} 0 \text{ for every } \square \circ 0;$$

which reduces to E $Z_{N;i}^2$ I $Z_{N;i} \circ \bigcap^2 \overline{N} \bigcap_{N=1}^{(\Box)} 0$ by exchangeability of the components. This corresponds to condition i) in Theorem 5.14 above.

Recall that Assumption 5.8 stipulates the envelope of class F is square-integrable function H and that under Assumption 5.1, there is some p_{\Box} ° 0 such that for all i PU_N , $p(W_i) \bullet p_{\Box}$. Hence, we have

$$\sum_{i=1}^{N} p(W_i) (1 \Box p(W_i) \bullet p_{\Box}(N \Box E(n)) = p_{\Box} N \prod_{i=1}^{D} \frac{E(n)}{N}$$

as well as

$$\square_{N;p}(f) \stackrel{\square}{\Longrightarrow} \square_{N;p}(H) \S \frac{1}{p_{\square}} \frac{1}{N \square E(n)} \prod_{i=1}^{[N]} H(X_i) \dagger 1.$$

We thus obtain:

i

$$\sup_{f \in F} \left[\frac{f(X_i)}{p(W_i)} \square_{N;p}(f) \right]^2 \S 2 \quad \sup_{f \in F} \left[\frac{f(X_i)}{p(W_i)} \right]^2 + \sup_{f \in F} \left[\frac{f(X_i)}{p(W_i)} \right]^2 + \frac{1}{p(H)} \left[\frac{1}{2} \right]^2 \\ \$ 2 \quad \left[\frac{H(X_i)}{p(W_i)} \right]^2 + \left[\frac{1}{N} \right]_{p(H)} \left[\frac{1}{2} \right]^2 .$$

Set

$$G_{1;i} = |\Box_i \Box_p(W_i)|^2 + H(X_i) + W_{N;p}(H)^2;$$

$$G_{2;i} = (\Box_i \Box_p(W_i))^2 + H(X_i)^2 + W_{N;p}(H)^2.$$

Observe that it is thus sufficient to check that

$$\begin{array}{c} \square & \square & P_{X;W} \\ E_{\mathsf{P}_{X;W}} & E_{\mathsf{T}_{\mathsf{N}}} & G_{2;i} \\ \mathbf{I} & G_{1;i} \\ ^{\circ} & \square \\ \end{array} \xrightarrow{?} \begin{array}{c} \square \\ \mathbb{N} \\ \end{array} \xrightarrow{} \begin{array}{c} \square \\ \mathbb{N} \\ \mathbb{N}$$

Condition ii) can be checked immediately by noticing that, in the case of the Poisson process, the equicontinuity condition becomes

In practice, condition iii) is checked in an easier manner by using the uniform entropy condition given here, see also Lemma 2.11.6 in van der Vaart and Wellner (1996).

Finally, condition iv) is a direct consequence of Lemma 5.13.

5.5.3 Approximation result

We shall prove Lemma 5.16, the statement of which is recalled below.

Lemma Let R_N and T_N be two sampling designs and assume that T_N is entirely characterized by its first order inclusion probabilities, $\Box(T_N)$. Then, the empirical processes $G_{T_N}^{\Box(T_N)}$ and $G_{R_N}^{\Box(T_N)}$ valued in `1 (F) satisfy the relationships:

$$d_{\mathsf{BL}} \overset{\Box}{\bullet} G_{\mathsf{T}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})}; G_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})} \overset{\Box}{\$} \ \} \mathsf{R}_{\mathsf{N}} \ \Box \ \mathsf{T}_{\mathsf{N}} \ \! \}_{\mathsf{1}} \ \$ \ \frac{\mathsf{a}}{\mathsf{2}\mathsf{D}(\mathsf{T}_{\mathsf{N}}\,;\mathsf{R}_{\mathsf{N}})}.$$

Consequently, if the sequences $(R_N)_{N \cdot 1}$ and $(T_N)_{N \cdot 1}$ are such that $R_N \square T_N_1 - 0$ or $D(T_N; R_N) - 0$ as N - 1 and if there exists a Gaussian process G such that

$$\mathsf{d}_{\mathsf{BL}}(\mathbf{G}_{\mathsf{T}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})};\mathsf{G}) \bigsqcup_{\mathsf{N}=1} 0;$$

then we also have

$$d_{\mathsf{BL}}(\mathbf{G}_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{T}_{\mathsf{N}})};\mathsf{G}) \underset{\mathsf{N}=1}{\Box} 0.$$

The same result holds true when replacing $G_{T_N}^{\square(T_N)}$ and $G_{R_N}^{\square(T_N)}$ by $G_{T_N}^{\square(T_N)}$ and $G_{R_N}^{\square(T_N)}$ respectively.

Proof Let
$$b P BL_1(^1(F))$$
. We have

$$E_{T_N} \stackrel{\circ}{b} G_{R_N}^{(T_N)} \stackrel{\circ}{=} E_{R_N} \stackrel{\circ}{b} G_{T_N}^{(T_N)} \stackrel{\circ}{=} = \prod_{\substack{s P P(U_N) \\ \Box}} R_N(s) b \stackrel{\circ}{G_{T_N}^{(T_N)}} \stackrel{\circ}{=} \prod_{\substack{s P P(U_N) \\ \Box}} T_N(s) b \stackrel{\circ}{G_{R_N}^{(T_N)}} \stackrel{\circ}{=} \prod_{\substack{s P P(U_N) \\ \Box}} T_N(s) b \stackrel{\circ}{G_{R_N}^{(T_N)}} \stackrel{\circ}{=} \prod_{\substack{s P P(U_N) \\ \Box}} T_N(s) a \stackrel{\circ}{=} R_N(s) a \stackrel{\circ}{=} n \xrightarrow{s P P(U_N)} \stackrel{\sim}{=} n \xrightarrow{s$$

because b is bounded by 1 and

$$b \stackrel{\square}{G}_{T_{N}(s)}^{\square(T_{N})} = b \stackrel{\square}{G}_{R_{N}(s)}^{\square(T_{N})};$$

their expressions depending on the first order inclusion probabilities $\Box(T_N)$ solely.

The last inequality follows from the usual inequality between the total variation metric and the entropy (Berger, 1998, Lemma 2 p.219).

5.5.4 FCLT in the rejective survey case

We shall prove Corollary 5.18, the statement of which is recalled below.

Corollary Suppose that Assumption 5.1, Assumption 5.8, Assumption 5.11 and conditions i) and ii) of Theorem 5.14 are satisfied. Then, there exists a \Box_P -equicontinuous Gaussian process G in 1 (F) with covariance operator \Box given by Equation (5.4) such that

$$G_{R_N}^{\square(R_N)}$$
 Ò G weakly in `¹ (F); as N — 1.

Proof Following in the footsteps of Hajek (1964), in the rejective sampling situation where $p(W_i) = p_i$ for i Pt1;:::; Nu, we have

$$\max_{\substack{1 \leq i \leq N \\ i \leq i \leq N}} \frac{\varphi_i}{\prod_{i=1}^R} \Box 1 \prod_{\substack{N \to 1 \\ N \to 1}} 0.$$

We thus have

$$\max_{\substack{i \in \mathbb{N} \\ i \in \mathbb{N}}} \prod_{i=1}^{R} \square 1 \prod_{i=1}^{n} 0$$

under the hypothesis that Assumption 5.1 is fulfilled by the pi's. Then, we can write

$$\begin{split} G_{R_{N}}^{p(W)}f & \Box G_{R_{N}}^{\Box(R_{N})}f = \frac{2}{\overline{N}} \frac{\left| \overrightarrow{N} \right|}{\overline{N}} \frac{\Box_{i}}{p(W_{i})} \Box \frac{\Box_{i}}{\Box_{i}^{R}} f(X_{i}) \\ &= \frac{2}{\overline{N}} \frac{\left| \overrightarrow{N} \right|}{\overline{N}} \frac{\Box_{i}}{p(W_{i})} 1 \Box \frac{p_{i}}{\Box_{i}^{R}} f(X_{i}) \end{split}$$

and

$$\sup_{b \in PBL_{1}(\sum_{i=1}^{1} (F))} E = b \sup_{f \in PF} G_{R_{N}}^{p(W)} f = b \sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\sup_{f \in PF} G_{R_{N}}^{D(W)} f = \sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\sup_{f \in PF} G_{R_{N}}^{D(W)} f = \sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\sup_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\lim_{f \in PF} G_{R_{N}}^{\Box(R_{N})} f = b$$

$$\frac{1}{N} \prod_{i=1}^{R} \prod_{j=1}^{R} \prod_{j=1}^{R}$$

which quantity vanishes asymptotically under Assumption 5.1, according to Theorem 5.1, Equations (5.7) and (5.26) in Hàjek (1964, p. 1508-1510).

The desired convergence is finally established by combining this result with Theorem 5.17 and the functional version of Slutsky's theorem (see Theorem 3.4 in Resnick, 2007 for instance).

5.5.5 CLT for Hadamard differentiable functionals

We shall prove Theorem 5.21, the statement of which is recalled below.

Theorem Suppose that the assumptions of Theorem 5.14 hold and that functional T : L Ä 1 (F) — R^q is Hadamard differentiable at P with differential dT_P and influence function T⁽¹⁾(x; P). Then, as N — +1 , we have:

where G is a Gaussian process with covariance operator \Box , as in Equation (5.4).

Proof The idea is essentially to apply the Hadamard differentiability property to the sequence $h_N = \frac{1}{N} (P_{R_N}^{\square(R_N)} \square P_N) =: G_{R_N}^{\square(R_N)}$, which converges to h = G in `1 (F) and $t_N = \frac{21}{N} - 0$. We thus have, as N - +1:

$$\widehat{\mathsf{N}}^{\mathsf{T}}(\mathsf{P}_{\mathsf{R}_{\mathsf{N}}}^{\mathsf{D}(\mathsf{R}_{\mathsf{N}})}) \Box \mathsf{T}(\mathsf{P}_{\mathsf{N}}) \stackrel{\mathsf{D}}{=} \widehat{\widetilde{\mathsf{N}}}^{\mathsf{T}}\mathsf{T}(\mathsf{P}_{\mathsf{N}} + \frac{21}{\overline{\mathsf{N}}}\mathsf{h}_{\mathsf{N}}) \Box \mathsf{T}(\mathsf{P}_{\mathsf{N}}) \stackrel{\mathsf{D}}{=} \mathsf{d}\mathsf{T}_{\mathsf{P}}.\mathsf{G}.$$

6

TAIL INDEX ESTIMATION BASED ON SURVEY DATA

In the previous chapter we investigated the asymptotic behavior of a variant of the Horvitz-Thompson type of the traditional empirical process for survey schemes of the Poisson type. The ensuing Functional Central Limit Theorems proved particularly convenient to derive asymptotic properties for numerous families of estimators, as was suggested in Section 5.4. Unfortunately, when interested in extreme phenomena like the very high exposure to some food chemical, these results are no longer sufficient and estimators need to be analyzed one by one. In extreme value theory, the survey design is usually ignored and the ensuing statistics, already suffering from the rarity of tail observations, possibly exhibit an additional bias due to the omission of the sampling phase. Whereas asymptotic analysis of the Horvitz-Thompson estimator (Horvitz and Thompson, 1951) has been the subject of much attention, in particular in the context of mean estimation and regression (see Hajek, 1964; Rosen, 1972; Robinson, 1982; Gourieroux, 1987; Deville and Särndal, 1992; Berger, 1998 for instance), and the last few years have witnessed significant progress towards a comprehensive functional limit theory for the assessment of distribution functions (Gill et al., 1988; Breslow and Wellner, 2007, 2008; Breslow et al., 2009; Saegusa and Wellner, 2011; Bertail et al., 2013), no result on tail estimation has been documented in the survey sampling literature yet. In a modest attempt to start filling this gap, we make our contribution to the elaboration of an extreme value theory for survey data by focusing on a modified (Horvitz-Thompson) version of the most celebrated Hill estimator. This new statistic assesses the extreme value index of univariate heavytailed distributions, which essentially indicates to what extent extreme events are rare, while accounting for the sampling design by means of which the data have been collected. Therefore, contrary to the standard Hill estimator, it corrects the bias induced by the survey plan. Its consistency is established for any type of sampling scheme that fulfills some adequate assumptions on the first and second order inclusion probabilities. Then, following in the footsteps of Hajek (1964) and along the lines of the previous chapter, its asymptotic normality is investigated for Poisson-like survey designs. The results presented in this chapter originate from a paper written in collaboration with P. Bertail (Université Paris X, France) and S. Clémençon (Télécom ParisTech, France) and has been submitted for publication. Regrettably, our analysis does not encompass complex designs such as that of the INCA2 database

yet, hence no concrete application to dietary risk analysis is provided in this chapter. Such refinements will hopefully be the object of further research in the near future.

The chapter is structured as follows. Basics about survey sampling and tail index estimation in the standard iid setup are briefly recalled and notations are set up in Section 6.1. In Section 6.2, we describe at length the proposed modification of the Hill estimator in the context of a general sampling plan and prove its consistency. Its asymptotic normality is investigated next in Section 6.3. Finally, practical issues such as the selection of an optimal number of largest observations on which to base the estimation are discussed in Section 6.4, together with illustrative numerical experiments. Technical proofs are deferred to Section 6.6.

6.1 background and preliminaries

We first recall the crucial notions in survey sampling that are extensively used in the subsequent analysis, as well as basic concepts of heavy-tail modeling, including Pareto-type distributions and standard strategies for statistical estimation of the related parameters. For the sake of clarity, most notations are kept identical to those in Chapter 5: the Dirac mass at x PR is denoted by \Box_x and the indicator function of any event E by I t Eu. We also denote by #E the cardinality of any finite set E, and by P(E) its power set. The (left-continuous) general inverse of any non-decreasing function $H: (a; b) - R, \Box 1 \S a \dagger b \S + 1$, is denoted by

$$H^{-}(x) := infty P(a;b) : H(y) \bullet xu;$$

x PR, with the convention that the infimum over an empty set is $\Box 1$. When dealing with some multivariate distribution function $H : \mathbb{R}^d \longrightarrow \mathbb{R}$ with marginals $H_1; \ldots; H_d$, we shall write $H^-(x) := H_1^-(x_1); \ldots; H_d^-(x_d)^{\Box}$ for any $x := (x_1; \ldots; x_d) \mathbb{P} \mathbb{R}^d$. Finally, the minimum (resp. maximum) of two real numbers x and y is denoted by x^{A} y (resp. x_{-} y).

6.1.1 Survey sampling

In this section we recall a few essential definitions and set out the notations relative to survey sampling. More details about the ins and outs of these concepts can be found in Section 5.1.1 of Chapter 5.

6.1.1.1 Population, sample, inclusion probabilities and indicators

Here and throughout, we consider a finite population of size N • 1, denoted by $U_N := t1; \ldots; Nu$. We call a sample of (possibly random) size n § N, any subset $s := ti_1; \ldots; i_{n(s)}u$ in $P(U_N)$ with cardinality n =: n(s) less than N. A sampling

scheme (design/plan) without replacement is determined by a probability distribution R_N on the set of all possible samples s $P P(U_N)$. For any i P t1; ...; Nu, the following quantity, generally called (first order) inclusion probability,

$$\Box_{i}(\mathsf{R}_{\mathsf{N}}) := \mathsf{P}_{\mathsf{R}_{\mathsf{N}}} (i \mathsf{P} \mathsf{S});$$

is the probability that the unit i belongs to a random sample S drawn from distribution R_N . In vectorial form, we shall write $\Box(R_N) := (\Box_1(R_N); \ldots; \Box_N(R_N))$. First order inclusion probabilities are assumed to be strictly positive in the subsequent analysis: @ Pt1; \ldots ; nu, $\Box_i(R_N) \circ 0$. Additionally, the second order inclusion probabilities are denoted by

$$\Box_{i;j}(\mathsf{R}_{\mathsf{N}}) \coloneqq \mathsf{P}_{\mathsf{R}_{\mathsf{N}}} \stackrel{\Box}{(i;j)} \mathsf{P} \mathsf{S}^{2} \stackrel{\Box}{;}$$

for any i \Box j in t1;:::; Nu². When no confusion is possible, we shall fail to mention the dependence in R_N when writing the first/ second order probabilities of inclusion. The information related to the observed sample S Ä t1;:::; Nu is encapsulated by the random vector $\Box := (\Box_1; :::; \Box_N)$, where

$$= \begin{array}{c} # \\ 1 & \text{if i PS} \\ 0 & \text{otherwise.} \end{array}$$

The distribution of the sampling scheme \Box has 1-dimensional marginals that correspond to the Bernoulli distributions B(\Box_i), 1 § i § N, and covariance matrix given by

$$\Box_{N} := \Box_{i;j} \Box \Box_{i} \Box_{j} {}^{(}_{1 \\ \$ i;j \\ \$ N} .$$

Notice incidentally that, equipped with these notations, we have $\sum_{i=1}^{\infty} |a_i| = n(S)$.

The superpopulation model we consider here stipulates that a real-valued random variable X with distribution P and cdf F is observable on the population U_N , i.e $X_1; \ldots; X_N$ are iid realizations drawn from P. In practice, it is customary to determine the first order inclusion probabilities as a function of an auxiliary variable, which is observed on the entire population. Here, it is denoted by W with distribution P_W . Hence, for all i P t1; \ldots ; Nu we can write $\Box_i = \Box(W_i)$ for some link function $\Box(.)$. When W and X are strongly correlated, thus proceeding helps select more informative samples and subsequently reduce the variance of estimators (we refer to Section 5.1.1.3 of Chapter 5 for a more detailed discussion on the use of auxiliary information in survey sampling). One may refer to Cochran (1977); Gourieroux (1981); Deville (1987) for accounts of survey sampling techniques.

6.1.1.2 A crucial example: the Poisson survey scheme

Though of extreme simplicity, the Poisson scheme (without replacement) plays a crucial role in sampling theory, insofar as it can be used to approximate a wide range of survey plans. This is indeed a key observation to establish general asymptotic results in the survey context, see Hajek (1964), Chapter 5 and Section 6.3 of the

present chapter. For such a plan, denoted by T_N , the \Box_i 's are independent Bernoulli random variables with parameters $p_1; \ldots; p_N$ in (0; 1). Thus, the first order inclusion probabilities fully characterize such a plan. Observe in addition that the size n(S) of a sample generated this way is random and goes to infinity as N - +1 with probability one, provided that $\min_{1 \le i \le N} p_i$ remains bounded away from zero.

6.1.1.3 The Horvitz-Thompson empirical measure

We recall that the Horvitz-Thompson estimator of the empirical measure

$$\mathsf{P}_{\mathsf{N}} \coloneqq \frac{1}{\mathsf{N}} \prod_{i=1}^{\mathsf{N}} \Box_{\mathsf{X}_{i}}$$

based on the survey data described above is defined as follows (Horvitz and Thompson, 1951):

$$\mathsf{P}_{\mathsf{R}_{\mathsf{N}}}^{\Box(\mathsf{R}_{\mathsf{N}})} \coloneqq \frac{1}{\mathsf{N}} \frac{\mathsf{P}}{\mathsf{I}_{\mathsf{i}}} \frac{\Box_{\mathsf{i}}}{\Box_{\mathsf{i}}} \Box_{\mathsf{X}_{\mathsf{i}}} = \frac{1}{\mathsf{N}} \frac{\mathsf{P}}{\mathsf{I}_{\mathsf{P}}} \frac{1}{\Box_{\mathsf{i}}} \Box_{\mathsf{X}_{\mathsf{i}}};$$

$$\overline{F}_{N}^{\Box}(x) := \frac{1}{N} \frac{\prod_{i=1}^{N}}{\prod_{i=1}^{i}} \frac{\prod_{i=1}^{i}}{\prod_{i=1}^{i}} I t X_{i} \circ xu = \frac{1}{N} \frac{\prod_{i=1}^{N}}{\prod_{i=1}^{i}} I t X_{i} \circ xu.$$
(6.1)

6.1.2 Tail index inference - the Hill estimator

In a wide variety of situations, it is appropriate to assume that a statistical population is described by a heavy-tailed probability distribution (the field of heavy-tail analysis is well depicted in Resnick, 2007). A distribution with Pareto-like right tail is any probability measure P on R with cdf F such that for all x PR,

$$1 \square F(x) = \overline{F}(x) = x^{\square 1 = \square} L(x); \qquad (6.2)$$

where \square ° 0 is the extreme value index (EVI) of distribution P and L(x) is a slowly varying function, i.e. a function such that L(t x)=L(x) — 1 as x — +1 for all t ° 0. Notice that instead of the EVI, focus is often on $\square := 1=\square$, the tail index of the distribution P. Functions of the form introduced in Equation (6.2) are said to be regularly varying with index $\square 1=\square$; the set of such functions is denoted by R $_{\square 1=\square}$. One may refer to Bingham et al. (1987) for an account of the theory of regularly varying functions. The Hill estimator (Hill, 1975) provides a popular way of estimating the EVI \square . Its asymptotic behavior and the practical issues related to its computation are well-documented in the literature, see for instance Resnick (2007, Chapter 4) and the references therein. Given an iid population X₁;...;X_N of size N • 1 drawn from P and K Pt1;...;N \square 1u largest observations, it is written

$$H_{K;N} := \frac{1}{K} \prod_{i=1}^{[N]} \log \left[\frac{X_{N \ i+1;N}}{X_{N \ K;N}} \right];$$
(6.3)

where $X_{1;N}$ § \blacksquare § $X_{N;N}$ denote the order statistics related to the population. Whereas the theory has been extensively developed in the case where the observations are independent and identically distributed (including questions related to the choice of K in Equation (6.3)), to the best of our knowledge the Hill procedure has received no attention when data arise from a general survey. We point out that there exist alternative methods for tail index or EVI estimation, refer for instance to Beirlant et al. (2004, Chapter 4) for further details. The argument of the subsequent analysis paves the way for studying extensions of such techniques in the context of survey sampling models.

6.2 the hill estimator in survey sampling

$$U(x) \coloneqq F^{-1} \Box \frac{1}{x}^{\Box}.$$

Its empirical and Horvitz-Thompson equivalents are respectively denoted by

$$U_{N}(x) \coloneqq F_{N}^{-} \stackrel{\square}{1} \square \frac{1}{x}^{\square} \quad \text{and} \quad U_{N}^{\square}(x) \coloneqq (F_{N}^{\square})^{-} \stackrel{\square}{1} \square \frac{1}{x}^{\square}$$

Notice that when F P R₁₌₀, the corresponding tail quantile function U is also regularly varying with index \Box (Beirlant et al., 2004, Sections 2.3.2 and 2.9.3). The goal pursued here is to estimate the tail parameter \Box based on the survey data $X_{i_1}; \ldots; X_{i_n}$ and the sampling plan R_N .

6.2.1 The Horvitz-Thompson variant of the Hill estimator

Notice first that, under the heavy-tail assumption above, we have:

$$\Box = \lim_{x \to 1} \int_{x}^{a+1} \frac{\overline{F}(u)}{\overline{F}(x)} \frac{du}{u}; \qquad (6.4)$$

see Beirlant et al. (2004, Section 2.6) for instance. In the case of the iid population $X_1; \ldots; X_N$ drawn from P, one classically recovers the celebrated Hill estimator by substituting F with the empirical cdf F_N in Equation (6.4) and taking

$$\mathbf{x} = \mathbf{U}_{\mathbf{N}} (\mathbf{N} = \mathbf{K}) = \mathbf{X}_{\mathbf{N} \square \mathbf{K}; \mathbf{N}}$$

for some number 1 § K § N \square 1 of largest observations, supposedly representative of the tail of the distribution. Indeed, we have:

$$\begin{array}{c} \stackrel{a}{\xrightarrow{}} +1 \\ \stackrel{}{\xrightarrow{}} X_{N \ \square \ K;N} \end{array} \frac{\overline{F}_{N}\left(u\right)}{\overline{F}_{N}\left(X_{N \ \square \ K;N}\right)} \frac{du}{u} = \frac{\left[\stackrel{f}{\xrightarrow{}} \right]^{a} X_{N \ \square \ i+1;N}}{i=1} \frac{\overline{F}_{N}\left(u\right)}{\overline{F}_{N}\left(X_{N \ \square \ K;N}\right)} \frac{du}{u} \\ = \frac{\left[\stackrel{f}{\xrightarrow{}} \right]^{i}}{i=1} \frac{i}{K} \left(\log X_{N \ \square \ i+1;N} \ \square \ \log X_{N \ \square \ i;N}\right) \\ = \frac{1}{K} \frac{\left[\stackrel{f}{\xrightarrow{}} \right]}{i=1} \log \frac{X_{N \ \square \ i+1;N}}{X_{N \ \square \ K;N}} \\ = H_{K;N} . \end{array}$$

Consistency of this estimator is classically guaranteed as soon as K - +1 and K = o(N) when N - +1, see Mason (1982). In this case, the empirical threshold $X_{N \square K;N}$ (equivalent in probability to U(N=K)) goes to infinity (again in probability) as N - +1.

Going back to the survey data situation, one may naturally replace U(N=K) by U_N^{\Box} (N=K) and build a plug-in estimate of the EVI \Box based on the Horvitz-Thompson estimator given in Equation (6.1) of the tail probability $\overline{F}(x)$. Observe that by definition U_N^{\Box} (N=K) corresponds to one of the observations in the sample, say X_i^{\Box} with rank `Pt1;:::;nu. To this ` obviously corresponds an index k Pt0;:::;n \Box 1u such that ` = n \Box k, implying $X_{n \Box k;n} = U_N^{\Box}$ (N=K), the Horvitz-Thompson estimator of

the quantile of order $1 \square K=N$. We denote by \square_N^\square the map linking k to K in U_N under the sampling scheme R_N :

$$\Box_{N}^{\square}: \begin{array}{c} t1; \ldots; N \Box 1u \Box t1; \ldots; n \Box 1u \\ K \quad fi k \coloneqq \Box_{N}^{\square}(K) \end{array}$$

where

$$\Box_{N}^{\mu}(\mathsf{K}) \coloneqq \mathsf{n} \Box \inf^{\#} \mathsf{i} \mathsf{Pt1}; \ldots; \mathsf{n} \Box \mathsf{1u} : \prod_{j=1}^{|\mathsf{I}|} \frac{1}{\Box_{j;n}} \bullet \mathsf{N} \Box \mathsf{K} \quad . \tag{6.5}$$

This leads to the quantity:

$$\vec{p} = \frac{\prod_{i=1}^{k+1} \frac{\vec{F}_{R_{N}}^{((R_{N})}(u)}{\prod_{i=1}^{k} (R_{N})(X_{n-k};n)} \frac{du}{u}}{\vec{F}_{R_{N}}^{((R_{N})}(X_{n-k};n)} \frac{du}{u}} = \frac{\prod_{i=1}^{k} x_{n-i+1;n}}{x_{n-i+1;n}} \frac{\vec{F}_{R_{N}}^{((R_{N})}(u)}{\vec{F}_{R_{N}}^{((R_{N})}(X_{n-k};n)} \frac{du}{u}} = \frac{\prod_{i=1}^{k} x_{n-i+1;n}}{\prod_{i=1}^{k} (R_{N})(X_{n-k};n)} \frac{\vec{F}_{R_{N}}^{((R_{N})}(u)}{\vec{F}_{R_{N}}^{((R_{N})}(X_{n-k};n)} \frac{du}{u}} = \frac{\prod_{i=1}^{k} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2}}{\prod_{i=1}^{k} \frac{1}{2}} \frac{1}{2} \left(\log(X_{n-i+1;n}) \log(X_{n-i+1;n})\right)} = \frac{\prod_{i=1}^{k} \frac{1}{2} \frac{1}{2$$

Hence, k is to the sample what K is to the population: the number of upper values on which the estimation should rely. Notice that we may also write

$$H_{K;n}^{\Box} = \frac{\prod_{j=1}^{n} \frac{\sum_{N = j+1;N}^{j+1;N}}{\sum_{N = j+1;N}^{j+1;N}} \prod_{i=1}^{j-1} \frac{\prod_{N = i+1;N}^{i-1} \sum_{N = i+1;N}^{j-1} \log \frac{X_{N = i+1;N}}{X_{N = K;N}} = : H_{K;N}^{\Box}; \quad (6.6)$$

where K is the chosen number of largest observations in the population from which k was constructed. Observe that \Box_N^{\Box} in Equation (6.5) is a surjective, non-injective random map, which suggests that the subsequent asymptotic analysis better rely on some appropriately chosen K and its random image k := \Box_N^{\Box} (K) rather than the contrary. Since in practice only k can be computed from the X_i's and \Box_i 's, i P S, the total population being partly unobserved, considerations about the choice of an appropriate k are discussed in detail in Section 6.4. From this point forward, the Horvitz-Thompson Hill estimator shall be written $H_{K;N}^{\Box}$ with K Pt1;:::;N \Box 1u held fixed.

6.2.2 Consistency of H^{_}_{K:N}

Here we investigate the limit properties of the estimator $H_{K,N}^{\Box}$ as N and n simultaneously go to infinity, with n § N. The following assumptions, related to the sample design, shall be involved in the asymptotic analysis.

Assumption 6.1 There exist \square_0° 0 and $N_0^\circ P N^\circ$ such that for all $N \bullet N_0^\circ$ and i $P U_N^\circ$,

□_i ° □_□.

Assumption 6.2 There exists + 1 such that we almost-surely have @N • 1,

 $\max_{\substack{1 \leq i; j \leq N}} \Phi_{i;j} \Box \Box \Box \overset{\sim}{\Box} \overset{\sim}{\underset{j}{\boxtimes}} \overset{\sim}{\underset{n}{\boxtimes}}.$

Assumption 6.1 guarantees that first order inclusion probabilities do not vanish asymptotically, while Assumption 6.2 corresponds to the situation where the second order inclusion probabilities are not too different from those in the case of independent sampling (it is thus fulfilled by the Poisson design, see Section 6.1.1.2).

Remark 6.3 – On Assumption 6.1 and Assumption 6.2. The two assumptions introduced herein-before are rather mild and are fulfilled in a wide variety of situations. Indeed, Assumption 6.2 is standard in asymptotic analysis of sampling techniques, see Hartley and Rao (1962) and Hàjek (1964) for instance. As for Assumption 6.1, recall that $\Box_i = \Box(W_i)$ with W an auxiliary variable and $\Box(.)$ a link function. Then, it is fulfilled as soon as $\Box(.)$ is continuous and the support of P_W is a compact subset of $(R_+^{\Box})^d$, where d denotes the dimension of the random vector W and R_+^{\Box} the set of positive real numbers.

Given this framework, following in the footsteps of Resnick (2007, Section 4.4.1), the consistency of $H_{K;N}^{\Box}$ can be handled by exploiting the properties of regularly varying distributions. Indeed, under the heavy-tail assumption in Equation (6.2), provided that K - +1 and K=N - 0 as N - +1, we have

$$\frac{N}{K} P \stackrel{\sqcup}{\longrightarrow} \frac{X}{U(N=K)} P \stackrel{\sqcup}{\longrightarrow} \frac{\nabla}{N-1} \square_{1=0}(.)$$

in the space of Radon measures on $(0; \pm 1]$. There, " $\stackrel{\vee}{-}$ " stands for the vague convergence of measures¹ and $\square_{1=\square}(.)$ is such that for all x ° 0, $\square_{1=\square}(x; 1] = x^{\square 1=\square}$ (see Resnick, 2007, Theorem 3.6 for instance). Its empirical counterpart in the population, usually called the tail empirical measure, is defined as follows:

$$\Box_{\mathsf{N}} := \frac{1}{\mathsf{K}} \prod_{i=1}^{[\mathsf{N}]} \Box_{\mathsf{X}_i = \mathsf{U}(\mathsf{N} = \mathsf{K})}.$$

^{1.} Recall that, in the space of non-negative Radon measures on (0; +1], a sequence $(\Box_m)_{m+1}$ is said to converge vaguely to \Box iff for any compactly supported continuous function h: (0; +1] - R, we have: $\bigcup_{0}^{+1} h(x) \Box_m (dx) - \bigcup_{0}^{+1} h(x) \Box(dx)$ as m - +1.

When replacing U(N=K) by its estimate U_N (N=K) in the expression above and assuming that K = K(N) - +1 where K=N - 0 as N - +1, it can be shown that it converges to $\Box_{\Box 1=\Box}$ in probability (Resnick, 2007, Equation (4.21)). Since we have

$$H_{K;N} = \int_{1}^{a} \Box_{N} \frac{X_{N \Box K;N}}{U(N=K)} (x; \pm 1] \frac{dx}{x}$$
$$\Box = \int_{1}^{a} \Box_{\Box 1=\Box} (x; \pm 1] \frac{dx}{x};$$

and

the asymptotic properties of \Box_N naturally convey the consitency of the Hill estimator. Generalizing this result to the Horvitz-Thompson tail empirical measure

$$\Box_{N}^{\Box} = \frac{1}{K} \prod_{i=1}^{[N]} \frac{\Box_{i}}{\Box_{i}} \Box_{X_{i}=U_{N}^{\Box}(N=K)} = \frac{1}{K} \prod_{i \in S} \frac{1}{\Box_{i}} \Box_{X_{i}=U_{N}^{\Box}(N=K)}$$
(6.7)

would then yield the theorem below (see the proof in Section 6.6). It reveals that, in regard to the asymptotic statistical estimation of the EVI \Box , the Horvitz-Thompson variant of the Hill estimator $H_{K\cdot N}^{\Box}$ is consistent.

Theorem 6.4 – Consistency. Let K = K(N) be a sequence of integers such that K - +1 and K=N - 0 as N;n - +1. Provided that Assumption 6.1 and Assumption 6.2 are fulfilled, we then have, as N and n tend to +1:

$$\mathsf{H}_{\mathsf{K}:\mathsf{N}}^{\Box} \sqsubseteq^{\mathsf{P}} \Box. \tag{6.8}$$

6.3 asymptotic normality of $H_{K:N}^{\square}$

Whereas the consistency of the standard Hill estimator in Equation (6.3) can be proved for any sequence K going to infinity at a reasonable rate, asymptotic normality cannot be guaranteed at such a level of generality. Higher-order regular variation properties of the heavy-tail model in Equation (6.2) are required (de Haan and Peng, 1998; de Haan and Stadtmüller, 1996). More specifically, consider the hypothesis below, referred to as the Von Mises condition (Goldie and Smith, 1987).

Assumption 6.5 The regularly varying tail quantile function U P R_{\Box} with \Box ° 0 is such that there is a real parameter \Box † 0, referred to as the second order parameter, and a positive or negative function A with $\lim_{x \to \pm 1} A(x) = 0$ such that for any t ° 0,

$$\frac{1}{A(x)} \frac{U(tx)}{U(x)} \Box t^{\Box} \frac{\Box}{x-1} t^{\Box} \frac{t^{\Box} \Box}{\Box};$$

or equivalently

$$\frac{1}{A \xrightarrow{\frac{1}{\overline{F}(x)}}} \xrightarrow{\Box} \overline{\overline{F}(t x)} = t^{\Box 1 = \Box} \xrightarrow{\Box} t^{\Box 1 = \Box} \frac{t^{\Box = \Box} = 1}{\Box}$$

This condition simply establishes some constraints about the slowly varying function L(.) in Equation (6.2) to ensure its influence vanishes quickly enough not to interfere with the Pareto form $x^{\Box 1=\Box}$ of \overline{F} .

The limit distribution of the standard Hill estimator in Equation (6.3) has been investigated by means of Rényi's exponential representation of log-spacings under Assumption 6.5. Of course, this condition can hardly be checked in practice and the choice of the number of extremal observations is generally selected so as to minimize an estimate of the asymptotic mean squared error (MSE), see Section 6.4 and the references therein.

Remar k 6.6 – On Hill and the tail empirical process. Other approaches than the Rényi decomposition in log-spacings have been developed to prove the asymptotic normality of the Hill estimator $H_{K;N}$ (see Resnick, 2007, Chapters 4 and 9 and the references therein). Along the lines of the study of empirical processes in Chapter 5, the version presented at length in Resnick (2007) involves the preliminary study of the tail empirical process

$$\mathsf{T}_{\mathsf{N}} \coloneqq \stackrel{?}{\overset{\square}{\underset{\scriptstyle \mathsf{K}}{\overset{\scriptstyle \square}{\underset{\scriptstyle \mathsf{N}}{\overset{\scriptstyle \square}{\underset{\scriptstyle \mathsf{N}}{\overset{\scriptstyle \square}{\underset{\scriptstyle \mathsf{\Pi}=0}{\overset{\scriptstyle \amalg}{\underset{\scriptstyle !}}}}}}}}}}}}}}}}}}}}$$

from which the asymptotic properties of the Hill estimator are later deduced. We refer to Theorem 9.1 and Section 9.1.2 in Resnick (2007) for more details on this seemingly simple but actually quite intricate procedure.

Unfortunately, contrary to the classical empirical process, the asymptotic properties of T_N cannot be extended to its Horvitz-Thompson equivalent. This is essentially due to the fact that the population U_N contains a finite number of observations. Hence, there is always a maximum $X_{N;N}$ † 1 bounding P_N and sampling fails to distinguish between a distribution with finite support such as those in the Weibull domain of attraction and a heavy-tailed distribution with no endpoint (see Chapter 2 for an introduction to extreme value theory and maximum domains of attractions). Actually, these arguments are exactly the same as those introduced when discussing the applicability of bootstrap in extreme value analysis. Indeed, as explained in Remark 5.7 of Chapter 5, survey sampling under a superpopulation model can be viewed as a generalization of weighted bootstrap. The interested reader may refer to Resnick (2007, Section 6.4) for a brief introduction to the difficulty of bootstrapping heavy-tailed phenomena. In this section, we shall aim at proving first that under Poisson survey schemes, the Horvitz-Thompson version of the Hill estimator computed on the sample is asymptotically close to its standard version calculated over the entire population. This result is next extended to rejective sampling plans through a coupling argument, similar to that used in Hàjek (1964) and in Chapter 5.

6.3.1 The case of the Poisson survey scheme

In this section we assume that the vector \Box corresponds to that of a Poisson survey scheme, such as depicted in Section 6.1.1.2. The distribution of this design is denoted by T_N and the first order inclusion probabilities by $p_1; \ldots; p_N$. Under this setting, the Horvitz-Thompson variant of the Hill estimator is naturally denoted by $H_{K;N}^p$ and the \Box_i 's in Equation (6.6) are to be replaced by the corresponding p_i 's. In addition, just like we previously set $\Box_i = \Box(W_i)$, we write $p_i = p(W_i)$ for all i P U_N and p the Poisson link function. In keeping with the results obtained in Chapter 5, so as to prove the asymptotic normality of $H_{K;N}^p$, we shall require the p_i 's to fulfill Assumption 6.1 with lower bound p_{\Box} and the auxiliary variable from which they are built to satisfy the condition below.

Assumption 6.7 The random vectors W_1 ;:::; W_N are iid with continuous distribution P_W on W Ä R^d, d-variate cdf F_W with marginals F_{W_1} ;:::; F_{W_d} and density f_W . The joint distribution of the entailed iid sequence $t(X_i; W_i)$; 1 § i § Nu is denoted by $P_{X;W}$ with corresponding cdf $F_{X;W}$.

The following result reveals that under the Poisson survey scheme, when based on the K largest values among the whole population $X_1; \ldots; X_N$, $H^p_{K;N}$ converges at the same rate (1= \overline{K} namely) to the same limit distribution as $H_{K;N}$, up to a multiplicative term in the asymptotic variance induced by the sampling scheme. Further details about the convergence of the classical Hill estimator can be found e.g. in de Haan and Peng (1998, Theorem 1) and Resnick (2007, Section 9).

Theorem 6.8 – Limit distribution in the Poisson survey case. Suppose that Assumption 6.5 is fulfilled by the underlying heavy-tailed model and that Assumption 6.1 is satisfied by the considered sequence of Poisson inclusion probabilities $p_1; \ldots; p_N, N \bullet 1$, constructed from some set of auxiliary variables as in Assumption 6.7. Further assume the conditions introduced herein-after.

i) The marginal cdf F is absolutely continuous with respect to the Lebesgue measure with density f.

ii) The joint cdf $F_{X;W}$ is absolutely continuous with Lebesgue-integrable density $f_{X;W}$ such that for all $(x;w) P(0;+1] \square W$, $w = (w_1; :::; w_d)$,

$$f_{X;W}(x;w) \coloneqq c_{X;W}(F(x);F_{W_1}(w_1);\ldots;F_{W_d}(w_d)) f(x) \prod_{j=1}^{TP} f_{W_j}(w_j);$$

for some copula density $c_{X;W} : R^{\Box}_{+} \Box R^{d} - R$ and $f_{W_{1}}; :::; f_{W_{d}}$ the marginal densities of the distribution of W.

iii) $c_{X;W}$ bounded in a neighborhood of $t 1u \square [0; 1]^d$.

Then, for

$$\Box_{p}^{2} \coloneqq \int_{W}^{a} \frac{1}{p(w)} c_{X;W}(1; F_{W_{1}}(w_{1}); \ldots; F_{W_{d}}(w_{d})) \int_{j=1}^{T^{d}} f_{W_{j}}(w_{j}) dw_{1} \ldots dw_{d}$$

and provided that K - +1 as N - +1 so that $\overrightarrow{KA}(N=K) - \Box$ for some constant $\Box PR$, we have the convergence in distribution as N - +1:

$$? \overline{\mathsf{K}} \stackrel{\square}{\mathsf{H}}{}^{\mathsf{p}}_{\mathsf{K};\mathsf{N}} \square \square \stackrel{\square}{\to} \mathsf{N} \stackrel{\square}{\xrightarrow{}}{\frac{\square}{1 \square \square}}; \square^2 \square_{\mathsf{p}}^2 . \tag{6.9}$$

As can be seen by examining the proof of this theorem in Section 6.6, the limit result in Proof 6.6.2 can be obtained using the following decomposition:

$${}^{?}\overline{K} \stackrel{P}{H}^{p}_{K;N} \stackrel{P}{=} = \frac{{}^{?}\overline{K} \stackrel{P}{K} \stackrel{P}{K_{N}} \stackrel{P}{H}^{p}_{K;N} \stackrel{P}{=} \stackrel{P}{H}^{p}_{K;N} \stackrel{P}$$

where

$$\mathbf{r}_{\mathrm{K};\mathrm{N}} \coloneqq \mathbf{r}_{\mathrm{K};\mathrm{N}} \left(\mathbf{T}_{\mathrm{N}}; \mathbf{p} \right) \coloneqq \frac{1}{\mathrm{K}} \frac{\prod_{i=1}^{\mathrm{N}} \frac{\prod_{i=1}^{\mathrm{N}} \prod_{i=1}^{\mathrm{I}} 1;\mathrm{N}}{\mathrm{P}_{\mathrm{N}} \prod_{i=1}^{\mathrm{I}} 1;\mathrm{N}}.$$
 (6.11)

These three quantities are studied independently under the hypotheses required in Theorem 6.8. First, we show that $r_{K;N}$ converges to 1 in probability as N tends to 1. Combined with Rényi's decomposition in log-spacings of the Hill estimator (refer for instance to Beirlant et al., 2004, Section 4.4), this establishes the asymptotic convergence, in probability, of $Q_N^{(1)}$ to 0. It also implies that $Q_N^{(2)}$ is equivalent to $\frac{1}{K}$ ($H_{K;N} \square$), a well-known quantity which tends to a Gaussian distribution with expectation \square =(1 \square) and variance \square^2 under the second order condition stipulated in Assumption 6.5 (De Haan and Ferreira, 2006, Theorem 3.2.5). As for $Q_N^{(3)}$, we calculate its expectation and variance conditionally on the full vector of observations

survey scheme has been controlled. Then, following in the lines of a Lindeberg-Feller theorem for independent and non-identically distributed variables (Feller, 1971, Theorem 3, p.262), we exhibit some sufficient conditions (namely i), ii) and iii) in Theorem 6.8) under which the conditional variance has a finite limit in probability relative to t(X_i;W_i); 1 § i § Nu. Provided that they are fulfilled, $Q_N^{(3)}$ converges weakly to a centered Gaussian distribution with variance $\Box^2 (\Box_p^2 \Box 1)$. This lets us consider $Q_N^{(2)}$ and $Q_N^{(3)}$ as independent random variables (one depends on the data and the other on the survey scheme). Thus, the limit distribution of their sum is simply the sum of their limit distributions, thereby yielding Proof 6.6.2.

Remark 6.9 – On the asymptotic variance. Looking at the variance term in Proof 6.6.2 of Theorem 6.8, we see that the influence of the survey scheme is encapsulated by the multiplicative term $\Box_p^2 \bullet 1$. Ideally, we would like to have at our disposal inclusion probabilities for which \Box_p^2 is as close to 1 as possible. In that case, the Horvitz-Thompson version of the Hill estimator would perform as well as its population equivalent. For some chosen expected sample size $n^{\Box} := E(n(S))$, they would solve the following optimization program:

$$\min_{p(w)} \frac{1}{w} c_{X;W}(1; F_{W_1}(w_1); \dots; F_{W_d}(w_d)) \prod_{j=1}^{n^d} f_{W_j}(w_j) dw_1 \dots dw_d \square 1$$
subject to
$$\sup_{i=1}^{n^d} p(W_i) = n^{\square};$$

provided the ensuing sequence p_1 ;...; p_N satisfies Assumption 6.1.

6.3.2 Extension to rejective sampling schemes

We now show how the result stated in Theorem 6.8 can be extended to an important class of survey plans, namely rejective sampling schemes. For the sake of clarity, we first provide a brief description of the latter, refer to Hajek (1964) and Berger (1998) for further details.

Fix n § N and consider a vector $(\Box_1; \ldots; \Box_N)$ of first order inclusion probability. The rejective sampling, sometimes referred to as conditional Poisson sampling (CPS in short), exponential design without replacement or maximum entropy design (Tillé, 2006), is the sampling plan R_N which picks samples of fixed size n(S) := n in order to maximize the entropy measure

 $H(R_{N}) = \Box \prod_{\substack{t \ s PP(U_{N}): \ \#s = n \ u}} R_{N}(s) \log R_{N}(s)$

subject to the constraint stipulating that its vector of first order inclusion probabilities coincides with $(\Box_1; \ldots; \Box_N)$. It can be implemented in two steps, as follows.

- Draw a sample S with a Poisson sampling plan (without replacement), with properly chosen first order inclusion probabilities (p₁;:::;p_N). The representation is called canonical if p_i = n. In that case the relationships between p_i and □_i, 1 § i § N, are established in Hàjek (1964).
- 2. If $n(S) \square n$, then reject it and go back to step one, otherwise stop.

The vector $(p_1; :::; p_N)$ must be chosen in a way that the resulting first order inclusion probabilities coincide with $\Box_1; :::; \Box_N$, by means of a dedicated optimization algorithm, see Tillé (2006). The corresponding probability distribution is given by: @s PP(U_N),

$$\mathsf{R}_{\mathsf{N}}(\mathsf{s}) = \underbrace{\sim}_{\mathsf{ts}^{1}\mathsf{PP}(\mathsf{U}_{\mathsf{N}}): \ \texttt{#s}^{1}=\mathsf{nu}}^{\mathsf{T}_{\mathsf{N}}^{\mathsf{p}}} \mathsf{(s}^{1})}_{\mathsf{i} \ \mathsf{Ps}} 9 \overset{\mathsf{T}}{\underset{\mathsf{i} \ \mathsf{Ps}}{\mathsf{Ps}}} \overset{\mathsf{T}}{\underset{\mathsf{i} \ \mathsf{Rs}}{\mathsf{n}}} \mathsf{(1 \square p_{\mathsf{i}})} \square \ \mathsf{I} \ \mathsf{t} \ \texttt{#s} = \mathsf{nu}.$$

Refer to Hajek (1964, p. 1496) for more details on the pi's.

Turning now to the extension of the result stated in Theorem 6.8 for the Poisson survey scheme to the case of rejective sampling, we introduce the following quantities: $@K \ S \ N$,

$$\mathsf{D}_{\mathsf{K};\mathsf{N}}\left(\mathsf{R}_{\mathsf{N}}\,;\mathsf{T}_{\mathsf{N}}\right)\coloneqq\mathsf{r}_{\mathsf{K};\mathsf{N}}\left(\mathsf{R}_{\mathsf{N}}\,;\,\Box\right)\mathsf{H}_{\mathsf{K};\mathsf{N}}^{\Box}\,\,\Box\,\mathsf{r}_{\mathsf{K};\mathsf{N}}\left(\mathsf{T}_{\mathsf{N}}\,;p\right)\mathsf{H}_{\mathsf{K};\mathsf{N}}^{\mathsf{p}}$$

and

$$\mathbf{r}_{\mathsf{K};\mathsf{N}}\left(\mathsf{R}_{\mathsf{N}}\,;\,\mathsf{T}_{\mathsf{N}}\right) \coloneqq \mathbf{r}_{\mathsf{K};\mathsf{N}}\left(\mathsf{R}_{\mathsf{N}}\,;\,\Box\right) \,\Box\,\mathbf{r}_{\mathsf{K};\mathsf{N}}\left(\mathsf{T}_{\mathsf{N}}\,;\,p\right);$$

where $H_{K;N}^{\Box}$ (respectively $r_{K;N}(R_N; \Box)$) refers to the Horvitz-Thompson version of the Hill estimator (resp. of $r_{K;N}$ in Equation (6.11)) under rejective sampling, and $H_{K;N}^{p}$ (resp. $r_{K;N}(T_N;p)$) to its Poisson counterpart. The corresponding inclusion probabilities are denoted by $\Box_1; \ldots; \Box_N$ and $p_1; \ldots; p_N$ respectively. The ensuing approach follows in the footsteps of Hàjek (1964) and relies more specifically on the results displayed in Theorem 5.1, p.1508. Let us start by defining the quantities

$$d_N = \prod_{i=1}^{[N]} p_i (1 \Box p_i) \text{ and } \bar{p}_N = \frac{1}{d_N} \prod_{i=1}^{[N]} p_i^2 (1 \Box p_i).$$

We assume that both R_N and T_N fulfill Assumption 6.1 for minoring constants \Box_{\Box} and p_{\Box} respectively and that the Poisson inclusion probabilities further satisfy the following condition.

Assumption 6.10

$$\limsup_{N \to +1} \frac{1}{N} \prod_{i=1}^{N} p_i(T_N) + 1.$$

Notice that, in this situation, $d_N = o(1=K)$ and \bar{p}_N is bounded. In addition, as shown in Hàjek (1964) (see p.1510 therein), the decomposition below holds for all i Pt1;:::;Nu:

$$p_{i} \square \square_{i} = \frac{\bar{p}_{N} \square p_{i}}{d_{N}} + o(1=d_{N})^{\square} p_{i}(1 \square \square_{i}).$$

This roughly means that the inclusion probabilities of the rejective sampling scheme are very close to those of the underlying Poisson design from which it was built. So close in fact that $D_{K;N}(R_N;T_N)$ and $r_{K;N}(R_N;T_N)$ asymptotically vanish. Therefore, as revealed by the following result, Theorem 6.8 also holds when the sample is constructed with a rejective plan.

Theorem 6.11 – Limit distribution in the rejective survey case. Suppose that all the conditions required in Theorem 6.8 hold together with Assumption 6.10. Then, for

$$= \sum_{p}^{2} := \int_{W}^{a} \frac{1}{p(w)} c_{X;W}(1; F_{W_{1}}(w_{1}); \dots; F_{W_{d}}(w_{d})) \int_{j=1}^{T^{d}} f_{W_{j}}(w_{j}) dw_{1} \dots dw_{d}$$

and provided that K - +1 as N - +1 so that $\overline{K}A(N=K) - \Box$ for some constant $\Box PR$, we have the convergence in distribution as N - +1:

$${}^{?}\overline{\mathsf{K}}(\mathsf{H}_{\mathsf{K};\mathsf{N}}^{\Box}\Box) \dot{\mathsf{O}} \mathsf{N} \stackrel{\Box}{\xrightarrow{}} {}^{\Box}_{1}\Box_{\Box}; \Box^{2}\Box^{2}_{p}\overset{\Box}{\xrightarrow{}}.$$
(6.12)

The proof of this theorem is available in Section 6.6. Notice that the limit variance does not depend on the inclusion probabilities $\Box_1; \ldots; \Box_N$, but on those of the underlying Poisson design, $p_1; \ldots; p_N$ namely. In that sense, this result is very similar to those obtained in Chapter 5: the asymptotic properties of the rejective sampling scheme are intricately linked to the Poisson plan with which it is coupled.

6.4 practical issues and illustrative experiments

6.4.1 On the choice of an optimal k

All results presented in the previous section depend on some appropriate number K of largest observations in the population $X_1; \ldots; X_N$. Unfortunately, the estimated tail quantile $X_{N \ \square K;N}$ from which $H_{K;N}^{\square}$ is computed may not be included in the sample. Hence, we need to choose a number k of largest values in the sample to which we may associate some K that respects the necessary conditions for consistency and asymptotic normality to hold (K = K(N) — +1 and $\overline{K}A(N=K)$ — $\Box \dagger 1$ as N — +1). Recall that we defined \Box_N^{\square} in Equation (6.5), a non-injective random map that assigns an index k in the sample to any index K in the population so that $X_{n \ \square k;n} = U_N^{\square}(N=K)$. Setting

$$\mathbf{\hat{R}}(\mathbf{k}) \coloneqq (\Box_{\mathbf{N}}^{\Box})^{-}(\mathbf{k}) \coloneqq \mathbf{N} \Box \prod_{i=1}^{\mathbf{N}} \frac{\mathbf{W}}{\Box_{i;n}};$$
where rs is the ceiling function, it is straightforward to show that the limit results stated in the sections above remains true for $\mathbf{R}(k) = \mathbf{H}_{k;n} = \mathbf{R}(k) = \mathbf{R}(k) = \mathbf{R}(k)$ of $\mathbf{R}(k) = \mathbf{R}(k)$ and $\mathbf{R}(k) = \mathbf{R}(k)$. This result can be naturally used to ground the construction of asymptotic Gaussian confidence intervals. The only work left is to find a suitable estimator \mathbf{p}_p^2 of \mathbf{R}_p^2 such that, by virtue of Sutsky's Lemma combined with Theorem 6.4, the quantity $\mathbf{R}(k) = \mathbf{R}(k) = \mathbf{R}_{k;n} = \mathbf{R}_{$

In practice, choosing an optimal threshold $X_{N \square K;N}$ is already complicated in the iid case. Many techniques have been proposed in the literature, often based on the minimization of the MSE (see Danielsson et al., 2001; Gomes and Oliveira, 2001; Goegebeur et al., 2008 and the references therein). Since they involve in general the estimation of the second order parameter \square , which goes beyond the scope of our analysis, we leave such considerations for future research. In the meantime, we propose to simply rely on heuristics such as the stability of the Horvitz-Thompson version of the Hill estimator around the appropriate k.

6.4.2 Numerical experiments

As a complement to the theoretical results established in the previous section, we provide here some illustrations based on simulations. In particular, we consider a model that does not fulfill condition ii) in Theorem 6.8, which requires the absolute continuity of $F_{X;W}$. The encouraging empirical results we obtain nonetheless give hope that this assumption may be relaxed. Such desirable extensions are left for future research.

6.4.2.1 Experiment setting

Simulations were based on the following model, chosen for its simplicity in terms of both computation and interpretation:

$$X = \frac{(1 \Box F_{W}(W))^{\Box} \Box 1}{\Box}; \Box^{\circ} 0;$$
$$W ; TN(\Box; \Box_{W}^{2}; w_{\Box}; w^{\Box});$$

where X is the variable of interest, W the auxiliary information with cdf F_W and $TN(\Box; \Box_W^2; w_{\Box}; w^{\Box})$ refers to the truncated Normal distribution over $[w_{\Box}; w^{\Box}]$, with expectation \Box and variance \Box_W^2 . Under such a representation, the distribution of X is a General Pareto with scale parameter 1 and EVI \Box , i.e. $F(x) = 1 \Box (1 + \Box x)^{\Box 1 = \Box}$. This is a well-known family of distributions, the second order properties of which are easily derived (De Haan and Ferreira, 2006, Section 3.2). In particular, we have

 \square = \square and A(x) = x^{\square} = \square . Following De Haan and Ferreira (2006, p.80), the optimal number of largest observations in the population is

$$\mathsf{K}^{\Box} = \mathsf{K}^{\Box}(\mathsf{N}) \Box \frac{\mathsf{N}^{2\Box}}{2} \Box^{3} (\mathsf{1} + \Box)^{2}$$

where txu is the integer part of x. It follows that $\frac{7}{K^{\Box}}A(N=K^{\Box}) - 0$ as N - +1. Concerning the joint distribution of X and W, it is straightforward to see that

$$F_{X;W}(x;w) = F(x) \wedge F_{W}(w);$$

which means that the copula linking both marginals is the well-known singular copula M (u; v) := $u \land v$, (u; v) P [0; 1]². Unfortunately, it is not derivable on it entire support and condition ii) in Theorem 6.8 is not fulfilled here. However, as we shall see in the next subsection, this does not impede tail estimation.

For a given population U_N of size N, where it is assumed that tW_i ; i P $U_N u$ are independent realizations of W, inclusion probabilities of the Poisson sampling scheme are defined as

$$p_i = p(W_i) = n \frac{W_i}{\sum_{i=1}^{N} W_i};$$
 (6.13)

with $n = \Box N$, $\Box P(0; 1)$, the desired expected sample size (Hàjek, 1964, Section 6, p.1512); this is the same formula as in Section 5.4 of Chapter 5. Thus defined, $p(W) P[p_{\Box}; p^{\Box}]$, where $p_{\Box} = \Box w_{\Box}=\Box$ and $p^{\Box} = \Box w^{\Box}=\Box$, which offers an easy way of ensuring Assumption 6.1 is fulfilled. Furthermore, given the formula used to compute X as a function of W, the more extreme the observations, the greater the probabilities of inclusion.

Numerical experiments were conducted on a set of populations with increasing sizes N = 10^3 , $5 \square 10^3$, 10^4 and $5 \square 10^4$. Several scenarios were investigated depending on the EVI \square ; they are summarized in Table 6.1. For each scenario, two sample sizes were considered: one small with n = $0.1 \square$ N and one relatively large with n = $0.5 \square$ N. Parameters of the distribution of W were chosen to ensure that for all i P U_N, p_i P [0.01; 1]. Specifically, we set $\square = 1$, $\square_W^2 = 0.09$, w $_\square = 0.1$ and w $\square = 2$, thereby implying that (p $_\square$; p \square) = (0.01; 0.02) when n = $0.1 \square$ N and (p $_\square$; p \square) = (0.05; 1) when n = $0.5 \square$ N.

Table 6.1 – List of scenarios depending on \Box and corresponding optimal $K^{\Box}(N)$

		K□(N)			
Scenario		$N = 10^3$	$N = 5 \square 10^3$	$N = 10^4$	$N = 5 \Box 10^4$
S ₁	1=2	11	26	37	83
S ₂	1	125	368	584	1709
S ₃	2	514	1863	3245	11760

For each scenario, we drew 1000 samples according to a rejective sampling scheme, following Algorithm 5.9 in Tillé (2006). The true inclusion probabilities, denoted by \Box_i , 1 § i § N, were deduced from their Poisson equivalents defined in Equation (6.13) using a Monte-Carlo approximation technique, based on the repetition (10⁵ times) of the basic algorithm stated in Section 6.3.2. Notice that since rejective sampling is a Poisson sampling conditioned upon its size, we have ($p_i = 1$) \dot{O} ($\Box_i = 1$).

The Horvitz-Thompson version of the Hill estimator was calculated using Equation (6.6) on each of the 1000 simulated samples. The ensuing results are presented herein-after.

6.4.2.2 Experiment results

Illustrations of the behavior of $H_{K;N}^{\Box}$ in a neighborhood of $K^{\Box}(N)$ as N grows are presented in Figure 6.1 for each scenario. As a complement, we display in Figure 6.2 the empirical estimator of $\Box^2 \Box_p^2$ in a neighborhood of $K^{\Box}(N)$ for each scenario and each sample size as N increases; for large populations, this gives some indication as to the form of the variance of $K H_{K;N}^{\Box} \Box$.

Since we only considered one fixed population, these results should be interpreted with caution: they only illustrate the behavior of $\overline{K} H_{K;N}^{\Box} \square H_{K;N}^{\Box}$ given the full vector $(X_1; W_1); \ldots; (X_N; W_N)$. We can see on Figure 6.1 that the Horvitz-Thompson version of the Hill estimator behaves perfectly well, even if the condition ii) in Theorem 6.8 is not satisfied. In particular, both its mean and variance decrease with N, more quickly when n = 0.5 \square N than when n = 0.1 \square N, and the distribution of the estimator appears to be symmetric around its classical version, which advocates normality. Scrutinizing Figure 6.2, we can see that the asymptotic variance of $\overline{K} H_{K;N}^{\Box} \square H_{K;N}^{\Box}$ seems indeed to be finite, depending on both the sample size and \square^2 (the smaller the EVI, the smaller the variance). This gives hope that the existence of a joint density may not be necessary for the asymptotic normality of our estimator to hold.



Figure 6.1 – Average values of $H_{K;N}^{\Box}$ (red line) and empirical 95% confidence band (pink area) computed on the 1000 simulated samples under scenario S_2 for $n = 0.1 \Box N$ (left hand plots) and $n = 0.5 \Box N$ (right hand plots), then compared to $H_{K;N}$ (black dotted line) for $N = 10^3$ (upper plots) and $N = 5 \Box 10^4$ (lower plots)



Figur e 6.2 – Estimation of $\Box \Box^2 \Box_p^2$ based on the 1000 simulated samples for $\Box = 0.1$ (dotted lines) and $\Box = 0.5$ (plain lines) under scenarios S₁ (grey lines), S₂ (black lines) and S₃ (red lines)

6.5 discussion

In an attempt to start adapting classical extreme value analysis to the case of survey data, we introduced in Section 6.2 a Horvitz-Thompson version of the widely celebrated Hill estimator of the extreme value index. After exhibiting some sufficient hypotheses on both the superpopulation model and the sampling scheme for the consistency of this novel statistic to hold in Section 6.2.2, we proved in Section 6.3 its asymptotic convergence to a Gaussian distribution when the survey design is of Poisson type. The exhibited rate of convergence appeared to be the same as the standard Hill estimator, namely \vec{K} , and the asymptotic variance was simply perturbed by a multiplicative constant depending solely on the sampling plan. In view of the empirical results presented in Section 6.4, hope is that the existence of a density copula linking those of the variable of interest and of the auxiliary information is not

necessary for the asymptotic normality to be true. This encourages further research to try and relax this assumption. Other improvements may be brought to these first results, for instance situations where the true inclusion probabilities are not available and replaced by an estimated version issued from post-calibration methods could be inspected. The assumptions made on the first and second order inclusion probabilities are also quite restrictive. Following in the lines of Boistard et al. (2012), higher order conditions could permit to get rid of Assumption 6.1. This remark is also true concerning the results obtained in Chapter 5. Sampling designs of other nature than the Poisson type may be considered as well, especially complex ones such as that used in the INCA2 database.

Though we could not directly apply our findings to dietary risk analysis, they could be of great interest in the context of big data management. Indeed, it is more and more frequent to meet databases that increase regularly (in finance, information about the markets is stocked every hour at least) and cannot be saved, thus analyzed, on a single computer. When accessing such huge files becomes a challenge, sampling is a natural solution. In this context, the superpopulation model and the asymptotic nature of our results are perfectly relevant. Moreover, the analyst has then complete control over the survey scheme they desire to adopt, which is typically rarely the case with institutional data. Hence, the Poisson and rejective schemes, which are not of frequent use in practice, are revealed as especially convenient for such types of analyses. With these potential assets in mind, we hope that this preliminary step towards the elaboration of a new extreme value theory for survey data will engender further research in the near future.

6.6 proofs and supplements

6.6.1 Consistency of the Horvitz-Thompson variant of the Hill estimator

We start by establishing the following intermediate results, in order to describe next the limit behavior of the Horvitz-Thompson tail empirical process.

First, we introduce the point measure:

$$\widetilde{\Box}_{N}^{\Box} \coloneqq \frac{1}{K} \prod_{i=1}^{|\Gamma|} \frac{\Box_{i}}{\Box_{i}} \Box_{X_{i}} = U(N = K).$$

Notice that the point measure \Box_N^{\Box} can be obtained from the latter by replacing the threshold U(N=K) by the empirical counterpart U_N^{\Box} (N=K).

Lemma 6.12 Under the assumptions of Theorem 6.4, as N;n and K tend to infinity so that K=N - 0, we have:

$$\widetilde{\Box}_{N}^{\square} \dot{O} \square_{\square 1 = \square}; \tag{6.14}$$

ī.

where "Ò" denotes weak convergence in the space of positive Radon measures on (0;+1].

Proof Consider first the tail empirical process

$$\Box_{\mathsf{N}} := \frac{1}{\mathsf{K}} \prod_{i=1}^{[\mathsf{N}]} \Box_{\mathsf{X}_i = \mathsf{U}(\mathsf{N} = \mathsf{K})}.$$

We shall prove that for any t $\,^\circ\,$ 0, as N; n and K tend to + 1 , provided K=N converges to 0,

$$D_{N}^{\Box}(t) := \widetilde{\Box}_{N}^{\Box}(t; +1] \Box \Box_{N}(t; +1] - 0 \text{ in } L_{2}.$$
(6.15)

Indeed, @ ° 0, provided Assumption 6.1 and Assumption 6.2 hold, we have

$$E(D_{N}^{\Box}(t)) = \frac{1}{K} \bigcap_{i=1}^{[N]} E \qquad \frac{E \Box_{i} \Box_{t}(X_{i}; W_{i})_{1 \le i \le N} u}{\Box_{i}} \Box 1 \quad I \ t X_{i} \circ t \ U(N=K) u$$
$$= 0;$$

together with

$$E D_{N}^{\circ}(t)^{2} = \frac{1}{K^{2}} \prod_{i=1}^{N} E E \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} (X_{i}; W_{i})_{1 \le i \le N} u = I X_{i}^{\circ} t U(N=K) u + \frac{2}{K^{2}} \prod_{1 \le i = j \le N} E E \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} (X_{i}; W_{i})_{1 \le i \le N} u = \frac{1}{K^{2}} \prod_{i=1}^{N} E \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} (X_{i} \circ t U(N=K))^{i} + \frac{2}{K^{2}} \prod_{1 \le i = j \le N} E \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} (X_{i} \circ t U(N=K)) u + \frac{2}{K^{2}} \prod_{1 \le i = j \le N} E \prod_{i=1}^{Q} \frac{1}{Q} \prod_{i=1}^{Q} (X_{i} \circ t U(N=K)) + \frac{1}{Q} \prod_{i=1}^{Q} \frac{1}{K^{2}} \prod_{i=1}^{N} P (X_{i} \circ t U(N=K)) + \frac{1}{Q^{2}} \prod_{i=1}^{Q} \frac{2}{K^{2}} \prod_{i=1}^{N} P (X_{i} \circ t U(N=K)) + \frac{1}{Q^{2}} \prod_{i=1}^{Q} \frac{2}{K^{2}} \prod_{1 \le i \le N} P (X_{i} \circ X_{j} \circ t U(N=K))^{i} = : N; K(t).$$

Since \overline{F} is supposed to be regularly varying with index $\Box 1=\Box$ and K=N — 0, we have $P(X_i \circ t \cup (N=K)) \square t^{\square 1 = \square} K=N$ for all i $Pt1; \dots; Nu$ as N and K go to infinity. It follows that as N; n; K - +1,

$$\hat{N}_{K}(t) \square \frac{1}{\square_{\square}} \square 1 t^{\square 1=\square} \frac{1}{K} + \frac{\hat{n}_{\square}}{\square_{\square}^{2}} t^{\square 2=\square} 1 \square \frac{1}{N} - 0.$$

Hence, the convergence in Equation (6.15) is proved and the desired convergence will then result from the fact that $\Box_N \dot{O} \Box_{\Box 1=\Box}$, see Resnick (2007, Theorem 4.1).

We next prove the lemma below, claiming that the threshold U_N^{\Box} (N=K) and U(N=K) are asymptotically equivalent in probability.

Lemma 6.14 Under the assumptions of Theorem 6.4, we have: as N;n and K tend to infinity,

$$\frac{J_{N}^{\perp}(N=K)}{U(N=K)} - 1 \text{ in probability.}$$
(6.16)

Proof This is a straightforward consequence of Lemma 6.12. Indeed, for all " ° 0, we have:

$$P \stackrel{[U_{N}]}{\longrightarrow} (N=K) = 1 \stackrel{[V_{N}]}{\longrightarrow} = P (U_{N}^{\Box} (N=K) \circ (1+") U(N=K)) + P (U_{N}^{\Box} (N=K) \dagger (1 \Box ") U(N=K)) \\ \$ P \stackrel{[1]}{\longrightarrow} \stackrel{[N]}{\longrightarrow} \underset{i=1}{\square} \stackrel{[i]}{\longrightarrow} \underset{i=1}{\square} ((1+") U(N=K); +1] \circ \frac{K}{N} + P \stackrel{[1]}{\longrightarrow} \stackrel{[N]}{\longrightarrow} \underset{i=1}{\square} \stackrel{[i]}{\longrightarrow} \underset{i=1}{\square} ((1 \Box ") U(N=K); +1] \dagger \frac{K}{N} \\ \$ P (\square_{N}^{\Box} (1+"; +1] \circ 1) + P (\square_{N}^{\Box} (1 \Box "; +1] \$ 1).$$

Therefore, by virtue of the lemma previously established, we asymptotically have: $\tilde{\Box}_{N}^{\Box}(1 + "; +1] - 1 = (1 + ")^{1=\Box} + 1$ and $\tilde{\Box}_{N}^{\Box}(1 + "; +1] - 1 = (1 \Box ")^{1=\Box} \circ 1$ in probability. Combined with the bound above, this proves the lemma.

Equipped with these preliminary results, we may now tackle the proof of Theorem 6.4, which is recalled below for convenience.

Theorem – Consistency. Let K = K(N) be a sequence of integers such that K - +1 and K=N - 0 as N; n - +1. Provided that Assumption 6.1 and Assumption 6.2 are fulfilled, we then have, as N and n tend to +1:

$$H_{K:N}^{\Box} \square \square \text{ in probability.}$$
(6.17)

Proof The consistency result can be established by following line by line the proof for the consistency of the Hill estimator in the iid situation given in Resnick (2007): by a continuous mapping theorem argument, one derives from Lemma 6.12 and Lemma 6.14 that the Horvitz-Thompson tail empirical process \Box_N^{\Box} converges in probability to $\Box_{\Box 1=\Box}$ in the space of positive Radon measures on (0; +1]. Then, it classically suffices to integrate the tail measures against dt=t (df. Equation (6.4) and Equation (6.6)) and apply the convergence previously mentioned. See Resnick (2007, Section 4.4.1) for further details.

6.6.2 Limit distribution of $H^{p}_{K:N}$ in the Poisson survey case

Before handling Theorem 6.8, we start by establishing three intermediate results, which are introduced herein-after. The first lemma claims that the quantity $r_{K;N}$ defined in Equation (6.11) converges to 1 in probability.

Lemma 6.17 Let $\Box_1; \ldots; \Box_N$ and $p_1; \ldots; p_N$ be respectively the inclusion indicators and probabilities of a Poisson survey plan in some population $U_N \coloneqq t1; \ldots; Nu$. Then, provided Assumption 6.1 holds, for any K P t1; \ldots; Nu such that $K \coloneqq K(N) \Box_{N-1} + 1$ we have

$$\mathbf{r}_{\mathbf{K};\mathbf{N}} := \frac{1}{\mathbf{K}} \frac{\prod_{i=1}^{\mathbf{M}} \frac{\prod_{i=1}^{N} \prod_{i=1}^{i+1;\mathbf{N}} \prod_{i=1}^{\mathbf{P}} 1}{p_{\mathbf{N} \prod_{i=1}^{i+1;\mathbf{N}} \prod_{i=1}^{\mathbf{P}} 1}.$$

Proof Recall that under a Poisson sampling plan, all \Box_1 ;:::; \Box_N are independent. We had set $E \Box \Box (X_i; W_i) \coloneqq p_i$ for all i Pt1;:::; Nu, hence

$$\mathsf{E}(\mathsf{r}_{\mathsf{K};\mathsf{N}}) = \frac{1}{\mathsf{K}} \frac{\mathsf{I}^{\mathsf{N}}}{\mathsf{I}_{i=1}} \mathsf{E} \quad \frac{\mathsf{E} \left[\mathsf{I}_{\mathsf{N} \ i+1;\mathsf{N}} \right] \left[(\mathsf{X}_{1};\mathsf{W}_{1}); \ldots; (\mathsf{X}_{\mathsf{N}};\mathsf{W}_{\mathsf{N}}) \right]^{\mathsf{I}}}{\mathsf{P}_{\mathsf{N} \ i+1;\mathsf{N}}} = 1$$

In addition, $V = (X_i; W_i) := p_i (1 = p_i)$ for all i Pt1;:::; Nu, therefore

$$V(\mathbf{r}_{K;N}) = \frac{1}{K^{2}} \prod_{i=1}^{[N]} E \frac{V \prod_{N=i+1;N} (X_{1}; W_{1}); \dots; (X_{N}; W_{N})}{p_{N=i+1;N}^{2}} + \frac{1}{K^{2}} \prod_{i=1}^{[N]} V \frac{E \prod_{N=i+1;N} (X_{1}; W_{1}); \dots; (X_{N}; W_{N})}{p_{N=i+1;N}} = \frac{1}{K^{2}} \prod_{i=1}^{[N]} E \prod_{N=i+1;N} \frac{1}{p_{N=i+1;N}} \prod_{i=1}^{[N]} \frac{1}{K}.$$

Under Assumption 6.1 it is clear that V $(r_{K;N}) = O(\frac{1}{K})$ when K $\underset{N-1}{\square} + 1$. This concludes the proof.

We now move to the quantity $\overrightarrow{K}^{r}_{K;N} H^{p}_{K;N} \square H_{K;N}^{r}$ appearing in Equation (6.10). The result below reveals that is vanishes asymptotically.

Lemma 6.19 Suppose that Assumption 6.5 is fulfilled by the underlying heavy-tailed model and that Assumption 6.1 is satisfied by the considered sequence of Poisson inclusion probabilities $p_1; \ldots; p_N$, N • 1. Assume also that K — +1 as N — +1 so that $\overline{KA}(N=K)$ — \Box for some constant \Box PR. Then we have

Proof Let us start by introducing the weighted versions of log-spacings in the population, given by

$$@ Pt1;:::;Nu; \square_i := i (log X_{N \square i+1;N} \square log X_{N \square i;N})$$

Let K Pt1;:::;Nu, these random variables are intrinsically linked to both $H_{K;N}$ and $H^p_{K;N}$, given by Equation (6.3) and Equation (6.6) respectively. Indeed, they can be expressed as

$$H_{K;N} = \frac{1}{K} \prod_{i=1}^{K} \Box_i$$

and

$$H_{K;N}^{p} = \frac{1}{K} \frac{\prod_{j=1}^{m} \frac{m_{j+1;N}}{p_{N-j+1;N}}}{\prod_{j=1}^{m} \frac{m_{j+1;N}}{p_{N-j+1;N}}} \frac{\prod_{j=1}^{l} \prod_{j=1}^{m} \frac{m_{j+1;N}}{p_{N-j+1;N}}}{\prod_{j=1}^{l} \frac{m_{j+1;N}}{p_{N-j+1;N}}} \square_{i}.$$

Combining the two we immediately obtain

$$r_{K;N} H^{p}_{K;N} \Box H_{K;N} = \frac{1}{K} \prod_{i=1}^{[n]} \frac{1}{i} \prod_{j=1}^{[n]} \frac{m_{n} \Box_{j+1;N}}{p_{N} \Box_{j+1;N}} \Box \prod_{i=1}^{!} \frac{1}{i} \prod_{j=1}^{!} \frac{m_{n} \Box_{j+1;N}}{p_{N} \Box_{j+1;N}} \Box \prod_{i=1}^{!} \frac{1}{i} \prod_{i=1}^{!} \frac{m_{n} \Box_{i+1;N}}{p_{N} \Box_{j+1;N}} \Box \prod_{i=1}^{!} \frac{1}{i} \prod_{i=1}^{!} \frac{m_{n} \Box_{i+1;N}}{p_{N} \Box_{i+1;N}} \Box \prod_{i=1}^{!} \frac{1}{i} \prod_{i=1}^{!} \frac{m_{n} \Box_{i+1;N}}{p_{N} \Box_{i+1;N}} \Box \prod_{i=1}^{!} \frac{1}{i} \prod_{i=1}^{!} \frac{m_{n} \Box_{i+1;N}}{p_{N} \Box_{i+1;N}}} \Box \prod_{i=1}^{!} \frac{m_{n} \Box_{i+1;N}}{p_{N} \Box_{i+1;N}}}$$

When Assumption 6.5 is fulfilled, it is possible to approximate the distribution of the \Box_i 's corresponding to the K + 1 largest values. Denoting by E_1 ;:::; E_K a collection of independent random variables with standard exponential distribution, the random variables \Box_i are approximately distributed as

$$\Box + \frac{i}{K+1} A \frac{N+1}{K+1} E_{i}; 1 \S i \S K.$$
 (6.19)

This property is at the basis of most of the asymptotic analyzes that were led concerning $H_{K;N}$, see de Haan and Resnick (1998) for more details. As mentioned in Remark 6.6, alternative approaches taking advantage of the Glivenko-Cantelli and Donsker theorems in the formulation of the Hill estimator in Equation (6.3) were also developed, see Resnick (2007) for instance.

Given the decomposition in Equation (6.18), just like $r_{K;N}$ in Lemma 6.17 the expectation and variance of $r_{K;N}$ $H^p_{K;N} \square H_{K;N}$ are easily derived by conditioning upon the full vector of observations $(X_1; W_1); \ldots; (X_N; W_N)$. In particular, it is straightforward to see that

$$\mathsf{E} \stackrel{\square}{\mathbf{r}}_{\mathsf{K};\mathsf{N}} \mathsf{H}^{\mathsf{p}}_{\mathsf{K};\mathsf{N}} \square \mathsf{H}_{\mathsf{K};\mathsf{N}} \stackrel{\square}{=} \mathsf{E} \quad \frac{1}{\mathsf{K}} \stackrel{\mathsf{I}^{\mathsf{q}}}{\underset{i=1}{\overset{\mathsf{I}}{\mathsf{I}}} \frac{1}{\mathsf{i}} \frac{\mathsf{I}^{\mathsf{q}}}{\underset{j=1}{\overset{\mathsf{I}}{\mathsf{I}}} \frac{\mathsf{\Box}_{\mathsf{N}} \square j+1;\mathsf{N}}{\mathsf{p}_{\mathsf{N}} \square j+1;\mathsf{N}} \square 1 \square_{\mathsf{I}} = 0.$$

Turning now to the variance, under Assumption 6.1, we have

$$V \stackrel{\Box}{\mathbf{r}}_{\mathbf{K};\mathbf{N}} \mathbf{H}^{\mathbf{p}}_{\mathbf{K};\mathbf{N}} \Box \mathbf{H}_{\mathbf{K};\mathbf{N}} \stackrel{\Box}{=} \frac{1}{\mathbf{K}^{2}} \frac{\prod_{i=1}^{n} \mathbf{E}}{\sum_{i=1}^{n} \sum_{j=1}^{n} \frac{1}{\mathbf{p}_{\mathbf{N}} \sum_{j=1}^{n} \sum_{i=1}^{n} 1}{\sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=$$

To simplify notations, set $c = p_{\Box}^{-1} \Box 1$ and consider the variable $A_{K;N} = A \frac{N+1}{K+1}$ involved in Equation (6.19). Using this particular equation, we are able to establish asymptotic properties of the right hand side in the inequality herein above: as N; K - +1 we have

ī.

$$\frac{1}{K^{2}} \prod_{i=1}^{M} E = \prod_{i=1}^{2} \frac{c}{i} = \frac{1}{K^{2}} \prod_{i=1}^{M} + \frac{i}{K+1} = \prod_{i=1}^{2} A_{K;N} = E_{i}^{2} \prod_{i=1}^{2} \frac{c}{i}$$

$$= \frac{2c}{K^{2}} \prod_{i=1}^{M} \frac{1}{i} + \frac{2c}{K^{2}(K+1)^{-2}} A_{K;N}^{2} \prod_{i=1}^{M} \frac{1}{i^{2-1}} \prod_{i=1}^{2} \prod_{i=1}^{2} \frac{1}{i^{2-1}} \prod_{i=1}^{2} \prod_$$

Because A PR₀ with \Box † 0, then A_{K;N} — 0 as N; K — +1 , and we can conclude that $\bigvee_{K;N}^{\Box} H_{K;N}^{p} \Box H_{K;N} = o(1=K).$

The last intermediate result concerns the quantity $Q_N^{(3)}$ in Equation (6.10). It claims that the latter converges weakly to a centered Normal distribution.

Lemma 6.21 Suppose that Assumption 6.1 is satisfied by the considered sequence of Poisson inclusion probabilities $p_1; \ldots; p_N$, N • 1, constructed from some set of auxiliary variables as in Assumption 6.7. Further assume the conditions i) \Box iii) introduced in Theorem 6.8. Then, for

$$\Box_p^2 \coloneqq \prod_{W=1}^{a} \frac{1}{p(w)} c_{X;W}(1; F_{W_1}(w_1); \ldots; F_{W_d}(w_d)) \prod_{j=1}^{T^d} f_{W_j}(w_j) dw_1 \ldots dw_d$$

and provided that K - + 1 as N - + 1 so that K = o(N), we have the convergence in distribution as N - + 1:

$$\overline{K}(1 \square r_{K;N}) O N \stackrel{\square}{0}; \square_p^2 \square 1^{\square}.$$

Proof This proof is based on the application of Feller (1971, Theorem 3, p.262) to the collection of random variables $tZ_{i;N}(\Box)$; 1 § i § Nu defined for all i Pt1;:::;Nu as follows:

$$Z_{i;N} (\Box) \coloneqq \frac{2^{1}}{\overline{K}} \ 1 \Box \frac{\Box_{i}}{p_{i}} \ I \ t \ X_{i} \ \circ \ X_{N \ \Box \ K;N} \ u;$$

with distribution P_i only depending on the survey scheme (they are conditioned upon the vectors $t(X_i; W_i); 1$ i Nu). Indeed, notice that we have

$$\widehat{K}(1 \Box r_{K;N}) = \prod_{i=1}^{N} Z_{i;N}(\Box).$$

In order to apply this theorem, we first have to check the three conditions below.

(C₁) For all i Pt1;:::; Nu, we have E $(Z_{i;N}(\Box)) = 0$ and V $(Z_{i;N}(\Box)) = \Box_i^2 \dagger 1$. (C₂) There exists some real constant $\Box^2 \dagger 1$ such that $S_N^2 := \sum_{i=1}^{\infty} \Box_i^2 \Box_{N-1}^P \Box^2$. (C₃) For each t ° 0, we have $\sum_{i=1}^{\infty} \sum_{|z| \bullet t S_N}^N z^2 P_i(dz) \Box_{N-1}^P 0$.

 $\acute{\mathbf{U}}$ Condition (C₁) Let us start by calculating the expectation of $Z_{i;N}$ (\Box) for all i P t1;:::; Nu. Because the \Box_1 ;:::; \Box_N are independent random variables with respective Bernoulli distributions B(p₁);:::; B(p_N), it is straightforward to see that

$$\mathsf{E}\left(\mathsf{Z}_{i\,;N}\left(\Box\right)\right)=\frac{2}{\overline{K}}\quad 1\Box\frac{\mathsf{E}\left[\Box\right]}{p_{i}}\frac{\Box}{\mathsf{t}}\left(\mathsf{X}_{i\,;W_{i}\,)_{1\,\S\,i\,\S\,N}\,u}\right]}{p_{i}}\quad\mathsf{I}\mathsf{t}\mathsf{X}_{i}\,\circ\,\mathsf{X}_{N\,\Box\,K;N}\,u=0.$$

As for the variance, we have

$$\begin{array}{l} V\left(Z_{i\,;N}\left(\Box\right)\right) = E \stackrel{\sqcup}{Z_{i\,;N}^{2}}\left(\Box\right)^{\sqcup} \\ = \frac{1}{K}E \stackrel{\Pi}{1} \Box \stackrel{\Box_{i}}{p_{i}} \stackrel{\Box_{2}}{=} t\left(X_{i\,;W_{i}}\right)_{1 \,\$\,i\,\$\,N} u \quad I \, t \, X_{i} \,\,^{\circ} \, X_{N \, \Box \, K;N} \, u \\ = \frac{1}{K} \frac{1 \, \Box \, p_{i}}{p_{i}} \, I \, t \, X_{i} \,\,^{\circ} \, X_{N \, \Box \, K;N} \, u =: \, \Box_{i}^{2} \, \dagger \, 1 \, . \end{array}$$

Therefore, condition (C_1) is fulfilled.

$$S_{N}^{2} := \prod_{i=1}^{N} \Box_{i}^{2} = \frac{1}{K} \prod_{i=1}^{N} \frac{1}{p_{i}} I t X_{i} \circ X_{N \square K; N} u \square 1;$$

where $X_{N \ \square K;N}$ is a consistent estimator of U(N=K) (Resnick, 2007, Section 4.4.1, p.81). With this remark in mind, we will proceed in two steps and successively prove that there exists a real constant $\ \square^2 \circ 0$ such that

$$\tilde{\mathbb{S}}_{N}^{2} \coloneqq \frac{1}{K} \prod_{i=1}^{|\Gamma|} \frac{1}{p_{i}} \operatorname{It} X_{i} \circ \operatorname{U}(N=K) u \Box 1 \prod_{N=1}^{P} \Box^{2};$$

then that $S_N^2 \square \tilde{S}_N^2 \square_{N-1}^P$ 0. Since $S_N^2 = S_N^2 \square \tilde{S}_N^2 + \tilde{S}_N^2$, this will yield the desired result.

Let us start with \tilde{S}_N^2 . By virtue of the law of large numbers, as N — +1 we have

$$\tilde{S}_{N}^{2} + 1 = \frac{N}{K} \int_{W}^{a} \frac{1}{U(N=K)} \frac{1}{p(w)} P_{X;W}(dx;dw) + o_{P}(1);$$

provided that this integral exists. Under Assumptions i) and ii) in Theorem 6.8, this yields

$$\tilde{S}_{N}^{2} + 1 = \frac{N}{K} \int_{W}^{a} \frac{1}{U(N=K)} \frac{1}{p(w)} c_{X;W} (F(x); F_{W_{1}}(w_{1}); \dots; F_{W_{d}}(w_{d}))$$
$$\Box f(x) \int_{j=1}^{T_{1}^{d}} f_{W_{j}}(w_{j}) dx dw_{1} \dots dw_{d} + o_{P}(1).$$

Further set $u := \frac{N}{K} \overline{F}(x)$, then for K := K(N) - +1 as N - +1 and K = o(N), under Assumption iii) we have:

$$\tilde{S}_{N}^{2} + 1 = \frac{\begin{bmatrix} a & a \\ W & 0 \end{bmatrix}}{\begin{bmatrix} T^{d} \\ W \\ 0 \end{bmatrix}} c_{X;W} \begin{bmatrix} T \\ W \\ C_{X;W} \end{bmatrix} \left[1 \\ \frac{K}{N} \\ U; F_{W_{1}}(w_{1}); \dots; F_{W_{d}}(w_{d}) \right]$$

$$= \frac{T^{d}}{\begin{bmatrix} F \\ W \\ W \end{bmatrix}} f_{W_{j}}(w_{j}) du dw_{1} \dots dw_{d} + o_{P}(1)$$

$$= \frac{P}{\begin{bmatrix} 0 \\ W \end{bmatrix}} \frac{1}{p(W)} c_{X;W}(1; F_{W_{1}}(w_{1}); \dots; F_{W_{d}}(w_{d}))$$

$$= \frac{T^{d}}{\begin{bmatrix} T^{d} \\ W \end{bmatrix}} f_{W_{j}}(w_{j}) dw_{1} \dots dw_{d} = : \square_{P}^{2}.$$

Provided that all the aforementioned hypotheses hold, we can conclude that

$$\tilde{S}_{N}^{2} \underset{N-1}{\square} \square_{p}^{2} \square 1.$$

There remains to control the quantity $S_N^2 \square \tilde{S}_N^2 \square$, which we denote by S_N for simplicity. For any fixed N P N \square and $\square \circ 0$, it can be decomposed as follows:

$$S_{N} := \frac{1}{K} \prod_{i=1}^{[N]} \frac{1}{p_{i}} (I \ t \ X_{i} \ \circ \ U(N=K) u \square I \ t \ X_{i} \ \circ \ X_{N \square K;N} \ u) = S_{N}^{(1)}(\square) + S_{N}^{(2)}(\square);$$

where

$$S_{N}^{(1)}(\Box) := S_{N} I \qquad \begin{array}{c} X_{N \ \Box \ K;N} \\ U(N=K) \end{array} \Box 1 \begin{array}{c} \\ \\ \end{array}$$

and

$$S_{N}^{(2)}(\Box) = S_{N} I \overset{"}{=} \frac{X_{N \Box K;N}}{U(N=K)} \Box 1 \overset{*}{\S} \Box$$

We shall use this decomposition to prove that $S_N \bigsqcup_{N=1}^{P} 0$. Referring to Assumption 6.1, it is easy to see that $S_N^{(1)}(\Box) \bigsqcup_{N=1}^{P} 0$. Indeed,

$$\mathbb{S}_{N}^{(1)}(\Box) \stackrel{P}{=} \frac{1}{p_{\Box}} \stackrel{P}{\longrightarrow} \frac{1}{K} \prod_{i=1}^{[N]} (I \ t \ X_{i} \ ^{\circ} \ U(N=K) u \ \Box \ I \ t \ X_{i} \ ^{\circ} \ X_{N \ \Box \ K; N} \ u)$$

The law of large numbers ensures that as N — 1 ,

$$\frac{1}{K} \prod_{i=1}^{[N]} |I tX_i \circ U(N=K)u \Box 1 = \frac{N}{K} P (X_i \circ U(N=K)) \Box 1 + o_P(1) - 0.$$

In addition, we know that for all \square ° 0,

$$I \stackrel{*}{=} \frac{X_{N \ \square K;N}}{U(N=K)} \square 1 \stackrel{*}{=} \stackrel{*}{\square} \frac{P}{N-1} 0$$

(Resnick, 2007, Section 4.4.1, p.81), hence $S_N^{(1)}(\Box) \Box_{N-1}^P = 0$.

As for $S_N^{(2)}(\Box)$, combining Assumption 6.1 with triangular inequalities provides an absolute bound:

$$\begin{split} & \overset{(2)}{\overset{(2)}}{\overset{(2)}}{\overset{(2)}{\overset{(2)}{\overset{(2)}{\overset{(2)}{\overset{(2)}}{\overset{(1}{\overset{(2)}{\overset{(2)}}{\overset{(1}{\overset{(2)}{\overset{(2)}}{\overset{(1}{\overset{(1}{\overset{(2)}}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}{\overset{(1}}{\overset{(1}{\atop(1}{\overset{(1}{\overset{(1}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{(1}}{\overset{$$

Denote by c_N () the quantity in the right hand part of the last inequality, i.e.

$$c_{N}(\Box) \coloneqq \frac{1}{p_{\Box}} \frac{1}{K} \prod_{i=1}^{|V|} I tX_{i} \circ (1 \Box \Box) U(N=K)u \Box I tX_{i} \circ (1+\Box) U(N=K)u.$$

Applying the law of large numbers and recalling that $\overline{F} \ P \ R_{\Box 1=\Box},$ we get that as N - +1 ,

$$c_{N} (\Box) = \frac{1}{p_{\Box}} \frac{N}{K} (1 \Box)^{\Box 1=\Box} (1 + \Box)^{\Box 1=\Box} \overline{F}(U(N=K)) + o_{P}(1)$$
$$\Box - \frac{P}{p_{\Box}} \frac{1}{p_{\Box}} (1 \Box)^{\Box 1=\Box} (1 + \Box)^{\Box 1=\Box} =: c_{\Box};$$

meaning that c_N () $\square = \frac{P}{N-1}$ c_{\square} . Furthermore, notice that $c_{\square} = \square = 0$, since

$$\begin{cases} \frac{2}{p_{\Box}} + o(\Box) = 0 \\ = 0 \end{cases}$$

We may now prove that $S_N \bigsqcup_{N=1}^{P} 0$ by going back to the definition.

Formally, we have shown that $|S_N| \S |S_N^{(1)}(\square)| + c_N(\square) =: \square_N(\square)$, for all $\square \circ 0$. We also know that $\square_N(\square) \square_{N-1}^P c_{\square}$ and $c_{\square} \square_{-0}^- 0$. Now for any fixed $\square \circ 0$, we have to verify that $P(|S_N| \circ 2\square) \square_{N-1}^- 0$. First choose $\square_0 \circ 0$ such that for all $0 \uparrow \square \S \square_0$, $|c_{\square}| \S \square$ then take some $\square \circ 0$. We need to prove that there exists N_0 PN such that for all $N \bullet N_0$, $P(|S_N| \circ 2\square) \S \square$ In order to construct this N_0 , first fix any $\square \circ 0$ such that $\square \S \square_0$. Since $\square_N(\square) \square_{N-1}^P c_{\square}$, there exists N_0 PN such that for all $N \bullet N_0$, we have $P(|\square_N(\square) \square c_{\square} \circ \square) \S \square$ In parallel, since $\square \S \square_0$, we have for all N PN :

 $|S_{N} | \S | \square_{N} (\square)| \S | \square_{N} (\square) \square c_{\square} + |c_{\square}| \S | \square_{N} \square c_{\square} + \square$

This implies that for all $N \bullet N_0$,

$$\mathsf{P}\left(|\mathsf{S}_{\mathsf{N}}|^{\circ} \ 2\Box\right) \ \mathsf{S} \ \mathsf{P}\left(|\Box_{\mathsf{N}}\left(\Box\right) \ \Box \ \mathsf{c}_{\Box}| + \ \Box^{\circ} \ 2\Box\right) = \ \mathsf{P}\left(|\Box_{\mathsf{N}}\left(\Box\right) \ \Box \ \mathsf{c}_{\Box}|^{\circ} \ \Box\right) \ \mathsf{S} \ \Box$$

Since this is true for any \square° 0, this means that $S_{N} \square_{N-1}^{P} 0$.

In fine, we can conclude that under all the hypotheses stated in Lemma 6.21,

$$S_N^2 \square \xrightarrow{P}_{N-1} \square_p^2 \square 1.$$

$$\begin{split} Z_{N}\left(t\right) &\coloneqq \left| \prod_{i=1}^{N} a \\ Z_{N}\left(t\right) &\coloneqq \left| \sum_{i=1}^{N} |z|^{\bullet} t S_{N} \\ &= \prod_{i=1}^{N} E \left| \sum_{i;N}^{2} (\Box) I t |Z_{i;N}\left(\Box\right)| \bullet t S_{N} u \\ &= \frac{1}{K} \prod_{i=1}^{N} I t X_{i} \circ X_{N \Box K;N} u \\ &= \frac{1}{K} \prod_{i=1}^{N} I t X_{i} \circ X_{N \Box K;N} u \\ &= E \left| 1 \Box \Box \Box \right|_{p_{i}}^{2} \left| \left| 1 \Box \Box \right|_{p_{i}}^{2} \right| = \frac{1}{K} \prod_{i=1}^{K} \frac{t S_{N}}{K I t X_{i} \circ X_{N \Box K;N} u} \right|_{t}^{*} \left| t (X_{i}; W_{i})_{1 \le i \le N} u \right|_{t}^{*} . \end{split}$$

Using Hölder's inequality, we obtain

$$Z_{N}(t) \S \frac{1}{K} \bigcap_{i=1}^{[N]} I t X_{i} \circ X_{N \ \cup \ K; N \ U} E \cap I \cap \frac{1}{p_{i}} \bigcap_{i=1}^{3} \int_{t}^{1} t (X_{i}; W_{i})_{1 \$ i \$ N \ U} \cap \frac{1}{p_{i}} \cap_{i=1}^{2} \cap_{i=1}^{2} \cap_{i=1}^{3} \cap_{i=1}^{2} \cap_{i=1}^{3} \cap_{i=1}^{2} \cap_{i=1}^{3} \cap_{i=$$

Observe that under Assumption 6.1, we have

Moreover, conditional on the vectors $t(X_i; W_i)$; 1 § i § Nu, the random variable $\Box = \frac{\Box_i}{p_i} = equals$ either $(p_i = 1) = p_i$ with probability p_i or 1 with probability $1 = p_i$. Therefore, by virtue of Markov's inequality, we can further bound $Z_N(t)$ from above:

$$Z_{N}(t) \begin{cases} \frac{1}{K} \bigcap_{i=1}^{N} |tX_{i} \circ X_{N \square K;N} u \ 3 \ \frac{1}{p_{\square}} \square 1 \\ \square \ \frac{1}{K} |tX_{i} \circ X_{N \square K;N} u \ \frac{2(1 \square p_{i})}{tS_{N}} \square^{1=3} \\ \end{cases} \\ \begin{cases} \frac{3^{2=3}}{K^{4=3}} \square \ \frac{1}{p_{\square}} \square 1 \\ \square \ \frac{1}{p_{\square}} \square 1 \\ \square \ \frac{1}{1} |tX_{i} \circ X_{N \square K;N} u \ \frac{2(1 \square p_{i})}{tS_{N}} \square^{1=3} \\ \end{cases} \\ \end{cases}$$

Using again Assumption 6.1, this yields

where we have shown that $S_N \bigsqcup_{N \to 1}^{P} \Box$. Consequently, the right hand part of this last inequality tends to 0 in probability as N tends to infinity for any t ° 0. Hence, condition (C₃) is fulfilled.

With all three conditions (C_1) , (C_2) and (C_3) satisfied, by virtue of Feller (1971, Theorem 3, p.262) we finally have

$$\stackrel{?}{\overline{\mathsf{K}}} (1 \Box \mathsf{r}_{\mathsf{K};\mathsf{N}}) \underset{\mathsf{N} \longrightarrow 1}{\overset{\mathsf{O}}{\longrightarrow}} \operatorname{N} \stackrel{\Box}{0}; \Box_{\mathsf{p}}^2 \Box 1$$

 \square

We are now fully equipped to prove Theorem 6.8, the statement of which is recalled below.

Theorem – Limit distribution in the Poisson survey case. Suppose that Assumption 6.5 is fulfilled by the underlying heavy-tailed model and that Assumption 6.1 is satisfied by the considered sequence of Poisson inclusion probabilities $p_1; \ldots; p_N$, N • 1, constructed from some set of auxiliary variables as in Assumption 6.7. Further assume the conditions introduced herein-after.

- i) The marginal cdf F is absolutely continuous with respect to the Lebesgue measure with density f.
- ii) The joint cdf $F_{X;W}$ is absolutely continuous with Lebesgue-integrable density $f_{X;W}$ such that for all $(x;w) P(0;+1] \square W$, $w = (w_1; :::; w_d)$,

$$f_{X;W}(x;w) \coloneqq c_{X;W}(F(x);F_{W_1}(w_1);\ldots;F_{W_d}(w_d)) f(x) \int_{j=1}^{T^d} f_{W_j}(w_j);$$

for some copula density $c_{X;W} : R^{\Box}_{+} \Box R^{d} - R$ and $f_{W_{1}}; \ldots; f_{W_{d}}$ the marginal densities of the distribution of W.

iii) $c_{X;W}$ bounded in a neighborhood of $t 1u \square [0; 1]^d$.

Then, for

$$= \sum_{p}^{a} := \int_{W}^{a} \frac{1}{p(w)} c_{X;W}(1; F_{W_{1}}(w_{1}); \ldots; F_{W_{d}}(w_{d})) \int_{j=1}^{T^{d}} f_{W_{j}}(w_{j}) dw_{1} \ldots dw_{d}$$

and provided that K - +1 as N - +1 so that $\widehat{KA}(N=K) - \Box$ for some constant $\Box PR$, we have the convergence in distribution as N - +1:

$$\stackrel{?}{\overline{\mathsf{K}}} \stackrel{\square}{\mathsf{H}}{}^{\mathsf{p}}_{\mathsf{K};\mathsf{N}} \square \square \stackrel{\square}{\bullet} \mathsf{O} \ \mathsf{N} \stackrel{\square}{\frac{\square}{1 \square \square}}; \square^2 \square^2_{\mathsf{p}}$$

Proof Recall the decomposition in Equation (6.10):

Combining Lemma 6.17 and Lemma 6.19, provided that Assumption 6.1 and Assumption 6.5 hold and that K = K(N) - +1, K = o(N) and $\overline{K}A(N=K) - \Box$ for some constant $\Box PR$, we have

$$Q_N^{(1)} \square P_{N-1} 0.$$

Lemma 6.17 also ensures that under Assumption 6.1, $Q_N^{(2)}$ is equivalent to

$$\widehat{\mathsf{K}} (\mathsf{H}_{\mathsf{K};\mathsf{N}} \Box \Box) .$$

Referring for instance to De Haan and Ferreira (2006, Theorem 3.2.5), this entails that provided Assumption 6.5 holds and that K = K(N) - +1, K = o(N) and

 $\overrightarrow{KA}(N=K)$ — \Box for some constant \Box PR, we have the convergence in distribution as N - +1:

$$Q_N^{(2)} \dot{O} N \stackrel{\square}{\frac{\square}{\square \square}}; \square^2$$
.

Finally, by virtue of Lemma 6.17 and Lemma 6.21, if Assumption 6.1, Assumption 6.7 and conditions i), ii), iii) in Theorem 6.8 hold together with K = K(N) - +1 and K = o(N), we have the convergence in distribution as N - +1:

$$Q_N^{(3)} \stackrel{\circ}{O} N \stackrel{\Box}{0}; \Box^2 \stackrel{\Box}{\Box_p^2} \Box 1 \stackrel{\Box}{1}.$$

Because the limit distribution of $Q_N^{(3)}$ was established conditionally on the set $t(X_i; W_i); 1 \S i \S Nu$ and in probability relative to this full vector of observations, we can consider $Q_N^{(2)}$ and $Q_N^{(3)}$ as independent random variables (one depends on the data and the other on the survey scheme). The limit distribution of their sum is then the sum of their limit distributions. This concludes the proof.

6.6.3 Limit distribution of $H_{K:N}^{\Box}$ in the rejective survey case

We shall prove Theorem 6.11, the statement of which is recalled below.

Theorem – Limit distribution in the rejective survey case. Suppose that all the conditions required in Theorem 6.8 hold together with Assumption 6.10. Then, for

$$\Box_p^2 \coloneqq \prod_{W=1}^{a} \frac{1}{p(W)} c_{X;W}(1; F_{W_1}(W_1); \ldots; F_{W_d}(W_d)) \prod_{j=1}^{T^d} f_{W_j}(W_j) dW_1 \ldots dW_d$$

and provided that K - +1 as N - +1 so that $\overline{KA}(N=K) - \Box$ for some constant $\Box PR$, we have the convergence in distribution as N - +1:

$$\widehat{\mathsf{K}}(\mathsf{H}_{\mathsf{K};\mathsf{N}}^{\Box} \Box) \stackrel{}{\bullet} \mathsf{N} \stackrel{\Box}{\xrightarrow{\Box}}; \Box^2 \Box_{\mathsf{P}}^2.$$

Proof We shall write $H_{K;N}^{\Box}(R_N)$ when the Horvitz-Thompson version of the Hill estimator involves the inclusion variables $\Box_1; \ldots; \Box_N$ drawn under the sampling plan R_N and the probabilities of inclusion $\Box_1; \ldots; \Box_N$. Consider a Poisson scheme T_N with probabilities $p_1; \ldots; p_N$ and a Rejective scheme R_N with probabilities $\Box_1; \ldots; \Box_N$. Since Theorem 6.8 establishes the asymptotic convergence of \overrightarrow{K} $H_{K;N}^p(T_N) \Box$ to a Normal distribution, we only have to control the quantity \overrightarrow{K} $H_{K;N}^p(R_N) \Box$ $H_{K;N}^p(T_N)$. Using triangular inequalities, we obtain

$$\stackrel{?}{\leftarrow} \stackrel{(A)}{\leftarrow} \stackrel{(A)$$

We shall successively prove that both $Q_N^{(4)}$ and $Q_N^{(5)}$ tend to 0 in probability as N-1 .

Let us start with $Q_N^{(4)}$. We use the bounded Lipschitz metric d_{BL} as defined in van der Vaart and Wellner (1996, p.73) and consider, conditionally on the full vector of observations t(X_i; W_i); 1 § i § Nu,

where $BL_1(R)$ is the set of Lipschitz real functions on R bounded by 1. By virtue of the results in van der Vaart and Wellner (1996, p.73), if this distance tends to 0 as N-1 then we have $Q_N^{(4)} \bigsqcup_{N=1}^{P} 0$. Take b PBL₁(R), then

$$\begin{array}{c} \mathsf{E}_{\mathsf{R}_{\mathsf{N}}} \stackrel{\text{\tiny O}}{\text{\tiny O}} \stackrel{\text{\tiny C}}{\overline{\mathsf{K}}} \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\text{\tiny O}}\left(\mathsf{R}_{\mathsf{N}}\right) \stackrel{\text{\tiny O}}{\text{\tiny O}} \mathsf{E}_{\mathsf{T}_{\mathsf{N}}} \stackrel{\text{\tiny O}}{\text{\tiny O}} \stackrel{\text{\tiny C}}{\overline{\mathsf{K}}} \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\text{\tiny O}}\left(\mathsf{T}_{\mathsf{N}}\right) \stackrel{\text{\tiny O}}{=} \\ \begin{array}{c} \Pi \\ \mathsf{D} \\ \mathsf{D} \\ \mathsf{F}_{\mathsf{K}}^{\text{\tiny O}}\left(\mathsf{R}_{\mathsf{N}}\left(\mathsf{s}\right)\right) \stackrel{\text{\tiny O}}{\text{\tiny R}_{\mathsf{N}}}\left(\mathsf{s}\right) \stackrel{\text{\tiny O}}{\text{\tiny O}} \stackrel{\text{\tiny O}}{\operatorname{\tiny O}} \stackrel{\text{\tiny O}}{\operatorname{\tiny O}} \stackrel{\text{\tiny O}}{\overline{\mathsf{K}}} \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\text{\tiny O}}\left(\mathsf{T}_{\mathsf{N}}\left(\mathsf{s}\right)\right) \stackrel{\text{\tiny O}}{\text{\scriptsize T}_{\mathsf{N}}}\left(\mathsf{s}\right) \\ \mathfrak{s}^{\mathsf{SPP}(\mathsf{U}_{\mathsf{N}})} \stackrel{\text{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}})}{\operatorname{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}})} \stackrel{\text{\tiny O}}{\operatorname{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}})} \stackrel{\text{\tiny O}}{\operatorname{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}}) \\ \mathfrak{s}^{\mathsf{SPP}(\mathsf{U}_{\mathsf{N}})} \stackrel{\text{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}})}{\operatorname{\scriptsize SPP}(\mathsf{U}_{\mathsf{N}})}$$

since b is bounded by 1 and for a fixed sample s,

As established in Berger (1998); Hàjek (1964), the right hand part of this last inequality converges to 0 as N — 1. Therefore, $Q_N^{(4)} \Box_{N-1}^P = 0$.

Let us now turn to $Q_N^{(5)}$. We have:

$$\begin{aligned} \mathsf{Q}_{\mathsf{N}}^{(5)} &= \frac{?}{\mathsf{K}} \stackrel{\mathsf{R}}{\overset{\mathsf{H}}_{\mathsf{K};\mathsf{N}}} (\mathsf{T}_{\mathsf{N}}) \boxdot \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\mathsf{p}} (\mathsf{T}_{\mathsf{N}}) \stackrel{\mathsf{H}}{\overset{\mathsf{P}}_{\mathsf{K};\mathsf{N}}} \frac{?}{\mathsf{K}} \stackrel{\mathsf{r}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}}) + \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\mathsf{p}} (\mathsf{T}_{\mathsf{N}}) \stackrel{?}{\overset{\mathsf{R}}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}})} \\ &= \frac{?}{\mathsf{K}} \stackrel{\mathsf{P}}{\overset{\mathsf{P}}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}}) + \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\mathsf{p}} (\mathsf{T}_{\mathsf{N}}) \stackrel{?}{\overset{\mathsf{R}}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}})} \\ &\frac{?}{\mathsf{K}} \frac{\mathsf{R}}{|\mathsf{D}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}})| + \mathsf{H}_{\mathsf{K};\mathsf{N}}^{\mathsf{p}} (\mathsf{T}_{\mathsf{N}}) \stackrel{?}{\overset{\mathsf{R}}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}})|} \\ &\frac{\mathsf{R}}{\mathsf{r}_{\mathsf{K};\mathsf{N}} (\mathsf{R}_{\mathsf{N}};\mathsf{T}_{\mathsf{N}}) + \mathsf{r}_{\mathsf{K};\mathsf{N}} (\mathsf{T}_{\mathsf{N}};\mathsf{p})} \end{aligned} \end{aligned}$$

We start by analyzing $D_{K;N}(R_N;T_N)$. Observe that it can be written as a function of d_N and \bar{p} , which have nice asymptotic properties:

$$D_{K;N}(R_{N};T_{N}) = \frac{1}{K} \begin{bmatrix} r_{1} \\ i=1 \end{bmatrix} \stackrel{(N \ i+1;N)}{\longrightarrow} \frac{1}{\square_{N \ i+1;N}} \stackrel{(P)}{\longrightarrow} \frac{1}{\square_{N \ i+1;N}} \stackrel{(P)}{\longrightarrow} \frac{1}{\square_{N \ i+1;N}} \log \frac{X_{N \ i+1;N}}{X_{N \ K;N}}$$
$$= \frac{1}{K} \begin{bmatrix} r_{1} \\ i=1 \end{bmatrix} \stackrel{(P)}{\longrightarrow} \frac{1}{\square_{N \ i+1;N}} \frac{P_{N \ i+1;N}}{\square_{N \ i+1;N}} \frac{P_{N \ i+1;N}}{\square_{N \ i+1;N}} \log \frac{X_{N \ i+1;N}}{X_{N \ K;N}}$$
$$= \frac{1}{K} \begin{bmatrix} r_{1} \\ i=1 \end{bmatrix} \stackrel{(N \ i+1;N)}{\square_{N \ i+1;N}} \frac{\bar{P}_{N \ i+1;N}}{\bar{P}_{N \ i+1;N}} + o(1=d_{N})$$

$$\square p_{N \ i + 1;N} \qquad \frac{1}{\square_{N \ i + 1;N}} \square 1 \qquad \frac{1}{\log} \frac{X_{N \ i + 1;N}}{X_{N \ K;N}} \square$$

Set $\Box_{\square} := (1 \Box p_{\square}) \xrightarrow{\Box_{\square}} 1$, then under Assumption 6.1 we have @K § N,

Recall that under Assumption 6.10 we have $d_N = o(1=K)$ and that $H_{K;N}$ is a consistent estimator of \Box (Mason, 1982). Therefore, by virtue of Lemma 6.19, we can conclude that $D_{K;N}(R_N;T_N) = o_P(1=\overline{K})$. Mimicking exactly this procedure and considering the same set of assumptions, we also obtain:

$$|r_{K;N}(R_N;T_N)| \ r_{K;N}(T_N;p) |\Box_0 + o(1)| \frac{1}{d_N};$$

leading to $r_{K;N}(R_N;T_N) = o_P(1=?\overline{K})$ by virtue of Lemma 6.17. Combining these two results with Theorem 6.4 yields

$$\stackrel{?}{\overset{\square}{K}}\stackrel{\square}{\overset{}{H}}_{K;N}^{}(T_{N}) \square \stackrel{H^{p}}{\overset{}{H}}_{K;N}^{}(T_{N})\stackrel{\square}{\overset{}{\overset{}{\Pi}}}\stackrel{P}{\overset{}{\overset{}{\Pi}}} 0;$$

provided that K = K(N) - +1 as N - +1 and that K = o(N). This concludes the proof.

BIBLIOGRAPHY

- Afssa. INCA 2 (2006-2007), Etude Individuelle Nationale des Consommations Alimentaires 2. Report of the Individual and the National Study on Food Consumption., 2009. URL www.afssa.fr/Documents/PASER-Sy-INCA2EN.pdf.
- N. Aires. Algorithms to find exact inclusion probabilities for conditional Poisson sampling and Pareto ps sampling designs. Methodology and Computing in Applied Probability, 1(4):457–469, 1999.
- O. Allais and J. Tressou. Using decomposed household food acquisitions as inputs of a kinetic dietary exposure model. Statistical Modelling, 9(1):27–50, 2009.
- M.I.F. Alves, M.I. Gomes, and L. de Haan. A new class of semi-parametric estimators of the second order parameter. Portugal. Math., 60(2):193–214, 2003. ISSN 0032-5155.
- Anses. Composition nutritionnelle des aliments, Table CIQUAL 2008, 2008. URL http://www.afssa.fr/TableCIQUAL/.
- Anses. Etude de l'alimentation totale française 2 (EAT 2), 2011. URL http://www. anses.fr/cgi-bin/countdocs.cgi?Documents/PASER2006sa0361Ra1.pdf.
- Anses. Rapport d'étude de l'Anses relatif à "Substitution du bisphénol A -L'identification des dangers des substituts potentiels au bisphénol A - Etat des lieux sur les alternatives au bisphénol A" en complément du rapport de l'Anses relatif à l'évaluation des risques liés au bisphénol A (BPA) pour la santé humaine et aux données toxicologiques et d'usage des bisphénols S, F, M, B, AP, AF, et BADGE, 25 March 2013. URL http://www.anses.fr/sites/default/files/ documents/CHIM2009sa0331Ra-3.pdf.
- M.A. Arcones and E. Giné. On the bootstrap of M-estimators and other statistical functionals. Exploring the Limits of Bootstrap, ed. by R. LePage and L. Billard, Wiley, pages 13–47, 1992.
- P. Barbe and P. Bertail. The weighted bootstrap, volume 98. Springer Verlag, 1995.
- S.M. Barlow, J.B. Greig, J.W. Bridges, A. Carere, A.J.M. Carpy, C.L. Galli, J. Kleiner, I. Knudsen, H.B.W.M. Koeter, L.S. Levy, et al. Hazard identification by methods of animal-based toxicology. Food and Chemical Toxicology, 40(2):145–191, 2002.
- C. Béchaux, M. Zetlaoui, J. Tressou, J.C. Leblanc, F. Héraud, and A. Crépet. Identification of pesticide mixtures and connection between combined exposure and diet. Submitted for publication in Food and Chemical Toxicology, 2013.

- J. Beirlant, Y. Goegebeur, J. Segers, and J. Teugels. Statistics of extremes: theory and applications. John Wiley & Sons Inc, 2004. ISBN 0471976474.
- J. Beirlant, T. de Wet, and Y. Goegebeur. A goodness-of-fit statistic for Pareto-type behaviour. J Comput. Appl. Math., 186(1):99–116, 2006. ISSN 0377-0427.
- J. Beirlant, F. Caeiro, and M.I. Gomes. An overview and open research topics in statistics of univariate extremes. Revstat, 10:1–31, 2012.
- J.C. Bemis and R.F. Seegal. Polychlorinated biphenyls and methylmercury act synergistically to reduce rat brain dopamine content in vitro. Environ. Health Persp., 107 (11):879, 1999.
- Y.G. Berger. Rate of convergence to normal distribution for the Horvitz-Thompson estimator. J. Stat. Plan. Inf, 67(2):209–226, 1998.
- Y.G. Berger. Pak. J. Statist. 2011 Vol. 27 (4), 407-426 Asymptotic consistency under large entropy sampling designs with unequal probabilities. Pak. J. Statist, 27(4): 407–426, 2011.
- P. Bertail and J. Tressou. Incomplete generalized U-statistics for food risk assessment. Biometrics, 62(1):66–74, 2006.
- P. Bertail, S. Clémençon, J. Tressou, et al. A storage model with random release rate for modeling exposure to food contaminants. Mathematical biosciences and engineering: MBE, 5(1):35, 2008.
- P. Bertail, S. Clémençon, and J. Tressou. Statistical analysis of a dynamic model for dietary contaminant exposure. Journal of Biological Dynamics, 4(2):212–234, 2010.
- P. Bertail, E. Chautru, and S. Clémençon. Empirical processes in survey sampling. Submitted to the Scandinavian Journal of Statistics, 2013.
- G. Biau and L. Bleakley. Statistical Inference on Graphs. Statistics & Decisions, 24: 209–232, 2006.
- N.H. Bingham, C.M. Goldie, and J.L. Teugels. Regular variation. Encyclopedia of Mathematics and its applications. Cambridge Univ Press, Cambridge, 1987.
- G. Blom. Some properties of incomplete U-statistics. Biometrika, 63(3):573–580, 1976.
- H. Boistard, H.P. Lopuhaä, and A. Ruiz-Gazen. Approximation of rejective sampling inclusion probabilities and application to high order correlations. Electronic Journal of Statistics, 6:1967–1983, 2012.
- M.O. Boldi and A.C. Davison. A mixture model for multivariate extremes. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 69(2):217–229, 2007.

- D. Bonnery. Propriétés asymptotiques de la distribution d'un échantillon dans le cas d'un plan de sondage informatif. PhD thesis, Université Rennes 1, 2011.
- D. Bonnéry, J. Breidt, and F. Coquet. Propriétés asymptotiques de l'échantillon dans le cas d'un plan de sondage informatif. Submitted for publication, 2011.
- P.E. Boon, M. Bonthuis, H. van der Voet, and JD. van Klaveren. Comparison of different exposure assessment methods to estimate the long-term dietary exposure to dioxins and ochratoxin A. Food and Chemical Toxicology, 49(9):1979–1988, 2011.
- S. Boucheron, O. Bousquet, and G. Lugosi. Theory of Classification: A Survey of Some Recent Advances. ESAIM: Probability and Statistics, 9:323–375, 2005.
- G.E.P. Box and D.R. Cox. An analysis of transformations. Journal of the Royal Statistical Society. Series B (Methodological), pages 211–252, 1964.
- N.E. Breslow and J.A. Wellner. Weighted likelihood for semiparametric models and two-phase stratified samples, with application to Cox regression. Scandinavian Journal of Statistics, 35:186–192, 2007.
- N.E. Breslow and J.A. Wellner. A Z-theorem with estimated nuisance parameters and correction note for "Weighted likelihood for semiparametric models and twophase stratified samples, with application to Cox regression". Scandinavian Journal of Statistics, 35:186–192, 2008.
- N.E. Breslow, T. Lumley, C. Ballantyne, L. Chambless, and M. Kulich. Improved Horvitz-Thompson estimation of model parameters from two-phase stratified samples: applications in epidemiology. Stat. Biosc., 1:32–49, 2009.
- B.M. Brown and D.G. Kildea. Reduced U-statistics and the Hodges-Lehmann estimator. The Annals of Statistics, 6:828–835, 1978.
- F. Caeiro and M.I. Gomes. Minimum-variance reduced-bias tail index and high quantile estimation. Revstat, 6(1):1–20, 2008.
- JJ. Cai, L. de Haan, and C. Zhou. Bias correction in extreme value statistics with index around zero. Extremes, pages 1–29, 2011.
- D.O. Carpenter, K. Arcaro, and D.C. Spink. Understanding the human health effects of chemical mixtures. Environ. Health Persp., 110(Suppl 1):25, 2002.
- R.B. Cattell. The scree test for the number of factors. Multivar. Behav. Res., 1(2):245–276, 1966.
- S. Clémençon. On U-processes and clustering performance. In J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F.C.N. Pereira, and K.Q. Weinberger, editors, Advances in Neural Information Processing Systems, pages 37–45, 2011.

- S. Clémençon and N. Vayatis. Overlaying Classifiers: a practical approach to optimal scoring. Constructive Approximation, 32(3):619–648, 2010.
- S. Clémençon, G. Lugosi, and N. Vayatis. Ranking and empirical risk minimization of U-statistics. The Annals of Statistics, 36(2):844–874, 2008.
- W.G. Cochran. Sampling techniques. Wiley, NY, 1977.
- U. Cormann and R.D. Reiss. Generalizing the Pareto to the log-Pareto model and statistical inference. Extremes, 12(1):93–105, 2009.
- E. Counil, P. Verger, and JL. Volatier. Fitness-for-purpose of dietary survey duration: A case-study with the assessment of exposure to ochratoxin A. Food and chemical toxicology, 44(4):499–509, 2006.
- A. Crépet and J. Tressou. Bayesian nonparametric model with clustering individual co-exposure to pesticides found in the French diet. Bayesian analysis, 6(1):127–144, 2011.
- A. Crépet, J. Tressou, V. Graillot, C. Béchaux, S. Pierlot, F. Héraud, and J.C. Leblanc. Accepted for publication in Environmental Research, 2012.
- A. Crépet, F. Héraud, C. Béchaux, M.E. Gouze, S. Pierlot, A. Fastier, JC. Leblanc, L. Le Hégarat, N. Takakura, V. Fessard, et al. The PERICLES research program: an integrated approach to characterize the combined effects of mixtures of pesticide residues to which the French population is exposed. Toxicology, 2013.
- S.P. Crispim, JH.M. de Vries, A. Geelen, O.W. Souverein, P.J.M. Hulshof, L. Lafay, A.S. Rousseau, I.T.L. Lillegaard, L.F. Andersen, I. Huybrechts, et al. Two nonconsecutive 24 h recalls using EPIC-Soft software are sufficiently valid for comparing protein and potassium intake between five European centres-results from the European Food Consumption Validation (EFCOVAL) study. British Journal of Nutrition, 105(03):447–458, 2011.
- J. Danielsson, L. De Haan, L. Peng, and C.G. De Vries. Using a bootstrap method to choose the sample fraction in tail index estimation. J. Multivariate Anal., 76(2): 226–248, 2001. ISSN 0047-259X.
- B. Das and S.I. Resnick. Detecting a conditional extreme value model. Extremes, 14 (1):29–61, 2011.
- B. Das, AbhimanyA. Mitra, and S.I. Resnick. Living on the multidimensional edge: seeking hidden risks using regular variation. Advances in Applied Probability, 45(1): 139–163, 2013.

- E.J. De Boer, N. Slimani, P. van'T Veer, H. Boeing, M. Feinberg, C. Leclercq, E. Trolle, P. Amiano, L.F. Andersen, H. Freisling, et al. Rationale and methods of the European Food Consumption Validation (EFCOVAL) Project. European journal of clinical nutrition, 65:S1–S4, 2011.
- W.J. de Boer, H. van der Voet, B.G.H. Bokkers, M.I. Bakker, and P.E. Boon. Comparison of two models for the estimation of usual intake addressing zero consumption and non-normality. Food Additives and Contaminants, 26(11):1433–1449, 2009.
- L. de Haan. Extremes in higher dimensions: the model and some statistics. In In Proceedings of 45th session international statistics institute, (paper 26.3). The Hague International Statistical Institute Z. Zhang de, pages 317–337, 1985.
- L. De Haan and A. Ferreira. Extreme value theory: an introduction. Springer Verlag, 2006.
- L. de Haan and L. Peng. Comparison of tail index estimators. Statist. Neerlandica, 52: 60–70, 1998.
- L. de Haan and S. Resnick. On asymptotic normality of the Hill estimator. Stochastic Models, 14:849–867, 1998.
- L. de Haan and S. Stadtmüller. Generalized regular variation of second order. J. Austral. Math. Soc. Ser. A, 61:381–295, 1996.
- A.L.M. Dekkers, J.H.J. Einmahl, and L. De Haan. A moment estimator for the index of an extreme-value distribution. Ann. Statist., 17(4):1833–1855, 1989.
- J.C. Deville. Réplications d'échantillons, demi-échantillons, Jackknife, bootstrap dans les sondages. Economica, Ed. Droesbeke, Tassi, Fichet, 1987.
- J.C. Deville and C.E. Särndal. Calibration estimators in survey sampling. JASA, 87: 376–382, 1992.
- L. Devroye, L. Györfi, and G. Lugosi. A Probabilistic Theory of Pattern Recognition. Springer, 1996.
- I.S. Dhillon, Y. Guan, and J. Kogan. Iterative clustering of high dimensional text data augmented by local search. In Data Mining, 2002. ICDM 2002. Proceedings. 2002 IEEE International Conference on, pages 131–138. IEEE, 2002.
- D. Dietrich, L.D. Haan, and J. Hüsler. Testing extreme value conditions. Extremes, 5 (1):71–85, 2002.
- K.W. Dodd, P.M. Guenther, L.S. Freedman, A.F. Subar, V. Kipnis, D. Midthune, J.A. Tooze, S.M. Krebs-Smith, et al. Statistical methods for estimating usual intake of nutrients and foods: a review of the theory. Journal of the American Dietetic Association, 106(10):1640–1650, 2006.

- D.L. Donoho. High-dimensional data analysis: The curses and blessings of dimensionality. AMS Math. Challenges Lecture, pages 1–32, 2000.
- J.J. Droesbeke, B. Fichet, and P. Tassi. Les sondages. Economica, 1987.
- R.M. Dudley. Nonlinear functionals of empirical measures and the bootstrap. In Probability in Banach Spaces 7, pages 63–82. Springer, 1990.
- R.M. Dudley. Uniform Central Limit Theorems. Cambridge University Press, 1999.
- E. Dybing, J. Doe, J. Groten, J. Kleiner, J. O'Brien, A.G. Renwick, J. Schlatter, P. Steinberg, A. Tritscher, R. Walker, et al. Hazard characterisation of chemicals in food and diet: dose response, mechanisms and extrapolation issues. Food and Chemical Toxicology, 40(2):237–282, 2002.
- L. Edler, K. Poirier, M. Dourson, J. Kleiner, B. Mileson, H. Nordmann, A. Renwick,
 W. Slob, K. Walton, and G. Würtzen. Mathematical modelling and quantitative methods. Food and Chemical Toxicology, 40(2):283–326, 2002.
- EFSA. Guidance of the Scientific Committee on a request from EFSA related to Uncertainties in Dietary Exposure Assessment. 14 December 2006. URL http: //www.efsa.europa.eu/en/efsajournal/pub/438.htm.
- JH.J. Einmahl and J. Segers. Maximum empirical likelihood estimation of the spectral measure of an extreme-value distribution. Ann. Statist., 37(5B):2953–2989, 2009.
- J.H.J. Einmahl, L. de Haan, and V.I. Piterbarg. Nonparametric estimation of the spectral measure of an extreme value distribution. Ann. Statist., pages 1401–1423, 2001. ISSN 0090-5364.
- P. Embrechts, C. Klüppelberg, and T. Mikosch. Modelling extremal events: for insurance and finance, volume 33. Springer, 2011.
- E. Enqvist. On sampling from sets of random variables with application to incomplete Ustatistics. PhD thesis, 1978.
- Max Feinberg, Patrice Bertail, Jessica Tressou, Philippe Verger, et al. Analysis of food risks. Editions Tec & Doc, 2006.
- W. Feller. An introduction to probability theory and its applications. Vol. II. . Second edition. John Wiley & Sons Inc., New York, 1971.
- C. Fischer, A. Fredriksson, and P. Eriksson. Neonatal co-exposure to low doses of an ortho-PCB (PCB 153) and methyl mercury exacerbate defective developmental neurobehavior in mice. Toxicology, 244(2-3):157–165, 2008.
- R.A. Fisher and L.H.C. Tippett. Limiting forms of the frequency distribution of the largest or smallest member of a sample. In Mathematical Proceedings of the Cambridge Philosophical Society, volume 24, pages 180–190. Cambridge Univ Press, 1928.

- P.T. Fletcher, C. Lu, S.M. Pizer, and S. Joshi. Principal geodesic analysis for the study of nonlinear statistics of shape. IEEE Trans. Med. Imag., 23(8):995–1005, 2004.
- J. Friedman, T. Hastie, and R. Tibshirani. The elements of statistical learning, volume 1. Springer Series in Statistics, 2001.
- M. Gessaman. A consistent nonparametric multivariate density estimator based on statistically equivalent blocks. Ann. Math. Stat., 41:1344–1346, 1970.
- R.D. Gill. Non- and semiparametric maximum likelihood estimators and the von Mises method. Scand. J. Statistics, 22:205–214, 1989.
- R.D. Gill, Y. Vardi, and J.A. Wellner. Large sample theory of empirical distributions in biased sampling models. The Annals of Statistics, 16(3):1069–1112, 1988.
- E. Giné and J. Zinn. Some limit theorems for empirical processes. The Annals of Probability, 12(4):929–989, 1984.
- B.V. Gnedenko. Sur la distribution limite du terme maximum d'une série aléatoire. Ann. Math., 44(3):423–453, 1943.
- Y. Goegebeur, J. Beirlant, and T. de Wet. Linking Pareto-tail kernel goodness-of-fit statistics with tail index at optimal threshold and second order estimation. Revstat, 6(1):51–69, 2008.
- C.M. Goldie and R.L. Smith. Slow variation with remainder: theory and applications. Quart. J. Math. Oxford, 38(1):45–71, 1987.
- M.I. Gomes and O. Oliveira. The bootstrap methodology in statistics of extremes choice of the optimal sample fraction. Extremes, 4(4):331–358, 2001. ISSN 1386-1999.
- C. Gourieroux. Théorie des sondages. Economica, 1981.
- C. Gourieroux. Effets d'un sondage: cas du □² et de la régression. Economica, Ed. Droesbeke, Tassi, Fichet, 1987.
- S. Guillotte, F. Perron, and J. Segers. Non-parametric Bayesian inference on bivariate extremes. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (3):377–406, 2011.
- J. Hajek. A symptotic theory of rejective sampling with varying probabilities from a finite population. The Annals of Mathematical Statistics, 35(4):1491–1523, 1964.
- F.R. Hampel, E.M. Ronchetti, P.J. Rousseeuw, and W.A. Stahel. Robust statistics: the approach based on influence functions. 1986.
- H.O. Hartley and J.N.K. Rao. Sampling with unequal probabilities and without replacement. Ann. Math. Statist., 33:350–374, 1962.

- S. Haug, C. Klüppelberg, and G. Kuhn. Dimension reduction based on extreme dependence. Submitted for publication, 12, 2010.
- J. Heffernan and S.I. Resnick. Hidden regular variation and the rank transform. Adv. in Appl. Probab., 37(2):393–414, 2005.
- S. Hercberg, K. Castetbon, S. Czernichow, A. Malon, C. Mejean, E. Kesse, M. Touvier, and P. Galan. The Nutrinet-Santé Study: a web-based prospective study on the relationship between nutrition and health and determinants of dietary patterns and nutritional status. BMC public health, 10(1):242, 2010.
- B.M. Hill. A simple general approach to inference about the tail of a distribution. Ann. Statist., 3(5):1163–1174, 1975.
- W. Hoeffding. A class of statistics with asymptotically normal distribution. Ann. Math. Stat., 19:293–325, 1948.
- J. Hoffmann-Jørgensen. Stochastic processes on Polish spaces. Aarhus Universitet, Denmark, 1991.
- D.G. Horvitz and D.J. Thompson. A generalization of sampling without replacement from a finite universe. JASA, 47:663–685, 1951.
- S. Huckemann and H. Ziezold. Principal component analysis for Riemannian manifolds, with an application to triangular shape spaces. Adv. in Appl. Probab., 38(2): 299–319, 2006.
- J. Hüsler and D. Li. On testing extreme value conditions. Extremes, 9(1):69–86, 2006.
- O.A.Y. Jackson. An analysis of departures from the exponential distribution. J. R. Stat. Soc. Ser. B Methodol., 29(3):540–549, 1967.
- S. Janson. The asymptotic distributions of Incomplete U-statistics. Z. Wahrsch. verw. Gebiete, 66:495–505, 1984.
- P. Jenner, A.H.V. Schapira, and C.D. Marsden. New insights into the cause of Parkinson's disease. Neurology, 42(12):2241, 1992.
- S. Jung, M. Foskey, and J.S. Marron. Principal arc analysis on direct product manifolds. Ann. Appl. Statist., 5(1):578–603, 2011.
- S. Jung, I.L. Dryden, and J.S. Marron. Analysis of Principal Nested Spheres. Biometrika, 99(3):551–568, 2012.
- M.C. Kennedy, V.J. Roelofs, C.W. Anderson, and J.D. Salazar. A hierarchical bayesian model for extreme pesticide residues. Food Chem. Toxicol., 49(1):222–232, 2011.

- R. Kroes, D. Muller, J. Lambe, M.R.H. Lowik, J. van Klaveren, J. Kleiner, R. Massey, S. Mayer, I. Urieta, P. Verger, et al. Assessment of intake from the diet. Food and Chemical Toxicology, 40:327–385, 2002.
- D.P. Kroese, T. Taimre, and Z.I. Botev. Handbook of Monte Carlo methods. Wiley, 2011.
- M Ledoux and M. Talagrand. Probability in Banach spaces: isoperimetry and processes. Springer-Verlag, 1991.
- A.J. Lee. U-statistics: Theory and practice. Marcel Dekker, Inc., New York, 1990.
- S. Maddipatla, R. Sreenivasan, and V. Rasbagh. On sums of independent random variables whose distributions belong to the max domain of attraction of max stable laws. Extremes, 14:267–283, 2011. ISSN 1386-1999. URL http://dx.doi.org/10.1007/s10687-010-0109-3. 10.1007/s10687-010-0109-3.
- R. Maitra and I.P. Ramler. A k-mean-directions algorithm for fast clustering of data on the sphere. J. Comput. Graph. Statist., 19(2):377–396, 2010.
- D.M. Mason. Laws of large numbers for sums of extreme values. Ann. Probab., 10: 756–764, 1982.
- P. Massart. Strong approximation for multivariate empirical and related processes, via KMT constructions. Ann. Probab., 17(1):266–291, 1989.
- C. McDiarmid. On the method of bounded differences, pages 148–188. Cambridge Univ. Press, 1989.
- G. McLachlan and D. Peel. Finite mixture models, volume 299. Wiley-Interscience, 2000.
- C. Neves and M.I.F. Alves. Testing extreme value conditions—an overview and recent approaches. Revstat, 6(1):83–100, 2008.
- V. Nichèle et al. L'évolution des achats alimentaires: 30 ans d'enquêtes auprès des ménages en France. Cahiers de Nutrition et de Diététique, 43(3):123–130, 2008.
- M.J. Paulo, H. Van der Voet, J.C. Wood, G.R. Marion, and J.D. Van Klaveren. Analysis of multivariate extreme intakes of food chemicals. Food Chem. Toxicol., 44(7):994– 1005, 2006. ISSN 0278-6915.
- D. Pollard. Convergence of stochastic processes. Springer-Verlag, New-York, 1984.
- W. Polonik. Minimum volume sets and generalized quantile processes. Stochastic Processes and their Applications, 69(1):1–24, 1997.
- O. Pons and E. de Turkheim. Von mises method, bootstrap and Hadamard differentiability for nonparametric general models. Statistics= A Journal of Theoretical and Applied Statistics, 22(2):205–214, 1991.

- R.D. Reiss and M. Thomas. Statistical analysis of extreme values: from insurance, finance, hydrology, and other fields. Boston: Basel, 1997.
- A.G. Renwick, S.M. Barlow, I. Hertz-Picciotto, A.R. Boobis, E. Dybing, L. Edler, G. Eisenbrand, J.B. Greig, J. Kleiner, J. Lambe, et al. Risk characterisation of chemicals in food and diet. Food and Chemical Toxicology, 41(9):1211–1271, 2003.
- S.I. Resnick. Hidden regular variation, second order regular variation and asymptotic independence. Extremes, 5(4):303–336, 2002.
- S.I. Resnick. Heavy-tail phenomena: probabilistic and statistical modeling. Springer Verlag, 2007. ISBN 0387242724.
- S.I. Resnick. Multivariate regular variation on cones: application to extreme values, hidden regular variation and conditioned limit laws. Stochastics, 80(2-3):269–298, 2008.
- C. Rigaux, S. Ancelet, F. Carlin, I. Albert, et al. Inferring an augmented Bayesian network to confront a complex quantitative microbial risk assessment model with durability studies: application to Bacillus Cereus on a courgette purée production chain. Risk Analysis, 2012.
- P.M. Robinson. On the convergence of the Horvitz-Thompson estimator. Australian Journal of Statistics, 24(2):234–238, 1982.
- P. Rosen. A symptotic theory for successive sampling. AMS, 43:373–397, 1972.
- M. Rueda, S. Martínez, H. Martínez, and A. Arcos. Estimation of the distribution function with calibration methods. Journal of statistical planning and inference, 137(2): 435–448, 2007.
- A. Sabourin and P. Naveau. Dirichlet Mixture model for multivariate extremes. 2012.
- T. Saegusa and J.A. Wellner. Weighted likelihood estimation under two-phase sampling. Preprint available at http://arxiv.org/abs/1112.4951v1, 2011.
- C. Scarrott and A. MacDonald. A review of extreme value threshold estimation and uncertainty quantification. REVSTAT–Statistical Journal, 10(1):33–60, 2012.
- C. Scott and R. Nowak. A Neyman-Pearson approach to statistical learning. IEEE Trans. Inf. Theory, 51(8):3806–3819, 2005.
- C. Scott and R. Nowak. Learning Minimum Volume Sets. Journal of Machine Learning Research, 7:665–704, 2006.
- G. Shorack and J.A. Wellner. Empirical processes with applications to statistics. Wiley, 1986.

- V. Sirot, T. Guérin, JL. Volatier, and JC. Leblanc. Dietary exposure and biomarkers of arsenic in consumers of fish and shellfish from France. Science of the Total Environment, 407(6):1875–1885, 2009.
- Steinwart, D. Hush, and C. Scovel. A classification framework for anomaly detection. J. Machine Learning Research, 6:211–232, 2005.
- A. Stephenson. Simulating multivariate extreme value distributions of logistic type. Extremes, 6(1):49–59, 2003.
- J.L. Teugels and G. Vanroelen. Box-Cox transformations and heavy-tailed distributions. Journal of Applied Probability, 41:213–227, 2004.
- Y. Tillé. Utilisation d'informations auxiliaires dans les enquêtes par sondage. Questiió: Quaderns d'Estadística, Sistemes, Informatica i Investigació Operativa, 23(3):491–505, 1999.
- Y. Tillé. Sampling algorithms. Springer Series in Statistics, 2006.
- JA. Tooze, D. Midthune, K.W. Dodd, L.S. Freedman, S.M. Krebs-Smith, A.F. Subar, P.M. Guenther, R.J. Carroll, and V. Kipnis. A new statistical method for estimating the usual intake of episodically consumed foods with application to their distribution. Journal of the American Dietetic Association, 106(10):1575–1587, 2006.
- JA. Tooze, V. Kipnis, D.W. Buckman, R.J. Carroll, L.S. Freedman, P.M. Guenther, S.M. Krebs-Smith, A.F. Subar, and K.W. Dodd. A mixed-effects model approach for estimating the distribution of usual intake of nutrients: The NCI method. Statist. Med., 2010.
- J. Tressou. Nonparametric modeling of the left censorship of analytical data in food risk assessment. J. Amer. Statist. Assoc., 101(476):1377–1386, 2006.
- J Tressou, A. Crépet, P. Bertail, M.H. Feinberg, and J.C. Leblanc. Probabilistic exposure assessment to food chemicals based on extreme value theory. Application to heavy metals from fish and sea products. Food Chem. Toxicol., 42(8):1349–1358, 2004a. ISSN 0278-6915.
- J Tressou, J.C. Leblanc, M.H. Feinberg, and P. Bertail. Statistical methodology to evaluate food exposure to a contaminant and influence of sanitary limits: application to Ochratoxin A. Regul. Toxicol. Pharm., 40(3):252–263, 2004b.
- A. Tsybakov. On nonparametric estimation of density level sets. Annals of Statistics, 25:948–969, 1997.
- Office of Pesticide Programs U.S. Environmental Protection Agency. Draft Guidance for Performing Aggregate Exposure and Risk Assessment. 1 February 1999.

S. van de Geer. Empirical processes in M-estimation. Cambridge University Press, 2000.

- A.W. Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.
- A.W. van der Vaart and J.A. Wellner. Weak convergence and empirical processes. Springer, 1996.
- JD. van Klaveren, P.W. Goedhart, D. Wapperom, and H. van der Voet. A European tool for usual intake distribution estimation in relation to data collection by EFSA. 2012. URL http://www.efsa.europa.eu/fr/supporting/doc/300e.pdf.
- R. Vert and J.P. Vert. Consistency and convergence rates of one-class SVM and related algorithms. J. Machine Learning Research, 17:817–854, 2006.
- JL. Wadsworth, J.A. Tawn, and P. Jonathan. Accounting for choice of measurement scale in extreme value modeling. The Annals of Applied Statistics, 4(3):1558–1578, 2010.
- P. Weihe, P. Grandjean, F. Debes, and R. White. Health implications for Faroe Islanders of heavy metals and PCBs from pilot whales. Sci. Total Environ., 186(1-2): 141–148, 1996.
- WHO. Joint FAO/ WHO Workshop on Methodology for Exposure Assessment of Contaminants and Toxins in Food. 7 8 June 2000. URL http://www.who.int/fsf.

RESUME : Véritable carrefour de problématiques économiques, biologiques, sociologiques, culturelles et sanitaires, l'alimentation suscite de nombreuses polémiques. Dans un contexte où les échanges mondiaux facilitent le transport de denrées alimentaires produites dans des conditions environnementales diverses, où la consommation de masse encourage les stratégies visant à réduire les coûts et maximiser le volume de production (OGM, pesticides, etc.) il devient nécessaire de quantifier les risques sanitaires que de tels procédés engendrent. Notre intérêt se place ici sur l'étude de l'exposition chronique, de l'ordre de l'année, à un ensemble de contaminants dont la nocivité à long terme est d'ores et déjà établie. Les dangers et bénéfices de l'alimentation ne se restreignant pas à l'ingestion ou non de substances toxiques, nous ajoutons à nos objectifs l'étude de certains apports nutritionnels. Nos travaux se centrent ainsi autour de trois axes principaux. Dans un premier temps, nous nous intéressons à l'analyse statistique des très fortes expositions chroniques à une ou plusieurs substances chimiques, en nous basant principalement sur des résultats issus de la théorie des valeurs extrêmes. Nous adaptons ensuite des méthodes d'apprentissage statistique de type ensembles de volume minimum pour l'identification de paniers de consommation réalisant un compromis entre risque toxicologique et bénéfice nutritionnel. Enfin, nous étudions les propriétés asymptotiques d'un certain nombre d'estimateurs permettant d'évaluer les caractéristiques de l'exposition, qui prennent en compte le plan de sondage utilisé pour collecter les données.

MOTS-CLEFS : Analyse des risques alimentaires - Apports nutritionnels de long terme - Théorie des valeurs extrêmes - Mesure spectrale - Théorie des sondages - Processus empiriques - Estimation de l'indice de valeurs extrêmes - Ensembles de volume minimum - U-statistiques - Risque-bénéfice





