

THÈSE DE DOCTORAT

de

L'UNIVERSITÉ PARIS-SACLAY

et de

L'UNIVERSITÉ GASTON BERGER

École doctorale de Mathématiques Hadamard (EDMH, ED 574) et École doctorale de
Sciences et technologies (EDST, UGB)

Établissements d'inscription : Université Paris-Sud et Université Gaston Berger

Laboratoire d'accueil : Laboratoire de Mathématiques d'Orsay, UMR 8626 CNRS

Spécialité de doctorat : Mathématiques Appliquées

Mor Absa LOUM

Modèle de mélange et modèles linéaires généralisés,
application aux données de co-infection
(arbovirus & paludisme)

Date de soutenance : 28 Août 2018

Après avis des rapporteurs : M. JEAN-MARC BARDET (Université de Paris 1 Panthéon-Sorbonne)
M. STÉPHANE GIRARD (Inria Grenoble Rhône-Alpes)

Jury de soutenance :

M. JEAN-MARC BARDET	(Université Paris 1 Panthéon-Sorbonne)	Rapporteur
M. ABDOU KÂ DIONGUE	(Université Gaston Berger)	Examineur
M. ALIOU DIOP	(Université Gaston Berger)	Co-directeur de thèse
MME CÉCILE DUROT	(Université Paris Nanterre)	Présidente
MME ELISABETH GASSIAT	(Université Paris Sud)	Directrice de thèse
M. CHRISTOPHE GIRAUD	(Université Paris Sud)	Examineur
M. CHEIKH LOUCOUBAR	(Institut Pasteur de Dakar)	Examineur

À El hadji Abdoulaye LOUM.

Remerciements

Il me sera très difficile de remercier tout le monde car c'est grâce à l'aide de nombreuses personnes que j'ai pu mener cette thèse à son terme.

En premier lieu, je tiens à exprimer ma profonde gratitude à l'égard de ma directrice de thèse Elisabeth Gassiat pour la confiance qu'elle m'a accordé en acceptant d'encadrer ma thèse et de l'avoir si bien fait. J'aimerais également lui dire à quel point j'ai apprécié ses compétences, sa rigueur et sa clairvoyance. J'ai beaucoup apprécié sa gentillesse, sa disponibilité et sa patience. Dans les périodes difficiles, elle m'a toujours offert du temps et de précieux conseils pour m'aider à avancer. Merci notamment pour l'exigence.

Je souhaite également remercier profondément mon co-directeur de thèse Aliou Diop qui malgré la distance a toujours été présent. Par ses conseils avisés, ses multiples relectures et remarques pertinentes, il m'a offert le meilleur encadrement possible afin de mener mes travaux de recherche. J'ai été extrêmement sensible à ses qualités humaines d'écoute, de disponibilité et de bienveillance.

Merci à tous les deux de m'avoir tant apporté aussi bien humainement que scientifiquement. Merci pour nos discussions si enrichissantes et d'avoir toujours été disponibles pour moi. J'ai pris un très grand plaisir à travailler avec vous. J'espère avoir été digne de la confiance que vous m'avez accordé et que nous continueront à travailler ensemble. J'ai beaucoup appris à vos côtés et je suis très honoré de vous avoir eu pour encadrants.

Je remercie Jean Marc Bardet et Stéphane Girard d'avoir accepté la tâche de rapporteur et pour leur relecture très attentive malgré un emploi du temps chargé. Je remercie aussi Cécile Durot, Christophe Giraud et Abdou Kâ Diongue de m'avoir fait l'honneur de participer au jury de ma thèse. Je leur suis sincèrement reconnaissant de s'être rendus disponibles pour cette soutenance. Je remercie également Cheikh Loucoubar de la collaboration depuis mon stage de master, mais aussi d'avoir accepté de participer au jury. Merci Cheikh pour les conseils et encouragements.

Je remercie très chaleureusement Marie-Anne Poursat pour son soutien, ses encouragements et sa collaboration scientifique. Merci Marie-Anne pour ton aide tout au long de ma thèse, en particulier dans la prise en main des données de co-infection. J'ai pu apprendre beaucoup de choses à tes côtés. Nos discussions autour d'un café ont toujours été enrichissantes.

Je souhaite remercier et également faire part de ma reconnaissance à Gilles Celeux et Benjamin Auder, avec qui j'ai eu à collaborer. Merci Gilles pour ton aide et tes conseils. Merci également à benjamin pour cette collaboration fructueuse. J'ai pu apprendre beaucoup de choses, sur le logiciel R, à tes côtés.

Je souhaite remercier et également faire part de ma reconnaissance à Amadou Alpha Sall pour m'avoir accueilli, depuis le stage de master, au sein de son unité à l'Institut Pasteur de Dakar. Merci également pour les données. Je remercie aussi tous les membres

de l'équipe de virologie et de G4 pour l'accueil et la bonne ambiance, en particulier Mamadou Diop et Fatou Sow Diop.

Je remercie très chaleureusement Damien Simon qui m'a encadrer en master (AIMS-Sénégal) et qui reste toujours à mes côtés pour m'encourager et me donner de précieux conseils. Merci mon professeur pour vos encouragements et les déjeuners à l'UPMC.

Je tiens à remercier mes collègues de bureau, François Delgove, Pierre Roux, Antony Poels, pour la bonne ambiance de travail qu'ils ont apportée. Je remercie également Jacques De Catelan, Martin Royer et Jeanne pour leurs pertinentes.

Je témoigne de mes remerciements à l'ensemble des membres du Laboratoire de Mathématiques d'Orsay, en particulier Abdelhadi Tout pour son aide précieux et sa sympathie, toute l'équipe informatique, l'équipe administrative, ceux qui sont encore présents comme ceux qui sont déjà partis, pour leur très grande sympathie. Je voudrais décerner une mention spéciale à l'équipe de foot du labo, emmenée par son capitaine emblématique (Rachid). Merci à tous les joueurs et à toutes les joueuses de l'équipe pour les très grands matchs, la sympathie et les "teufs" chez les chasseurs.

Je remercie les membres du Laboratoire d'Études et Recherches en Statistiques et Développement (LERSTAD) et de école doctorale Sciences et Techniques de l'université Gaston Berger, plus particulièrement Aminata Wagué, Mbaye Faye, Aissatou Conté, et tous autres les doctorants. Merci également à El hadji Dème pour les conseils et encouragements.

Je remercie mes amis Mamadou Dione, Cheikh Mbaye, Moustapha Thiam, Modou Bèye, Sidy Ndime, Sidy Bouya, Babacar Sakho, Moussa Thiam, Monsieur et Madame Mboup. Je remercie également mes amis de Massy Lamine Sene, Baye Cheikh, Modou Faye, Ibnou Seck, Dame Seck, Fodé, Cheikhouna, Moussou et Seynabou. Les longs débats autour d'un bon plat sénégalais, discutant de la politique et de la vie en général m'ont beaucoup plu. Je remercie aussi tous les membres du centre islamique de Taverny, en particulier le "dieuwrigne" Souleymane Diouf.

Enfin, je ne saurais terminer cette partie sans exprimer ma gratitude aux membres de ma famille, à ma femme Seynabou Diop pour son amour et son soutien inconditionnel, à mes très chers parents Ibrahima Loum et Ndeye Fatou Diop, à mes oncles Cheikh Ibnou Diop et Abdou Khadir Diop, à mes frères El Hadji Abdoulaye Thiam et Cheikh Ibrahima Biteye, à mes cousins Lat Déguène Loum et Abdoulaye Diop, qui m'ont toujours soutenu, encouragé et stimulé pendant mes études. Je n'aurais jamais pu arriver ici, sans leur soutien et encouragements. Merci !

Table des matières

Remerciements	v
1 Introduction	1
1.1 Contexte	2
1.1.1 Problématique	2
1.1.2 Objectifs	3
1.2 Modèles linéaires généralisés	3
1.2.1 Régression logistique binaire	4
1.2.2 Régression logistique multinomiale	5
1.3 Modèle de mélange fini	5
1.3.1 Mélange de modèles linéaires généralisés	6
1.3.2 Mélange de régression logistique	6
1.3.3 Estimation des paramètres	6
1.4 Méthode des moments	6
1.4.1 Méthode des moments classique	6
1.4.2 Méthode tensorielle	7
1.5 Modèle de mélange et théorie des valeurs extrêmes	8
1.5.1 Rappels sur la théorie des valeurs extrêmes	8
1.5.2 Notion de censure	10
1.5.3 Mélange d'extrêmes en présence de censure	10
1.6 Présentation générale des résultats	10
1.6.1 Chapitre 2 Modèle de régression logistique multinomiale pour le diagnostic de la co-infection	11
1.6.2 Chapitre 3 Mélange de modèles linéaires généralisés et méthode des moments : identifiabilité & applications	12
1.6.3 Chapitre 4 Extensions des mélanges de modèles linéaires généralisés	13
1.6.4 Chapitre 5 Mélange de valeurs extrêmes en présence de censure	14
2 Modèle de régression logistique multinomiale pour le diagnostic et la recommandation de traitement en cas de co-infection entre deux maladies (<i>Arbovirus</i> & <i>Paludisme à Kédougou</i>)	17
2.1 Introduction	20
2.2 Description des données	21
2.2.1 Réponse multinomiale et jeux de données	22
2.2.2 Covariables	24
2.3 Modèle Statistique	26

2.3.1	La régression logistique multinomiale	27
2.3.2	La régression logistique multinomiale dans le cas Arbovirus-Paludisme	30
2.3.3	Sélection de variables	30
2.3.4	Analyse des facteurs influents des différentes maladies	42
2.4	Analyse prédictive	47
2.4.1	Test d'indépendance entre arbovirus et paludisme	47
2.4.2	Diagnostic de la co-infection	51
2.5	Discussion	54
3	Mélange de modèles linéaires généralisés et méthode des moments : identifiabilité & applications	55
3.1	Introduction	58
3.2	Notation & modèle	59
3.2.1	Notations et définitions	59
3.2.2	Modèles	60
3.3	Algorithme	61
3.3.1	Estimation des directions	61
3.3.2	Estimation de tous les paramètres du modèle	68
3.4	Résultats théoriques	69
3.4.1	Identifiabilité	69
3.4.2	Consistance	73
3.4.3	Normalité Asymptotique	75
3.5	Applications	83
3.5.1	Package R	83
3.5.2	Simulations	84
3.5.3	Sélection de variables	94
4	Extensions des mélanges de modèles linéaires généralisés	99
4.1	Introduction	100
4.2	Mélanges de modèles linéaires généralisés	100
4.2.1	Covariables continues	100
4.2.2	Covariables continues et catégorielles	102
4.3	Données longitudinales	103
5	Mélange de valeurs extrêmes en présence de censure	105
5.1	Introduction	108
5.2	Modèle de mélange et valeurs extrêmes	110
5.3	Estimation des paramètres	111
5.3.1	Vraisemblance en dessous du seuil u	111
5.3.2	Vraisemblance au dessus du seuil u	112
5.3.3	Estimation	113
5.4	Estimation des quantiles extrêmes	115
5.4.1	Par la fonction de répartition du modèle extrême	115
5.4.2	Par la méthode de reparamétrisation	116
5.5	Étude de simulation	116
5.5.1	Pour un seuil u fixé	117

5.5.2	Pour un seuil u inconnu	119
5.5.3	Conclusion	121
5.6	Discussion & conclusion	122
6	Conclusion et perspectives	127
A	Multinomial logistic model for coinfection diagnosis between arbovirus and malaria in Kedougou	129
A.1	Introduction	129
A.2	Data description	131
A.2.1	Data set	131
A.2.2	Covariates	132
A.3	Statistical analysis of the coinfection influential factors	134
A.3.1	Multinomial logit model	135
A.3.2	Variable selection using random forests	136
A.3.3	Influence of selected covariates on disease status	138
A.4	Predictive analysis	143
A.4.1	Testing independence between arbovirus and malaria	143
A.4.2	Diagnosis of arboviral disease	144
A.5	Discussion	147
B	Mixture of generalized linear models : identifiability and applications	149
B.1	Introduction	149
B.2	Model and notations	150
B.3	Moment identifiability and estimation	152
B.3.1	Identifiability results	152
B.3.2	The least squares moment estimator	154
B.3.3	Algorithm	154
B.4	Simulations	155
B.4.1	R package	155
B.4.2	Experiments	156
B.5	Some other identifiability results	160
B.5.1	Continuous covariates	161
B.5.2	Continuous and categorical covariates	163
B.5.3	Longitudinal observations	164
B.5.4	Some perspectives	165
B.6	Proofs	165
B.6.1	Proof of Theorem B.3.1.1	165
B.6.2	Proof of Theorem B.3.1.2	165
B.6.3	Proof of Theorem B.3.2.1	167
C	Morpheus-package : Estimate Parameters of Mixtures of Logistic Regressions	173
C.1	R topics documented :	174
C.1.1	morpheus-package	174
C.1.2	alignMatrices	175
C.1.3	computeMoments	175

C.1.4	computeMu	176
C.1.5	generateSampleIO	177
C.1.6	multiRun	177
C.1.7	normalize	179
C.1.8	optimParams	180
C.1.9	plotBox	181
C.1.10	plotCoefs	181
C.1.11	plotHist	182
C.1.12	plotQn	182

Bibliographie	183
----------------------	------------

Chapitre 1

Introduction

Sommaire

1.1	Contexte	2
1.1.1	Problématique	2
1.1.2	Objectifs	3
1.2	Modèles linéaires généralisés	3
1.2.1	Régression logistique binaire	4
1.2.2	Régression logistique multinomiale	5
1.3	Modèle de mélange fini	5
1.3.1	Mélange de modèles linéaires généralisés	6
1.3.2	Mélange de régression logistique	6
1.3.3	Estimation des paramètres	6
1.4	Méthode des moments	6
1.4.1	Méthode des moments classique	6
1.4.2	Méthode tensorielle	7
1.5	Modèle de mélange et théorie des valeurs extrêmes	8
1.5.1	Rappels sur la théorie des valeurs extrêmes	8
1.5.2	Notion de censure	10
1.5.3	Mélange d'extrêmes en présence de censure	10
1.6	Présentation générale des résultats	10
1.6.1	Chapitre 2 Modèle de régression logistique multinomiale pour le diagnostic de la co-infection	11
1.6.2	Chapitre 3 Mélange de modèles linéaires généralisés et méthode des moments : identifiabilité & applications	12
1.6.3	Chapitre 4 Extensions des mélanges de modèles linéaires généralisés	13
1.6.4	Chapitre 5 Mélange de valeurs extrêmes en présence de censure	14

1.1 Contexte

1.1.1 Problématique

En Afrique, en Amérique du sud et en Asie (en zone tropical généralement), plusieurs pays continuent de souffrir de la mortalité due au paludisme. La forte ressemblance clinique entre le paludisme et d'autres maladies infectieuses pose un réel problème de diagnostic clinique ([1, 2, 3, 4, 5, 6, 7, 8, 9]). En effet, il existe une forte ressemblance entre les symptômes du paludisme et d'autres maladies infectieuses ([10]) telles que les arboviroses (maladies causées par des arbovirus), les infections bactériennes, . . . De plus, le paludisme et certaines arboviroses sont pratiquement endémiques sur les mêmes zones (par exemple la Dengue et le Paludisme, Figure 1.1).

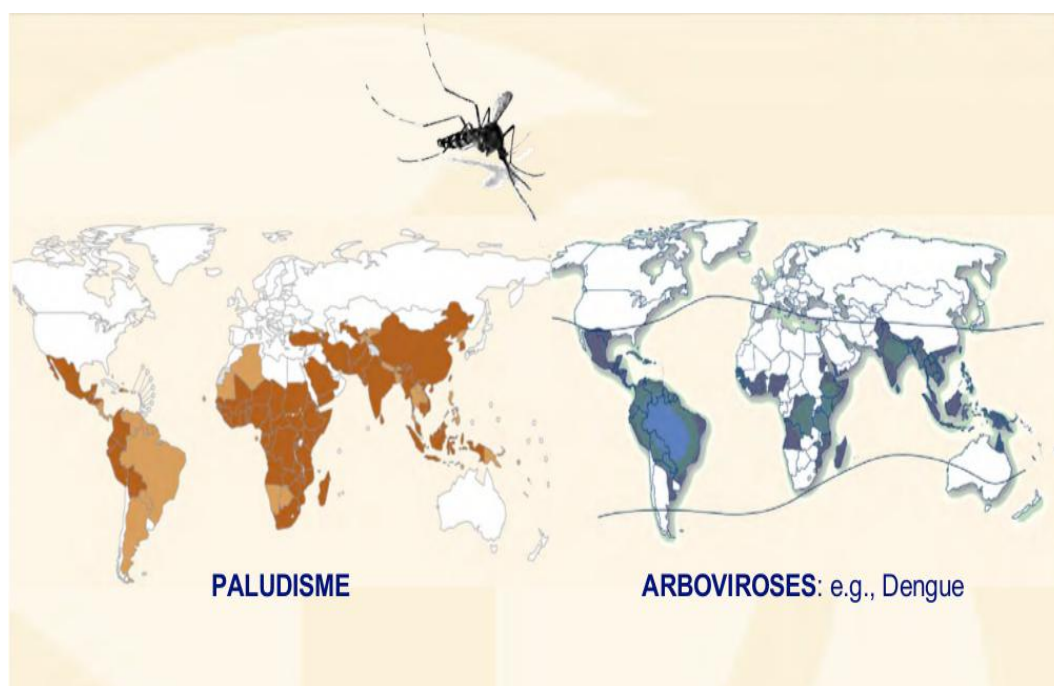


FIGURE 1.1 – Pays ou zones à risque de transmission du Paludisme et de la Dengue. Source : adaptée de l’OMS, 2013.

Au Sénégal, depuis l’introduction du test de diagnostic rapide (TDR) en 2007, une baisse importante de la prévalence du paludisme a été constatée. Ce qui laisse penser que le paludisme a été surestimé dans le passé au détriment d’autres maladies ([11, 12]) (exemple : arboviroses, infections bactériennes, . . .). Dans la région de Kédougou (Sud-Est du Sénégal, voir Figure 1.2), le paludisme et les arboviroses sont toutes deux endémiques à cause du climat et de la position géographique. La co-circulation des parasites du paludisme et des arbovirus peut expliquer l’observation des cas de co-infection (des individus positifs à un ou plusieurs arbovirus et au paludisme) signalés dans cette région ([13]). Ainsi, la co-infection a besoin d’être mieux diagnostiquée dans ces zones. De plus, en cas de co-infection, il serait nécessaire de savoir lequel des virus est à l’origine de la maladie.

Dans ce contexte, deux axes ont été définis dans cette thèse : une partie pratique, pour l’étude et la prise en main des données ; une partie méthodologique, pour la mise



FIGURE 1.2 – Région de Kédougou.

en place de méthodes pouvant aider au diagnostic et au traitement de la co-infection.

1.1.2 Objectifs

Dans le but de fournir un outil de diagnostic à l'Institut Pasteur de Dakar (IPD), nous utilisons les données collectées à Kédougou entre 2009 et 2013. En effet, depuis 2009, l'IPD a mis en place à Kédougou une surveillance des maladies fébriles aiguës appelée AFI (acute febrile illnesses). L'objectif était de pouvoir détecter de manière précoce les épidémies des arboviroses et du paludisme afin de réduire la mortalité et la morbidité due à ces maladies dans cette région. Ainsi, à partir des données, deux objectifs ont été définis : (1) fournir un outil de diagnostic de la co-infection qui, en cas d'absence de test des arbovirus, pourra orienter les décideurs sur la co-infection ou non du patient ; (2) fournir des recommandations de traitement en cas de co-infection, en se basant sur les symptômes présentés par le patient.

Dans l'étude pratique, nous avons utilisé un modèle multinomial classique bien connu dans la littérature. Suite aux résultats, nous avons pensé à les améliorer en utilisant d'autres modèles beaucoup plus adaptés. Par exemple, poser le problème comme un mélange de modèles. Dans ce cas, on n'utilisera pas la méthode de l'algorithme EM habituelle (Expectation Maximisation). L'étude se fera par une méthode des moments basée sur une méthode spectrale.

1.2 Modèles linéaires généralisés

Un modèle de régression linéaire est un modèle de régression qui cherche à établir une relation linéaire entre une variable réponse et une ou plusieurs covariables. Le modèle

linéaire généralisé, communément appelé “Generalized Linear Model” (GLM) en anglais, est une généralisation naturelle du modèle de régression linéaire. Cette généralisation permet de relier la variable réponse à une ou plusieurs covariables via une fonction lien. Le modèle linéaire généralisé est initialement développé en 1972 par Nelder et Wedderburn ([14]). Un exposé détaillé de ces travaux est présenté dans les ouvrages de Nelder et McCullagh (1989, [15]), Agresti (2013, [16]) ou Antoniadis et al. (1992, [17]). Le modèle linéaire généralisé peut être vu sous forme de trois composantes :

1. Distribution et densité de la variable à expliquer

Soit Y_1, \dots, Y_n une suite de variables aléatoires indépendantes et identiquement distribuées de loi appartenant à la famille exponentielle. C’est-à-dire que la densité par rapport à la mesure de Lebesgue ou à une mesure de comptage, s’écrit sous la forme

$$f(y_i; \theta_i, \phi) = \exp\left(\frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i; \phi)\right) \quad (1.1)$$

où θ_i est le paramètre de position et ϕ le paramètre de dispersion. Comme expliqué par McCullagh et Nelder ([15]), l’espérance et la variance de Y sont données par $\mathbb{E}(Y) = b'(\theta)$ et $Var(Y) = b''(\theta)a(\theta)$ respectivement.

2. Prédicteur linéaire

Le prédicteur linéaire noté η est la composante déterministe du modèle. Il est donné par

$$\eta = X\beta \in \mathbb{R}^n \quad (1.2)$$

avec $X \in \mathbb{R}^{n \times d}$ la matrice des vecteurs de covariables et β le vecteur de paramètre de taille d .

3. La fonction lien

La fonction lien permet d’exprimer Y_i en fonction du prédicteur linéaire. Notons par μ_i , l’espérance de Y_i sachant X_i ,

$$\mu_i = \mathbb{E}[Y_i | X_i], \quad i = 1 \dots, n.$$

On a la fonction lien g est telle que

$$g(\eta_i) = \mu_i, \quad i = 1, \dots, n \quad \text{où } \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \\ \dots \\ \eta_n \end{pmatrix}.$$

On suppose que la fonction lien g est monotone et différentiable. Si la réponse est discrète, on peut parler de régression logistique (binaire ou multinomiale) ou régression probit.

1.2.1 Régression logistique binaire

La régression logistique binaire est un cas particulier du modèle linéaire généralisé où la variable réponse ne prend que deux entrées (exemple 0 et 1). On dit dans ce cas que le lien est logit. C’est-à-dire que la fonction lien g est telle que

$$g(x) = \frac{e^x}{(1 + e^x)}.$$

Notons par $\pi(x)$, la probabilité que $Y = 1$ sachant $X = x$, ie

$$\pi = \mathbb{P}(Y = 1|X).$$

On a alors $\mathbb{E}(Y|X) = \pi$ et la densité est donnée par

$$f(y|x) = \pi^y(1 - \pi)^{1-y}. \quad (1.3)$$

On a dans ce cas

$$\mathbb{E}(Y|X = x) = \frac{e^{\langle x, \beta \rangle}}{1 + e^{\langle x, \beta \rangle}}$$

et le modèle dite de régression logistique binaire s'écrit :

$$\log\left(\frac{\pi}{1 - \pi}\right) = \langle x, \beta \rangle \quad (1.4)$$

où x est le vecteur de covariable. Pour plus de détails sur la régression logistique binaire, le lecteur pourra consulter le livre de Hosmer et Lemeshow ([18]).

1.2.2 Régression logistique multinomiale

La régression logistique multinomiale est une généralisation de la régression logistique binaire. En effet, la variable réponse peut prendre K modalités ($k = 0, \dots, K - 1$) avec $K > 2$. Cette partie est détaillée à la section 2.3.1. Pour plus de détails, le lecteur pourra consulter [18] ou [15].

1.3 Modèle de mélange fini

Le modèle de mélange fini est un modèle statistique utilisé pour prendre en compte l'hétérogénéité d'une population. Il est utilisé depuis plus d'un siècle (Newcomb (1886, [19]), Pearson (1894, [20])). Le modèle de mélange fini consiste à supposer que les données proviennent d'une source contenant plusieurs sous-populations. Chaque population est modélisée de manière séparée. La population totale est un mélange de ces différentes sous-populations. Soit Y une variable aléatoire de densité $f(\cdot)$. Si la loi de Y est un mélange de K sous-populations, alors sa densité s'écrit

$$f(y|\Theta) = \sum_{k=1}^K p_k f_k(y|\Theta_k) \quad (1.5)$$

où p_k représente la probabilité a priori de la composante k . Les probabilités p_k , $k = 1, \dots, K$ vérifient $\sum_{k=1}^K p_k = 1$ et $0 \leq p_k \leq 1$. Θ désigne l'ensemble des paramètres du modèle et Θ_k les paramètres du sous-modèle k .

Pour plus de détails, sur le modèle de mélange fini, le lecteur pourra consulter le livre de Sylvia Frühwirth-Schnatter ([21]).

Dans cette thèse, nous accordons un intérêt particulier à la classe des mélanges de modèles linéaires généralisés.

1.3.1 Mélange de modèles linéaires généralisés

On parle de mélange de modèles linéaires généralisés si les composantes du mélange sont des modèles linéaires généralisés. C'est-à-dire, $\forall k = 1, \dots, K$, $f_k(y|\Theta_k)$ de l'équation (1.5) est sous la forme de l'équation (1.1).

Comme présenté par Grün ([22]), les mélanges de modèles linéaires sont utilisés dans beaucoup de domaines. On note par exemple des applications en biologie ou en médecine (Aitkin (1999, [23]), Follmann et Lambert (1989, [24]), Wang et al. (1996, [25]), Wang et Puterman (1998, [26])). Un exemple d'application biologique est illustré à travers les données "Alphis" dans les travaux de Boiteau et al. (1998, [27]). Pour plus de détails sur les mélanges de modèles linéaires généralisés, le lecteur pourra consulter les travaux de Bettina Grün ([22]).

1.3.2 Mélange de régression logistique

Dans la classe des mélanges de modèles linéaires généralisés, nous nous intéressons à la partie concernant les observations binaires. C'est-à-dire que, si Y est la variable réponse et x le vecteur covariables,

$$\mathbb{P}(Y = 1|X = x) = \sum_{k=1}^K p_k g(\langle \beta_k, x \rangle + b_k) \quad (1.6)$$

où g est la fonction lien, et pour tout $k = 1, \dots, K$, p_k est la proportion de la sous-population k , β_k est le vecteur de paramètre de régression et b_k l'intercept.

1.3.3 Estimation des paramètres

Dans le passé, le modèle de mélange fini était étudié à l'aide des méthodes basées sur la maximisation de la vraisemblance du modèle (algorithme EM, Dempster et al. [28]) ou des méthodes variationnelles de Bayes (Bishop et Svensen [29], [21]). Ces méthodes ont eu beaucoup de succès dans différents domaines ([21], Jordan et al. [30]; Xu et al. [31]; Grün [32]). Mais elles peuvent converger vers des optimums locaux et peuvent présenter des vitesses de convergence faibles en grande dimension. De plus, elles peuvent présenter des temps de calcul assez longs. Dans cette thèse, l'idée est d'utiliser une méthode des moments pour estimer les paramètres du mélange.

1.4 Méthode des moments

1.4.1 Méthode des moments classique

La méthode des moments est très ancienne et remonte à Pearson en 1894 ([20]). L'idée de base de la méthode des moments est d'estimer (ou approcher) un moment théorique par un moment empirique. Par exemple estimer une espérance mathématique par une moyenne empirique, une variance par une variance empirique, ...

Soit X une variable aléatoire de loi \mathbb{P}_θ ($\theta \in \Theta$). Soit $\psi(\theta)$ le vecteur des s premiers moments donné par

$$\psi(\theta) = [\mathbb{E}_\theta(X), \dots, \mathbb{E}_\theta(X^s)] = [\alpha_1(\theta), \dots, \alpha_s(\theta)]. \quad (1.7)$$

Soient $[\hat{\alpha}_k]_{1 \leq k \leq s}$, les s premiers moments empiriques. On estime θ par une solution du système d'équations suivant :

$$\begin{cases} \hat{\alpha}_1 = \alpha_1(\theta) \\ \dots \\ \dots \\ \hat{\alpha}_s = \alpha_s(\theta) \end{cases} \quad (1.8)$$

1.4.2 Méthode tensorielle

La méthode tensorielle (tensor method [33]) est une méthode statistique d'estimation basée sur la définition et la décomposition des tenseurs en fonction des paramètres du modèle. Ces dernières années, son utilisation dans le cadre de l'estimation des paramètres du mélange a eu beaucoup de succès ([34], [35],[36, 37], [38]). Cette estimation est basée sur l'écriture des moments sous forme de tenseur symétrique. Soit (x, y) un couple de variable aléatoire de loi $\mathbb{P}_\theta = \mathcal{L}(x, y)$, avec y une variable réponse ($y \in \{0, 1\}$) et $x \in \mathbb{R}^d$ un vecteur de covariables. Supposons que conditionnellement à x , y provient d'un mélange de K populations modélisé par

$$\mathbb{P}(Y = 1|X = x) = \sum_{k=1}^K \omega_k g(\langle \beta_k, x \rangle + b_k) \quad (1.9)$$

avec $\theta = (\omega, \beta, b) \in \Theta$, $\beta = [|\beta_1|, \dots, |\beta_K|] \in \mathbb{R}^{d \times K}$ la matrice des K vecteurs de paramètres, $b = (b_1, \dots, b_K)$ le vecteur des intercepts et g la fonction lien. On note $\omega = (\omega_1, \dots, \omega_K)$ le vecteur des poids du mélange, sous les contraintes $\sum_{k=1}^K \omega_k = 1$ et $\omega_k > 0$

On peut définir les moments croisés d'ordre j , $j \geq 1$ entre y et x par

$$M_j = \mathbb{E} [y \cdot x^{\otimes j}] \in \mathbb{R}^{d^j}, \quad (1.10)$$

avec $x^{\otimes j}$ le produit tensoriel de x d'ordre j . Par exemple, si $j = 2$, $x^{\otimes 2} = x \otimes x \in \mathbb{R}^{d \times d}$ est une matrice dont les coordonnées (i, j) sont données par les $x_i x_j$.

L'idée ici est d'écrire les moments M_j sous la forme

$$M_j = \sum_{k=1}^K \lambda_k \beta_k^{\otimes j} \quad j \geq 1 \quad (1.11)$$

avec $\lambda_k \in \mathbb{R}$ et β_k les vecteurs de paramètres. Pour plus de détails, le lecteur pourra consulter les travaux de Anima Anandkumar et al. ([36]).

Une manière d'estimer les paramètres du mélange en utilisant la méthodes des tenseurs est d'utiliser le moment M_3 d'ordre 3

$$M_3 = \mathbb{E} [y \cdot x^{\otimes 3}]$$

pour construire un ensemble de matrice à diagonaliser simultanément (diagonalisation jointe [39], [40]). Cette diagonalisation jointe permet d'estimer les vecteurs de paramètres normalisés. Pour retrouver entièrement les paramètres du modèle, nous proposons dans cette thèse une estimation en deux étapes : (1) une première étape de diagonalisation jointe pour estimer les vecteurs de paramètres normalisés et (2) une étape d'estimation où on minimise un critère du type moindres carrés.

1.5 Modèle de mélange et théorie des valeurs extrêmes

Dans cette partie, nous ferons quelques rappels sur la théorie des valeurs extrêmes et sur la censure. On évoquera aussi les mélanges d'extrêmes en présence de censure.

1.5.1 Rappels sur la théorie des valeurs extrêmes

La théorie des valeurs extrêmes est une vaste théorie dont le but est d'étudier les événements rares. C'est-à-dire les événements dont la probabilité d'apparition est faible. Autrement dit, on cherche à calculer les probabilités des événements qui ont peu de chances de se réaliser. Pour cela, considérons X_1, \dots, X_n une suite de variables aléatoires indépendantes et identiquement distribuées de même loi que X . Notons par F la fonction de répartition de X . Soit $M_n = \max(X_1, \dots, X_n)$. La loi du maximum M_n peut être donnée par

$$\mathbb{P}(M_n \leq x) = \prod_{i=1}^n \mathbb{P}(X_i \leq x) = F(x)^n. \quad (1.12)$$

Si la distribution F est connue, on peut entièrement déterminer la distribution de M_n par F^n . Dans le cas où la distribution F est inconnue, la distribution de M_n pourra être approchée en utilisant la théorie asymptotique sur M_n . On utilise dans ce cas, la forme normalisée de M_n donnée par

$$M_n^* = \frac{M_n - b_n}{a_n},$$

pour des suites (a_n) et (b_n) telles que pour tout n , $a_n > 0$ et $b_n \in \mathbb{R}$. Dans ce cas, la distribution asymptotique de M_n est donnée par le théorème suivant :

Théorème 1.5.1.1 (Coles (2001, [41])) *S'il existe des suites $\{a_n > 0\}_{n \geq 1}$ et $\{b_n\}_{n \geq 1}$ telles que si $n \rightarrow +\infty$,*

$$\mathbb{P}(M_n^* \leq x) = \mathbb{P}\left(\frac{M_n - b_n}{a_n} \leq x\right) \rightarrow G(x) \quad (1.13)$$

où G est une distribution non-dégénérée, alors G appartient à une des familles de distributions suivantes :

I. Gumbel :

$$G(x) = \exp\left\{-\exp\left[-\left(\frac{x-\mu}{\sigma}\right)\right]\right\}, \quad x \in \mathbb{R} \quad (1.14)$$

II. Fréchet :

$$G(x) = \begin{cases} 0 & \text{si } x \leq \mu \\ \exp\left[-\left(\frac{x-\mu}{\sigma}\right)^{-\xi}\right] & \text{si } x > \mu \end{cases} \quad (1.15)$$

III. Weibull :

$$G(x) = \begin{cases} \exp\left\{-\left[-\left(\frac{x-\mu}{\sigma}\right)^\xi\right]\right\} & \text{si } x \leq \mu \\ 1 & \text{si } x > \mu \end{cases} \quad (1.16)$$

où $\sigma > 0$, $\mu \in \mathbb{R}$ et $\xi > 0$.

Preuve 1.5.1.1 *Consulter Coles (2001, [41]) ou Leadbetter (1983, [42]) pour la preuve.*

Generalized Extreme Value (GEV)

Les trois familles de distributions (Gumbel, Fréchet et Weibull) peuvent être combinées en une seule famille de distributions donnée par

$$G(x) = \exp \left\{ - \left[1 + \xi \left(\frac{x - \mu}{\sigma} \right) \right]^{-\frac{1}{\xi}} \right\} \quad (1.17)$$

définie sur $\{x : 1 + \xi(x - \mu)/\sigma > 0\}$, où $\mu \in \mathbb{R}$, $\sigma > 0$ et $\xi \in \mathbb{R}$. Cette distribution est appelée *GEV*(μ, σ, ξ). Les distributions de type *II* et *III* correspondent aux cas $\xi > 0$ et $\xi < 0$ respectivement. Le cas $\xi = 0$ correspond aux distributions de **Gumbel**. De manière générale, la fonction de répartition d'une *GEV*(μ, σ, ξ) est donnée par :

$$F(x|\mu, \sigma, \xi) = \begin{cases} \exp \left\{ - \left[1 + \xi \left(\frac{x - \mu}{\sigma} \right) \right]^{-\frac{1}{\xi}} \right\} & \text{si } \xi \neq 0 \\ \exp \left\{ - \exp \left[- \left(\frac{x - \mu}{\sigma} \right) \right] \right\} & \text{si } \xi = 0 \end{cases} \quad (1.18)$$

On peut déterminer dans ce cas le quantile x_p d'ordre $1 - p$ d'une GEV par :

$$x_p = \begin{cases} \mu - \frac{\sigma}{\xi} \left\{ 1 - [-\log(1 - p)]^{-\xi} \right\} & \text{si } \xi \neq 0 \\ \mu - \sigma \log [-\log(1 - p)] & \text{si } \xi = 0. \end{cases} \quad (1.19)$$

Generalized Pareto Distribution (GPD)

L'approche basée sur les distributions GEV peut être réductrice du fait que l'utilisation d'un seul maxima conduit à une perte d'information contenue dans les autres grandes valeurs de l'échantillon. La solution est de considérer plusieurs grandes valeurs au lieu de la plus grande (voir la section 3.5 de Coles (2001, [41])). Une autre approche appelée POT (Peaks Over Threshold) consiste à utiliser les observations qui dépassent un certain seuil, plus particulièrement les différences entre ces observations et le seuil u , appelées excès ($X - u$). La loi de ces excès sachant $X > u$ peut être approchée par la loi de Pareto généralisée est communément appelée "Generalized Pareto Distribution (GPD)" en anglais.

Théorème 1.5.1.2 (Pickands [43], [41]) *Soit X_1, \dots, X_n une suite de variables aléatoires indépendantes et identiquement distribuées de distribution F satisfaisant le théorème 1.5.1.1. C'est-à-dire pour $n \rightarrow +\infty$,*

$$\mathbb{P}(M_n^* \leq x) \rightarrow G(x)$$

avec $G(x)$ donnée par l'équation (1.17). Alors pour u assez grand, la distribution de $(X - u)$ conditionnellement à $X > u$ est approchée par

$$H(y) = 1 - \left(1 + \frac{\xi y}{\tilde{\sigma}} \right)^{-\frac{1}{\xi}} \quad (1.20)$$

où H est définie sur $\{y : y > 0 \text{ et } (1 + \xi y)/\tilde{\sigma}\}$ et $\tilde{\sigma} = \sigma + \xi(u - \mu)$.

Preuve 1.5.1.2 Consulter les travaux de Pickands (1975, [43]), Coles (2001, [41]) ou Leadbetter (1983, [42]) pour la preuve.

La famille de distributions donnée par l'équation (1.20) est appelée famille de distributions de Pareto généralisée (GPD). De manière plus précise, la fonction de répartition d'une GPD est donnée par

$$G(x|u, \sigma_u, \xi) = \begin{cases} 1 - \left[1 + \xi \left(\frac{x-u}{\sigma_u}\right)\right]_+^{-\frac{1}{\xi}} & \text{si } \xi \neq 0 \\ 1 - \exp\left[-\left(\frac{x-u}{\sigma_u}\right)\right]_+ & \text{si } \xi = 0 \end{cases} \quad (1.21)$$

où $x > 0$, $\sigma_u > 0$ et $\left[1 + \xi \left(\frac{x-u}{\sigma_u}\right)\right]_+ > 0$. Pour plus de détails sur la théorie des valeurs extrêmes, le lecteur pourra aussi consulter les thèses de Toulemonde ([44]), de Ndao [45], de Garrido [46], de Lekina [47], de Stupfler [48], de Gardes [49] ou de El Methni [50].

1.5.2 Notion de censure

Soit X une variable aléatoire positive ou nulle. On parle de censure ou de troncature si on n'arrive pas à observer les données de manière complète. En cas de censure, il existe une variable de censure C qui permet de modéliser les données non observées. Il existe différents types de censures : si au lieu d'observer X , on observe C et que $X > C$ (respectivement $X < C$, $C_1 < X < C_2$), on dit qu'on a une censure à droite (respectivement à gauche, censure par intervalle). Pour plus de détails sur la censure, consulter la thèse de Ndao ([45]).

1.5.3 Mélange d'extrêmes en présence de censure

L'un des défis dans la théorie des valeurs extrêmes est d'estimer un seuil u au dessus duquel les données sont considérées extrêmes. Si on considère que le seuil est un paramètre à estimer, on peut voir le modèle comme un mélange à deux composantes : (1) une composante en dessous du seuil u , appelée "bulk model" et (2) une composante au dessus de u appelée "tail model". Hu détaille dans sa thèse ([51]) la bibliographie et les méthodes utilisées pour étudier ce genre de modèle.

On suppose dans cette thèse qu'au delà du seuil u , les données ne sont pas complètement observées. Elles sont censurées aléatoirement à gauche. Pour des raisons que nous avons expliquées au chapitre 5, la méthode habituelle du maximum de vraisemblance ne marche pas très bien. Nous proposons ici, une méthode d'estimation à deux étapes (voir le chapitre 5 pour les détails).

Pour plus de détails sur les modèles de mélanges extrêmes, voir la thèse de Hu ([51]) ou celle de MacDonald ([52]).

1.6 Présentation générale des résultats

L'analyse des données de co-infection a donné lieu au chapitre 2. Dans ce chapitre, nous présentons le jeu de données. Ensuite, nous décrirons un modèle logistique multinomial

en utilisant des stratégies de sélection de variables. Nous avons aussi analysé les facteurs influents des maladies et présenté une analyse prédictive permettant de diagnostiquer la co-infection et de donner une recommandation de traitement. Ce chapitre fait l'objet d'un article, Loum et al.(2017, [53], voir annexe A), soumis et révisé à la revue "The International Journal of Biostatistics".

Dans le chapitre 3, nous avons étudié les mélanges de régression logistique. Nous avons présenté une méthode d'estimation à deux étapes : (1) par diagonalisation jointe, puis (2) par minimisation d'un critère du type moindres carrés. Deux résultats d'identifiabilité ont été montrés. La consistance et la normalité asymptotique des estimateurs sont aussi présentées. Dans le chapitre 4 nous avons présenté une extension des mélanges étudiés dans le chapitre 3. Des résultats d'identifiabilité ont été montrés dans le cas où le vecteur de covariable est continu mais aussi dans le cas où le vecteur de covariable est composé d'une partie continue et d'une partie catégorielle. Nous avons aussi montré un résultat d'identifiabilité dans le cas des données longitudinales. Les résultats de ces deux chapitres font l'objet d'un article (Loum et al. (2018a, [54])) soumis à la revue "Journal of Classification" (voir Annexe B) et d'un Package R *Morpheus* (Loum et Auder, [55]) disponible sur CRAN (voir annexe C).

Le chapitre 5, présente les mélanges de modèles extrêmes en présence de censure. Nous avons présenté dans ce chapitre une méthode d'estimation à deux étapes basées sur le maximum de vraisemblance. Ce chapitre fait l'objet d'un article (Loum et al. (2018c, [56])) à soumettre dans une revue internationale à comité de lecture.

1.6.1 Chapitre 2 Modèle de régression logistique multinomiale pour le diagnostic de la co-infection

L'objectif de l'analyse des données de co-infection était de fournir un outil de diagnostic. De plus, en cas de co-infection, une recommandation de traitement était nécessaire.

Données de co-infection

La base de données initiale contient 15523 individus avec 25 covariables et 3 variables réponses. Les 3 variables réponses ont servi à construire une réponse multinomiale permettant de prendre en compte toutes les maladies. Une analyse descriptive a permis de réduire le nombre de covariables à 15. Deux jeux de données ont été considérés ici selon la définition de l'infection arbovirale (*IgM* ou *IgM/IgG*).

Modèle statistique

Nous avons mis en place la régression logistique multinomiale dans le cas de la co-infection. Nous avons effectué une sélection de variables en utilisant les forêts aléatoires et le test du rapport de vraisemblance. Les variables retenues ont été utilisées dans l'analyse des facteurs influents des différents profils d'infection (co-infection, paludisme seul et arbovirus seul).

Analyse prédictive

Un test d'indépendance entre arbovirus et paludisme a été mis en place. Nous avons retenu que l'arbovirose et le paludisme sont liés. Ce qui nous a permis de mettre en place une méthode de classification basée sur la probabilité d'être co-infecté sachant qu'on a le paludisme.

Résultats

Nous avons trouvé que si cette probabilité de co-infection est supérieure à un seuil (à calibrer sur les données), que la durée de la maladie est supérieure à trois jours et que l'âge du patient est supérieur à 10 ans, on peut alors penser que la personne est malade d'arbovirus. Si la température corporelle est supérieure à $40^{\circ}C$ (C'est-à-dire que le patient présente de la fièvre) et que le patient présente des nausées/vomissements pendant la saison des pluies, on peut alors penser que la personne est malade du paludisme.

1.6.2 Chapitre 3 Mélange de modèles linéaires généralisés et méthode des moments : identifiabilité & applications

Pour $Y \in \{0, 1\}$ et $X \in \mathbb{R}^d$, on considère (X, Y) est de loi \mathbb{P}_θ telle que

$$\mathbb{P}(Y = 1|X = x) = \sum_{k=1}^K \omega_k g(\langle x, \beta_k \rangle + b_k) \quad (1.22)$$

et X suivant une loi normale. On note ici $\theta = (\omega, \beta, b) \in \Theta$ avec $\beta = [\beta_1, \dots, \beta_K] \in \mathbb{R}^{d \times K}$ la matrice des K vecteurs de paramètres, $b = (b_1, \dots, b_K)$ le vecteur des intercepts et g la fonction lien. On note $\omega = (\omega_1, \dots, \omega_K)$ le vecteur des poids du mélange, sous les contraintes $\sum_{k=1}^K \omega_k = 1$ et $\omega_k > 0$.

Algorithme

On estime d'abord les directions des vecteurs de β . C'est-à-dire les $\mu_k = \beta_k / \|\beta_k\|$, $k = 1, \dots, K$. Pour cela, nous utilisons la diagonalisation jointe ([40]). Nous estimons ensuite tous les paramètres du modèle en minimisant

$$\begin{aligned} Q_n(\theta) &= \sum_{j \in [d]} \left\{ \hat{M}_1[j] - M_1(\theta)[j] \right\}^2 \\ &+ \sum_{j, k \in [d]} \left\{ \hat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 \\ &+ \sum_{j, k, l \in [d]} \left\{ \hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2. \end{aligned} \quad (1.23)$$

avec M_j , le moment (croisé entre x et y) d'ordre j et \hat{M}_j son estimateur, pour $j = 1, 2, 3$.

Résultats

Pour montrer qu'on peut bien retrouver les paramètres à partir des moments, deux résultats d'identifiabilité ont été montrés. La consistance et la normalité asymptotique aussi ont été montrées.

1. Identifiabilité Probit

Sous l'hypothèse (H1) du chapitre 3, on a que si la fonction lien est probit, alors on peut retrouver K (le nombre de composante du mélange) et le paramètre $\theta = (\omega, \beta, b)$ à partir de $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$.

2. Identifiabilité générale

Si les hypothèses (H1), (H2) et (H3) du chapitre 3 sont vraies et que la fonction lien vérifie $g^{(3)}(0) \neq 0$, alors il existe $L > 0$ et $B > 0$ tels que si $\|\beta_k\| < L$ et $|b_k| < B$, on peut retrouver K et θ à partir de $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$.

3. Consistance

Supposons que le modèle est identifiable et que Θ compact. Si les hypothèses de (H1) à (H5) du chapitre 3 sont vraies, alors $\hat{\theta}_n$ converge en probabilité vers θ^* , où θ^* est la vraie valeur de θ .

4. Normalité Asymptotique

Supposons que les hypothèses de (H1) à (H4) du chapitre 3 sont vraies et que $\hat{\theta}_n$ est consistant pour θ^* . Alors $\sqrt{n}(\hat{\theta}_n - \theta^*)$ converge en loi vers une loi normale centrée dont on notera Σ la matrice de covariance.

1.6.3 Chapitre 4 Extensions des mélanges de modèles linéaires généralisés

Modèle

Soit (X, Y) un vecteur aléatoire de loi \mathbb{P}_θ telle que $X \in \mathbb{R}^d$, $Y \in \mathbb{R}$. On suppose que la loi de Y sachant X est un mélange donné par

$$E(Y|X) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + b_k),$$

avec $\omega = (\omega_1, \dots, \omega_K)$, $\beta = [|\beta_1|, \dots, |\beta_K|] \in \mathbb{R}^{d \times K}$, et $b = (b_1, \dots, b_K) \in \mathbb{R}^K$. On suppose que pour tout k , $\omega_k \geq 0$, que $\sum_{k=1}^K \omega_k = 1$, et que g est une fonction à valeurs dans $(0, 1)$.

Si $\theta = (K, g, \omega, \beta, b)$, on cherche à retrouver les paramètres inconnus K , g , ω , β et b . On notera $\mu_k = \beta_k / \|\beta_k\|$ et $\lambda_k = \|\beta_k\|$, $k = 1, \dots, K$, de sorte que $\beta_k = \lambda_k \mu_k$.

Résultats d'identifiabilité

On montre que ce modèle est identifiable dans les cas suivants : (1) si le vecteur de covariable est continue, (2) si le vecteur de covariable est composé d'une partie continue et d'une partie catégorielle, (3) les données sont longitudinales.

1. *Covariables continues : identifiabilité des vect* On considère que le vecteur des covariables X est continu. Dans ce cas, on peut retrouver le nombre de composantes K et les vecteurs normalisés μ_k , sous les hypothèses (H1), (H2) et (H3) du chapitre 4. Comme dans le chapitre 3, le but est de retrouver tous les paramètres K , g , ω ,

β et b . Si on remplace l'hypothèse $H2$ par l'hypothèse (H2bis) du chapitre 4, on peut retrouver les paramètres K, g, ω, β et b .

2. *covariables continues et covariables catégorielles*

On considère maintenant qu'une partie $X \in \mathbb{R}^d$ des covariables est continue et une autre partie Z , de dimension d' , est constituée de variables catégorielles à valeurs dans $\{z_1, \dots, z_m\} \subset \mathbb{R}^{d'}$. Sous ce modèle, on a l'identifiabilité en ajoutant l'hypothèse (H4) du chapitre 4.

3. *Données longitudinales* On suppose ici qu'on est dans le cas des données longitudinales. C'est-à-dire, on a des répétitions indépendantes (conditionnellement à la population) avec des covariables différentes. L'observation Y est m -dimensionnelle (si on a m répétitions) et on a m covariables. Le modèle est donné dans ce cas par

$$E(Y|X, Z) = \sum_{k=1}^K \omega_k \otimes_{j=1}^m g(\langle \beta_k, X_j \rangle + \langle \gamma_k, Z_j \rangle + b_k).$$

On a alors l'identifiabilité sous les hypothèses (H1)-(H4) du chapitre 4.

1.6.4 Chapitre 5 Mélange de valeurs extrêmes en présence de censure

Nous considérons ici une variable aléatoire X de loi donnée par (MacDonald et al. (2011,[57]))

$$F_X(x) = \begin{cases} (1 - \phi_u) \frac{H(x|\beta, \lambda)}{H(u|\beta, \lambda)} & \text{si } x \leq u \\ (1 - \phi_u) + \phi_u G(x|u, \sigma, \xi) & \text{si } x > u \end{cases} \quad (1.24)$$

avec $\phi_u = \mathbb{P}(X > u)$ et $G(\cdot|u, \sigma, \xi)$ la fonction de répartition de la loi de Pareto généralisée (GPD). On suppose que la variable X n'est pas complètement observée mais elle est censurée. On observe alors Z et δ avec

$$Z = \begin{cases} X & \text{si } x \leq u \\ X \wedge C & \text{si } x > u \end{cases} \quad \text{et } \delta = \mathbf{1}_{\{X \leq C\}} \quad (1.25)$$

où C est une variable de censure de loi $GPD(\cdot|u, \sigma', \xi')$.

Estimation des paramètres $\theta = (\beta, \lambda, u, \sigma, \xi)$

Nous avons estimé dans ce chapitre le paramètre θ par maximum de vraisemblance en deux étapes :

1. On suppose que le seuil u fixé et connu. On cherche alors à estimer $\theta'_u = (\beta, \lambda, \sigma, \xi)$ par maximum de vraisemblance.
2. On choisit ensuite plusieurs valeurs de u (u_1, \dots, u_L) et pour chaque valeur u_l , on note $\hat{\theta}_n^l$ l'estimateur du maximum de vraisemblance de θ'_{u_l} , $l = 1, \dots, L$. On obtient ensuite $\hat{\theta}_n$ de θ en maximisant les log-vraisemblances sur toutes les valeurs de u .

Estimation des quantiles extrêmes

nous avons ensuite estimé les quantiles extrêmes, en utilisant deux méthodes :

1. Méthode classique

D'abord, en utilisant la méthode classique, c'est-à-dire inverser la fonction de répartition.

2. Méthode de reparamétrisation

On utilise la méthode de reparamétrisation ; c'est-à-dire, on écrit le paramètre u en fonction du quantile x_p d'ordre $1 - p$. Ensuite, on ré-estime les paramètres par maximum de vraisemblance.

Chapitre 2

Modèle de régression logistique multinomiale pour le diagnostic et la recommandation de traitement en cas de co-infection entre deux maladies (*Arbovirus & Paludisme à Kédougou*)

Sommaire

2.1	Introduction	20
2.2	Description des données	21
2.2.1	Réponse multinomiale et jeux de données	22
2.2.2	Covariables	24
2.3	Modèle Statistique	26
2.3.1	La régression logistique multinomiale	27
2.3.2	La régression logistique multinomiale dans le cas Arbovirus-Paludisme	30
2.3.3	Sélection de variables	30
2.3.4	Analyse des facteurs influents des différentes maladies	42
2.4	Analyse prédictive	47
2.4.1	Test d'indépendance entre arbovirus et paludisme	47
2.4.2	Diagnostic de la co-infection	51
2.5	Discussion	54

Résumé

Dans les régions tropicales, les populations continuent de souffrir de la mortalité due au paludisme et aux maladies arbovirales. Dans la région de Kédougou (Sud-Est du Sénégal), ces maladies sont toutes endémiques à cause du climat et de la position géographique. La co-circulation des parasites du paludisme et des arbovirus peut expliquer l'observation des cas de co-infection signalés dans cette région. En effet, il existe une forte ressemblance entre les symptômes du paludisme et ceux des maladies arbovirales. Cette ressemblance pose un vrai problème de diagnostic médical de la co-infection ; du fait que l'origine de ces maladies n'est pas complètement connue. Certaines personnes peuvent être immunisées à l'un ou l'autre agent pathogène. Ces immunités sont souvent obtenues avec l'âge ou à l'exposition dans des zones endémiques. Ainsi, la co-infection a besoin d'être mieux diagnostiquée dans ces zones. En utilisant les données collectées à Kédougou par l'Institut Pasteur de Dakar, entre 2009 et 2013, nous sélectionnons les variables importantes pouvant expliquer la co-infection et les autres formes simples d'infections. Nous ajustons ensuite un modèle logistique multinomial en utilisant les variables sélectionnées. Nous observons un ensemble de variables relatif à chaque maladie et à la co-infection. Nous fournissons une analyse prédictive en commençant par tester l'indépendance entre arbovirus et paludisme. Ce qui nous permet de calculer la probabilité d'être co-infecté sachant qu'on a le paludisme. Si cette probabilité est supérieure à un seuil à calibrer sur les données, que la durée de la maladie est supérieure à 3 jours et que l'âge du patient est supérieur à 10 ans, alors on peut penser que la personne est malade d'arbovirus. Si la température corporelle est supérieure à 40°C et que le patient présente des nausées et/ou vomissements pendant la saison des pluies, on peut penser que la personne est malade du paludisme.

Mots clés : Arbovirus, co-infection, Classification, Forêts aléatoires, Paludisme, Régression logistique multinomiale, Sélection de variable, Stepwise.

Abstract

In tropical regions, populations continue to suffer morbidity and mortality from malaria and arboviral diseases. In Kedougou (Senegal), these illnesses are all endemic due to the climate and its geographical position. The co-circulation of malaria parasites and arboviruses can explain the observation of co-infected cases. Indeed there is strong resemblance in symptoms between these diseases making problematic targeted medical care of co-infected cases. This is due to the fact that the origin of illness is not obviously known. Some cases could be immunized against one or the other of the pathogens, immunity typically acquired with factors like age and exposure as usual for endemic area. Then, co-infection needs to be better diagnosed. Using data collected from patients in Kedougou region, from 2009 to 2013, we selected relevant variables in explaining co-infection status and adjusted a multinomial logistic model using these relevant variables. We observed specific sets of variables explaining each of the diseases exclusively and the co-infection. We tested the independence between arboviral and malaria infections and derived co-infection probabilities from the model fitting. In case of a co-infection probability greater than a threshold value to be calibrated on the data, duration of illness above 3 days and age above 10 years-old are mostly indicative of arboviral disease while body temperature higher than 40°C and presence of nausea or vomiting symptoms during the rainy season are mostly indicative of malaria disease.

Keywords : Arbovirus, co-infection, malaria, multinomial logistic regression, random forest classification, variable selection, stepwise.

2.1 Introduction

Les infections simultanées sont souvent observées, chez les maladies à transmission vectorielles telles que le paludisme et les maladies virales à base d'arbovirus, dans les régions tropicales ([1], [2]). C'est le cas du paludisme et de la dengue dans les régions tropicales américaines, africaines et asiatiques ([3, 4, 5, 6, 7, 8, 9]). Le paludisme peut être facilement attribué à d'autres maladies fébriles ; car ses symptômes cliniques sont souvent indiscernables de ceux qui ont été observés dans la dengue ou le chikungunya ([10]). Depuis l'introduction du test de diagnostic rapide (TDR) au Sénégal (en 2007), le paludisme a été mieux diagnostiqué et une baisse importante des chiffres (estimations antérieures de la prévalence) a été observée. Ceci laisse penser que le paludisme a été surestimé dans le passé, au détriment d'autres maladies infectieuses comme celles causées par les arbovirus ou les maladies bactériennes ([11, 12]). Le traitement de la fièvre avec des médicaments antipaludiques était largement pratiqué pour réduire la mortalité attribuable au paludisme. Mais cette pratique signifie que les patients malades peuvent être traités de manière inappropriée, en particulier lorsque le test de diagnostic rapide n'est pas disponible ou lorsqu'on ne peut pas tester les infections arbovirales. Ainsi, un diagnostic erroné des co-infections d'arbovirus comme infections de paludisme peut être une cause de sous-estimation des infections à arbovirus émergentes. En 2009, la surveillance de la maladie fébrile aiguë (AFI) a été mise en place à Kédougou (sud-est du Sénégal) pour la détection précoce des épidémies d'arbovirus et du paludisme et afin de mesurer avec précision la morbidité et la mortalité dans cette région. En raison de la co-circulation des parasites du paludisme et des arbovirus, principalement la dengue (DEN), la Chikungunya (CHIK), la Zika (ZIK), la fièvre jaune (YF) et les virus de la fièvre de la Vallée du Rift (RVFV) dans cette région (négligeant la prévalence des autres infections arbovirales), des cas d'infections simultanées ont été observés et posent un défi pour le diagnostic médical ([13]). Nous comparons ici les profils cliniques des patients co-infectés aux profils cliniques des patients mono-infectés par l'analyse statistique d'un ensemble de données recueillies auprès de patients fébriles dans la région de Kédougou au Sénégal de 2009 à 2013. Notre étude vise à caractériser les facteurs de risque des co-infections et à fournir des indicateurs statistiques qui améliorent le diagnostic de l'arbovirus.

Les données de notre étude ont été fournies par l'Institut Pasteur de Dakar (IPD) à Kédougou. Dans cette région, le paludisme et les arbovirus sont endémiques en raison du climat et des mouvements de la population. Les données ont été collectées dans sept centres de santé de la région : L'hôpital rural de Ninfesha, les centres de santé de Kédougou et de Saraya, le poste de santé de Bandafassi et de Khossanto, le centre de santé militaire de Kédougou et l'équipe mobile de la mission catholique. Les critères d'inclusion dans la collecte de données étaient :

- (i) avoir au moins un (01) an à la date de visite,
- (ii) avoir de la fièvre (i.e, Température $\geq 38^{\circ}C$) et
- (iii) présenter au moins l'un des symptômes cliniques décrits à la section 2.2.

Les patients satisfaisant les critères d'inclusion ont été sélectionnés après avoir signé un accord de participation dans l'étude.

Dans ce chapitre, nous proposons un modèle logistique multinomial pour étudier la co-infection entre le paludisme et les arboviroses. La réponse multinomiale est construite à partir des trois réponses enregistrées dans le jeu de données initial. Elle est constituée

de quatre catégories associées aux quatre groupes de patients qui sont : le groupe des patients positifs aux arbovirus seuls (positif à au moins l'un des 5 arbovirus testés), le groupe des patients positifs au paludisme seul, le groupe des co-infectés (positif à au moins un des arbovirus et au paludisme) et le groupe des autres maladies (négatifs à tous les tests). Les symptômes de ce dernier groupe témoin sont probablement dus à d'autres pathogènes pour lesquels tous les groupes étaient supposés être également exposés. Après une description des données d'étude, nous avons fait une sélection de variables. Nous avons utilisé la méthode pas à pas (stepwise) et celle du test du rapport de vraisemblance pour sélectionner les variables importantes. La robustesse de cette sélection de variables a été étudiée en utilisant les forêts aléatoires, qui est une méthode basée sur la mesure de l'importance des variables ([58]). Ensuite, nous avons ajusté un modèle paramétrique logistique multinomial en utilisant les variables sélectionnées précédemment. Ainsi, à partir de l'analyse des facteurs influents des différentes réponses, nous pourrions discuter des questions suivantes : le paludisme est-il vraiment lié à l'infection aux arbovirus ? Quels sont les facteurs qui peuvent expliquer la co-infection ? Quels sont les facteurs risque permettant de distinguer le paludisme de l'infection aux arbovirus ? Une analyse prédictive basée sur la probabilité d'être co-infecté sachant que l'individu ait le paludisme a été faite. Pour calculer cette probabilité, nous avons testé l'association ou la dépendance entre l'arbovirus et le paludisme, en proposant un test statistique du type Wald. S'il y a dépendance, nous pourrions donc calculer pour chaque individu la probabilité d'être co-infecté sachant qu'il a le paludisme.

Ce chapitre est organisé comme suit. D'abord, à la section 2.2, on présentera le jeu de données. Ensuite, on décrira le modèle statistique à la section 2.3, ainsi que la sélection de variables et l'analyse des facteurs influents des maladies. A la section 2.4, nous présenterons l'analyse prédictive. On donnera enfin une discussion à la section 2.5.

2.2 Description des données

La base de données initiale contient 15523 individus avec 25 covariables et 3 réponses. Les réponses sont constituées de 3 variables binaires indiquant (i) la présence ou l'absence de parasites du paludisme dans le sang, (ii) la détection d'anticorps *IgM* contre le virus et (iii) la détection de virus ou d'anticorps *IgG* contre le virus. Les covariables sont constituées de variables quantitatives telles que l'âge (*age*), la température (*temperature*), le nombre de jours de maladie (*number of sick days*) et la pluviométrie (*pluviometry*), et des variables binaires : le sexe (*sex*), la gorge irritée (*irithated throat*), l'éruption (*eruption*), la fièvre (*fever*), les douleurs articulaires (*joint pain*), les douleurs musculaires (*muscl pain*), les douleurs aux yeux (*eye pain*), la céphalée (*cephalalgia*), les toux (*cough*), les nausées/vomissements (*nausea/vomiting*), les frissons (*chills*), la diarrhée (*diarrhea*), la congestion nasale (*nasal congestion*), l'ictère/jaunice (*icterus/jaudice*), sang dans les selles (*blood in stool*), les vomissements de sang (*vomiting blood*), les pertes vaginales sanguines (*blood vaginal discharge*), du sang dans les gencives (*blood gums*), les pétéchies (*petechea*), les pertes nasales sanguines (*blood nose*) et du sang dans les urines (*blood in the urine*). La figure 2.1 représente un résumé de ces 25 covariables et des 3 variables réponses.

Designation	Type	For categorical variables		For quantitative variables		
		# levels	mean	median	min	max
Age	quantitative		19.5	16.5	1	90
Temperature	quantitative		38.97	39	38	42
Number sick days	quantitative		3.039	3	0	19
Pluviometry	quantitative		147.5	76.1	0	500.2
Sex	categorical	2				
Fever	categorical	2				
Cephalalgia	categorical	2				
Eye pain	categorical	2				
Muscl pain	categorical	2				
Joint pain	categorical	2				
Eruption	categorical	2				
Blood nose	categorical	2				
Nausea/vomiting	categorical	2				
Diarrhea	categorical	2				
Chills	categorical	2				
Cough	categorical	2				
Vomiting blood	categorical	2				
Blood in stool	categorical	2				
Blood in the urine	categorical	2				
Petechea	categorical	2				
Blood gums	categorical	2				
Blood vaginal discharge	categorical	2				
Nasal congestion	categorical	2				
Irritated throat	categorical	2				
Icterus/Jaundice	categorical	2				
Malaria	categorical	2				
IgM	categorical	2				
IgG	categorical	2				

FIGURE 2.1 – Résumé des variables du jeu de données initial.

2.2.1 Réponse multinomiale et jeux de données

Nous considérons comme cas positifs au paludisme, les individus positifs au test par la goutte épaisse. Pour l'infection à l'arbovirus, nous considérons comme cas positifs, les individus positifs par test ELisa (IgM ou IgG), PCR à temps réel et isolement, à l'un des arbovirus suivants : dengue (DEN), chikungunya (CHIK), zika (ZIK), fièvre jaune (YF), fièvre de la vallée du rift (FVR). Ainsi on a trois variables réponses à partir des cas considérés (paludisme, IgM–arbovirus, IgG–arbovirus). En se basant sur ces données, nous construisons une nouvelle réponse multinomiale avec 4 catégories, en considérant les combinaisons des différentes infections, pour définir et coder les quatre profils d'infection :

$$y = \begin{cases} 0 & \text{“Autres maladies (O)”} \\ 1 & \text{“Arbovirus Seul (A)”} \\ 2 & \text{“Paludisme Seul (M)”} \\ 3 & \text{“co-infection (C)”} \end{cases}$$

1. *“Arbovirus Seul”* codé 1 : qui correspond aux individus positifs au test d'au moins l'un des arbovirus listé plus haut et négatifs au test du paludisme.
2. *“Paludisme Seul”* codé 2 : correspond aux individus positifs au test du paludisme et négatifs aux tests des arbovirus listés plus haut.
3. *“co-infection”* codé 3 : correspondant aux individus positifs au test du paludisme et à celui d'au moins un des arbovirus listé plus haut.
4. *“Autres maladies”* codé 0 : c'est la classe de référence qui correspond aux individus négatifs à tous les tests (arbovirus et paludisme). C'est le groupe de contrôle.

Dans cette étude, les cas d'arbovirus sont diagnostiqués par la détection d'anticorps

IgM ou IgG. Nous pouvons avoir deux manières de définir un cas d'arbovirus : (1) en considérant seulement la détection d'anticorps IgM ou (2) en considérant la détection des deux anticorps IgM ou IgG. La détection d'anticorps IgM chez un patient signifie, biologiquement, qu'il a eu une infection récente à l'arbovirus. Ainsi, on considère que les cas positifs d'infection aux arbovirus sont les individus positifs à la détection d'anticorps *IgM* seulement. En ignorant les individus avec au moins une donnée manquante (974 données manquantes pour la réponse paludisme et 803 données manquantes pour les covariables), on obtient un jeu de données de taille $n = 12288$, appelé *IgM - data*. Notons que les distributions des différentes variables avec et sans données manquantes restent similaires. Un résumé du jeu de données *IgM - data* est donné par le tableau 2.1. On peut voir, à partir de ce tableau, que ce jeu de données est très déséquilibré (3 cas d'arbovirus ou de co-infection sur 1000 patients) et nécessite une analyse statistique spécifique.

Arbovirus \ Paludisme	+	-	Total
+	21 (<i>C</i>)	18 (<i>A</i>)	39 (<i>A+</i>)
-	7069 (<i>M</i>)	5180 (<i>O</i>)	12305 (<i>A-</i>)
Total	7087 (<i>M+</i>)	5201 (<i>M-</i>)	12288

TABLE 2.1 – Résumé du jeu de données *IgM - data*. *A+* pour les individus positifs à l'arbovirus, *A-* pour les individus négatifs à l'arbovirus, *M+* pour les individus positifs au paludisme et *M-* pour les individus négatifs au paludisme.

Figure 2.2 présente la cinétique des infections aux arbovirus et les tests de diagnostic appropriés. Ainsi, on voit à partir de cette figure que le diagnostic d'une infection aux arbovirus est idéalement basé sur la détection d'anticorps IgM chez le patient. En effet, les anticorps IgG peuvent persister plusieurs années ou toute une vie entière.

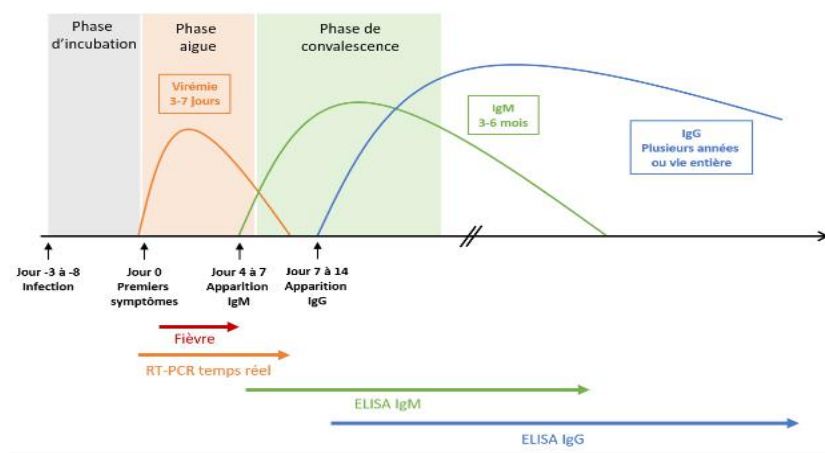


FIGURE 2.2 – Cinétique des infections aux arbovirus et les tests appropriés.

Cependant, afin d'obtenir un jeu de données plus équilibré, nous avons décidé d'en construire un autre en considérant comme cas positif d'arbovirus, les patients positifs aux tests de détection d'anticorps IgM ou IgG. Comme 13412 données manquantes ont

été enregistrées sur la variable IgG, la taille du jeu a été considérablement réduite et nous avons obtenu un jeu de données de taille $n = 1976$, appelé *IgM/IgG – data* et résumé dans le tableau 2.2. Pour ce jeu de données, nous avons comparé les distributions de chaque covariable avec et sans données manquantes sur la réponse IgG. À l’exception de la variable congestion nasale variable qui est surreprésentée (60% des cas positifs dans l’échantillon comparé à 40% dans le jeu de données initial, voir tableau 2.3), les distributions des autres variables sont similaires. Nous avons donc considéré qu’ignorer les individus avec des données manquantes n’affecterait pas considérablement l’analyse prédictive.

Paludisme			
Arbovirus	+	–	Total
+	397 (<i>C</i>)	263 (<i>A</i>)	633 (<i>A+</i>)
–	751 (<i>M</i>)	565 (<i>O</i>)	1318 (<i>A–</i>)
Total	1148 (<i>M+</i>)	828 (<i>M–</i>)	1976

TABLE 2.2 – Résumé du jeu de données *IgM/IgG*. *A+* pour les individus positifs à l’arbovirus, *A–* pour les individus négatifs à l’arbovirus, *M+* pour les individus positifs au paludisme et *M–* pour les individus négatifs au paludisme.

Réponses			
Données	Paludisme	IgM	IgG
Avec Données manquantes	31%	0.9%	40%
Sans données manquantes	41%	46%	60%

TABLE 2.3 – Un résumé des distributions des réponses pour la variable Congestion Nasale. Avec données manquantes, pour chaque réponse, correspond au taux de cas positifs à la congestion nasale sur le jeu de données initial. Sans données manquantes, pour chaque réponse, correspond au taux de cas positifs à la congestion nasale sur le jeu de données sans données manquantes dans la réponse.

Pour la suite, nous considérerons deux jeux de données qui sont dérivés d’un même jeu de données initial en utilisant deux codages différents : 1. le jeu de données *IgM/IgG – data*, assez équilibré pour appliquer notre méthodologie ; 2. le jeu de données *IgM – data* contenant les vrais cas d’arbovirus (d’un point de vue biologique) mais très déséquilibré. Nous utiliserons dans la section 2.3.3 une stratégie de ré-échantillonnage pour régler ce problème.

2.2.2 Covariables

Après une analyse descriptive, nous avons constaté qu’il y avait des variables redondantes, des variables corrélées et des variables déséquilibrées. Ainsi certaines décisions ont été prises : il y avait une redondance entre la variable *fever* et *temperature*. Nous avons choisi de garder *temperature* aux dépens de *fever* qui était catégorielle. Les deux variables *irithated throat* et *eruption* sont très déséquilibrées avec moins de 2% de positifs et nous avons décidé de ne pas les garder dans l’étude. Les variables (symptômes) hémorragiques

(*blood in stool, vomiting blood, blood vaginal discharge, blood gums, petechea, blood nose et blood in the urine*) ne sont pas observées chez les individus qui sont positifs aux arbovirus. En effet, ces symptômes ne sont observés qu'à un état très avancé de la maladie. Par suite, nous avons aussi décidé de ne pas les garder dans l'étude.

Nous allons garder 15 variables sur les 25 de la base de données. Dans les deux jeux de données, on a quatre (04) variables quantitatives : la température, mesurée en degré Celsius le jour de la consultation (*temperature*), le nombre de jours de maladie défini par le nombre de jours entre la date de la consultation et la date d'apparition des symptômes (*number of sick days*), l'âge du patient enregistré en année (*age*). Pour chaque patient, la pluviométrie, en millimètre, correspondant au mois de consultation est aussi donnée (*pluviometry*). En plus de ces quatre variables quantitatives, on a 11 variables qualitatives : le sexe (*sex*) codé par M (Male) et F (Female) et 10 autres variables binaires, codées par 0 et 1, qui représentent les symptômes chez les patients. Pour chacune de ces variables, on note 0 l'absence et 1 la présence du symptôme correspondant. Les variables retenues dans l'étude sont résumées dans le tableau de la Figure 2.3.

Designation	For categorical variables			quantitative variables			
	# levels	0 (%)	1 (%)	mean	median	min	max
Age				19.5	16.5	1	90
Temperature				38.97	39	38	42
Number of sick days				3.039	3	0	19
Rainfall				147.5	76.1	0	500.2
Sex (F=0 and H=1)	2	42	58				
Cephalalgia	2	6	94				
Nausea/vomiting	2	50	50				
Diarrhea	2	83	17				
Chills	2	45	55				
Cough	2	64	36				
Eye pain	2	95	5				
Joint pain	2	77	23				
Muscl pain	2	71	29				
Nasal congestion	2	54	46				
Ictere/jaundice	2	95	5				
Malaria	2	42	58				
IgM	2	99	1				
IgG	2	95	5				

FIGURE 2.3 – Résumé des variables retenues dans l'étude.

Dans les deux jeux de données, le groupe de sexe masculin représente environ 58% des individus tandis que les femmes constituent 42%. Notons que les variables qualitatives considérées dans l'étude sont assez équilibrées (en termes de proportion de positives et de négatives). Dans le jeu de données *IgM – data*, les deux catégories “*co-infection*” et “*arbovirus seul*” sont sous-représentées, ce qui explique qu'elles ne sont pas bien représentées dans la figure 2.5. Une analyse du jeu de données *IgM/IgG – data* montre que l'âge est positivement corrélé à l'infection aux arbovirus alors que la température, les nausées/vomissements et la pluviométrie sont associés au paludisme. Par exemple,

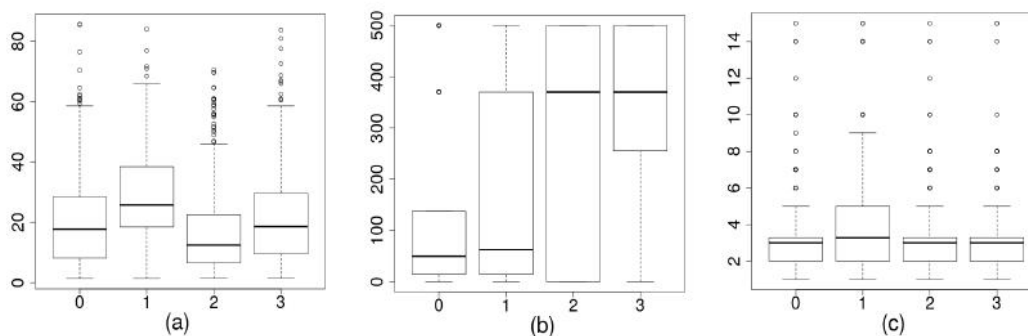


FIGURE 2.4 – *IgM/IgG-data* ; (a) *age*, (b) *pluviometry* and (c) *number of sick days* distribution empirique en fonction de la réponse Y .

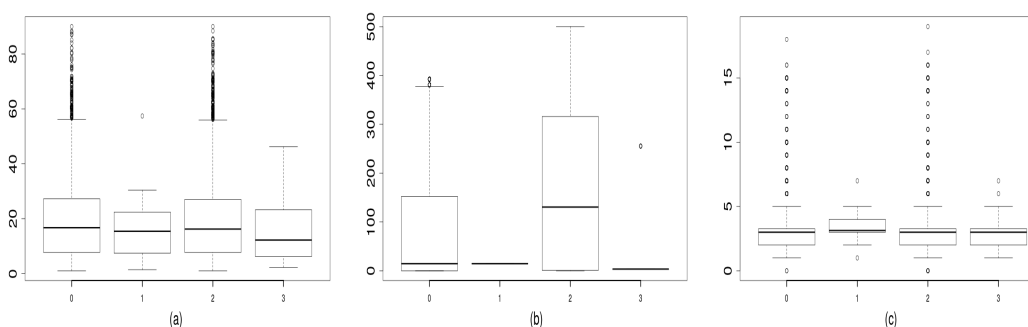


FIGURE 2.5 – *IgM-data* ; (a) *age*, (b) *pluviometry* and (c) *number of sick days* distribution empirique en fonction de la réponse Y .

sur les patients présentant des nausées/vomissements, 45% ont le paludisme seul, 10% ont l'arbovirus seul et 23% sont co-infectés. Sur les patients présentant le symptôme de congestion nasale, 31% d'entre eux sont positifs au paludisme seul, 14% ont d'entre eux sont positifs à au moins un des arbovirus seul et 21% sont co-infectés. La figure 2.4 présente la distribution de l'âge, la pluviométrie et le nombre de jours de maladie selon les quatre catégories de maladies du jeu de données *IgM/IgG – data*. On observe que les patients infectés à l'arbovirus sont généralement plus âgés que les patients positifs au paludisme. Le nombre de jours de maladie est aussi plus grand pour les cas d'arbovirus. Les fortes fièvres ont été observées chez les patients ayant le paludisme ou la co-infection. On peut aussi voir dans la figure 2.4.(b) que les cas de paludisme et co-infection sont plus fréquents durant la saison des pluies.

2.3 Modèle Statistique

L'objectif de cette section est de proposer une méthodologie qui pourra identifier les symptômes importants pour le diagnostic de l'arbovirus en cas d'absence de confirmation par laboratoire et qui pourra aider à la prise de décision pour le traitement en cas de co-infection.

La sélection de variables est appréciée en cas d'analyse de données médicales car cela permettrait de faire le diagnostic de la maladie sur un nombre minimal de variables. En plus, réduire le nombre de variables peut également permettre d'avoir des résultats de classification plus précis. Ainsi, dans une première étape, nous présentons le modèle logistique multinomial. Dans une deuxième étape, nous sélectionnons les variables importantes qui peuvent expliquer la réponse via un modèle logistique multinomial. L'analyse statistique est difficile à cause du nombre faible de cas positifs d'arbovirus (39) par rapport au nombre total d'observations (12288). Les cas les plus importants pour l'étude sont rares et il en existe peu sur le jeu de données disponible. Ce qui nous place dans le cas de ce qu'on appelle, dans la littérature, un problème de jeu de données déséquilibré (imbalanced data). Pour résoudre ce problème, nous avons proposé deux approches de pré-traitement de données. La première est basée sur la considération biologique et l'extension des cas d'arbovirus de 39 à 633 cas (voir tableaux 2.1 et 2.2). C'est-à-dire, en considérant comme cas d'arbovirus les individus positifs à l'IgM arbovirus ou à l'IgG arbovirus. On obtient dans ce cas le jeu de données *IgM/IgG – data* décrit dans la section précédente. La deuxième consiste à appliquer une technique de sous-échantillonnage, en retirant aléatoirement des observations dans la classe majoritaire, afin d'équilibrer les données. Comme la distribution des données est changée, on s'attend à ce que les modèles ajustés soient biaisés par rapport aux objectifs. Dans une troisième étape, nous étudions la robustesse de la sélection de variables en utilisant les forêts aléatoires. Nous quantifions l'effet des variables sélectionnées, dans une quatrième étape, en utilisant les rapports de cotes (odds ratios). Nous calculons les odds ratio d'une variable pour une catégorie de maladie par rapport à une autre.

2.3.1 La régression logistique multinomiale

La régression logistique multinomiale est une généralisation de la régression logistique binaire (où l'on n'a que 2 modalités), où la variable d'intérêt prend K modalités ($k = 0, \dots, K - 1, K > 2$). Elle est *ordinaire* s'il y a une relation d'ordre sur les modalités et *nominale* dans le cas contraire.

Pour $k = 0, \dots, K - 1$, la probabilité d'appartenance à une modalité k pour un individu de covariable x est donnée par :

$$\pi_k(x) = \mathbb{P}(Y = k | X = x) \quad \text{avec} \quad \sum_{k=0}^{K-1} \pi_k(x) = 1 \quad (2.1)$$

- **Modalité de référence**

Dans le cas de la régression logistique binomiale, la modalité de référence (l'échec : $Y = 0$) est utilisée pour définir le logit de la réponse $Y = 1$ contre la réponse $Y = 0$. C'est-à-dire, définir

$$\log \frac{\mathbb{P}(Y = 1 | X = x)}{\mathbb{P}(Y = 0 | X = x)} = \langle x, \beta \rangle.$$

Dans le cas de la régression logistique multinomiale, il est nécessaire de choisir une modalité de référence pour définir les $K - 1$ logits. Supposons ici que la modalité de référence est $k = 0$, alors il existe $\beta_1, \beta_2, \dots, \beta_{K-1} \in \mathbb{R}^{p+1}$ tels que, pour tout $k = 1, \dots, K - 1$ et tout vecteur de covariable x ,

$$\log \left(\frac{\pi_k(x)}{\pi_0(x)} \right) = \langle x, \beta_k \rangle \quad (2.2)$$

avec

$$\langle x, \beta_k \rangle = \sum_{j=0}^p x_j \beta_{jk}$$

et $x_0 = 1$ pour inclure le paramètre d'intercept β_{0k} , $k = 1, \dots, K - 1$. Par conséquent, les probabilités $\pi_k(x)$, pour $k = 1, \dots, K - 1$ et x covariable, sont données par :

$$\pi_k(x) = \frac{e^{\langle x, \beta_k \rangle}}{1 + \sum_{j=1}^{K-1} e^{\langle x, \beta_j \rangle}}, \quad (2.3)$$

et

$$\pi_0(x) = 1 - \sum_{j=1}^{K-1} \pi_k(x) = \frac{1}{1 + \sum_{j=1}^{K-1} e^{\langle x, \beta_j \rangle}} \quad (2.4)$$

Il est important de préciser ici que $\beta = [\beta_1, \dots, \beta_{K-1}]$ est une matrice de $\mathbb{R}^{(p+1) \times (K-1)}$ dont les colonnes sont les β_k . Les éléments de la matrice β sont notés β_{kl} pour la k -ième ligne et la l -ième colonne.

- **La vraisemblance du modèle**

Soit $(X_n, Y_n) = ((x_1, y_1), \dots, (x_n, y_n))$ un échantillon indépendant et identiquement distribué de même loi que (X, Y) . La vraisemblance du modèle est donnée, dans ce cas, par

$$L(\beta) = \prod_{i=1}^n \left[\prod_{k=0}^{K-1} (\pi_k(x_i))^{y_k(i)} \right] \quad (2.5)$$

avec

$$y_k(i) = \begin{cases} 1 & \text{si } y_i = k \\ 0 & \text{si non} \end{cases}$$

Par suite, la Log-vraisemblance est donnée par

$$LL(\beta) = \sum_{i=1}^n \sum_{k=0}^{K-1} y_k(i) \log[\pi_k(x_i)] \quad (2.6)$$

- **Estimation des paramètres**

L'estimateur $\hat{\beta}$ du paramètre β , contenant tous les éléments des vecteurs β_k , est obtenu par maximum de vraisemblance. Mais il n'existe pas de solution analytique. On utilise l'algorithme de *Fisher Scoring* (Newton-Raphson) pour maximiser la log-vraisemblance du modèle. Cet algorithme est basé sur le score et la matrice Hessienne.

→ **Le score** $\nabla LL(\beta)$

Le score $\nabla LL(\beta)$ est défini comme la dérivée de $LL(\beta)$ par rapport à β :

$$\nabla LL(\beta) = \begin{pmatrix} \nabla_1 LL(\beta) \\ \vdots \\ \nabla_{K-1} LL(\beta) \end{pmatrix} \text{ avec } \nabla_k LL(\beta) = \begin{pmatrix} \nabla_{k0} LL(\beta) \\ \nabla_{k1} LL(\beta) \\ \vdots \\ \nabla_{kp} LL(\beta) \end{pmatrix} \quad (2.7)$$

Il est défini par bloc ; et chaque bloc $\nabla_k LL(\beta)$ est relatif à une modalité k . Les $\nabla_{kl} LL(\beta)$ sont les dérivées par rapport aux β_{kl} :

$$\nabla_{kl} LL(\beta) = \frac{\partial LL(\beta)}{\partial \beta_{kl}} = \sum_{i=1}^n X_{il} [y_k(i) - \pi_k(x_i)] \quad (2.8)$$

avec X_{il} les coordonnées de la matrice X associés à la i -ème ligne et à la l -ième colonne.

→ **La matrice hessienne** $H(\beta)$

La matrice hessienne est la matrice des dérivées secondes de la log-vraisemblance. Elle est définie comme suit :

$$H(\beta) = \begin{pmatrix} H_{11}(\beta) & H_{12}(\beta) & \cdots & H_{1,K-1}(\beta) \\ H_{21}(\beta) & H_{22}(\beta) & \cdots & H_{2,K-1}(\beta) \\ \vdots & \vdots & \cdots & \vdots \\ H_{K-1,1}(\beta) & H_{K-1,2}(\beta) & \cdots & H_{K-1,K-1}(\beta) \end{pmatrix} \quad (2.9)$$

Avec

$$H_{ij|_{i \neq j}}(\beta) = \begin{pmatrix} \nabla_{i0} \nabla_{j0} LL(\beta) & \nabla_{i0} \nabla_{j1} LL(\beta) & \cdots & \nabla_{i0} \nabla_{jp} LL(\beta) \\ \nabla_{i1} \nabla_{j0} LL(\beta) & \nabla_{i1} \nabla_{j1} LL(\beta) & \cdots & \nabla_{i1} \nabla_{jp} LL(\beta) \\ \vdots & \vdots & \cdots & \vdots \\ \nabla_{ip} \nabla_{j0} LL(\beta) & \nabla_{ip} \nabla_{j1} LL(\beta) & \cdots & \nabla_{ip} \nabla_{jp} LL(\beta) \end{pmatrix} \quad (2.10)$$

et

$$H_{ii} = \begin{pmatrix} \nabla_{i0}^2 LL(\beta) & \nabla_{i0} \nabla_{i1} LL(\beta) & \cdots & \nabla_{i0} \nabla_{ip} LL(\beta) \\ \nabla_{i1} \nabla_{i0} LL(\beta) & \nabla_{i1}^2 LL(\beta) & \cdots & \nabla_{i1} \nabla_{ip} LL(\beta) \\ \vdots & \vdots & \cdots & \vdots \\ \nabla_{ip} \nabla_{i0} LL(\beta) & \nabla_{ip} \nabla_{i1} LL(\beta) & \cdots & \nabla_{ip}^2 LL(\beta) \end{pmatrix} \quad (2.11)$$

Les dérivées secondes sont définies par :

$$\nabla_{kl}^2 LL(\beta) = \frac{\partial^2 LL(\beta)}{\partial \beta_{kl}^2} = \sum_{i=1}^n X_{il}^2 \pi_k(x_i) [\pi_k(x_i) - 1] \quad (2.12)$$

et

$$\nabla_{kl} \nabla_{k'l'} LL(\beta) = \begin{cases} \sum_{i=1}^n X_{il'} X_{il} \pi_k(x_i) [\pi_k(x_i) - 1] & \text{si } k = k' \\ \sum_{i=1}^n X_{il'} X_{il} \pi_k(x_i) \pi_{k'}(x_i) & \text{si } k \neq k' \end{cases} \quad (2.13)$$

A partir du score $\nabla LL(\beta)$ et de la matrice hessienne $H(\beta)$, on peut estimer β par l'algorithme suivant : pour chaque itération ($m+1$),

$$\beta_{m+1} = \beta_m - H(\beta_m)^{-1} \nabla LL(\beta_m). \quad (2.14)$$

Pour plus de détails sur cette partie, le lecteur pourra consulter ([18, 16]).

2.3.2 La régression logistique multinomiale dans le cas Arbovirus-Paludisme

Soit $X_n = (X_1, \dots, X_p) \in \mathbb{R}^{n \times p}$ la matrice des p vecteurs de covariables. Pour chaque individu i , posons $x_i = (1, X_{i1}, \dots, X_{ip})$ le vecteur de covariable (avec l'intercept) associé et y_i la réponse, $i = 1, \dots, n$. Les modalités de notre variable d'intérêt sont définies dans le tableau ci-dessous :

k (Modalité)	P (Paludisme)	A (Arboviruse)	y (Réponse)
0	0	0	(O) “Autres maladies”
1	0	1	(A) “Infection seule aux arbovirus”
2	1	0	(M) “Paludisme seul”
3	1	1	(C) “co-infection”

TABLE 2.4 – Modalités de la variable réponse

Pour un individu avec un vecteur de covariable x_i , on cherche à prédire la probabilité d'appartenance à la modalité k sachant x_i :

$$\pi_k(x_i) = \mathbb{P}(y_i = k | x_i), \quad k = 1, 2, 3, \quad i = 1, \dots, n. \quad (2.15)$$

En estimant les $\hat{\beta}_k$ par maximum de vraisemblance, nous pourrions estimer les $\pi_k(i)$ par

$$\hat{\pi}_k(x_i) = \frac{e^{\langle x_i, \hat{\beta}_k \rangle}}{1 + \sum_{j=1}^3 e^{\langle x_i, \hat{\beta}_j \rangle}}. \quad (2.16)$$

2.3.3 Sélection de variables

En utilisant le modèle logistique multinomial, on peut sélectionner les variables importantes pour expliquer la réponse y . Pour ce faire, la méthode du test du rapport de vraisemblance est souvent utilisée pour comparer les modèles deux à deux.

Par test du rapport de vraisemblance

La sélection de variables peut être faite via un test du rapport de vraisemblance, qui est une méthode paramétrique dépendant du modèle multinomial défini à la section 2.3.1. Le test du rapport de vraisemblance est un test qui permet d'évaluer l'apport de variables explicatives (covariables) supplémentaires dans l'ajustement du modèle. Le principe consiste à comparer les vraisemblances de deux modèles emboîtés \mathcal{M}_1 et \mathcal{M}_2 (Hosmer et al.[18]). Supposons que \mathcal{M}_1 comporte p_1 covariables et que \mathcal{M}_2 comporte p_2 covariables (avec $p_2 < p_1$). De plus toutes les variables du modèle \mathcal{M}_2 sont dans le modèle \mathcal{M}_1 . On cherche à tester l'hypothèse :

$$H_0 : \text{le modèle } \mathcal{M}_2 \text{ est adéquat vs } H_1 : \text{le modèle } \mathcal{M}_1 \text{ est adéquat}$$

Pour ce faire, on définit la statistique de test suivant :

$$TRV = -2 \log \left(\frac{L(\mathcal{M}_2)}{L(\mathcal{M}_1)} \right) \quad (2.17)$$

où $L(\mathcal{M}_r)$ est le log du maximum de vraisemblance dans le modèle \mathcal{M}_r , $r = 1, 2$. Sous l'hypothèse H_0 la statistique TRV converge en loi vers une loi de khi-deux $\mathcal{X}^2(p_1 - p_2)$. Le modèle \mathcal{M}_r , $r = 1, 2$, comporte $p_r + 1$ paramètres à estimer.

Ce test est aussi appelé le test de la différence de déviance. En effet, la déviance du modèle \mathcal{M}_r est définie par :

$$D_r = 2 \log L(\mathcal{M}_{sat}) - 2 \log L(\mathcal{M}_r) \quad (2.18)$$

avec \mathcal{M}_{sat} : le modèle saturé (contenant toutes les variables). La statistique de test TRV peut être vue comme différence de déviances des deux modèles :

$$\begin{aligned} D_2 - D_1 &= 2 \log L(\mathcal{M}_{sat}) - 2 \log L(\mathcal{M}_2) - 2 \log L(\mathcal{M}_{sat}) + 2 \log L(\mathcal{M}_1) \\ &= -2 \log L(\mathcal{M}_2) - (-2 \log L(\mathcal{M}_1)) \\ &= TRV \end{aligned}$$

Pour tester la pertinence d'une variable X_j , $j \in \{1, \dots, p\}$, on considère deux modèles : l'un (\mathcal{M}_1) contenant certaines variables en plus de X_j et l'autre (\mathcal{M}_2) contenant les mêmes variables mais sans X_j . Si H_0 est rejetée, alors la variable X_j est pertinente. Si non la variable X_j n'est pas pertinente. Cette procédure peut être faite dans les deux directions : direction ascendante ("Forward direction"), où on commence par le modèle nul (ne contenant aucune variable) en ajoutant les variables une à une et en faisant à chaque étape le test ; la direction descendante ("Backward direction"), où on commence par le modèle saturé (contenant toutes les variables) en enlevant les variables une à une tout en testant la significativité de la variable enlevée à chaque étape. Si la variable est significative, on la garde et on retire une autre, ainsi de suite.

— **Application au jeu de données *IgM/IgG - data***

Nous avons effectué cette méthode de sélection de variable sur le jeu de données *IgM/IgG - data*, d'abord dans la direction ascendante, c'est-à-dire en partant du modèle ne contenant pas de variable, en ajoutant les variables une à une et en testant à chaque étape la significativité de la variable ajoutée. Si cette dernière est significative, on la garde si non on l'enlève.

La même procédure est appliquée dans la direction descendante. C'est-à-dire en partant du modèle contenant les 15 variables du jeu de données, on enlève les variables une à une en testant à chaque fois si la variable enlevée est significative. Si c'est le cas, on la garde si non on l'enlève pour de bon. Le tableau 2.5 nous donne un résumé avec les p -values associées à chaque variable.

On voit à partir du tableau 2.5 que si on prend la direction ascendante, au niveau $\alpha = 0.05$, on sélectionne toutes les variables sauf *Sex*, *Chills* et *Jaudice*. C'est-à-dire, on retient le modèle composé de douze variables. Si on prend la direction descendante, on sélectionne huit variables (*Age*, *Temperature*, *Number of sick day*, *Pluviometry*, *Joint pain*, *Nausea/vomiting*, *Cough* et *Nasal congestion*). Suivant les directions, on obtient alors des modèles différents. De plus, suivant l'ordre d'ajout ou de retrait des variables, on obtient aussi d'autres modèles. Ce qui laisse penser que la sélection de variables par test du rapport de vraisemblance n'est pas assez robuste pour sélectionner les variables importantes de notre étude.

<i>Variable</i> \ <i>Direction</i>	<i>Ascendante</i>	<i>Descendante</i>
Age	$< 10^{-16}$	$< 10^{-16}$
Temperature	$2.19.10^{-8}$	$2.30.10^{-9}$
Number of sick days	$3.28.10^{-12}$	$5.25.10^{-12}$
Pluviometry	$< 10^{-16}$	$< 10^{-16}$
Sex	0.67	0.68
Cephalalgia	0.04	0.13
Eye pain	0.01	0.27
Muscle pain	0.02	0.47
Joint pain	0.04	0.01
Nause/vomiting	$1.7.10^{-13}$	$6.2.10^{-13}$
Diarrhea	0.04	0.22
Chills	0.57	0.76
Cough	$7.43.10^{-7}$	$4.38.10^{-6}$
Nasal congestion	$< 10^{-16}$	$< 10^{-16}$
Jaudice	0.74	0.88

TABLE 2.5 – Sélection de variable par test du rapport de vraisemblance : p-values associées à chaque variable dans les deux directions

Par stepwise

Comme expliqué à la section précédente, le test du rapport de vraisemblance est souvent utilisé pour comparer des modèles deux à deux. Il ne s'applique qu'à des modèles emboîtés (dérivant l'un de l'autre par ajout ou suppression de variable). Quand de nombreux modèles doivent être comparés entre eux, le risque de rejeter l'hypothèse nulle alors qu'elle est vraie augmente. Pour résoudre ce problème, une solution possible consiste à comparer les modèles en utilisant le critère d'information d'Akaike (Akaike, 1974 [59]), défini par

$$AIC = -2 \log L(\mathcal{M}) + 2k \quad (2.19)$$

avec $L(\mathcal{M})$ la vraisemblance et k le nombre de paramètres du modèle \mathcal{M} . C'est-à-dire la déviance du modèle est pénalisée par deux fois le nombre de paramètres.

La sélection des variables importantes se fait ici par étape (stepwise selection). Soient \mathcal{M}_0 le modèle nul, ne contenant aucune variable et \mathcal{M}_p le modèle saturé contenant les p variables.

1. Dans la direction ascendante (c'est-à-dire, en partant du modèle \mathcal{M}_0), à chaque étape k , $k = 0, 1, \dots, p - 1$, on considère $p - k$ modèles. Ces modèles sont obtenus en ajoutant à chaque fois une des $p - k$ variables dans le modèle \mathcal{M}_k . On choisit ensuite le meilleur des $p - k$ modèles en terme de déviance (modèle avec la plus petite déviance) et on l'appelle \mathcal{M}_{k+1} . En répétant la procédure sur les p étapes, on a $p + 1$ modèles, qui sont $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$. Le critère AIC est alors calculé pour chacun de ces modèles et on sélectionne le modèle \mathcal{M}_\star ayant le plus petit AIC. Dans ce cas, les variables sélectionnées sont les variables qui constituent le modèle \mathcal{M}_\star . La sélection de variables dans la direction ascendante est donnée par l'algorithme 1, tableau 2.6.

Algorithme 1 : Direction ascendante (Forward direction)

1. Soit \mathcal{M}_0 le modèle nul, ne contenant aucune variable.
 2. Pour $k = 0, 1, \dots, p - 1$:
 - (a) Considérer $p - k$ modèles qui sont obtenus en ajoutant à chaque fois une des $p - k$ variables dans le modèles \mathcal{M}_k
 - (b) Choisir le meilleur des $p - k$ modèles en terme de déviance en l'appelant \mathcal{M}_{k+1} .
 3. Sélectionner le meilleur modèle parmi les $p + 1$ modèles $\mathcal{M}_0, \mathcal{M}_1 \dots, \mathcal{M}_p$ en terme d'AIC.
-
-

TABLE 2.6 – Algorithme de la sélection de variables par critère AIC : direction ascendante

2. Dans la direction descendante (c'est-à-dire, en partant du modèle saturé \mathcal{M}_p), à chaque étape k , $k = p, p - 1, \dots, 1$, on considère k modèles. Ces modèles sont obtenus en retirant à chaque fois une des k variables dans le modèle \mathcal{M}_k . On choisit ensuite le meilleur des k modèles en terme de déviance (modèle avec la plus petite déviance) et on l'appelle \mathcal{M}_{k-1} . En répétant la procédure sur les p étapes, on a $p + 1$ modèles, qui sont $\mathcal{M}_0, \mathcal{M}_1 \dots, \mathcal{M}_p$. Le critère AIC est alors calculé pour chacun de ces modèles et on sélectionne le modèle \mathcal{M}_* ayant le plus petit AIC. Dans ce cas, les variables sélectionnées sont les variables qui constituent le modèle \mathcal{M}_* . La sélection de variables dans la direction descendante est donnée par l'algorithme 2, tableau 2.7.

Algorithme 2 : Direction descendante (Backward direction)

1. Soit \mathcal{M}_{sat} le modèle saturé, contenant toutes les p variables.
 2. Pour $k = p, p - 1, \dots, 1$:
 - (a) Considérer k modèles qui sont obtenus en retirant à chaque fois une des k variables dans le modèles \mathcal{M}_k
 - (b) Choisir le meilleur des k modèles en terme de déviance en l'appelant \mathcal{M}_{k-1} .
 3. Sélectionner le meilleur modèle parmi les $p + 1$ modèles $\mathcal{M}_0, \mathcal{M}_1 \dots, \mathcal{M}_p$ en terme d'AIC.
-
-

TABLE 2.7 – Algorithme de la sélection de variables par critère AIC : direction descendante

Cette méthode de sélection de variable peut être appliquée sur les deux différents jeux de données. Pour chaque jeu de données, nous combinons les deux directions (Backward et Forward). Pour plus de détails sur la sélection de variable par stepwise, le lecteur pourra consulter le livre de James et al. [60].

a. **Application au jeu de données $IgM/IgG - data$**

Comme le jeu de données $IgM/IgG - data$ est assez équilibré, la sélection de variable, dans ce cas, est juste basée sur l'analyse standard du modèle logistique multinomial. Ainsi, nous commençons par donner les résultats pour ce jeu de données. Le modèle logistique multinomial a été ajusté sur les données $IgM/IgG - data$ en utilisant la fonction `multinom` ou la fonction `vglm` des packages R `nnet` et `VGAM`. La sélection par stepwise est faite suivant les deux directions et le même modèle a été retenu. Ce modèle est composé des variables suivantes : *age*, *temperature*, *number*

of sick days, rainfall, nausea or vomiting, cough, nasal congestion and joint pain. Dans le but de tester la robustesse de la sélection de variable par stepwise, on utilise d’abord, le test de rapport de vraisemblance entre le modèle contenant les huit variables sélectionnées et les sous-modèles obtenus en enlevant une des variables à chaque fois. On obtient que tous les variables sélectionnées sont significatives. Toutes les p -values obtenues sont inférieures à 10^{-9} sauf pour la variable *joint pain* qui sort avec une p -value égale à $7.44 \cdot 10^{-3}$.

Ensuite, on effectue le test du rapport de vraisemblance entre le modèle contenant les huit variables sélectionnées et les modèles obtenus en ajoutant les variables non retenues une à une. On obtient que toutes les variables non sélectionnées par stepwise ne sont pas significatives avec des p -values supérieures à 5% (voir Tableau 2.8).

Variables	p-values
Sexe	0.78
Céphalée	0.18
Diarrhée	0.2
Frissons	0.65
Douleur aux yeux	0.07
Douleur musculaire	0.48
Ictère/Jaunisse	0.76

TABLE 2.8 – Jeu de données *IgM/IgG*. Pvalues du test du rapport de vraisemblance associé aux variables non retenues par stepwise.

b. **Application au jeu de données *IgM – data***

Le jeu de données *IgM – data* contient 18 cas de mono-infection aux arbovirus, 21 cas de co-infection, 7069 cas de mono-infection au paludisme et 5180 cas d’autres maladies. En utilisant ce jeu de données, le modèle logistique multinomial ajusté ne prédit que les classes 0 et 2. Ce qui signifie qu’il ignore les deux classes minoritaires 1 et 3 aux dépens des classes majoritaires.

Pour résoudre ce problème, nous proposons d’utiliser une stratégie de ré-échantillonnage (voir [61] pour un aperçu sur les méthodes existantes). Deux des approches les plus simples sont le sous-échantillonnage et le sur-échantillonnage. Comme le jeu de données *IgM – data* est très déséquilibré avec un très grand nombre d’observations pour les deux classes majoritaires, nous utilisons une stratégie de sous-échantillonnage qui enlève des observations dans les classes majoritaires et réduit la taille de l’échantillon. Le sous-échantillonnage se fait par tirage sans remise de 50 cas dans chacune des classes majoritaires (classes 0 et 2) pour construire un échantillon de taille $18 + 21 + 50 + 50 = 139$. En faisant l’apprentissage sur ce jeu de données, le modèle prédit quatre classes.

Dans cette stratégie de sous-échantillonnage, on risque de perdre des informations par le fait d’enlever des observations dans les classes majoritaires. Pour éviter ce problème, on répète 1000 fois l’étape de construction d’un sous-échantillon de taille 139 et on utilise les 1000 sous-échantillons équilibrés pour faire la sélection de variables sur le jeu de données *IgM – data*. Le modèle multinomial a été ajusté pour chaque sous-échantillon, et la sélection de variables par stepwise a

été appliquée. Le pourcentage d'apparition (ranking) de chaque variable sur les 1000 sélections de variables effectuées est calculé. Ce "ranking" est présenté à la figure 2.6. Contrairement à ce qu'on avait observé sur le jeu de données *IgM/IgG-data*, les variables *Nasal congestion* et *Number of sick days* font parties des moins importantes et que la variable *Eye pain* est l'une des plus importantes. On observe ainsi une certaine variabilité qui est, peut être, due à la perte d'informations liée à la procédure de sous-échantillonnage.

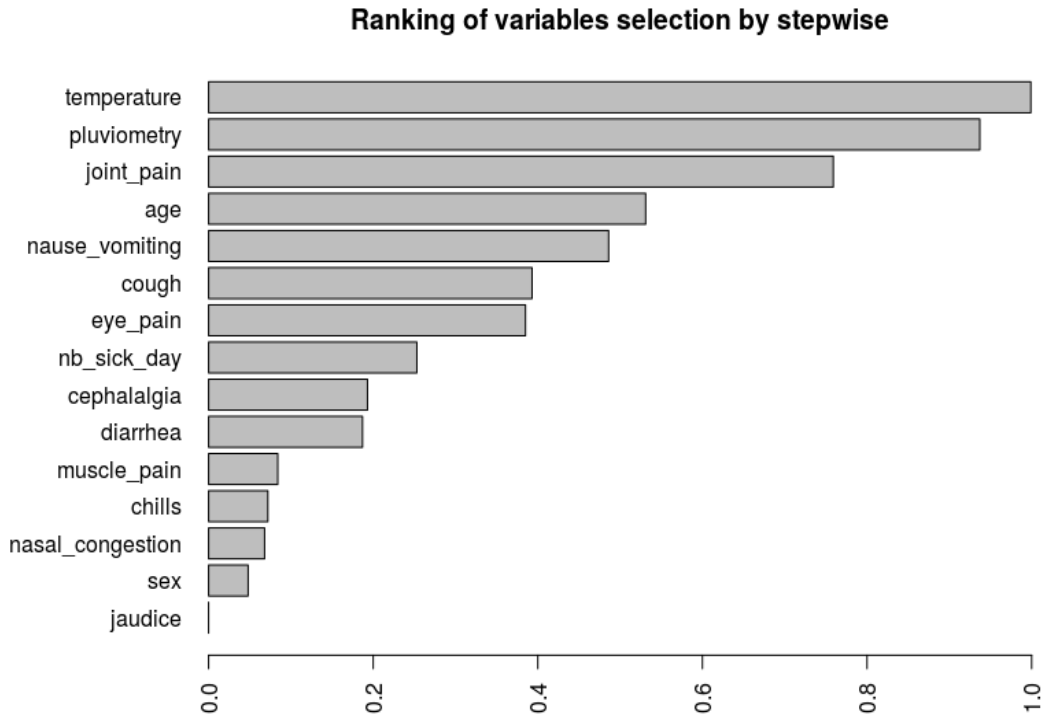


FIGURE 2.6 – *IgM – data*. Ranking des variables pour la sélection de variables par stepwise sur les 1000 sous-échantillons.

Cette variabilité observée dans les 1000 sélections de variables soulève des questions de robustesse et pour répondre à ces questions, nous proposons une analyse non-paramétrique basée sur l'algorithme des forêts aléatoires. Ces dernières années, des méthodes combinant le ré-échantillonnage et l'apprentissage sont de plus en plus utilisées ([61]). Nous trouvons, ici, que la mesure d'importance basée sur des forêts aléatoires peut être un outil pratique pour résumer les informations obtenues pour la sélection de variables sur les 1000 sous-échantillons. En utilisant la mesure d'importance fournie par les forêts aléatoires, on arrive à étudier la robustesse de la sélection de variables sur les deux jeux de données. On verra que la sélection de variables reste stable et qu'on aura le même modèle sur les deux jeux de données.

Par forêts aléatoires

La méthode des forêts aléatoires est une méthode statistique non-paramétrique proposée par Léo Breiman en 2001 [62]. C'est une méthode basée sur des arbres de décision

obtenus par la méthode CART. Elle s'avère très performante dans de nombreuses applications.

a. ***Arbre de décision (CART)***

Il existe plusieurs algorithmes de construction des arbres aléatoires dont le plus connu est l'algorithme CART (Classification And Regression Trees). L'algorithme CART est une méthode introduite par Breiman et al. ([63]) pour construire des prédicteurs par arbre en régression et en classification. L'idée est de partitionner l'ensemble des données d'entrées (X) en deux sous-parties, puis de déterminer une sous partie pour la prédiction. A chaque étape de l'algorithme, on découpe une partie en deux sous-parties. La racine de l'arbre est constituée de l'ensemble des données et les fils correspondent aux deux sous parties. Supposons que l'on a p variables quantitatives. En partant de la racine, on veut découper de la manière suivante : pour un $j \in \{1, \dots, p\}$ et un $d \in \mathbb{R}$,

$$\{X^j \leq d\} \cup \{X^j > d\}.$$

C'est-à-dire : les observations avec la valeur de la j -ième variable plus petite que d sont dans le nœud fils gauche, et celles avec une valeur plus grande que d sont dans le nœud fils droite. On cherche le meilleur couple (j, d) qui minimise un critère :

— En régression, on cherche à minimiser la variance des nœuds. Pour un nœud t , on minimise donc

$$\sum_{i: X_i \in t} (Y_i - \bar{Y}_t)^2, \text{ où } \bar{Y}_t \text{ est la moyenne des observations au nœud } t.$$

— En classification, on cherche à minimiser l'indice de Gini des nœuds. Pour un nœud t , on minimise donc

$$\sum_{c=1}^L \hat{p}_t^c (1 - \hat{p}_t^c),$$

avec \hat{p}_t^c proportion de la classe c au nœud t .

Après le découpage de la racine, la même procédure est répétée sur les nœuds fils et ainsi de suite, jusqu'à ce que chaque nœud ne contienne qu'un seul élément ou des observations de même classe. L'arbre ainsi construit est appelé arbre maximal. A chaque feuille est associée une prédiction définie par la classe majoritaire des observations qu'elle contient.

La deuxième étape de la construction de l'arbre CART est l'élagage, qui consiste à trouver le meilleur sous-arbre. C'est en fait une sorte de sélection de modèles (où les modèles sont les sous-arbres élagués de l'arbre maximal). Cette étape n'est pas utilisée dans la construction de la forêt. En effet, l'objectif de cette étape est de diminuer la variance. Alors que dans la construction de la forêt, cet objectif est atteint grâce à l'agrégation.

b. ***Forêts aléatoires***

Le principe des forêts aléatoires est de construire des arbres CART (q arbres) à partir de q échantillons bootstrap $\mathcal{L}_n^1, \dots, \mathcal{L}_n^q$. Pour chaque arbre, le découpage d'un nœud se fait de la manière suivante : on tire aléatoirement m variables et la meilleure coupure est cherchée sur les m variables sélectionnées. Le prédicteur est ainsi obtenu en agrégeant la collection d'arbres construite (c'est-à-dire la moyenne

en régression et le vote majoritaire en classification). A chaque nœud, le tirage des m variables est sans remise et uniforme parmi toutes les variables. Le nombre de variables $m \leq p$ (*mtry*) est très important. Il est fixé au début et est identique pour tous les arbres. Dans le paquet *randomForest* de R, il a comme valeur par défaut $m = \sqrt{p}$ en classification et $m = p/3$ en régression. Dans cette thèse, on est bien dans le cas d'une classification, ce qui fait que la valeur par défaut de m est égale à \sqrt{p} . L'autre paramètre important est le nombre d'arbres q de la forêt, noté *nntree* dans le paquet *randomForest* de R avec une valeur par défaut de $q = 500$. Avec un bon choix des paramètres m et q , on peut faire la sélection de variables par forêts aléatoires en se basant sur la mesure de l'importance des variables notée *MDA* (Mean Decrease Accuracy). La *MDA* est calculée à partir d'une mesure appelée erreur OOB (*Out-of-Bag*) qui signifie ici "en dehors du bootstrap".

•L'erreur OOB

Soit (X_i, Y_i) une observation de l'échantillon \mathcal{L}_n . Pour fabriquer une prédiction \hat{Y}_i de Y_i , l'agrégation se fait uniquement sur les prédictions des arbres construites sur les échantillons bootstrap ne contenant pas cette observation. Cette procédure est répétée sur toutes les observations afin de calculer l'erreur OOB notée *errOOB*. Cette erreur est donnée par l'erreur quadratique moyenne en régression

$$errOOB = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2, \quad (2.20)$$

et le taux de mal classés en classification

$$errOOB = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\hat{Y}_i \neq Y_i\}}. \quad (2.21)$$

•La mesure d'importance MDA

Pour $j \in \{1, \dots, p\}$ fixé, on cherche à calculer l'importance de la variable X_j . Soient \mathcal{L}_n^l un échantillon bootstrap et $echOOB_l$ l'échantillon OOB associé, c'est-à-dire l'ensemble des observations qui n'apparaissent pas dans \mathcal{L}_n^l . Notons par $errOOB_l$ l'erreur OOB commise sur l'échantillon $echOOB_l$ par l'arbre construit sur \mathcal{L}_n^l . Si on permute les valeurs de la j -ième variable dans l'échantillon $echOOB_l$, on obtient un échantillon perturbé $echOOB_l^j$. Ainsi on peut calculer l'erreur $errOOB_l^j$, commise sur l'échantillon $echOOB_l^j$. Cette procédure est répétée sur tous les échantillons bootstrap. Par suite l'importance de variable X_j , $MDA(X_j)$ est donnée par la différence entre l'erreur moyenne d'un arbre sur l'échantillon OOB perturbé et celle sur l'échantillon OOB :

$$MDA(X_j) = \frac{1}{q} \sum_{l=1}^q (errOOB_l^j - errOOB_l). \quad (2.22)$$

En utilisant le paquet **randomForest** de R, où le calcul est fait comme décrit ci-dessus, nous pourrions classer les variables par ordre d'importance. Rappelons que plus cette importance (MDA) est grande plus la variable est importante.

Pour plus de détails dans cette partie, le lecteur pourra consulter la thèse de Robin Genuer ([64]) ou les travaux récents de Genuer et Poggi ([65]).

Application au jeu de données *IgM/IgG – data*

Pour une bonne classification, nous avons donc besoin de faire un bon choix des paramètres `mtry` et `ntree` qui doivent être calibrés sur le jeu de données. La figure 2.7 nous donne, pour différentes valeurs de `mtry`, l'erreur OOB en fonction du nombre d'arbres (`ntree`). On peut voir sur cette figure qu'on peut choisir comme paramètres `mtry= 5` et `ntree= 500`. En effet, on cherche à choisir le couple de paramètres (`mtry`, `ntree`) à partir duquel l'erreur OOB se stabilise.

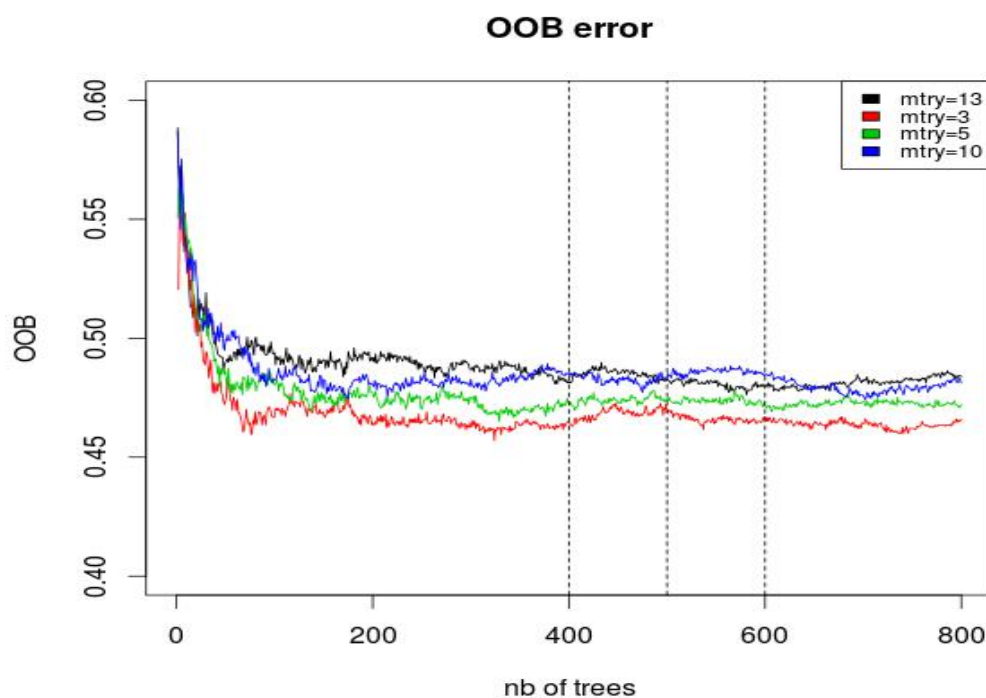


FIGURE 2.7 – Calibrage des paramètres `ntree` et `mtry` du package `randomForest`.

Avec ce choix de paramètres, nous utilisons le paquet `randomForest` (version 4.6-12) de R pour faire la sélection de variables sur le jeu de données *IgM/IgG – data*. Nous utilisons dans ce chapitre la version 3.2.2 (2015-08-14) de R. La figure 2.8 représente les 15 variables par ordre croissante d'importance.

On peut noter que la variable la plus importante pour expliquer la réponse est *pluviometry*, qui est une variable indicative pour la saison des pluies. Rappelons que le parasite du paludisme se développe beaucoup en saison des pluies. Un second groupe de variables suivantes est composé de *nasal congestion*, *age* et *number of sick days*. Ensuite viennent les variables *cough* et *temperature*. Ainsi de suite, jusqu'à la variable *eye pain* qui est la moins importante. Notons que le paquet `randomForest` donne les variables par ordre d'importance mais pas une coupure pour identifier exactement les variables qui expliquent la réponse. Pour ce faire, nous utilisons le paquet `VSURF` (Variable Se-

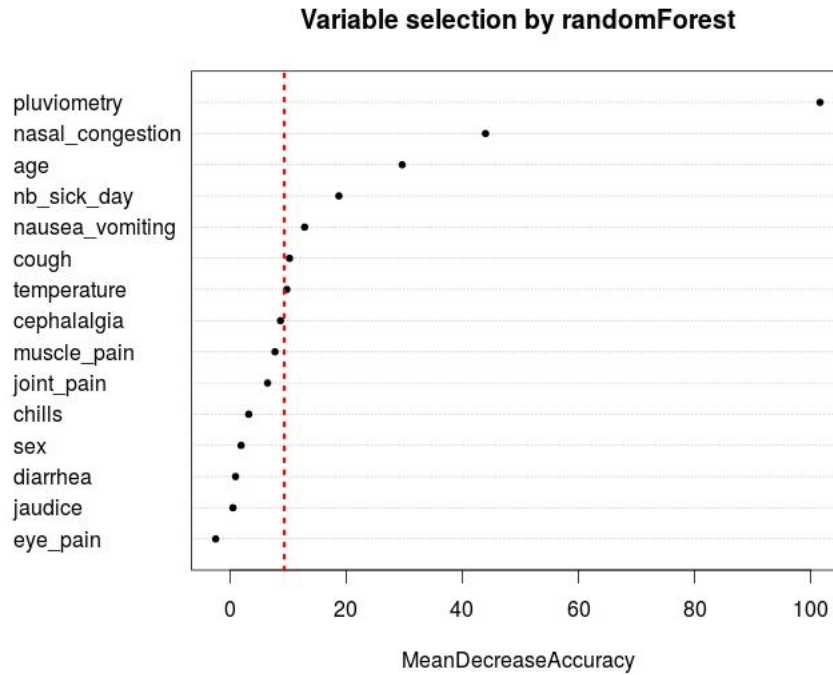


FIGURE 2.8 – *IgM/IgG – data*. graphe d’importance des variables : “mean decrease of accuracy” (MDA) des covariables par ordre croissante. Les variables avec des *MDA* au delà de la ligne verticale sont sélectionnées par la procédure VSURF.

lection Using Random Forest, version 1.0.2) lié à `randomForest` pour faire le seuillage et sélectionner les variables pertinentes ([66]). Le découpe est fait par VSURF suivant deux étapes. (1) la première est une étape d’élimination et de classement. Les variables sont classées par ordre d’importance décroissante. Les variables avec une importance plus petite qu’un seuil (la valeur minimale des écarts types des mesures d’importance) sont éliminées. Soit m_1 le nombre de variables conservées. (2) la deuxième étape est dite de sélection de variables. on construit la collection de modèles emboîtés constitués par les forêts construites sur les k premières variables, pour $k = 1, \dots, m_1$ et on sélectionne les variables du modèle conduisant à la plus faible erreur OOB. Ceci conduit à considérer m_2 variables, avec $m_2 \leq m_1$. Pour plus de détails sur le seuillage utilisé par VSURF, voir les récents travaux de Genuer et Poggi ([65]). En utilisant le seuillage de VSURF, on peut voir sur la figure 2.8 que les variables les plus pertinentes pour expliquer la réponse multinomiale y sont : *pluviometry*, *nasal congestion*, *age*, *number of sick days*, *nausea or vomiting*, *cough* et *temperature*.

En comparant avec la sélection de variables par stepwise présentée à la section 2.3.3, on obtient que les deux méthodes sélectionnent les mêmes variables à la différence de la variable *joint pain*. Nous avons testé la significativité de cette variable par un test du rapport de vraisemblance (entre le modèle contenant les variables sélectionnées par forêts aléatoires et le modèle contenant ces sept variables en plus de la variable *joint pain*), et nous l’avons trouvé significative avec une *p-value* égale à $7.44 \cdot 10^{-3}$. Par suite, le modèle final retenu sur le jeu de données *IgM/IgG – data* est composé de huit variables qui sont : *pluviometry*, *nasal congestion*, *age*, *number of sick days*, *nausea/vomiting*, *cough*,

temperature et joint pain.

Application au jeu de données *IgM – data*

La figure 2.9 classe les 15 variables par ordre d'importance sur les 1000 sous-échantillons. Comme sur le jeu de données *IgM/IgG – data*, la variable la plus importante est *pluviometry*. Un second groupe de variables moins importantes que la première, est composé de *cough*, *age* et *joint pain*. Un autre groupe, formé de *number of sick days*, *temperature*, *nausea/vomiting*, *eye pain* et *nasal congestion*, s'en suit. Enfin le groupe des 6 dernières variables les moins importantes est composé de *muscle pain*, *chills*, *cephalalgia*, *jaudice*, *diarrhea* et *sex*.

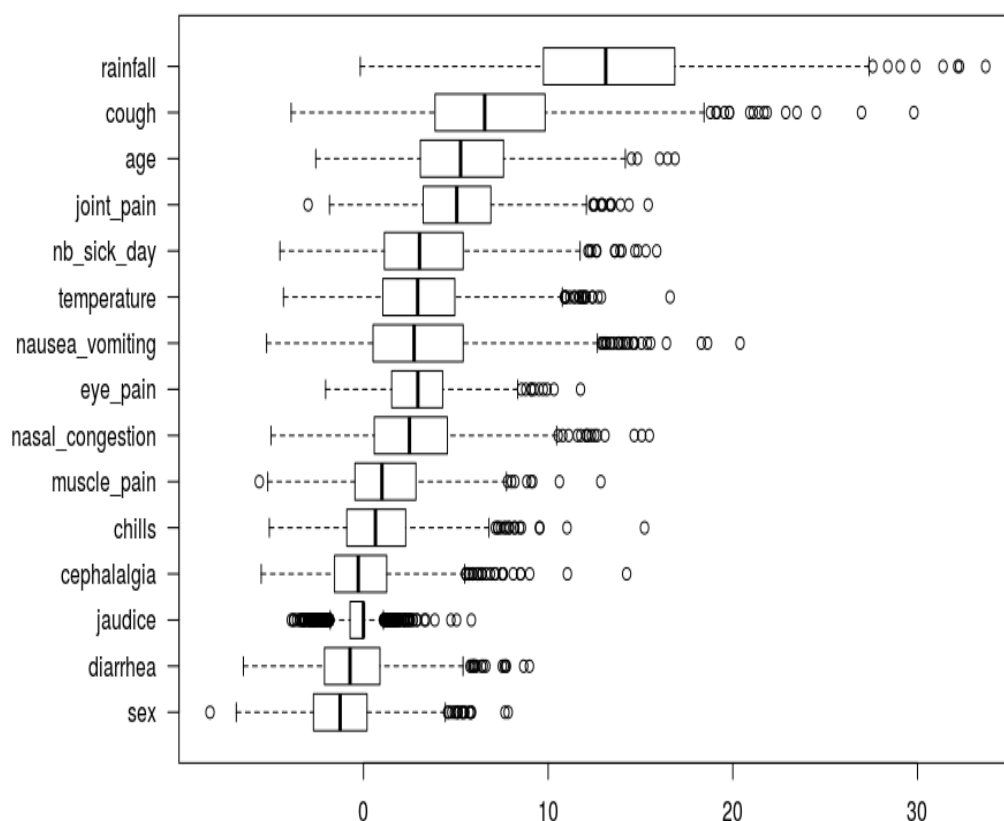


FIGURE 2.9 – *IgM – data*. Graphe d'importance des variables. Chaque boxplot résume la distribution de l'importance de la variable associée à travers les 1000 sous échantillons.

La frontière entre les deux derniers groupes de variables n'est pas claire et comme précisé lors de l'analyse du jeu de données *IgM/IgG – data*, le paquet `randomForest` ne donne pas de seuil pour séparer les variables pertinentes des variables non-pertinentes. Nous utilisons `VSURF` pour faire ce seuillage. La figure 2.10 résume le résultat de sélection de variables par `VSURF` sur les 1000 sous-échantillons. La variable *pluviometry* (92.2%) est toujours la plus importante; ensuite les variables les plus importantes sont *cough* (29.1%), *age*(28.3%), *joint pain* (19.8%), *nausea/vomiting* (16.4%), *number of sick days*

(16.1%), *temperature* (16.1%) et *nasal congestion* (11%) dans l'ordre décroissant. Les autres variables sont sélectionnées dans moins de 10% des sous-échantillons.

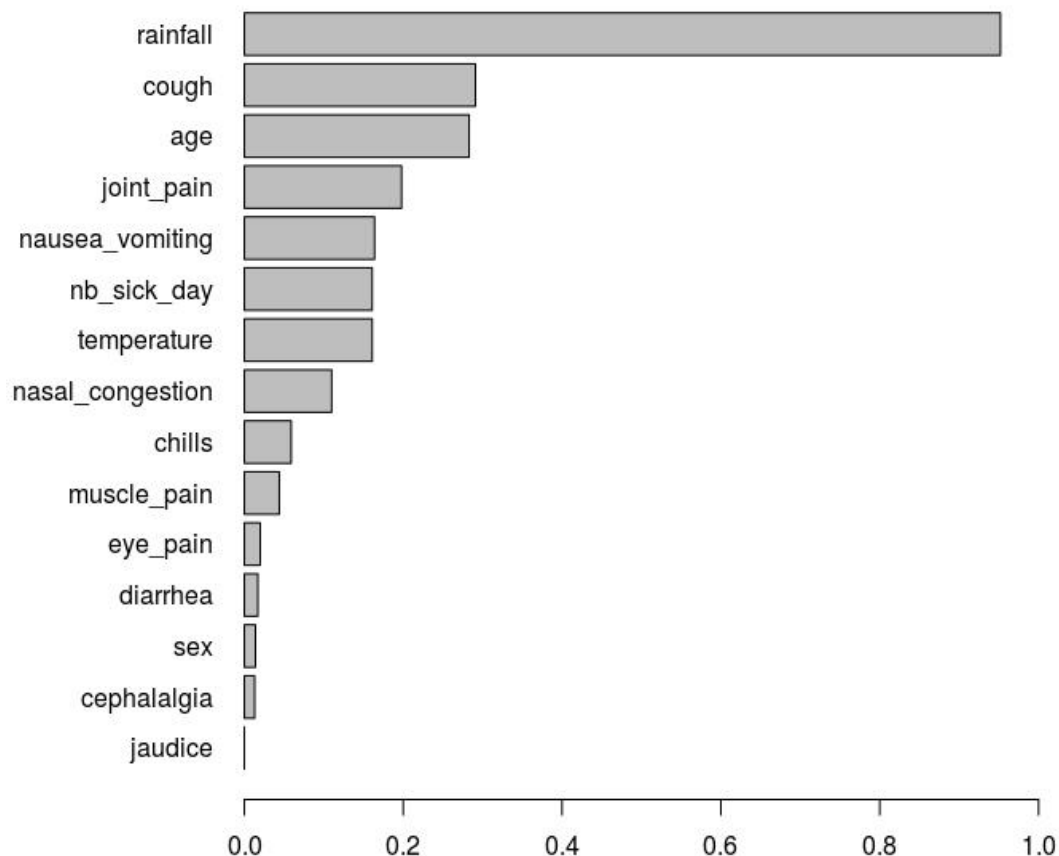


FIGURE 2.10 – *IgM – data*. Ranking par la procédure VSURF : pour chaque variable, la longueur de la barre correspond à la probabilité empirique d’être sélectionnée par la procédure VSURF à travers les 1000 *IgM* sous-échantillons.

Notons, à partir des deux figures 2.8 et 2.10, qu’on sélectionne les mêmes variables pour les deux jeux de données (*IgM – data* et *IgG/IgG – data*). Pour d’autres méthodes de classification des données déséquilibrées, le lecteur pourra consulter les travaux de Chen et al. [67] ou ceux de Chawal [68]). En pratique, nous trouvons pour notre cas de sélection de variable, que l’agrégation des résultats des 1000 sous-échantillons reste stable pour sélectionner les variables les plus pertinentes. Par suite, le modèle final retenu sur le jeu de données *IgM – data* est composé de huit variables qui sont : *pluviometry*, *nasal congestion*, *age*, *number of sick days*, *nausea/vomiting*, *cough*, *temperature* et *joint pain*.

Nous avons comparé la sélection de variables par stepwise et par forêts aléatoires sur le jeu de données *IgM/IgG – data* et nous pouvons conclure que les deux méthodes sont en accord sur les variables sélectionnées. En effet, on peut voir à partir de la Figure 2.8

que la méthode des forêts aléatoires a retenu sept des huit variables sélectionnées par stepwise sur le jeu de données *IgM/IgG – data*. Une variable supplémentaire, *joint pain* a été rajoutée par le test du rapport de vraisemblance. Sur le jeu de données *IgM – data*, les résultats des forêts aléatoires sur l'importance des variables restent assez robustes et le même groupe de huit variables à été retenu (Figures 2.9 et 2.10). En conclusion, nous décidons d'ajuster le même modèle multinomial, avec huit variables (*pluviometry, nasal congestion, age, number of sick days, nausea/vomiting, cough, temperature et joint pain.*), sur les deux jeux de données de notre analyse. C'est-à-dire, sur chaque jeu de données, on a un échantillon $(X_n, Y_n) = ((x_1, y_1), \dots, (x_n, y_n))$, avec

$$x_j \in \mathbb{R}^p, \quad j = 1, \dots, n;$$

où $p = 8$, $n = 139$ pour *IgM – data* et $n = 1976$ pour *IgM/IgG – data*.

2.3.4 Analyse des facteurs influents des différentes maladies

En utilisant le modèle final retenu à la section 2.3.3, nous analysons les facteurs influents des différentes maladies (paludisme, arbovirus seul et co-infection). Nous commençons par présenter l'analyse par odd ratio pour quantifier l'effet des covariables. Ensuite nous présentons les résultats sur les deux jeux de données. Enfin, nous donnons une petite conclusion sur cette analyse.

Odds ratio

La régression logistique multinomiale permet de quantifier l'effet d'une variable en terme d'odd ratio (OR) ou logarithme odd ratio (log OR).

Dans le cas de la régression logistique binomial (y binaire), notons par p la probabilité que $y = 1$ conditionnellement à x ,

$$p = \mathbb{P}(y = 1|x).$$

Dans ce cas, l'odd (ou "cote") associé à cette probabilité p est défini par

$$Odd_p = \frac{p}{1-p}. \quad (2.23)$$

Ce qui peut se généraliser dans le cas de la régression logistique multinomiale. Comme la modalité $k = 0$ est la référence, on peut calculer, pour les probabilités

$$\pi_k(x) = \mathbb{P}(Y = k|X = x), \quad k = 1, \dots, K-1,$$

les odds associés par

$$Odd_{\pi_k(x)} = \frac{\pi_k(x)}{\pi_0(x)} = \frac{\mathbb{P}(Y = k|X = x)}{\mathbb{P}(Y = 0|X = x)}, \quad k = 1, \dots, K-1.$$

L'Odds Ratio (ou "rapport de cote") associé à la probabilité d'appartenance à une modalité k , pour deux valeurs différentes d'une variable X_j ($X_j = a$ vs $X_j = b$) est le rapport des cotes associés à la probabilité d'avoir la maladie k entre les patients qui ont

la valeur de la covariable $X_j = a$ et ceux qui ont la valeur de la variable $X_j = b$. Il est défini par

$$OR_k(a|b) = \frac{Odds_{\pi_k(a)}}{Odds_{\pi_k(b)}} = \frac{\pi_k(a)/\pi_0(a)}{\pi_k(b)/\pi_0(b)}. \quad (2.24)$$

On suppose ici que toutes les autres covariables sont à valeurs égales sauf la variable X_j et que β_{kj} est la composante de β_k qui correspond à la variable X_j . Dans notre cas, l'odds est donné par

$$Odds_{\pi_k(a)} = \frac{\pi_k(a)}{\pi_0(a)} = e^{a\beta_{kj}}. \quad (2.25)$$

Par suite, l'odds ratio est donnée par

$$OR_k(X_j = a|X_j = b) = e^{(a-b)\beta_{kj}}. \quad (2.26)$$

On peut aussi définir l'Odds Ratio associé à la probabilité d'appartenance à une modalité k_1 , pour deux valeurs différentes d'une variable X_j ($X_j = a$ vs $X_j = b$) par rapport à une autre modalité k_2 . On a dans ce cas un rapport d'odds ratio entre les modalités k_1 et k_2 . Il est défini par

$$OR_{k_1|k_2}(a|b) = \frac{OR_{k_1}(X_j = a|X_j = b)}{OR_{k_2}(X_j = a|X_j = b)} = e^{(a-b)(\beta_{k_1j} - \beta_{k_2j})}. \quad (2.27)$$

L'OR peut être interprété comme l'effet d'un certain changement de la variable, en supposant toutes les autres variables constantes.

Interprétation des ORs par rapport à la modalité de référence

- Pour une variable quantitative
 - $OR_k(x = x_1|x = x_2) \simeq 1$, l'augmentation de $x_1 - x_2$ de la variable x est indépendante de la maladie k (elle est non significative).
 - $OR_k(x = x_1|x = x_2) > 1$, l'augmentation de $x_1 - x_2$ de la variable x est significative pour la maladie k ; elle augmente la probabilité d'avoir la maladie k .
 - $OR_k(x = x_1|x = x_2) < 1$, l'augmentation de $x_1 - x_2$ de la variable x est significative pour la maladie k ; elle diminue la probabilité d'avoir la maladie k .
- Pour une variable qualitative (Présence (1) ou Absence (0) d'un symptôme)
 - $OR_k(x = 1|x = 0) \simeq 1$, la maladie k est indépendante du symptôme x (elle n'est pas significative).
 - $OR_k(x = 1|x = 0) > 1$, la présence du symptôme x est significative pour la maladie k ; elle augmente la probabilité d'avoir la maladie k .
 - $OR_k(x = 1|x = 0) < 1$, la présence du symptôme x est significative pour la maladie k ; elle diminue la probabilité d'avoir la maladie k .

Interprétation des ORs entre deux modalités

- Pour une variable quantitative
 - $OR_{k_1|k_2}(x = x_1|x = x_2) \simeq 1$, l'augmentation de $x_1 - x_2$ de la variable x ne différencie pas les maladies k_1 et k_2 (elle est non significative).
 - $OR_{k_1|k_2}(x = x_1|x = x_2) > 1$, l'augmentation de $x_1 - x_2$ de la variable x différencie les maladies k_1 et k_2 ; elle augmente la probabilité d'avoir la maladie k_1 plutôt que la maladie k_2 .

- $OR_{k_1|k_2}(x = x_1|x = x_2) < 1$, l'augmentation de $x_1 - x_2$ de la variable x différencie les maladies k_1 et k_2 ; elle augmente la probabilité d'avoir la maladie k_2 plutôt que la maladie k_1 .
- Pour une variable qualitative (Présence ou Absence d'un symptôme)
 - $OR_{k_1|k_2}(x = 1|x = 0) \simeq 1$, le symptôme x ne différencie pas les maladies k_1 et k_2 (il n'est pas significatif).
 - $OR_{k_1|k_2}(x = 1|x = 0) > 1$, la présence du symptôme x est significative pour différencier les maladies k_1 et k_2 ; elle augmente la probabilité d'avoir la maladie k_1 plutôt que la maladie k_2 .
 - $OR_{k_1|k_2}(x = 1|x = 0) < 1$, la présence du symptôme x est significative pour différencier les maladies k_1 et k_2 ; elle augmente la probabilité d'avoir la maladie k_2 plutôt que la maladie k_1 .

Pour chaque variable sélectionnée, nous avons calculé son OR et son intervalle de confiance dans chaque catégorie de la variable réponse y . Par défaut, les ORs sont relatifs à la catégorie de référence (Table 2.9 et Figure 2.12). Nous avons calculé les ORs en faisant évoluer la température de 38 à 40 degrés Celsius et le nombre de jours de maladie de 2 à 6 jours. Les quantiles extérieurs (premier et troisième quantiles) de l'âge sont 8 et 28 respectivement. Ce qui permet de calculer l'OR, sur le demi-échantillon, pour l'âge. De façon similaire, nous avons calculé les ORs sur le demi-échantillon pour la variable pluviométrie, c'est-à-dire entre 14 mm et 370 mm. Les ORs sur les variables binaires sont calculés en comparant les deux modalités : 0 pour absence et 1 pour la présence du symptôme. Les différences entre les 3 groupes cliniques, arbovirus vs paludisme, co-infection vs arbovirus et co-infection vs paludisme sont représentés à travers le calcul des ORs pour une catégorie de maladie par rapport à une autre catégorie de maladie. Ces résultats sont représentés graphiquement dans les Figures 2.11 et 2.13.

Sur le jeu de données *IgM/IgG – data*

A partir du tableau 2.9, nous pouvons dire que si la température augmente de 38 à 40 degrés Celsius, la probabilité d'être co-infecté est doublée et que la probabilité d'avoir le paludisme est multipliée par 2.5. La probabilité d'avoir l'infection aux arbovirus seul est multipliée par 1.71 pour un adulte comparé à un enfant, alors que la probabilité d'avoir le paludisme seul est diminuée par un facteur 0.61. Elle est multipliée par 2.54 si le nombre de jours de maladie augmente de 2 à 6. Les nausées/vomissements augmentent significativement la probabilité d'être co-infecté ou d'avoir le paludisme seul.

La Figure 2.11 présente les odds ratio de chaque variable pour une maladie comparée à une autre. Figure 2.11(a) présente le log odds ratio entre le paludisme seul et les infections aux arbovirus seul. Nous pouvons dire que les variables douleurs articulaires et toux ne sont pas significatives pour distinguer le paludisme seul des infections arbovirales; alors que les autres variables sont toutes significatives. Congestion nasale, nombre de jours de maladie et âge sont associés aux arbovirus. *Temperature*, *pluviometry* et *nausea/vomiting* sont associées au paludisme seul. Ce constat est confirmé par la figure 2.11(b). La figure 2.11(c) montre que quatre variables restent significatives pour distinguer la co-infection du paludisme seul. En particulier, *age*, *number of sick days* et *nasal congestion* sont indicatives pour l'infection aux arbovirus comparé aux cas de co-infection.

Variables \ Diseases	Arbovirus	co-infection	Malaria
<i>Age</i>	1.71 [1.42; 2.07]	1.12 [0.92; 1.36]	0.61 [0.50; 0.73]
<i>Temperature</i>	1.02 [0.69; 1.49]	2.16 [1.52; 3.07]	2.47 [1.82; 3.35]
<i>Number of sick days</i>	2.54 [1.91; 3.37]	1.43 [1.04; 1.96]	1.04 [0.77; 1.39]
<i>Rainfall</i>	2.19 [1.53; 3.14]	17.0 [12.0; 24.0]	9.81 [7.18; 13.4]
<i>Nausea /vomiting</i>	0.83 [0.60; 1.13]	2.07 [1.55; 2.78]	2.15 [1.67; 2.77]
<i>Cough</i>	0.79 [0.58; 1.10]	0.46 [0.33; 0.63]	0.57 [0.44; 0.74]
<i>Nasal congestion</i>	0.52 [0.35; 0.75]	0.13 [0.09; 0.2]	0.10 [0.07; 0.13]
<i>Joint pain</i>	1.52 [0.99; 2.32]	1.90 [1.26; 2.83]	1.74 [1.21; 2.50]

TABLE 2.9 – *IgM/IgG – data* : Les ORs (par rapport à la modalité de référence) calculés pour toutes les variables du modèle et les intervalles de confiance (95%) associés. Résumé sur les différentes maladies.

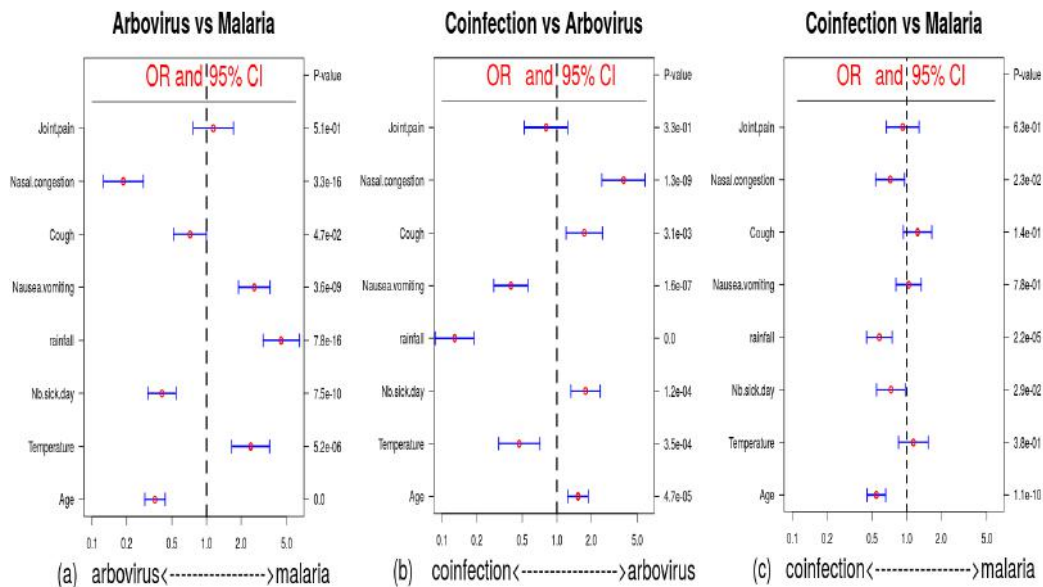


FIGURE 2.11 – *IgM/IgG – data* : odds ratios et intervalles de confiance à 95% représentés par les points et les barres respectivement. (a) Arbovirus *vs* Malaria (b) co-infection *vs* Arbovirus (c) co-infection *vs* Malaria.

Sur le jeu de données *IgM – data*

La figure 2.12 représente la distribution des log ORs, pour chaque variable, sur les 1000 sous-échantillons du jeu de données *IgM*–arbovirus. A partir de la figure 2.12(a), on peut dire que la probabilité d’avoir l’infection aux arbovirus seul augmente avec l’âge et le nombre de jours de maladie, alors que les fortes températures et pluviométrie réduit la probabilité d’avoir l’infection aux arbovirus seul. La comparaison de la co-infection et du paludisme seul, aux figures 2.12(b) et 2.12(c), montre que la température, la pluviométrie, la présence de nausées/vomissement et les douleurs articulaires augmentent la probabilité d’avoir ces maladies, alors que l’âge, nombre de jours de maladie, toux et congestion nasale

diminuent la probabilité d’avoir ces deux maladies.

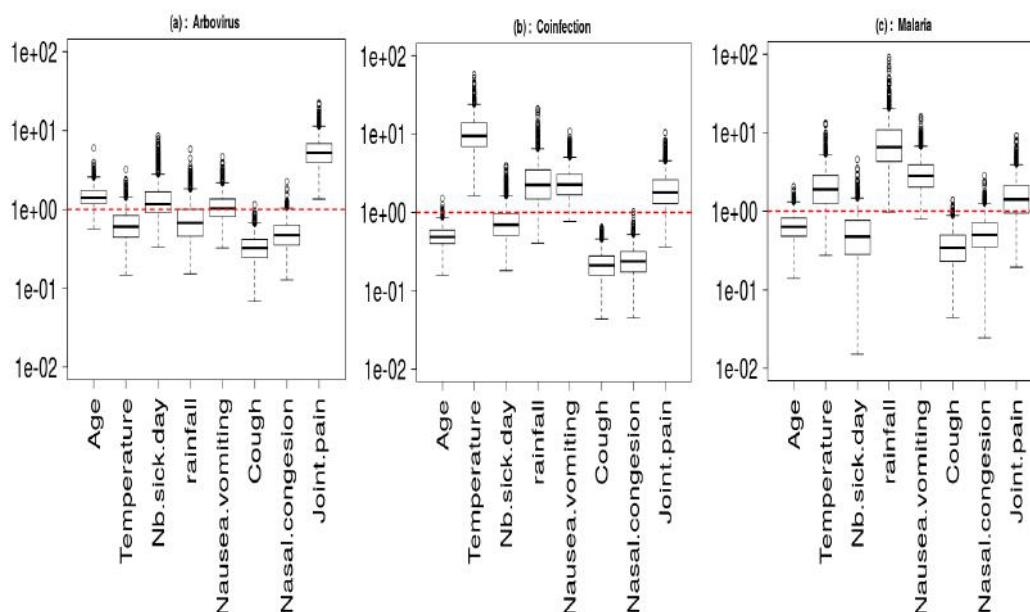


FIGURE 2.12 – *IgM* – data : odds ratios par rapport à la modalité de référence (“autres maladies”) : (a) Arbovirus (b) co-infection (c) Malaria. Distributions à sur les 1000 sous-échantillons.

Nous avons aussi évalué les log ORs pour chaque variable et dans chacun des 1000 sous-échantillons, pour une maladie donnée comparée à une autre. La figure 2.13 présente la distribution des log ORs sur les 1000 sous-échantillons. Cette figure montre que l’âge, et nombre de jours de maladie sont associés aux arbovirus alors que la température et la pluviométrie sont associées au paludisme. La figure 2.13(c) montre qu’une augmentation du nombre de jours de maladie est en faveur de la co-infection et que les fortes températures sont plus observées pour des cas de co-infection que des cas de paludisme seul.

Conclusion

Les résultats basés sur les deux jeux de données montrent que les fortes températures et la présence des nausées/vomissements lors de la saison des pluies sont plus indicatives à l’infection aux parasites du paludisme alors que le nombre de jour de maladies et l’âge sont indicatives à l’infection aux arbovirus. La congestion nasale et les douleurs articulaires sont des symptômes ne montrant pas de résultats assez clairs pour être interprétés. La question principale de notre étude était d’identifier les facteurs risques qui pourront aider le docteur à diagnostiquer la co-infection entre l’infection aux arbovirus et le paludisme. A partir de ces résultats, la température et la pluviométrie sont les seules facteurs risques permettant de différencier la co-infection au paludisme seul ou à l’arbovirus seul.

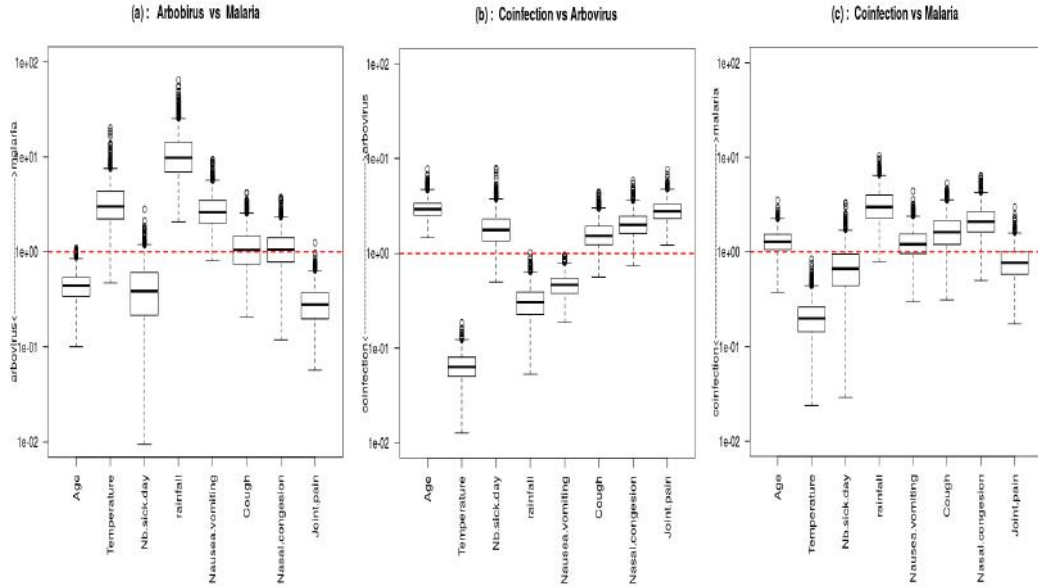


FIGURE 2.13 – *IgM* – data : odds ratios entre deux catégories : (a) Arbovirus *vs* Malaria (b) co-infection *vs* Arbovirus (c) co-infection *vs* Malaria

2.4 Analyse prédictive

Dans cette section, nous essayerons de proposer une méthodologie qui pourra aider au diagnostic de la co-infection et à la recommandation de traitement en cas de co-infection. Nous commencerons par présenter le test d'indépendance entre arbovirus et paludisme. Nous donnerons par la suite une analyse prédictive permettant de discriminer les patients positifs à l'arbovirus aux patients négatifs à l'arbovirus.

2.4.1 Test d'indépendance entre arbovirus et paludisme

A partir du modèle multinomial, on peut construire un test d'indépendance entre deux infections. Nous nous intéressons ici au cas arbovirus (A^+) versus paludisme (M^+). Comme dans les tableaux 2.1 et 2.2, rappelons que A^+ correspond aux individus des catégories 1 et 3, c'est-à-dire aux individus qui sont au moins positifs à l'arbovirus (Arbovirus/co-infection) et que M^+ correspond aux individus des catégories 2 et 3, c'est-à-dire aux individus qui sont au moins positifs au paludisme (Paludisme/co-infection). La loi jointe entre infections ($M_i^+, A_i^+, i = 1, \dots, n$) est donnée par le tableau 2.10. Ainsi indépendance entre infections arbovirales et le paludisme s'écrit comme suit : $\forall (l_1, l_2) \in \{0, 1\}$

$$\mathbb{P}(M^+ = l_1, A^+ = l_2) = \mathbb{P}(M^+ = l_1) \times \mathbb{P}(A^+ = l_2).$$

Selon les valeurs de l_1 et de l_2 , on distingue quatre cas :

- (cas 1) : Si $l_1 = 0$ et $l_2 = 0$, alors

Loi de M \ Loi de A	A = 0	A = 1	Loi de M ⁺
M = 0	π_0	π_1	$\mathbb{P}(M^+ = 0) = \pi_0 + \pi_1$
M = 1	π_2	π_3	$\mathbb{P}(M^+ = 1) = \pi_2 + \pi_3$
Loi de A ⁺	$\mathbb{P}(A^+ = 0) = \pi_0 + \pi_2$	$\mathbb{P}(A^+ = 1) = \pi_1 + \pi_3$	1

TABLE 2.10 – Loi jointe entre l’infection aux arbovirus et le paludisme

$$\begin{aligned}
\mathbb{P}(M^+ = 0, A^+ = 0) &= \mathbb{P}(M^+ = 0) \times \mathbb{P}(A^+ = 0) \\
\pi_0 &= (\pi_0 + \pi_1) \times (\pi_0 + \pi_2) \\
\frac{1}{D} &= \frac{1 + e^{\langle x_i, \beta_1 \rangle}}{D} \times \frac{1 + e^{\langle x_i, \beta_2 \rangle}}{D}, \text{ avec } D = 1 + \sum_{j=1}^3 e^{\langle x_i, \beta_j \rangle} \\
e^{\langle x_i, \beta_3 \rangle} &= e^{\langle x_i, \beta_1 \rangle} e^{\langle x_i, \beta_2 \rangle} \\
\beta_3 &= \beta_1 + \beta_2
\end{aligned}$$

Le calcul est similaire pour les autres cas suivants :

— (cas 2) : Si $l_1 = 0$ et $l_2 = 1$, alors

$$\begin{aligned}
\mathbb{P}(M^+ = 0, A^+ = 1) &= \mathbb{P}(M^+ = 0) \times \mathbb{P}(A^+ = 1) \\
\pi_1 &= (\pi_0 + \pi_1) \times (\pi_1 + \pi_3) \\
\beta_3 &= \beta_1 + \beta_2
\end{aligned}$$

— (cas 3) : Si $l_1 = 1$ et $l_2 = 0$, alors

$$\begin{aligned}
\mathbb{P}(M^+ = 1, A^+ = 0) &= \mathbb{P}(M^+ = 1) \times \mathbb{P}(A^+ = 0) \\
\pi_2 &= (\pi_2 + \pi_3) \times (\pi_0 + \pi_2) \\
\beta_3 &= \beta_1 + \beta_2
\end{aligned}$$

— (cas 4) : Si $l_1 = 1$ et $l_2 = 1$, alors

$$\begin{aligned}
\mathbb{P}(M^+ = 1, A^+ = 1) &= \mathbb{P}(M^+ = 1) \times \mathbb{P}(A^+ = 1) \\
\pi_3 &= (\pi_2 + \pi_3) \times (\pi_1 + \pi_3) \\
\beta_3 &= \beta_1 + \beta_2
\end{aligned}$$

Ainsi $\forall (l_1, l_2) \in \{0, 1\}$

$$\mathbb{P}(M^+ = l_1, A^+ = l_2) = \mathbb{P}(M = l_1) \times \mathbb{P}(A = l_2) \iff e^{X_i \beta_3} = e^{X_i \beta_1} e^{X_i \beta_2}.$$

Par suite, l’hypothèse H_0 est donnée par

$$\beta_3 = \beta_1 + \beta_2 \text{ avec } \beta_k = (\beta_{k0}, \beta_{k1}, \dots, \beta_{kp})^T.$$

Statistique de test

Notons par $\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{pmatrix}$ l'estimateur du maximum de vraisemblance (EMV) et $\beta^* \in \mathbb{R}^q$ la vraie valeur de β , avec $q = 3(p + 1)$. On a :

$$\hat{\beta} - \beta^* = \left(-nH(\hat{\beta})^{-1}\right)^{-\frac{1}{2}} \frac{U_n}{\sqrt{n}} + o_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right),$$

avec $U_n \rightsquigarrow U$ où $U \sim \mathcal{N}(0, I_q)$, et $H(\hat{\beta})$ la matrice hessienne. Définissons

$$h : \beta \mapsto h(\beta) = \beta_3 - \beta_2 - \beta_1$$

On a, par développement de Taylor, que :

$$h(\hat{\beta}) = h(\beta^*) + Dh \left(-nH(\hat{\beta})^{-1}\right)^{-\frac{1}{2}} \frac{U_n}{\sqrt{n}} + o_{\mathbb{P}}\left(\frac{1}{\sqrt{n}}\right),$$

avec Dh la matrice jacobienne de h définie par : $Dh = (-Id_{p+1}, -Id_{p+1}, Id_{p+1})$.

$$\sqrt{n} \left(h(\hat{\beta}) - h(\beta^*)\right) = Dh \left(-nH(\hat{\beta})^{-1}\right)^{-\frac{1}{2}} U_n + o_{\mathbb{P}}(1).$$

Par suite, on a

$$\sqrt{n}\Sigma_n^{-1/2} \left(h(\hat{\beta}) - h(\beta^*)\right) \rightsquigarrow \mathcal{N}(0, I_q),$$

avec $\Sigma_n = -nDhH(\hat{\beta})^{-1}Dh^T$.

Nous voulons tester l'hypothèse :

$$H_0 : h(\beta^*) = 0 \text{ vs } H_1 : h(\beta^*) \neq 0.$$

Nous définissons alors la statistique de test W_n , donnée par :

$$W_n = h(\hat{\beta})^T \Sigma_n^{-1} h(\hat{\beta}).$$

Sous l'hypothèse H_0 , W_n converge en loi vers $\chi^2(p + 1)$. La zone de rejet est donnée par :

$$\mathcal{R} = \left\{W_n > \chi_{1-\frac{\alpha}{2}}^2(p + 1)\right\} \cup \left\{W_n < \chi_{\frac{\alpha}{2}}^2(p + 1)\right\},$$

Ainsi en utilisant l'EMV donné par le modèle obtenu après la sélection de variables, on pourra tester l'indépendance sur les différents jeux de données.

Résultats de test

Pour les deux jeux de données, nous considérons le même modèle final composé des huit variables suivantes : *pluviometry*, *age*, *temperature*, *joint pain*, *nasal congestion*, *number of sick days*, *nausea/vomiting* et *cough*. En utilisant ce modèle, nous testons l'indépendance entre arbovirus et parasites du paludisme sur le jeu de données *IgM/IgG - data*. Dans ce cas, l'hypothèse H_0 est rejetée avec une *p-value* égale à

Variable	p-value du test
joint pain	$1.8.10^{-3}$
nasal congestion	0.41
cough	$5.3.10^{-8}$
nausea/vomiting	$1.71.10^{-8}$
number of sick days	$1.1.10^{-5}$
pluviometry	$6.25.10^{-6}$
temperature	$5.99.10^{-5}$
age	$1.64.10^{-5}$

TABLE 2.11 – *IgM/IgG – data*. P-values du test en enlevant une des variables importantes.

Variable	p-value du test
sex	$2.18.10^{-4}$
cephalalgia	0.18
eye pain	0.35
muscle pain	$1.09.10^{-3}$
chills	$3.53.10^{-4}$
diarrhea	$9.2.10^{-3}$
jaudice	0.88

TABLE 2.12 – *IgM/IgG – data*. P-values du test en ajoutant une des variables non-importantes.

$1.46.10^{-6}$. Nous avons aussi étudié la robustesse de la décision du test par rapport à la sélection de variables. C'est-à-dire, en oubliant une variable importante ou en ajoutant une variable non importante et en regardant la décision du test. Hormis les variables *nasal congestion*, *cephalalgia*, *eye pain* et *jaudice*, nous obtenons toujours des *pvalues* inférieures ou égales à 10^{-3} (voir tableaux 2.11 et 2.12).

Pour appliquer le test sur le jeu de données *IgM – data*, nous calculons les 1000 *pvalues* de test correspondant aux 1000 sous-échantillons construits lors de la sélection de variables. Nous trouvons que 42,5% d'entre elles sont inférieures ou égales à 0.05. Ce faible taux de rejet peut être expliqué par le fait que la taille des sous-échantillons est petite (139), ce qui fait que l'approximation asymptotique de la statistique de test ne marche pas bien. On peut aussi penser au fait qu'on utilise le même modèle sur tous les sous-échantillons alors qu'on voit que le test était sensible à la sélection de certaines variables. Sur chacun des 1000 sous-échantillons une sélection de variables a été effectuée. Les variables retenues ont été considérées dans le modèle final du sous-échantillon associé. Le test a été fait sous ces conditions pour chaque sous-échantillon. On obtient que dans 73,8% des cas, le test est rejeté si on considère un taux d'erreur de 5%. Les p-values de test sont données par la figure 2.14. Ainsi, dans plus de 15% des cas, le test n'est pas rejeté. On peut alors penser que le jeu de données *IgM – data* peut ne pas contenir assez d'informations pour mettre en évidence la dépendance entre arbovirus et paludisme. Ce

qui signifie que le test d'indépendance perd en puissance.

Dans la suite, on considérera le jeu de données *IgM/IgG – data* pour proposer une analyse prédictive.

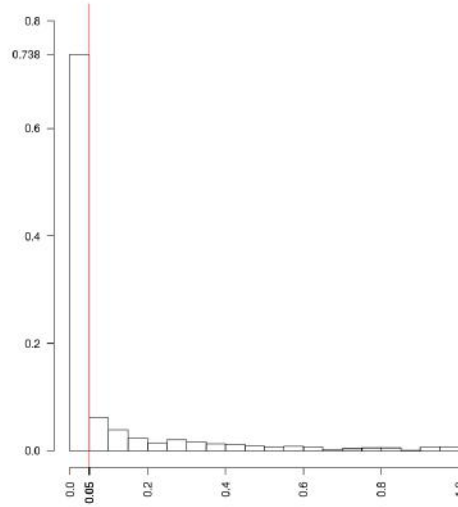


FIGURE 2.14 – *IgM – data* : histogramme des p-values de test sur les 1000 sous-échantillons.

2.4.2 Diagnostic de la co-infection

En se basant sur les résultats précédents, nous proposons une méthodologie de diagnostic qui pourra aider le médecin à mieux diagnostiquer les infections arbovirus sachant que le patient a le paludisme.

Le test d'indépendance décrit à la section 2.4.1 montre une association entre paludisme et infections aux arbovirus. Alors nous pouvons évaluer la probabilité d'être co-infecté sachant que le paludisme est observé. Cette probabilité peut être calculée en fonction des probabilités π_k estimées à partir de la régression logistique multinomiale. Pour chaque individu i ,

$$\mathbb{P}(\widehat{C/P})_i = \frac{\widehat{\pi}_3(i)}{\widehat{\pi}_3(i) + \widehat{\pi}_2(i)} = \frac{e^{X_i \widehat{\beta}_3}}{e^{X_i \widehat{\beta}_3} + e^{X_i \widehat{\beta}_2}}.$$

La probabilité conditionnelle de co-infection peut être utilisée pour différencier laquelle des maladies doit-on traiter. Nous proposons une classification binaire et nous prédisons un cas de co-infection si la probabilité conditionnelle de co-infection dépasse un certain seuil γ :

$$\begin{cases} \text{Si } \mathbb{P}(C|M) \geq \gamma : & \text{Cas positif aux arbovirus,} \\ \text{Si } \mathbb{P}(C|M) < \gamma : & \text{Cas négatif aux arbovirus.} \end{cases}$$

L'évaluation de la classification est basée sur la matrice de confusion et la précision globale de la classification. La matrice de confusion est utilisée pour calculer les vrais positifs aux arbovirus (TP), les faux positifs aux arbovirus (FP), les vrais négatifs aux arbovirus (TN) et les faux négatifs aux arbovirus (FN). Une mesure de performance

globale est donnée par le taux de mal classés (MCR) défini par :

$$\text{MCR} = \frac{FP + FN}{N},$$

avec $N = TP + FP + TN + FN$.

L'analyse binaire présentée dans cette partie est basée sur 1148 individus du jeu de données *IgM/IgG*–arbovirus qui correspondent aux patients infectés aux parasites du paludisme (c'est-à-dire co-infection et paludisme). L'apprentissage de la régression logistique multinomiale est fait sur 70% du jeu de données *IgM/IgG*–arbovirus, soit 1317 individus, et le test est fait sur un échantillon de 377 individus positifs au paludisme. Pour choisir le seuil de classification γ , une pratique standard est de minimiser le taux de mal classés (MCR). Nous calculons l'estimateur du MCR par validation croisée sur 5-échantillons. Nous pouvons voir sur la figure 2.15 que le seuil optimal est autour de $\gamma = 0.5$.

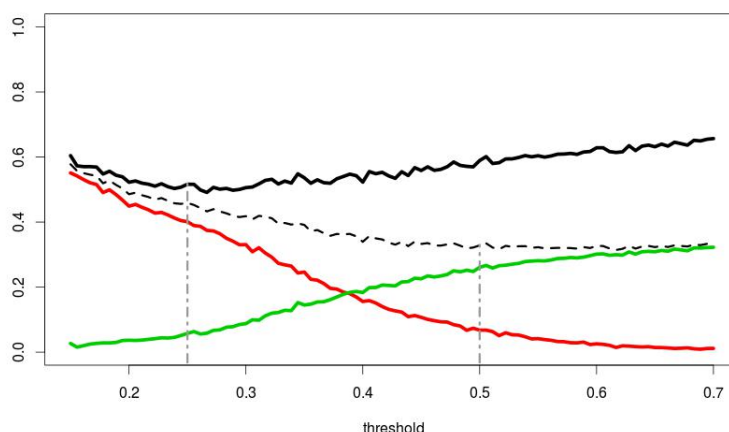


FIGURE 2.15 – *IgM/IgG* – data : Taux de mal-calssés estimé par validation croisée. La ligne pleine noire représente le WMCR. Le MCR est donné par la ligne noire pointillée. Une augmentation de γ augmente le nombre de FN (ligne verte) et diminue le nombre de FP (ligne rouge).

La validation croisée a été effectuée plusieurs fois avec des 5-échantillons différents et la valeur optimale du seuil reste stable. Alors une classification avec $\gamma = 0.5$ a été utilisée pour prédire le type de maladie d'un patient en se basant sur ses symptômes cliniques. Les prédictions et les actuels cas d'arbovirus ont été comparés en utilisant l'échantillon test (sur 377 individus), comme présenté au tableau 2.13.

Les lignes de la matrice sont les classes actuelles et les colonnes correspondent au prédictions. On observe que le MCR est de 38%, et que le nombre de faux négatifs (FN) est très grand. Dans le cas de diagnostic de maladie, il est préférable d'avoir une classification qui réduit le nombre de FN, parce que les FN peuvent être plus dangereux si on en rate beaucoup lors du traitement. Différentes stratégies peuvent être adoptées. Une possibilité est de réduire le nombre de FN en minimisant une version pondérée du MCR :

$$\text{WMCR} = \frac{FP + 2FN}{N}.$$

<i>True</i> \ <i>Predicted</i>	0	1
0	211	29
1	114	23

TABLE 2.13 – Table de confusion avec $\gamma = 0.5$.

<i>True</i> \ <i>Predicted</i>	0	1
0	88	152
1	24	113

TABLE 2.14 – Table de confusion avec $\gamma = 0.25$.

La valeur du seuil qui minimise le WMCR est de 0.25. Avec ce choix de γ , on observe sur le tableau 2.14 que le nombre de FN est réduit mais le taux de mal classés a augmenté. Dans une autre étape, nous proposons de sélectionner, sur les patients prédits positifs, ceux qui ont un age supérieur à 10 et un nombre de jours de maladie supérieur à 3. En effet, nous avons conclu (pour le jeu de données *IgM/IgG – data*) à la section 2.3.4 que ces deux variables sont plus indicatives aux arbovirus. Le tableau 2.15 donne le résultat correspondant : le MCR décroît jusqu’à 36% et que le nombre de FN reste plus petit que le nombre de FN sur le tableau 2.13. De même le nombre de TP est doublé.

<i>True</i> \ <i>Predicted</i>	0	1
0	190	50
1	85	52

TABLE 2.15 – Table de confusion avec $\gamma = 0.25$, *Age* = 10 and Number of sick days = 3.

L’objectif de ces prédictions était d’affecter un patient au groupe des “paludisme” ou au groupe des “arbovirus”, et à traiter des cas de co-infection en fonction de la similitude des symptômes avec ceux de ces deux maladies. La procédure de classification est basée sur le calcul de la probabilité conditionnelle $\mathbb{P}(C|M)$. Le paramètre de seuil est calibré sur les données en minimisant le taux de mal classés pondéré (WMCR). Pour plus de précision dans la classification, nous proposons d’utiliser les deux variables considérées comme indicatives à l’arbovirus.

La performance de la procédure de classification est fortement affectée par la qualité des données. Notre analyse est basée sur deux jeux de données. Nous nous basons sur le jeu de données *IgM/IgG – data* construit à partir des patients positifs à l’*IgM* ou à l’*IgG* pour fournir une analyse prédictive. Malheureusement, être positif à l’*IgG* ne veut pas forcément dire qu’on a eu une infection récente aux arbovirus, parce que les anticorps auront été peut-être développés depuis longtemps. Ce qui minimise la possibilité de trouver une vraie corrélation avec les symptômes enregistrés initialement. Ces limites réduisent la capacité de prédiction dans la procédure de classification. Les faux positifs

et faux négatifs dûs aux tests biologiques peuvent impacter les résultats. Cependant, les tests de diagnostic utilisés dans cette étude présentent de forts paramètres de sensibilité et de sensibilité. Leurs impacts peuvent être considérés négligeables.

2.5 Discussion

Un mauvais diagnostic de la co-infection entre arbovirus et paludisme peut augmenter la propagation des maladies arbovirales dans les zones où les tests ne sont pas accessibles. Cette étude propose une méthodologie statistique qui pourra aider le médecin à élaborer un bon diagnostic de la co-infection en cas de paludisme, mais aussi lui offrir une recommandation de traitement en cas de co-infection. L'un des principaux objectifs de ce chapitre était de disposer d'une procédure qui peut pré-traiter les données, appliquer des techniques statistiques appropriées afin de fournir des prédictions raisonnablement précises et cliniquement utiles. Notre analyse est basée sur un jeu de données réelles. Le jeu de données *IgM – data*, les individus positifs aux arbovirus sont identifiés dans les premiers stades de la maladie. Cependant, les cas positifs constituent une très petite partie (39 cas sur 12288 individus). Plusieurs stratégies d'échantillonnage sont développées pour étudier des données déséquilibrées ([67]) et faire une bonne classification. Branco et al. [61] propose de classer ces approches en deux catégories : pré-traitement de données et modifications d'algorithmes d'apprentissage. Des stratégies de traitement incluant les solutions via forêts aléatoires y sont discutées ([68], [69]).

Pour analyser les données de co-infection, nous avons proposé une méthodologie avec trois étapes : 1. une sélection de variables par forêts aléatoires, par test du rapport de vraisemblance et par stepwise, 2. une analyse des facteurs influents par le calcul des odds ratio à partir de la régression logistique multinomiale, 3. une analyse prédictive basée sur la probabilité de co-infection.

A partir de notre analyse, on peut dire que la combinaison des forêts aléatoire et du test du rapport de vraisemblance est une méthode robuste pour sélectionner les variables importantes pour les différentes maladies. L'analyse des odds ratio permet d'identifier les facteurs de risque qui caractérisent chaque maladie. Nous avons observé que les valeurs élevées du nombre de jours de maladie et de l'âge sont principalement révélatrices d'une maladie arbovirale alors que des fortes températures et la présence de nausée et/ou vomissements pendant la saison des pluies sont principalement révélatrices d'une maladie du paludisme. La règle de classification basée sur la probabilité de co-infection, l'âge et le nombre de jours de maladie identifie les patients co-infectés à traiter pour des maladies arbovirales avec une précision globale de 65%. Les résultats pourraient être améliorés sur un autre jeu de données beaucoup plus approprié. Une future étude appliquera cette méthodologie aux données de co-infection entre le paludisme et d'autres agents pathogènes plus facilement détectables, au début de l'infection, que les arbovirus.

Chapitre 3

Mélange de modèles linéaires généralisés et méthode des moments : identifiabilité & applications

Sommaire

3.1	Introduction	58
3.2	Notation & modèle	59
3.2.1	Notations et définitions	59
3.2.2	Modèles	60
3.3	Algorithme	61
3.3.1	Estimation des directions	61
3.3.2	Estimation de tous les paramètres du modèle	68
3.4	Résultats théoriques	69
3.4.1	Identifiabilité	69
3.4.2	Consistance	73
3.4.3	Normalité Asymptotique	75
3.5	Applications	83
3.5.1	Package R	83
3.5.2	Simulations	84
3.5.3	Sélection de variables	94

Résumé

Le modèle de mélange fini consiste à supposer que les données proviennent d'une source contenant plusieurs sous-populations. Chaque sous-population est modélisée de manière séparée. La population totale est un mélange de ces différentes sous-populations. Les méthodes de vraisemblance pour les modèles de mélange peuvent utiliser l'algorithme EM (Dempster et al. [28]). Des méthodes variationnelles de Bayes ont aussi été développées pour l'étude des modèles de mélange fini. Mais ces méthodes peuvent converger vers des optimums locaux et peuvent présenter des vitesses de convergence faibles en grande dimension. En plus, elles peuvent présenter des temps de calcul assez longs.

Nous présentons ici une méthode de moindres carrés qui est une combinaison d'une méthode spectrale et d'une méthode de diagonalisation jointe. Sous certaines conditions, notre méthode garantit de bien retrouver les paramètres. L'idée de base est d'écrire les moments croisés entre les entrées x et la réponse y comme des tenseurs symétriques et d'utiliser l'identité de Stein.

Dans cette étude, nous nous intéressons aux modèles de mélanges linéaires généralisés pour des observations binaires. C'est-à-dire que, si y est la réponse et x le vecteur de covariables,

$$\mathbb{P}(Y = 1|X) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + b_k)$$

où g est la fonction lien, par exemple la fonction logistique, et où pour $k = 1, \dots, K$, et pour la k -ième sous-population, ω_k est la proportion de la sous-population, β_k est le vecteur de régression et b_k est l'intercept.

Avec ce modèle, nous présentons l'algorithme spectral en deux étapes : (1) une étape d'estimation des vecteurs de régression normalisés et (2) une étape d'estimation de tous les paramètres du modèle. Nous prouvons ensuite les résultats théoriques, sous des hypothèses raisonnables : identifiabilité du modèle, consistance et normalité asymptotique des estimateurs. Nous montrons dans une étape d'application (par simulation) qu'on arrive à bien estimer les paramètres. Nous montrons aussi que les estimateurs par la méthode spectrale peuvent être meilleurs pour des échantillons de taille modérée que ceux obtenus par maximum de vraisemblance quand la dimension augmente. Le temps de calcul aussi reste assez faible comparé à la méthode du maximum de vraisemblance. Toutes ces applications sont présentées via un package R qui pourra être utilisé pour une éventuelle futur étude.

Mots clés : Diagonalisation jointe, Modèle de mélange, Modèle linéaire généralisé, Méthode spectrale, Méthode des moments.

Abstract

The finite mixture model assumes that the data come from a source containing several subpopulations. Each subpopulation is modeled separately. The total population is a mixture of these different subpopulations. Likelihood methods for the finite mixture model can use the EM algorithm (Dempster et al., [28]). Bayesian variational methods have also been developed to deal with such models. But both of these methods can converge to spurious local optima and can have low convergence rate in high dimensional models. In addition, they can have quite long computing times.

We present here a least squares method which is a combination of a moment method and of a joint diagonalization method. Under some weak conditions, our method is proved to recover the parameters. The basic idea is to write cross moments between the entries x and the response y as symmetric tensors and to use Stein's identity.

In this study, we will focus on finite mixtures of regression models for binary output. That means that, if y is the output and x the vector of covariates :

$$\mathbb{P}(y = 1|x) = \sum_{k=1}^K \omega_k g(\langle \beta_k, x \rangle + b_k)$$

where g is the link function, for instance the logistic function, and where for $k = 1, \dots, K$, and for the k -th subpopulation, ω_k is the proportion of the subpopulation, β_k is the regression vector and b_k is the intercept. With this model, we present our spectral algorithm in two steps : (1) first, we present the estimation of the normalized regression vectors and (2) secondly we present the estimation of all the model parameters. We then prove our theoretical results, which hold under mild assumptions : identifiability of the model, consistency and asymptotic normality of the estimators. We show in an application step (using simulation studies) that we can recover all the model parameters. We also show that the estimators provided by spectral method may be better for finite samples than those obtained by maximum likelihood when the dimension increases. The computation time also remains rather low compared to the maximum likelihood method. All these applications are presented using an R package which can be used for a possible future study.

Keywords : Generalized linear model, Joint diagonalization, Mixture Model, Moments method, Spectral method.

3.1 Introduction

Le modèle de mélange fini est utilisé depuis plus d'un siècle (Newcomb (1886) [19], Pearson (1894) [20]) mais depuis plusieurs années, l'utilisation du modèle de mélange s'est considérablement développée avec la parution de l'article de Dempster et al. [28]. Le modèle de mélange fini consiste à supposer que les données proviennent d'une source contenant plusieurs sous-populations. Chaque sous-population est modélisée de manière séparée. La population totale est un mélange de ces différentes sous-populations. Le modèle résultant est un modèle de mélange fini de densité f définie par :

$$f(\cdot) = \sum_{k=1}^K \pi_k f_k(\cdot),$$

avec, π_k : proportions du mélange, f_k : densités des composantes du mélange.

Le modèle linéaire généralisé est une extension de la régression linéaire dans lequel la réponse peut être écrite comme une fonction non linéaire des entrées à travers une fonction lien (McCullagh et Nelder (1989) [15]). Le modèle linéaire généralisé peut être étudié à l'aide de plusieurs méthodes présentes dans la littérature (Kakade et al. [70]).

Dans certaines situations, le modèle linéaire généralisé ne suffit pas pour étudier les données disponibles du fait qu'elles proviennent de plusieurs groupes qui ont des caractéristiques différentes. Dans ce cas, le modèle de mélange de modèles linéaires généralisés est plus approprié. Le modèle de mélange constitue aujourd'hui un outil fondamental en statistique appliquée et en apprentissage dans différents domaines. Le modèle de mélange de modèles linéaires généralisés est utilisé aujourd'hui dans beaucoup de domaines d'applications (comme souligné dans les travaux de Sedghi et al. [34]) tels que la reconnaissance d'objets (Quattoni et al. [71]), la reconnaissance d'action humaine (Wang et Mori [72]), l'analyse syntaxique (Petrov et Klein [73]) et la traduction automatique (Liang et al. [74]).

Dans le passé, le modèle de mélange fini était étudié à l'aide des méthodes basées sur le maximum de vraisemblance telles que l'algorithme EM (Dempster et al. [28], Jordan et al. [30]; Xu et al. [31]; Grün [32]) ou des méthodes variationnelles de Bayes (Bistrop et Svensen [29]). Mais ces méthodes peuvent converger vers des optimums locaux et peuvent présenter des vitesses de convergence faibles en grande dimension. En plus, elles peuvent présenter des temps de calcul assez longs.

Nous présentons ici une méthode d'estimation qui est une combinaison d'une méthode des moments (Pearson, [20]) et d'une méthode spectrale (diagonalisation jointe, [40]). L'idée de base est d'écrire les moments croisés entre les entrées x et la réponse y comme des tenseurs symétriques, et d'utiliser l'identité de Stein. La méthode des moments est très ancienne et remonte à Pearson (Pearson 1894, [20]) mais a refait surface ces dernières années avec l'utilisation des tenseurs ([36]). Ces méthodes ont eu beaucoup de succès ces dernières années ([35, 37, 38]). Sous certaines conditions, ces méthodes garantissent de bien retrouver les vecteurs de paramètres normalisés. Ces vecteurs de paramètres pourront servir dans l'initialisation de l'algorithme EM mais aussi à sélectionner les variables importantes. En plus de cela, on peut retrouver entièrement les paramètres du modèle de mélange en un temps raisonnable.

Dans cette étude, nous nous intéresserons à une partie de modèles de mélanges linéaires généralisés (Grün [32]) qui est constituée des mélanges de modèles de type régression

logistique. C'est-à-dire que, si y est la réponse (binaire) et x le vecteur de covariables,

$$\mathbb{P}(Y = 1|X) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + b_k)$$

où g est la fonction lien, par exemple la fonction logistique, et où pour $k = 1, \dots, K$, et pour la k -ième sous-population, ω_k est la proportion de la sous-population, β_k est le vecteur de régression et b_k est l'intercept.

Avec ce modèle, l'algorithme spectral permet d'estimer les vecteurs de régression à un facteur près ([75]). Nous utilisons par la suite une méthode des moindres carrés pour estimer tous les paramètres du modèle.

Ce chapitre est organisé comme suit. La section 3.2 présentera le modèle à étudier et les notations utilisées. La section 3.3 décrit l'algorithme d'estimation des paramètres. On commencera par estimer d'abord les vecteurs de régression normalisés pour ensuite estimer globalement tous les paramètres du modèle. La section 3.4 présente les résultats théoriques de l'étude, à savoir l'identifiabilité du modèle, ainsi que la consistance et la normalité asymptotique des estimateurs présentés à la section 3.2. La section 3.5 présente des applications liées au modèle présenté à la section 3.2. On présentera le package **R** prévu à cet effet et quelques simulations.

3.2 Notation & modèle

Nous présenterons dans cette partie la notation et les modèles qui seront utilisés dans ce chapitre.

3.2.1 Notations et définitions

Notons par $[n]$ l'ensemble des valeurs $\{1, 2, \dots, n\}$ et $e_i \in \mathbb{R}^d$, le i -ième vecteur de la base canonique de \mathbb{R}^d . Notons aussi par $I_d \in \mathbb{R}^{d \times d}$ la matrice identité dans \mathbb{R}^d . Le produit tensoriel de p espaces euclidiens \mathbb{R}^{d_i} , $i \in [p]$ est noté $\bigotimes_{i=1}^p \mathbb{R}^{d_i}$. T est un tenseur réel d'ordre p si $T \in \bigotimes_{i=1}^p \mathbb{R}^{d_i}$. Si $p = 1$, T est un vecteur de \mathbb{R}^d et une matrice de $\mathbb{R}^{d \times d}$ pour $p = 2$. Les coordonnées (i_1, i_2, \dots, i_p) de T par rapport à la base canonique sont des réels notés $T[i_1, i_2, \dots, i_p] : i_1, i_2, \dots, i_p \in [d]$.

Définition 3.2.1.1 Soit $T \in \mathbb{R}^{d \times d \times d}$ un tenseur d'ordre 3 et $M_r \in \mathbb{R}^{d \times d_r}$, $r = 1, 2, 3$. Le tenseur d'ordre 3 $T(M_1, M_2, M_3) \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ est défini par

$$T(M_1, M_2, M_3)[i_1, i_2, i_3] := \sum_{j_1, j_2, j_3 \in [d]} T[j_1, j_2, j_3] M_1[j_1, i_1] M_2[j_2, i_2] M_3[j_3, i_3].$$

En particulier, si T est un tenseur d'ordre 3 tel que, pour des vecteurs v_1, \dots, v_K de \mathbb{R}^d

$$T = \sum_{k \in [K]} v_k^{\otimes 3} \in \mathbb{R}^{d \times d \times d}$$

et si $W \in \mathbb{R}^{d \times K}$ est une matrice, on peut alors définir les tenseurs suivants :

- $T(I, I, W) := \sum_{k \in [K]} v_k \otimes v_k \otimes (W^T v_k) \in \mathbb{R}^{K \times d \times d}$,
- $T(I, W, W) := \sum_{k \in [K]} v_k \otimes (W^T v_k) \otimes (W^T v_k) \in \mathbb{R}^{K \times K \times d}$,
- $T(W, W, W) := \sum_{k \in [K]} (W^T v_k)^{\otimes 3} \in \mathbb{R}^{K \times K \times K}$.

Pour des vecteurs $u, v, w \in \mathbb{R}^d$, on note

$$T(u, v, w) = \sum_{i, j, l \in [d]} u_i v_j w_l T[i, j, l] \in \mathbb{R}, \quad (3.1)$$

$$T(I_d, v, w) = \sum_{j, l \in [d]} v_j w_l T[:, j, l] \in \mathbb{R}^d, \quad (3.2)$$

$$T(I_d, I_d, w) = \sum_{l \in [d]} w_l T[:, :, l] \in \mathbb{R}^{d \times d}. \quad (3.3)$$

3.2.2 Modèles

Soit $(X_n, Y_n) = ((x_1, y_1), \dots, (x_n, y_n))$ un échantillon indépendant identiquement distribué de loi $\mathbb{P}_\theta = \mathcal{L}(X, Y)$, avec Y la variable réponse ($Y \in \{0, 1\}$) et $X \in \mathbb{R}^d$ le vecteur de covariables. Supposons que conditionnellement à $X = x$, Y provient d'un mélange de K populations modélisé par

$$\mathbb{P}(Y = 1 | X = x) = \sum_{k=1}^K \omega_k g(\langle \beta_k, x \rangle + b_k) \quad (3.4)$$

avec $\theta = (\omega, \beta, b) \in \Theta$, $\beta = [\beta_1 | \dots | \beta_K] \in \mathbb{R}^{d \times K}$ la matrice des K vecteurs de paramètres, $b = (b_1, \dots, b_K)$ le vecteur des intercepts et g la fonction lien. On note $\omega = (\omega_1, \dots, \omega_K)$ le vecteur des poids du mélange, sous les contraintes $\sum_{k=1}^K \omega_k = 1$ et $\omega_k > 0$.

Modèle 3.2.2.1 *On suppose que*

$$x \sim N(0, I_d) \quad (3.5)$$

Modèle 3.2.2.2 *On suppose que*

$$x \sim N(m, \Sigma) \quad (3.6)$$

avec m la moyenne et Σ la matrice covariance.

Notons par $\lambda_k = \|\beta_k\|_2$ la norme du vecteur β_k pour le modèle 3.2.2.1 et $\lambda_k = \|\Sigma \beta_k\|_2$ la norme du vecteur $\Sigma \beta_k$ pour le modèle 3.2.2.2, $k = 1, \dots, K$.

Pour tout individu i , $i = 1, \dots, n$ et pour tout k , $k = 1, \dots, K$, on note

$$\pi_{ik} = \mathbb{P}(y_i = 1 | x_i, k) = g(\langle \beta_k, x_i \rangle + b_k). \quad (3.7)$$

— Si le lien est logit, alors la fonction lien est

$$g(z) = \frac{e^z}{1 + e^z}$$

— Si le lien est probit, alors la fonction lien est

$$g(z) = \Phi(z),$$

avec Φ la fonction de répartition d'une loi normale centrée réduite.

— De manière générale, on suppose que g est connue, que g est une fonction strictement croissante de \mathbb{R} vers $]0, 1[$ et indéfiniment dérivable.

3.3 Algorithme

Nous présenterons dans cette partie, les outils permettant d'écrire l'algorithme d'estimation. D'abord nous présenterons l'estimation des directions ensuite nous essayerons d'estimer les paramètres du modèle.

3.3.1 Estimation des directions

L'estimation de directions des vecteurs de paramètres consiste à estimer les $\mu_k = \beta_k / \|\beta_k\|_2$ pour le modèle 3.2.2.1 et les $\mu_k = \beta'_k / \|\beta'_k\|_2$ pour le modèle 3.2.2.2, avec $\beta'_k = \Sigma \beta_k \in \mathbb{R}^d$, $k = 1, \dots, K$. Rappelons que Σ est la matrice de variance des entrées x dans le modèle 3.2.2.2.

Définition 3.3.1.1 *Sous le modèle 3.2.2.1, nous pouvons définir les moments croisés, entre la réponse y et le vecteur de covariable, jusqu'à l'ordre 3 par :*

- $M_1(\theta) := \mathbb{E}_\theta[y.x]$, moment d'ordre 1,
- $M_2(\theta) := \mathbb{E}_\theta[y.(x \otimes x - I_d)]$, moment d'ordre 2 et
- $M_3(\theta) := \mathbb{E}_\theta[y.x \otimes x \otimes x] - \sum_{j \in [d]} \mathbb{E}_\theta \left[y(e_j \otimes x \otimes e_j + e_j \otimes e_j \otimes x + x \otimes e_j \otimes e_j) \right]$ moment d'ordre 3.

Définition 3.3.1.2 *Sous le modèle 3.2.2.2, nous pouvons définir les moments croisés, entre la réponse y et le vecteur de covariable, jusqu'à l'ordre 3 par :*

- $M_1(\theta) := \mathbb{E}_\theta[y.(x - m)]$, moment d'ordre 1,
- $M_2(\theta) := \mathbb{E}_\theta \left[y.(x \otimes x - x \otimes m - m \otimes x + m \otimes m - \Sigma) \right]$, moment d'ordre 2, et
- $M_3(\theta) := \mathbb{E}_\theta \left[y(x \otimes x \otimes x - x \otimes x \otimes m - m \otimes x \otimes x - x \otimes m \otimes x + x \otimes m \otimes m + m \otimes m \otimes x + m \otimes x \otimes m - m \otimes m \otimes m) \right] - T(\theta)$ le moment d'ordre 3, avec

$$T(\theta)[a, b, c] = \Sigma[a, b]M_1(\theta)[c] + \Sigma[a, c]M_1(\theta)[b] + \Sigma[b, c]M_1(\theta)[a]$$

Dans les différents modèles, on peut écrire les moments croisés entre x et y comme des tenseurs symétriques en fonction des vecteurs de paramètres β_k , $k = 1, \dots, K$.

Lemme 3.3.1.1 *(Stein, 1972 [76])*

Soit $X \sim \mathcal{N}(0, 1)$ une variable aléatoire de loi normale centrée réduite. Soit g une fonction réelle continûment dérivable telle que $\mathbb{E}(|g(X).X|) < \infty$ et $\mathbb{E}(|g'(X)|) < \infty$. Alors

$$\mathbb{E}[g(X).X] = \mathbb{E}[g'(X)].$$

Preuve 3.3.1.1 *Voir [76] pour la preuve.*

Lemme 3.3.1.2 ([37])

Sous le modèle 3.2.2.1, on peut écrire les moments donnés à la définition 3.3.1.1 comme suit

$$M_1(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g'(\langle \beta_k, x \rangle + b_k)] \cdot \beta_k, \quad (3.8)$$

$$M_2(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g''(\langle \beta_k, x \rangle + b_k)] \cdot \beta_k \otimes \beta_k, \quad (3.9)$$

$$M_3(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g^{(3)}(\langle \beta_k, x \rangle + b_k)] \cdot \beta_k \otimes \beta_k \otimes \beta_k, \quad (3.10)$$

Preuve 3.3.1.2 Voir la preuve dans les travaux de Anandkumar et al. [37].

Lemme 3.3.1.3 Sous le modèle 3.2.2.2, on peut écrire les moments donnés à la définition 3.3.1.2 comme suit

$$M_1(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g'(\langle \beta_k, x \rangle + b_k)] \cdot \beta'_k, \quad (3.11)$$

$$M_2(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g''(\langle \beta_k, x \rangle + b_k)] \cdot \beta'_k \otimes \beta'_k, \quad (3.12)$$

$$M_3(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g^{(3)}(\langle \beta_k, x \rangle + b_k)] \cdot \beta'_k \otimes \beta'_k \otimes \beta'_k, \quad (3.13)$$

avec $\beta'_k = \Sigma \beta_k$.

Preuve du Lemme 3.3.1.3.

La preuve est similaire à la preuve du lemme 3.3.1.2 ([37]).

1. Le moment d'ordre 1 :

Par définition,

$$\mathbb{E} [y \cdot x] = \sum_{k \in [K]} \omega_k \mathbb{E} [g(\langle x, \beta_k \rangle + b_k) \cdot x] \quad (3.14)$$

Comme $x = m + \Sigma^{\frac{1}{2}} v$, avec $v \sim N(0, I_d)$, on a

$$\mathbb{E} [y \cdot x] = \mathbb{E} [y \cdot m] + \sum_{k \in [K]} \omega_k \mathbb{E} \left[g(\langle x, \beta_k \rangle + b_k) \cdot \Sigma^{\frac{1}{2}} v \right] \quad (3.15)$$

et

$$\langle x, \beta_k \rangle + b_k = \langle \Sigma^{1/2} v, \beta_k \rangle + b'_k, \quad \text{avec } b'_k = \langle m, \beta_k \rangle + b_k.$$

Posons

$$J^1 = \mathbb{E} [g(\langle x, \beta_k \rangle + b_k) \cdot \Sigma^{1/2} v].$$

Les coordonnées i , $i \in [d]$ sont données par

$$\begin{aligned} J_i^1 &= \mathbb{E} \left[g \left(\langle v, \Sigma^{\frac{1}{2}} \beta_k \rangle + b'_k \right) \cdot \left(\Sigma^{\frac{1}{2}} v \right)_i \right] \\ &= \sum_{r=1}^d \Sigma_{i,r}^{\frac{1}{2}} \mathbb{E} \left[g \left(\langle v, \Sigma^{\frac{1}{2}} \beta_k \rangle + b'_k \right) \cdot v_r \right] \end{aligned}$$

Par le lemme de Stein (Lemme 3.3.1.1), on a

$$\begin{aligned} J_i^1 &= \sum_{r=1}^d \Sigma_{i,r}^{\frac{1}{2}} \left(\Sigma^{\frac{1}{2}} \beta_k \right)_r \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] \\ &= \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] (\Sigma \beta_k)_i \end{aligned} \quad (3.16)$$

En remplaçant dans l'équation (3.14), on a

$$M_1(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g'(\langle \beta_k, x \rangle + b_k)] \cdot \beta'_k \text{ avec } \beta'_k = \Sigma \beta_k$$

2. Le moment d'ordre 2 :

Par définition

$$\mathbb{E}[y.x \otimes x] = \sum_{k \in [K]} \omega_k \mathbb{E} [g(\langle x, \beta_k \rangle + b_k) . x \otimes x]. \quad (3.17)$$

Comme

$$x \otimes x = x \otimes m + m \otimes x - m \otimes m + (\Sigma^{\frac{1}{2}} v) \otimes (\Sigma^{\frac{1}{2}} v),$$

l'équation (3.17) s'écrit

$$\begin{aligned} \mathbb{E}[y.x \otimes x] &= \mathbb{E} [y.(x \otimes m + m \otimes x - m \otimes m)] \\ &+ \sum_{k \in [K]} \omega_k \mathbb{E} \left[g(\langle x, \beta_k \rangle + b_k) . \left(\Sigma^{\frac{1}{2}} v \right) \otimes \left(\Sigma^{\frac{1}{2}} v \right) \right]. \end{aligned} \quad (3.18)$$

Posons

$$J^2 = \mathbb{E} \left[g(\langle x, \beta_k \rangle + b_k) . \left(\Sigma^{\frac{1}{2}} v \right) \otimes \left(\Sigma^{\frac{1}{2}} v \right) \right].$$

Les coordonnées $i, j \in [d]$ sont donnés par

$$\begin{aligned} J_{i,j}^2 &= \mathbb{E} \left[g(\langle v, \Sigma^{\frac{1}{2}} \beta_k \rangle + b'_k) . \left(\Sigma^{\frac{1}{2}} v \right)_i \left(\Sigma^{\frac{1}{2}} v \right)_j \right] \\ &= \sum_{r,s=1}^d \Sigma_{i,r}^{\frac{1}{2}} \Sigma_{j,s}^{\frac{1}{2}} \mathbb{E} \left[g(\langle v, \Sigma^{\frac{1}{2}} \beta_k \rangle + b'_k) . v_r v_s \right]. \end{aligned}$$

Par le lemme de Stein (lemme 3.3.1.1), on a

$$\begin{aligned} J_{i,j}^2 &= \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] . \left(\sum_{r=1}^d \Sigma_{i,r}^{\frac{1}{2}} \left(\Sigma^{\frac{1}{2}} \beta_k \right)_r \right) \left(\sum_{s=1}^d \Sigma_{j,s}^{\frac{1}{2}} \left(\Sigma^{\frac{1}{2}} \beta_k \right)_s \right) \\ &+ \sum_{r=1}^d \Sigma_{i,r}^{\frac{1}{2}} \Sigma_{r,j}^{\frac{1}{2}} \mathbb{E} [g(\langle x, \beta_k \rangle + b_k)] \\ &= \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] . (\Sigma \beta_k)_i (\Sigma \beta_k)_j + \mathbb{E} [g(\langle x, \beta_k \rangle + b_k)] . \Sigma_{i,j} \end{aligned}$$

En remplaçant dans l'équation (3.18), on a

$$M_2(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E} [g''(\langle \beta_k, x \rangle + b_k)] \cdot \beta'_k \otimes \beta'_k \text{ avec } \beta'_k = \Sigma \beta_k$$

3. Le moment d'ordre 3 :

Par définition, on a

$$\mathbb{E}[y.x \otimes x \otimes x] = \sum_{k \in [K]} \omega_k \mathbb{E}[g(\langle x, \beta_k \rangle + b_k) .x \otimes x \otimes x] \quad (3.19)$$

Commençons par évaluer

$$\begin{aligned} x \otimes x \otimes x &= x \otimes x \otimes m + m \otimes x \otimes x + x \otimes m \otimes x + m \otimes m \otimes m \\ &- x \otimes m \otimes m - m \otimes m \otimes x - m \otimes x \otimes m + (\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v) \end{aligned}$$

En prenant l'espérance, on a

$$\begin{aligned} \mathbb{E}[g(\langle x, \beta_k \rangle + b_k) .x \otimes x \otimes x] &= \mathbb{E}\left[g(\langle x, \beta_k \rangle + b_k) \left(x \otimes x \otimes m + m \otimes x \otimes x + x \otimes m \otimes x \right. \right. \\ &+ m \otimes m \otimes m - x \otimes m \otimes m - m \otimes m \otimes x - m \otimes x \otimes m \left. \left. \right) \right. \\ &+ \left. \mathbb{E}\left[g(\langle x, \beta_k \rangle + b_k) .(\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v)\right] \right. \end{aligned}$$

Posons

$$J^3 = \mathbb{E}\left[g(x'_k) .(\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v) \otimes (\Sigma^{\frac{1}{2}}v)\right]. \quad (3.20)$$

on a alors

$$J_{a,b,c}^3 = \sum_{r,s,t} \Sigma^{\frac{1}{2}}[a, r] \Sigma^{\frac{1}{2}}[b, s] \Sigma^{\frac{1}{2}}[c, t] \mathbb{E}\left[g(x'_k) .v_r v_s v_t\right] \quad (3.21)$$

Par le lemme de Stein (lemme 3.3.1.1), on a

$$\mathbb{E}\left[g(x'_k) .v_r v_s v_t\right] = \alpha_k[r] \alpha_k[s] \alpha_k[t] \mathbb{E}\left[g^{(3)}(x'_k)\right] + \alpha_k[t] \mathbb{E}\left[g'(x'_k)\right] \mathbf{1}_{\{r=s\}} \quad (3.22)$$

$$+ \alpha_k[r] \mathbb{E}\left[g'(x'_k)\right] \mathbf{1}_{\{s=t\}} + \alpha_k[s] \mathbb{E}\left[g'(x'_k)\right] \mathbf{1}_{\{r=t\}} \quad (3.23)$$

avec $\alpha_k = \Sigma^{\frac{1}{2}}\beta$. En remplaçant dans l'équation (3.21), on a

$$Z[a, b, c] = \mathbb{E}\left[g^{(3)}(x'_k) .(\Sigma^{\frac{1}{2}}\alpha_k)^{\otimes 3}\right][a, b, c] + A + B + C \quad (3.24)$$

où

$$A = \Sigma[b, c] \mathbb{E}\left[g'(x'_k)\right] \beta'_k[a], \quad (3.25)$$

$$B = \Sigma[a, c] \mathbb{E}\left[g'(x'_k)\right] \beta'_k[b], \quad (3.26)$$

et

$$C = \Sigma[a, b] \mathbb{E}\left[g'(x'_k)\right] \beta'_k[c] \quad (3.27)$$

En remplaçant dans l'équation (3.19), on a

$$M_3(\theta) = \sum_{k \in [K]} \omega_k \mathbb{E}[g^{(3)}(\langle \beta_k, x \rangle + b_k) .\beta'_k \otimes \beta'_k \otimes \beta'_k].$$

Dans les modèles 3.2.2.1 et 3.2.2.2, on peut réécrire les moments symétriques en fonction des β_k normalisés comme suit :

$$M_1(\theta) = \sum_{k=1}^K r_k^{(1)} \cdot \mu_k \quad (3.28)$$

$$M_2(\theta) = \sum_{k=1}^K r_k^{(2)} \cdot \mu_k \otimes \mu_k \quad (3.29)$$

$$M_3(\theta) = \sum_{k=1}^K r_k^{(3)} \cdot \mu_k \otimes \mu_k \otimes \mu_k \quad (3.30)$$

avec

$$\begin{aligned} r_k^{(1)} &= \omega_k \lambda_k \mathbb{E}[g'(\lambda_k \langle x, \mu_k \rangle + b_k)], \\ r_k^{(2)} &= \omega_k \lambda_k^2 \mathbb{E}[g''(\lambda_k \langle x, \mu_k \rangle + b_k)], \end{aligned}$$

et

$$r_k^{(3)} = \omega_k \lambda_k^3 \mathbb{E}[g^{(3)}(\lambda_k \langle x, \mu_k \rangle + b_k)].$$

Rappelons que $\lambda_k = \|\beta_k\|_2$ pour le modèle 3.2.2.1 et $\lambda_k = \|\beta'_k\|_2$ pour le modèle 3.2.2.2, $k = 1, \dots, K$.

Double diagonalisation

La double diagonalisation consiste à passer par deux étapes de diagonalisation à partir du moment d'ordre 3 pour estimer les vecteurs normalisés de β . L'idée est de construire un tenseur orthogonal à partir d'un tenseur symétrique d'ordre 3 présenté au début de la section 3.3.1. C'est-à-dire, écrire le tenseur symétrique donné par le lemme 3.3.1.2 sous forme d'un tenseur orthogonal pour construire les vecteurs propres du tenseur en question. On commence par construire une matrice permettant d'orthogonaliser le tenseur M_3 .

• Construction de la matrice d'orthogonalisation

A partir du tenseur M_3 , on construit des matrices de la façon suivante. Soit $\xi \in \mathbb{R}^d$,

$$M_2^\xi = M_3(I, I, \xi) = \sum_{k \in [K]} r_k^{(3)} \langle \xi, \mu_k \rangle \mu_k^{\otimes 2}.$$

Ainsi

$$M_2^\xi = \sum_{k \in [K]} \alpha_k(\xi) \mu_k^{\otimes 2}, \text{ avec } \alpha_k(\xi) = r_k^{(3)} \langle \xi, \mu_k \rangle.$$

La matrice M_2^ξ est symétrique, donc par diagonalisation, on a

$$M_2^\xi = U D U^T$$

avec

- $U \in \mathbb{R}^{d \times K}$ matrice orthogonale des vecteurs propres de M_2^ξ ,
- $D \in \mathbb{R}^{K \times K}$ matrice diagonale des valeurs propres de M_2^ξ .

Supposons que les éléments de la diagonale de la matrice D sont strictement positifs. Définissons la matrice d'orthogonalisation $W := UD^{-\frac{1}{2}} \in \mathbb{R}^{d \times K}$. Alors, on a

$$M_2^\xi(W, W) = W^T M_2^\xi W = D^{-\frac{1}{2}} U^T U D U^T U D^{-\frac{1}{2}} = I_K.$$

Or par définition, on a

$$M_2^\xi(W, W) = \sum_{k \in [K]} W^T \left(\sqrt{\alpha_k(\xi)} \mu_k \right) \left(\sqrt{\alpha_k(\xi)} \mu_k \right)^T W = \sum_{k \in [K]} \tilde{\mu}_k \tilde{\mu}_k^T, \text{ avec } \tilde{\mu}_k = \sqrt{\alpha_k(\xi)} W^T \mu_k.$$

Par suite, les $\tilde{\mu}_k$ définis sont orthonormaux. A partir de W et M_3 , on peut définir un tenseur orthogonal

$$\tilde{M}_3 := M_3(W, W, W) \in \mathbb{R}^{K \times K \times K}.$$

En effet, on a

$$\tilde{M}_3 = \sum_{k \in [K]} r_k^{(3)} (W^T \mu_k)^{\otimes 3} = \sum_{k \in [K]} \frac{1}{\sqrt{\alpha_k(\xi)} \langle \mu_k, \xi \rangle} (\sqrt{\alpha_k(\xi)} W^T \mu_k)^{\otimes 3}.$$

Par suite, on a un tenseur orthogonal

$$\tilde{M}_3 = \sum_{k \in [K]} \tilde{\omega}_k \tilde{\mu}_k^{\otimes 3}, \text{ avec } \tilde{\omega}_k = \frac{1}{\sqrt{\alpha_k(\xi)} \langle \mu_k, \xi \rangle} \text{ et } \tilde{\mu}_k = \sqrt{\alpha_k(\xi)} W^T \mu_k.$$

A partir du tenseur orthogonal construit, on peut estimer les vecteurs de paramètres normalisés. En effet, pour $\rho \in \mathbb{R}^K$, on peut construire une matrice orthogonale \tilde{M}_2 :

$$\tilde{M}_2 := \tilde{M}_3(I, I, \rho) = \sum_{k \in [K]} \tilde{\omega}_k \langle \tilde{\mu}_k, \rho \rangle \tilde{\mu}_k^{\otimes 2} \quad (3.31)$$

$$= \sum_{k \in [K]} \tilde{\alpha}_k \tilde{\mu}_k^{\otimes 2} \text{ avec } \tilde{\alpha}_k = \tilde{\omega}_k \langle \tilde{\mu}_k, \rho \rangle \quad (3.32)$$

On a $(\tilde{\alpha}_k, \tilde{\mu}_k) = (\text{valeur propre, vecteur propre})$ de \tilde{M}_2 qui est une matrice symétrique. Donc par diagonalisation, on peut obtenir les K vecteurs propres de \tilde{M}_2 dont les valeurs propres associées sont non nulles. C'est-à-dire

$$\tilde{M}_2 = \tilde{U} \tilde{D} \tilde{U}^T$$

avec

- $\tilde{U} \in \mathbb{R}^{K \times K}$ matrice orthogonale des vecteurs propres de \tilde{M}_2 ,
- $\tilde{D} \in \mathbb{R}^{K \times K}$ matrice diagonale des valeurs propres de \tilde{M}_2 .
- **Retour vers les vecteurs de paramètres normalisés**

On a

$$\tilde{\mu}_k = \sqrt{\alpha_k(\xi)} W^T \mu_k.$$

Si B est l'inverse de Moore-Penrose de W^T , on peut retrouver les paramètres μ_k par

$$\tilde{\omega} B \tilde{\mu}_k = \frac{\langle \tilde{\mu}_k, \rho \rangle}{\langle \mu_k, \xi \rangle} \mu_k, \quad k = 1, \dots, K$$

pour les mêmes raisons que le point 3 du théorème 4.3 de [36].

L'étape de construction de la matrice de d'orthogonalisation est appelée, dans la littérature, l'étape de blanchissement (ou Whitening) et la matrice W est appelé *Whitening matrix*. Notons que pour la mettre en oeuvre, il faut trouver un vecteur ξ tel que les valeurs propres de M_2^ξ soient strictement positives. Même si cela est possible, en pratique, cette méthode reste assez longue à exécuter du fait de la recherche d'un tel ξ .

Diagonalisation jointe

Soient z_1, z_2, \dots, z_P , P vecteurs de \mathbb{R}^d . On peut définir un ensemble de matrices $\{B(z_p)\}_{p=1}^P$, par

$$B(z_p) := M_3(I, I, z_p), \quad p = 1, \dots, P,$$

c'est-à-dire

$$B(z_p) = \sum_{k=1}^K r_k^{(3)} \cdot \langle \mu_k, z_p \rangle \mu_k^{\otimes 2}.$$

La méthode de "diagonalisation jointe" consiste à trouver, pour un ensemble de matrices symétriques données $\{B(z_p)\}_{p=1}^P$, une matrice V telle que les matrices $VB(z_p)V^T$ soient le plus diagonales possible, c'est-à-dire minimisant la somme des distances de Frobenius des $VB(z_p)V^T$ aux matrices $diag(a_p)$, $p = 1, \dots, P$, pour des $a_p \in \mathbb{R}^d$.

Même si cette méthode est appelée "diagonalisation jointe" dans la littérature, on ne fait pas une diagonalisation proprement dite des $B(z_p)$. En effet, les vecteurs de V^{-1} ne sont pas forcément des vecteurs propres associées aux matrices $B(z_p)$. C'est-à-dire les vecteurs de V^{-1} ne forment pas forcément une base orthonormée. Cette méthode de "diagonalisation jointe" nous permet de retrouver, aux signes et aux permutations près, les vecteurs des paramètres normalisés. Il suffit de partir sur un nombre suffisant de matrice symétriques $\{B(z_p)\}_{p=1}^P$. Le nombre minimal de matrices nécessaire et génériquement suffisant est de 2 avec $P \leq d$; c'est-à-dire $2 \leq P \leq d$. Pour plus de détails, voir [40].

On obtient les vecteurs de paramètres normalisés, au signe et à permutation près, en prenant les K premiers vecteurs de la matrice V^{-1} . Soit U la matrice constituée de ces K vecteurs. Alors U est de dimension $d \times K$. Ainsi, on peut retrouver les signes de tous les vecteurs de paramètres normalisés. En effet, on a

$$B(z_p) = V^{-1}Diag(a_p)(V^{-1})^T \quad (3.33)$$

$$= \sum_{k=1}^d (a_p)_k (V^{-1})_k (V^{-1})_k^T \quad (3.34)$$

Remarquons que $U_k = (V^{-1})_k$ pour $k = 1, \dots, K$ et que les $r_k^{(1)} = \omega_k \lambda_k \mathbb{E}[g'(\langle x, \mu_k \rangle + b_k)]$ sont positifs dans l'équation (3.28). Posons $O = U_*^{-1}M_1(\theta) \in \mathbb{R}^K$, avec $M_1(\theta)$ donnée par l'équation (3.28) et U_*^{-1} l'inverse généralisé de la matrice U . Les éléments de O sont égaux aux $r_k^{(1)}$, $k = 1, \dots, K$ aux signes près. On obtient alors les signes des U_k en multipliant par -1 tous les vecteurs U_k associés aux valeurs négatives de O .

En pratique, on obtient des estimateurs des U_k en mettant en oeuvre la méthode de diagonalisation jointe à partir de \hat{M}_1 , \hat{M}_2 et \hat{M}_3 les estimateurs respectifs des moments d'ordre 1, 2 et 3, calculés à partir des entrées X et de la réponse Y . C'est-à-dire

$$- \hat{M}_1 = \frac{1}{n} \sum_{i=1}^n y_i \cdot x_i,$$

- $\hat{M}_2 = \frac{1}{n} \sum_{i=1}^n y_i \cdot (x_i \otimes x_i - I_d)$, et
- $\hat{M}_3 = \frac{1}{n} \sum_{i=1}^n y_i \cdot x_i \otimes x_i \otimes x_i - \sum_{j \in [d]} \frac{1}{n} \sum_{i=1}^n \left[y_i (e_j \otimes x_i \otimes e_j + e_j \otimes e_j \otimes x_i + x_i \otimes e_j \otimes e_j) \right]$.

Cette méthode reste parfois instable en pratique et ne permet pas de retrouver tous les paramètres du modèle. Mais elle pourra servir d'initialisation dans l'estimation globale des paramètres du modèle.

3.3.2 Estimation de tous les paramètres du modèle

Dans cette partie, nous cherchons à estimer tous les paramètres du modèle. La question que nous posons est alors : peut-on retrouver tous les paramètres du modèle à partir des moments $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$? La réponse est oui, et on obtiendra un estimateur consistant en minimisant un critère du type moindres carrés, défini par

$$\begin{aligned}
Q_n(\theta) &= \sum_{j \in [d]} \left\{ \hat{M}_1[j] - M_1(\theta)[j] \right\}^2 \\
&+ \sum_{j, k \in [d]} \left\{ \hat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 \\
&+ \sum_{j, k, l \in [d]} \left\{ \hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2.
\end{aligned} \tag{3.35}$$

Sous certaines hypothèses, on obtiendra la consistance et la normalité asymptotique de l'estimateur. On notera $G_r \in \mathbb{R}^{K \times K}$, $r = 1, \dots, 5$, les matrices diagonales dont les éléments de la diagonale sont donnés par les $\mathbb{E}[g^{(r)}(\langle x, \beta_k \rangle + b_k)]$, $k = 1, \dots, K$.

- (H1) Les vecteurs β_1, \dots, β_K sont linéairement indépendants et les poids sont positifs : $\omega_k > 0$, $k = 1, \dots, K$.
- (H2) La fonction lien g est strictement croissante, infiniment dérivable avec une dérivée décroissante sur $[0, +\infty[$, et satisfait

$$\forall z \in \mathbb{R}, g(z) + g(-z) = 1.$$

- (H3) Tous les coefficients de la diagonale de G_3 sont non nuls.
- (H4) Tous les coefficients de la diagonale de $G_1 G_3 - G_2^2$ sont non nuls.
- (H5) Il existe un voisinage \mathcal{O} de $(0, 0)$ dans $\mathbb{R}_+^* \times \mathbb{R}$ et des fonctions L_s , $s = 1, 2, 3$, telles que $\forall z \in \mathbb{R}, \forall (\lambda, b) \in \mathcal{O}$, on a

$$(|z| + 1) \left| \frac{\partial g^{(s+1)}}{\partial \lambda}(\lambda z + b) \right| \leq L_s(z)$$

et de plus pour $s = 1, 2, 3$,

$$\int_{\mathbb{R}} L_s(z) e^{-z^2/2} dz < +\infty.$$

Si l'hypothèse (H2) est vraie, on a les propriétés suivantes :

- (P1) la fonction g' est positive et satisfait $g'(x) = g'(-x)$ pour tout $x \in \mathbb{R}$,
- (P2) la fonction g'' satisfait $g''(x) = -g''(-x)$ pour tout $x \in \mathbb{R}$ et $g''(x) < 0$ pour $x > 0$.

Notons que l'hypothèse (H1) implique que $d \geq K$ et que (H2) est vraie dans le cas d'un lien probit où la fonction g est telle que $g(z) = \Phi(z)$, avec Φ la fonction de répartition d'une loi normale centrée réduite. L'hypothèse (H2) est aussi vraie dans le cas d'un lien logistique où la fonction g est telle que $g(z) = e^z/(1 + e^z)$.

3.4 Résultats théoriques

Nous donnerons dans cette partie les résultats théoriques qui permettent de répondre à la question posée précédemment, à savoir l'identifiabilité du modèle, la consistance et la normalité asymptotique de l'estimateur. Dans la suite, nous présenterons les résultats pour le modèle 3.2.2.1. Les résultats du modèle 3.2.2.2 peuvent être retrouvés de manière similaire.

3.4.1 Identifiabilité

Le but est de prouver qu'à partir des moments $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$, qui peuvent être estimés en utilisant les données X et Y , nous pouvons retrouver les paramètres du modèle ω , β , b . Rappelons que $\theta = (\omega, \beta, b)$ et que les moments sont donnés par les équations (3.28), (3.29) et (3.30).

En utilisant $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$ on peut retrouver les 3-uplets

$$(\omega_k \mathbb{E}[g'(\lambda_k \langle \mu_k, x \rangle + b_k)] \lambda_k; \omega_k \mathbb{E}[g''(\lambda_k \langle \mu_k, x \rangle + b_k)] \lambda_k^2; \omega_k \mathbb{E}[g^{(3)}(\lambda_k \langle \mu_k, x \rangle + b_k)] \lambda_k^3), \quad k = 1, \dots, K.$$

En effet, la diagonalisation jointe, décrite à la section 3.3.2, permet de retrouver entièrement les μ_k et les $r_k^{(1)}$ aux permutations près, à partir des moments $M_1(\theta)$ et $M_3(\theta)$. Par suite, il suffit d'utiliser les μ_k , les moments $M_2(\theta)$ et $M_3(\theta)$ pour retrouver les $r_k^{(2)}$ et les $r_k^{(3)}$.

Comme $z = \langle \mu, x \rangle \sim \mathcal{N}(0, 1)$ et a donc pour densité $\frac{1}{\sqrt{2\pi}} e^{-z^2/2}$, on a l'identifiabilité du modèle si la fonction qui à tout (ω, λ, b) dans $]0, +\infty[\times]0, +\infty[\times \mathbb{R}$ associe

$$\left(\omega \lambda \int g'(\lambda z + b) e^{-z^2/2} dz; \omega \lambda^2 \int g''(\lambda z + b) e^{-z^2/2} dz; \omega \lambda^3 \int g^{(3)}(\lambda z + b) e^{-z^2/2} dz \right)$$

est bijective.

En utilisant une intégration par partie, la bijectivité de la fonction est équivalent à montrer que la fonction définie dans $]0, +\infty[\times]0, +\infty[\times \mathbb{R}$ qui à tout (ω, λ, b) associe

$$\lambda \left(\omega \int g'(\lambda z + b) e^{-z^2/2} dz; \omega \int z g'(\lambda z + b) e^{-z^2/2} dz; \omega \int z^2 g'(\lambda z + b) e^{-z^2/2} dz \right)$$

est bijective. Ce qui est encore équivalent au fait que la fonction définie dans $]0, +\infty[\times \mathbb{R}$ qui à tout (λ, b) associe

$$\left(\frac{\int z g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz}; \frac{\int z^2 g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz} \right)$$

est bijective. Pour tout $(b, \lambda) \in \mathbb{R} \times]0, +\infty[$, définissons

$$dQ_{(b, \lambda)}(z) = \frac{g'(\lambda z + b) e^{-z^2/2}}{\int g'(\lambda z + b) e^{-z^2/2} dz} dz. \quad (3.36)$$

Alors l'identifiabilité revient à montrer que la fonction définie par

$$(\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) := \left(\int z dQ_{(b,\lambda)}(z); \int z^2 dQ_{(b,\lambda)}(z) \right) \quad (3.37)$$

est bijective.

Théorème 3.4.1.1 (*Identifiabilité probit*) *Supposons que l'hypothèse (H1) est vraie et que la fonction lien est probit, alors on peut retrouver le nombre de composants du mélange K et le paramètre $\theta = (\omega, \beta, b)$ à partir des moments $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$.*

Preuve du Théorème 3.4.1.1.

Comme l'hypothèse (H1) est vraie, alors le rang de $M_2(\theta)$ est égal à K . La fonction lien g est probit, on a alors

$$g(z) = \Phi(z)$$

et

$$g'(z) = \varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2},$$

avec Φ et φ les fonctions de répartition et de densité de la loi normale centrée réduite respectivement. En remplaçant dans l'équation (3.36), on a

$$\begin{aligned} dQ_{(b,\lambda)}(z) &= \frac{\varphi(\lambda z + b) e^{-z^2/2}}{\int \varphi(\lambda z + b) e^{-z^2/2} dz} dz \\ &= \frac{e^{-\frac{1}{2}((\lambda z + b)^2 + z^2)}}{\int e^{-\frac{1}{2}((\lambda z + b)^2 + z^2)} dz} dz. \end{aligned} \quad (3.38)$$

Comme

$$\begin{aligned} (\lambda z + b)^2 + z^2 &= (\lambda^2 + 1)z^2 + 2\lambda b z + b^2 \\ &= (\lambda^2 + 1) \left(z^2 + \frac{2\lambda b}{\lambda^2 + 1} z + \frac{b^2}{\lambda^2 + 1} \right) \\ &= (\lambda^2 + 1) \left\{ \left(z + \frac{\lambda b}{\lambda^2 + 1} \right)^2 + \frac{b^2}{(\lambda^2 + 1)^2} \right\}, \end{aligned}$$

on a, à partir de l'équation (3.38), que

$$dQ_{(b,\lambda)}(z) = \frac{e^{-\frac{1}{2}(\lambda^2+1)\left(z+\frac{\lambda b}{\lambda^2+1}\right)^2}}{\int e^{-\frac{1}{2}(\lambda^2+1)\left(z+\frac{\lambda b}{\lambda^2+1}\right)^2} dz} dz.$$

Par suite, on a

$$dQ_{(b,\lambda)}(z) = \frac{\sqrt{\lambda^2 + 1}}{\sqrt{2\pi}} e^{-\frac{1}{2}(\lambda^2+1)\left(z+\frac{\lambda b}{\lambda^2+1}\right)^2} dz.$$

Ainsi, $Q_{(b,\lambda)} = \mathcal{N}\left(-\frac{\lambda b}{\lambda^2+1}; \frac{1}{\lambda^2+1}\right)$ et les moments d'ordre 1 et 2 de Z sont donnés par

$$(\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) = \left(-\frac{\lambda b}{\lambda^2 + 1}; \frac{\lambda^2 b^2 + \lambda^2 + 1}{(\lambda^2 + 1)^2} \right).$$

En posant

$$\begin{cases} -\frac{\lambda b}{\lambda^2+1} = \alpha_1 \\ \frac{\lambda^2 b^2 + \lambda^2 + 1}{(\lambda^2+1)^2} = \alpha_2, \end{cases} \quad (3.39)$$

on a $\alpha_2 - \alpha_1^2 = \frac{1}{\lambda^2+1}$ ce qui implique que

$$\lambda = ((\alpha_2 - \alpha_1^2)^{-1} - 1)^{\frac{1}{2}}.$$

A partir de la première ligne de l'équation (3.39), on a

$$b = -\alpha_1 \frac{(\lambda^2 + 1)}{\lambda}.$$

Comme

$$\lambda^2 + 1 = (\alpha_2 - \alpha_1^2)^{-1},$$

on a par suite

$$b = -\alpha_1 (\alpha_2 - \alpha_1^2)^{-1} [(\alpha_2 - \alpha_1^2)^{-1} - 1]^{\frac{1}{2}}.$$

Ainsi, on retrouve b et λ si $\lambda > 0$.

Dans le cas général, l'identifiabilité est obtenue au moins dans un ouvert. Pour le montrer, il suffit de prouver que pour un certain $B > 0$ et $L > 0$, si l'hypothèse (H2) est vérifiée, alors, la fonction qui à tout $(b, \lambda) \in]-B, B[\times]0, L[$ associe $(E_{(b,\lambda)}(Z); E_{(b,\lambda)}(Z^2))$ est bijective.

Théorème 3.4.1.2 *Supposons que les hypothèses (H1), (H2) et (H5) sont vraies. Si $g^{(3)}(0) \neq 0$, alors il existe $L > 0$ et $B > 0$ tels que si $\|\beta_k\| < L$ et $|b_k| < B$, on peut retrouver K et θ à partir de $M_1(\theta)$, $M_2(\theta)$ et $M_3(\theta)$.*

Preuve du Théorème 3.4.1.2.

La fonction définie dans l'équation (3.37) est donnée par

$$\begin{aligned} (\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) &= \left(\int z dQ_{(b,\lambda)}(z); \int z^2 dQ_{(b,\lambda)}(z) \right) \\ &= \left(\frac{\int z g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz}; \frac{\int z^2 g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz} \right). \end{aligned} \quad (3.40)$$

En utilisant une intégration par parties, on a

$$(\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) = \left(\frac{\lambda \int g''(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz}; 1 + \frac{\lambda^2 \int g^{(3)}(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz} \right). \quad (3.41)$$

Comme l'hypothèse (H2) est vraie, on a

- (P3) $g'(0) > 0$
- (P4) $g''(0) = g^{(4)}(0) = 0$

Définissons pour tout z les fonctions K_s , $s = 1, 2, 3$ telles

$$K_s : \mathbb{R}_+^* \times \mathbb{R} \rightarrow \mathbb{R} \\ (\lambda, b) \mapsto K_s(\lambda, b)$$

avec

$$K_s(\lambda, b) = \int g^{(s)}(\lambda z + b) e^{-z^2/2} dz.$$

Un développement de Taylor autour de $(0, 0)$ des fonctions K_s donne, pour tout z :

$$K_s(\lambda, b) = K_s(0, 0) + \langle \nabla K_s(0, 0), (\lambda, b) \rangle + o(\lambda^2 + b^2) \quad (3.42)$$

On a

$$\nabla K_s(\lambda, b) = \left(\frac{\partial K_s}{\partial \lambda}(\lambda, b); \frac{\partial K_s}{\partial b}(\lambda, b) \right)$$

avec

$$\frac{\partial K_s}{\partial \lambda}(\lambda, b) = \frac{\partial}{\partial \lambda} \int g^{(s)}(\lambda z + b) e^{-z^2/2} dz \quad (3.43)$$

et

$$\frac{\partial K_s}{\partial b}(\lambda, b) = \frac{\partial}{\partial b} \int g^{(s)}(\lambda z + b) e^{-z^2/2} dz. \quad (3.44)$$

Comme $(H5)$ est vraie, les équations (3.43) et (3.44) sont données au point $(0, 0)$ par

$$\frac{\partial K_s}{\partial \lambda}(0, 0) = \int z g^{(s+1)}(0) e^{-z^2/2} dz \quad (3.45)$$

et

$$\frac{\partial K_s}{\partial b}(0, 0) = \int g^{(s+1)}(0) e^{-z^2/2} dz. \quad (3.46)$$

Par suite

$$K_s(\lambda, b) = g^{(s)}(0) \int e^{-z^2/2} dz + g^{(s+1)}(0) \int (\lambda z + b) e^{-z^2/2} dz + o(\lambda^2 + b^2). \quad (3.47)$$

En utilisant $(P4)$ et l'équation (3.47), on a

$$\int g'(\lambda z + b) e^{-z^2/2} dz = \sqrt{2\pi} g'(0) + o(\lambda^2 + b^2), \quad (3.48)$$

$$\int g''(\lambda z + b) e^{-z^2/2} dz = \sqrt{2\pi} g^{(3)}(0) b + o(\lambda^2 + b^2) \quad (3.49)$$

et

$$\int g^{(3)}(\lambda z + b)e^{-z^2/2} dz = \sqrt{2\pi}g^{(3)}(0) + o(\lambda^2 + b^2). \quad (3.50)$$

Par suite, en remplaçant (3.48),(3.49) et (3.50), dans l'équation (3.41), on a

$$E_{(b,\lambda)}(Z) = \frac{g^{(3)}(0)}{g'(0)}\lambda b + o(\lambda^2 + b^2)$$

et

$$E_{(b,\lambda)}(Z^2) = 1 + \frac{g^{(3)}(0)}{g'(0)}\lambda^2 + o(\lambda^2 + b^2).$$

Ainsi, en utilisant le théorème d'inversion local, on a

$$\lambda^2 = \frac{g'(0)}{g^{(3)}(0)} (E_{(b,\lambda)}(Z)^2 - 1) + o(|E_{(b,\lambda)}(Z)^2 - 1| + |E_{(b,\lambda)}(Z)|)$$

et

$$\lambda b = \frac{g'(0)}{g^{(3)}(0)} E_{(b,\lambda)}(Z) + o(|E_{(b,\lambda)}(Z)^2 - 1| + |E_{(b,\lambda)}(Z)|).$$

Ce qui nous donne le résultat général de l'identifiabilité.

3.4.2 Consistance

Le paramètre θ pourra être estimé par

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} Q_n(\theta) \quad (3.51)$$

avec

$$\begin{aligned} Q_n(\theta) &= \sum_{j \in [d]} \left\{ \hat{M}_1[j] - M_1(\theta)[j] \right\}^2 + \sum_{j,k \in [d]} \left\{ \hat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 \\ &+ \sum_{j,k,l \in [d]} \left\{ \hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2 \end{aligned}$$

En utilisant la loi des grands nombres, on a que \hat{M}_r converge en probabilité vers $M_r(\theta^*)$, $r = 1, 2, 3$.

Théorème 3.4.2.1 *Supposons que le modèle est identifiable et Θ compact. Si les hypothèses de (H1) à (H5) sont vraies, alors $\hat{\theta}_n$ converge en probabilité vers θ^* .*

Preuve du Théorème 3.4.2.1.

Définissons Q par :

$$\begin{aligned}
Q(\theta) &= \sum_{j \in [d]} \left\{ M_1(\theta^*)[j] - M_1(\theta)[j] \right\}^2 + \sum_{j, k \in [d]} \left\{ M_2(\theta^*)[j, k] - M_2(\theta)[j, k] \right\}^2 \\
&+ \sum_{j, k, l \in [d]} \left\{ M_3(\theta^*)[j, k, l] - M_3(\theta)[j, k, l] \right\}^2
\end{aligned}$$

Pour $\hat{\theta}_n$, Q_n et Q définis plus haut, on a

$$\begin{aligned}
Q_n(\theta) - Q(\theta) &= \sum_{j \in [d]} \left\{ \hat{M}_1[j] - M_1(\theta)[j] \right\}^2 - \left\{ M_1(\theta^*)[j] - M_1(\theta)[j] \right\}^2 \\
&+ \sum_{j, k \in [d]} \left\{ \hat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 - \left\{ M_2(\theta^*)[j, k] - M_2(\theta)[j, k] \right\}^2 \\
&+ \sum_{j, k, l \in [d]} \left\{ \hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2 - \left\{ M_3(\theta^*)[j, k, l] - M_3(\theta)[j, k, l] \right\}^2 \\
&= \sum_{j \in [d]} \left(\hat{M}_1[j] - M_1(\theta^*)[j] \right) \left(\hat{M}_1[j] + M_1(\theta^*)[j] - 2M_1(\theta)[j] \right) \\
&+ \sum_{j, k \in [d]} \left(\hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right) \left(\hat{M}_2[j, k] + M_2(\theta^*)[j, k] - 2M_2(\theta)[j, k] \right) \\
&+ \sum_{j, k, l \in [d]} \left(\hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right) \left(\hat{M}_3[j, k, l] + M_3(\theta^*)[j, k, l] - 2M_3(\theta)[j, k, l] \right)
\end{aligned}$$

$$\begin{aligned}
|Q_n(\theta) - Q(\theta)| &\leq \sum_{j \in [d]} \left(\left| \hat{M}_1[j] - M_1(\theta^*)[j] \right| \right) \left(\left| \hat{M}_1[j] \right| + \left| M_1(\theta^*)[j] \right| + 2\left| M_1(\theta)[j] \right| \right) \\
&+ \sum_{j, k \in [d]} \left(\left| \hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right| \right) \left(\left| \hat{M}_2[j, k] \right| + \left| M_2(\theta^*)[j, k] \right| + 2\left| M_2(\theta)[j, k] \right| \right) \\
&+ \sum_{j, k, l \in [d]} \left(\left| \hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right| \right) \left(\left| \hat{M}_3[j, k, l] \right| \right. \\
&\quad \left. + \left| M_3(\theta^*)[j, k, l] \right| + 2\left| M_3(\theta)[j, k, l] \right| \right)
\end{aligned}$$

Posons

$$S = \sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)|$$

Ainsi, on a

$$\begin{aligned}
S &\leq \sum_{j \in [d]} \left(\left| \hat{M}_1[j] - M_1(\theta^*)[j] \right| \right) \left(\left| \hat{M}_1[j] \right| + \left| M_1(\theta^*)[j] \right| + 2 \sup_{\theta \in \Theta} \left| M_1(\theta)[j] \right| \right) \\
&+ \sum_{j, k \in [d]} \left(\left| \hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right| \right) \left(\left| \hat{M}_2[j, k] \right| + \left| M_2(\theta^*)[j, k] \right| + 2 \sup_{\theta \in \Theta} \left| M_2(\theta)[j, k] \right| \right) \\
&+ \sum_{j, k, l \in [d]} \left(\left| \hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right| \right) \left(\left| \hat{M}_3[j, k, l] \right| + \left| M_3(\theta^*)[j, k, l] \right| + 2 \sup_{\theta \in \Theta} \left| M_3(\theta)[j, k, l] \right| \right)
\end{aligned}$$

Comme $\theta \mapsto M_r(\theta)$, $r = 1, 2, 3$ est continue et Θ est compact, donc il existe c_1 , c_2 et c_3 tels que

$$\begin{aligned}
S &\leq \sum_{j \in [d]} \left(c_1 + \left| \hat{M}_1[j] \right| \right) \left(\left| \hat{M}_1[j] - M_1(\theta^*)[j] \right| \right) + \sum_{j, k \in [d]} \left(c_2 + \left| \hat{M}_2[j, k] \right| \right) \left(\left| \hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right| \right) \\
&+ \sum_{j, k, l \in [d]} \left(c_3 + \left| \hat{M}_3[j, k, l] \right| \right) \left(\left| \hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right| \right)
\end{aligned}$$

De plus, comme $\left| \hat{M}_r - M_r(\theta^*) \right| = o_{\mathbb{P}}(1)$ par la loi des grands nombres, alors

$$\sup_{\theta \in \Theta} \left| Q_n(\theta) - Q(\theta) \right| = o_{\mathbb{P}}(1). \quad (\star)$$

Soit $\epsilon > 0$ et $A_\epsilon = \{\theta \in \Theta : d(\theta, \theta^*) \geq \epsilon\}$. Posons $\alpha_\epsilon = \inf_{\theta \in A_\epsilon} Q(\theta)$.

La continuité de la distance et la compacité de Θ nous donnent la compacité de A_ϵ . Ainsi il existe $\theta_\epsilon \in A_\epsilon$ tel que $\alpha_\epsilon = Q(\theta_\epsilon)$. Puis

$$\theta_\epsilon \in A_\epsilon \implies d(\theta_\epsilon, \theta^*) \geq \epsilon$$

Ce qui implique que $\theta_\epsilon \neq \theta^*$. Par identifiabilité, on a $Q(\theta_\epsilon) \neq Q(\theta^*) = 0$. Ce qui prouve que

$$\forall \epsilon > 0 \quad \inf_{\theta \in \Theta : d(\theta, \theta^*) \geq \epsilon} Q(\theta) > Q(\theta^*). \quad (\star\star)$$

Par définition de Q_n et $\hat{\theta}_n$ (Equation (3.51)), on a

$$Q_n(\hat{\theta}_n) \leq \inf_{\theta \in \Theta} Q_n(\theta) + u_n, \quad \text{avec } u_n \xrightarrow{n \rightarrow \infty} 0. (\star\star\star)$$

En utilisant (\star) , $(\star\star)$, $(\star\star\star)$ et le théorème 5.7 (page 45 du livre de Van der Vaart [77]), on obtient la consistance de $\hat{\theta}_n$.

3.4.3 Normalité Asymptotique

En utilisant les hypothèses de (H1) à (H4), on peut montrer la normalité asymptotique de notre estimateur $\hat{\theta}_n$ de θ^* .

Théorème 3.4.3.1 *Supposons que les hypothèses de (H1) à (H4) sont vraies et que $\hat{\theta}_n$ est consistant pour θ^* . Alors $\sqrt{n}(\hat{\theta}_n - \theta^*)$ converge en loi vers une loi normale centrée dont on notera Σ la variance.*

Rappelons que la variance Σ dépend du paramètre inconnu θ^* . En pratique, on peut l'estimer en utilisant la méthode du Bootstrap.

Preuve du Théorème 3.4.3.1.

Définissons Z_n par

$$Z_n(\theta) = \nabla_{\theta} Q_n(\theta), \quad \theta \in \Theta, \quad (3.52)$$

avec

$$\begin{aligned} Q_n(\theta) &= \sum_{j \in [d]} \left\{ \hat{M}_1[j] - M_1(\theta)[j] \right\}^2 + \sum_{j,k \in [d]} \left\{ \hat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 \\ &+ \sum_{j,k,l \in [d]} \left\{ \hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2. \end{aligned}$$

Les composantes de $Z_n(\theta)$ peuvent être obtenues par

$$\begin{aligned} \frac{\partial Q_n(\theta)}{\partial \theta_r} &= -2 \left\{ \sum_{j \in [d]} \frac{\partial M_1(\theta)[j]}{\partial \theta_r} \left[\hat{M}_1[j] - M_1(\theta)[j] \right] + \sum_{j,k \in [d]} \frac{\partial M_2(\theta)[j, k]}{\partial \theta_r} \left[\hat{M}_2[j, k] - M_2(\theta)[j, k] \right] \right. \\ &\left. + \sum_{j,k,l \in [d]} \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_r} \left[\hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right] \right\} \end{aligned}$$

En utilisant un développement de Taylor, on a

$$Z_n(\hat{\theta}_n) = Z_n(\theta^*) + \int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] (\hat{\theta}_n - \theta^*) dt \quad (3.53)$$

où $D_1 Z_n$ est la matrice dérivée de Z_n . Comme $Z_n(\hat{\theta}_n) = 0$, on a

$$-\sqrt{n} Z_n(\theta^*) = \left[\int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] dt \right] \sqrt{n} (\hat{\theta}_n - \theta^*) \quad (3.54)$$

Posons

$$\hat{M} = \left(\hat{M}_1[j], \hat{M}_2[j, k], \hat{M}_3[j, k, l] \right)_{1 \leq j, k, l \leq d}$$

et

$$M(\theta^*) = \left(M_1(\theta^*)[j], M_2(\theta^*)[j, k], M_3(\theta^*)[j, k, l] \right)_{1 \leq j, k, l \leq d}$$

On a, par le théorème central limite, que

$$\sqrt{n}(\hat{M} - M(\theta^*)) \text{ converge en loi vers } \mathcal{N}(0, B) \quad (3.55)$$

avec B la matrice de covariance de \hat{M} . En utilisant la méthode delta et les équations (3.52), (3.55),

$$\sqrt{n} Z_n(\theta^*) \text{ converge en loi} \quad (3.56)$$

vers une loi normale de moyenne nulle, dont on notera L la variance.

On a $D_1 Z_n(\theta) = \nabla_{\theta}^2 Q_n(\theta)$. Les composantes de cette matrice sont données par

$$\begin{aligned} \frac{\partial^2 Q_n(\theta)}{\partial \theta_{r_1} \partial \theta_{r_2}} &= -2 \left\{ \sum_{j \in [d]} \left[\frac{\partial^2 M_1(\theta)[j]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_1[j] - M_1(\theta)[j]] - \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_1}} \times \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_2}} \right] \right. \\ &+ \sum_{j, k \in [d]} \left[\frac{\partial^2 M_2(\theta)[j, k]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_2[j, k] - M_2(\theta)[j, k]] - \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_1}} \times \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_2}} \right] \\ &+ \sum_{j, k, l \in [d]} \left[\frac{\partial^2 M_3(\theta)[j, k, l]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_3[j, k, l] - M_3(\theta)[j, k, l]] - \right. \\ &\left. \left. \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_1}} \times \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_2}} \right] \right\} \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 Q_n(\theta)}{\partial \theta_{r_1} \partial \theta_{r_2}} &= -2 \sum_{j \in [d]} \frac{\partial^2 M_1(\theta)[j]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_1[j] - M_1(\theta)[j]] \\ &- 2 \sum_{j, k \in [d]} \frac{\partial^2 M_2(\theta)[j, k]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_2[j, k] - M_2(\theta)[j, k]] \\ &- 2 \sum_{j, k, l \in [d]} \frac{\partial^2 M_3(\theta)[j, k, l]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_3[j, k, l] - M_3(\theta)[j, k, l]] + V_{r_1 r_2}(\theta) \end{aligned}$$

avec

$$\begin{aligned} V_{r_1 r_2}(\theta) &= 2 \sum_{j \in [d]} \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_1}} \times \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_2}} + 2 \sum_{j, k \in [d]} \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_1}} \times \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_2}} \\ &+ 2 \sum_{j, k, l \in [d]} \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_1}} \times \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_2}} \end{aligned}$$

Posons $\theta_t = \theta^* + t(\hat{\theta}_n - \theta^*)$, alors $\theta_t = \theta^* + o_{\mathbb{P}}(1)$.

On doit montrer que

$$V_n = \int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] dt$$

converge en probabilité ver la matrice $V(\theta^*)$. On a

$$\|V_n - V(\theta^*)\| \leq \int_0^1 \|D_1 Z_n(\theta_t) - V(\theta^*)\| dt \quad (3.57)$$

avec

$$\|v\| = \sum_{i_1 i_2=1}^d |v_{i_1 i_2}|, \text{ pour } v \in \mathbb{R}^{d \times d}.$$

Par suite,

$$\left\| V_n - V(\theta^*) \right\| \leq \int_0^1 \left\| D_1 Z_n(\theta_t) - V(\theta_t) \right\| dt + \int_0^1 \left\| V(\theta_t) - V(\theta^*) \right\| dt \quad (3.58)$$

Pour tout t , $\left\| V(\theta_t) - V(\theta^*) \right\| = o_{\mathbb{P}}(1)$ comme $\theta_t = \theta^* + o_{\mathbb{P}}(1)$ et que V est continue. On a aussi $\left\| V(\theta_t) - V(\theta^*) \right\| \leq 2 \sup_{\theta \in \Theta} V(\theta)$ avec Θ compact. On a donc par le théorème de convergence dominée, que le second terme de (3.58) est égal à $o_{\mathbb{P}}(1)$.

Evaluons le premier terme de (3.58) :

$$\begin{aligned} \left\| D_1 Z_n(\theta_t) - V(\theta_t) \right\| &\leq 2 \sum_{j \in [d]} \left| \frac{\partial^2 M_1(\theta_t)[j]}{\partial \theta_{r_1} \partial \theta_{r_2}} \right| \times \left| \hat{M}_1[j] - M_1(\theta_t)[j] \right| \\ &+ 2 \sum_{j, k \in [d]} \left| \frac{\partial^2 M_2(\theta_t)[j, k]}{\partial \theta_{r_1} \partial \theta_{r_2}} \right| \times \left| \hat{M}_2[j, k] - M_2(\theta_t)[j, k] \right| \\ &+ 2 \sum_{j, k, l \in [d]} \left| \frac{\partial^2 M_3(\theta_t)[j, k, l]}{\partial \theta_{r_1} \partial \theta_{r_2}} \right| \times \left| \hat{M}_3[j, k, l] - M_3(\theta_t)[j, k, l] \right|. \end{aligned}$$

Ainsi

$$\left| \hat{M}_r[\cdot] - M_r(\theta_t)[\cdot] \right| \leq \left| \hat{M}_r[\cdot] - M_r(\theta^*)[\cdot] \right| + \left| M_r(\theta_t)[\cdot] - M_r(\theta^*)[\cdot] \right| \quad (3.59)$$

En utilisant le fait que $\theta_t = \theta^* + o_{\mathbb{P}}(1)$ uniformément pour $t \in [0, 1]$ et la continuité de M_r , on a $\left| M_r(\theta_t)[\cdot] - M_r(\theta^*)[\cdot] \right| = o_{\mathbb{P}}(1)$ uniformément pour $t \in [0, 1]$. Puis par la loi des grands nombres, $\left| \hat{M}_r[\cdot] - M_r(\theta^*)[\cdot] \right| = o_{\mathbb{P}}(1)$ et donc $\int_0^1 \left| M_r(\theta_t)[\cdot] - M_r(\theta^*)[\cdot] \right| dt = o_{\mathbb{P}}(1)$ et $\int_0^1 \left| \hat{M}_r[\cdot] - M_r(\theta^*)[\cdot] \right| dt = o_{\mathbb{P}}(1)$. Par le lemme de Slutsky avec $\Sigma = V^{-1}L(V^{-1})^T$, on a la normalité asymptotique si V est inversible, avec $V = V(\theta^*)$. On finit la preuve en montrant que V est inversible.

Rappelons que V est de dimension $q \times q$ avec $q = K(2 + d) - 1$. Soit $U \in \mathbb{R}^q$. Nous pouvons noter les coordonnées de U en fonction des paramètres. En utilisant la forme de V nous pouvons dire que

$$U^T V U = \sum_{r_1, r_2} U_{r_1} V_{r_1 r_2} U_{r_2} = 0$$

si et seulement si :

$$U^T D M_1(\theta)[j] = 0, \quad j = 1, \dots, q, \quad (3.60)$$

et

$$U^T D M_2(\theta)[j, l] = 0, \quad j, l = 1, \dots, q, \quad (3.61)$$

et

$$U^T D M_3(\theta)[j, l, m] = 0, \quad j, l, m = 1, \dots, q. \quad (3.62)$$

Ici, $DM[\cdot]$ est le vecteur gradient des coordonnées de M . Notons par $U(\beta_{kl})$ les coordonnées de U associés aux paramètres β_{kl} , $U(b_k)$ les coordonnées de U associés aux paramètres b_k et $U(\omega_k)$ les coordonnées de U associés aux paramètres ω_k . Notons aussi par $\bar{0}$ le vecteur zéro de dimension d , $\bar{0} \otimes \bar{0}$ le tenseur zéro d'ordre 2 de dimension $d \times d$ et $\bar{0} \otimes \bar{0} \otimes \bar{0}$ le tenseur zéro d'ordre 3 de dimension $d \times d \times d$. Pour tout $j, l, m = 1, \dots, d$ les coordonnées dans les équations (3.60)-(3.62) sont données par :

$$\begin{aligned} U^T DM_1(\theta)[j] &= \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_1(\theta)[j]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_1(\theta)[j]}{\partial b_k} \\ &+ \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_1(\theta)[j]}{\partial \beta_{mk}}, \end{aligned} \quad (3.63)$$

$$\begin{aligned} U^T DM_2(\theta)[j, l] &= \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_2(\theta)[j, l]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_2(\theta)[j, l]}{\partial b_k} \\ &+ \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_2(\theta)[j, l]}{\partial \beta_{mk}} \end{aligned} \quad (3.64)$$

et

$$\begin{aligned} U^T DM_3(\theta)[j, l, m] &= \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial b_k} \\ &+ \sum_{k=1}^K \sum_{s=1}^d U(\beta_{sk}) \frac{\partial M_3(\theta)[j, l, m]}{\partial \beta_{sk}}. \end{aligned} \quad (3.65)$$

En utilisant le fait que $\sum_{k=1}^d \omega_k = 1$, les premiers termes des équations (3.63)-(3.65) s'écrivent comme :

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_1(\theta)[j]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \right. \\ &\left. - \mathbb{E} [g'(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \right\}, \end{aligned} \quad (3.66)$$

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_2(\theta)[j, l]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \right. \\ &\left. - \mathbb{E} [g''(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \beta_K(l) \right\} \end{aligned} \quad (3.67)$$

et

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \beta_k(m) \right. \\ &\left. - \mathbb{E} [g^{(3)}(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \beta_K(l) \beta_K(m) \right\} \end{aligned} \quad (3.68)$$

respectivement. De même les seconds termes des équations (3.63)-(3.65) s'écrivent comme :

$$\sum_{k=1}^K U(b_k) \frac{\partial M_1(\theta)[j]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] \cdot \beta(j), \quad (3.69)$$

$$\sum_{k=1}^K U(b_k) \frac{\partial M_2(\theta)[j, l]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \cdot \beta(j) \beta(l) \quad (3.70)$$

et

$$\sum_{k=1}^K U(b_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g^{(4)}(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \beta_k(m) \quad (3.71)$$

respectivement. En dérivant par rapport aux coordonnées des β_k et en utilisant le lemme de Stein (lemme 3.3.1.1), les derniers termes des équations (3.63)-(3.65) s'écrivent comme :

$$\begin{aligned} \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_1(\theta)[j]}{\partial \beta_{mk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] U(\beta_k(j)), \end{aligned} \quad (3.72)$$

$$\begin{aligned} \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_2(\theta)[j, l]}{\partial \beta_{mk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(4)}(\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \beta_k(l) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] \left\{ \beta_k(j) U(\beta_k(l)) \right. \\ &+ \left. \beta_k(l) U(\beta_k(j)) \right\}, \end{aligned} \quad (3.73)$$

et

$$\begin{aligned} \sum_{k=1}^K \sum_{s=1}^d U(\beta_{sk}) \frac{\partial M_3(\theta)[j, l, m]}{\partial \beta_{sk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(5)}(\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \beta_k(l) \beta_k(s) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \left\{ \beta_k(j) \beta_k(l) U(\beta_k(s)) \right. \\ &+ \left. \beta_k(j) U(\beta_k(l)) \beta_k(s) + U(\beta_k(j)) \beta_k(l) \beta_k(s) \right\} \end{aligned} \quad (3.74)$$

respectivement. Alors, en utilisant les équations (3.63)-(3.74), on peut réécrire l'équation

(3.60) comme :

$$\begin{aligned}
\bar{0} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E}[g'(\langle x, \beta_k \rangle + b_k)] \beta_k - \mathbb{E}[g'(\langle x, \beta_K \rangle + b_K)] \beta_K \right\} \\
&+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E}[g''(\langle x, \beta_k \rangle + b_k)] \beta_k + \sum_{k=1}^K \omega_k \mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g'(\langle x, \beta_k \rangle + b_k)] U(\beta_k), \tag{3.75}
\end{aligned}$$

réécrire l'équation (3.61) comme :

$$\begin{aligned}
\bar{0} \otimes \bar{0} &= \sum_{k=1}^{K-1} U(\omega_k) \left[\mathbb{E}[g''(\langle x, \beta_k \rangle + b_k)] \beta_k \otimes \beta_k - \mathbb{E}[g''(\langle x, \beta_K \rangle + b_K)] \beta_K \otimes \beta_K \right] \\
&+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \beta_k \otimes \beta_k + \sum_{k=1}^K \omega_k \mathbb{E}[g^{(4)}(\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k \otimes \beta_k \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g''(\langle x, \beta_k \rangle + b_k)] \left(U(\beta_k) \otimes \beta_k + \beta_k \otimes U(\beta_k) \right), \tag{3.76}
\end{aligned}$$

et réécrire l'équation (3.62) comme :

$$\begin{aligned}
\bar{0}^{\otimes 3} &= \sum_{k=1}^{K-1} U(\omega_k) \left[\mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \beta_k^{\otimes 3} - \mathbb{E}[g^{(3)}(\langle x, \beta_K \rangle + b_K)] \beta_K^{\otimes 3} \right] \tag{3.77} \\
&+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E}[g^{(4)}(\langle x, \beta_k \rangle + b_k)] \beta_k^{\otimes 3} + \sum_{k=1}^K \omega_k \mathbb{E}[g^{(5)}(\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k^{\otimes 3} \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \left(U(\beta_k) \otimes \beta_k \otimes \beta_k + \beta_k \otimes U(\beta_k) \otimes \beta_k + \beta_k \otimes \beta_k \otimes U(\beta_k) \right).
\end{aligned}$$

Nous allons montrer que les vecteurs $U(\beta_1), \dots, U(\beta_K)$ sont tous dans l'espace linéaire engendré par β_1, \dots, β_K .

Soit W un vecteur qui est orthogonal à cet espace linéaire. En multipliant l'équation (3.77) à droite par W , et en utilisant le fait que les β_1, \dots, β_K sont linéairement indépendants par (H1), nous obtenons

$$\forall k = 1, \dots, K, \quad \omega_k (G_3)_k \langle U(\beta_k), W \rangle = 0.$$

En utilisant (H1), nous avons $\omega_k > 0$, $k = 1, \dots, K$ de sorte que

$$\forall k = 1, \dots, K, \quad (G_3)_k \langle U(\beta_k), W \rangle = 0.$$

Alors, si (H3) est vraie, on a, pour tout k et pour tout W , $\langle U(\beta_k), W \rangle = 0$, ce qui prouve que sous (H3), les vecteurs $U(\beta_1), \dots, U(\beta_K)$ sont tous dans l'espace linéaire engendré par β_1, \dots, β_K .

Soit $B \in \mathbb{R}^{d \times K}$ la matrice dont les colonnes sont les β_1, \dots, β_K . Soit $U(\beta) \in \mathbb{R}^{d \times K} d \times K$ la matrice dont les colonnes sont les $U(\beta_1), \dots, U(\beta_K)$. Si (H3) est vraie, alors il existe une matrice, de dimension $K \times K$, $A = (A_1, \dots, A_K)$ telle que $U(\beta) = BA$.

Si R est un vecteur de dimension K , notons par $Diag(R) \in \mathbb{R}^{K \times K}$ la matrice diagonale dont les éléments de la diagonale sont donnés par R . Posons

$$U(\omega) = \left(U(\omega_1), \dots, U(\omega_{K-1}), -\sum_{k=1}^{K-1} U(\omega_k) \right),$$

$$U(b) = (U(b_1), \dots, U(b_K))$$

et rappelons que

$$\omega = \left(\omega_1, \dots, \omega_{K-1}, 1 - \sum_{k=1}^{K-1} \omega_k \right).$$

Soient P , Q et Δ des matrices diagonales telles que $P = Diag(U(\omega))$, $Q = Diag(U(b))$ et $\Delta = Diag(\omega)$. Pour $W \in \mathbb{R}^d$, posons aussi, $D = Diag(\langle \beta_1, W \rangle, \dots, \langle \beta_K, W \rangle)$. Alors, en utilisant le fait que B est de rang plein, on a à partir de l'équation (3.77), que

$$G_3PD + G_4\Delta QD + G_5\Delta Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle), BA_K)D$$

$$+ AG_3\Delta D + G_3\Delta DA^T + G_3Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle) = \bar{0} \otimes \bar{0}. \quad (3.78)$$

Comme $U(\beta) = BA$, alors $U(\beta_k) = \sum_{r=1}^K \beta_r A_{rk} = BA_k$. Ce qui implique que

$$Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle) = Diag(\langle BA_1, \beta_1 \rangle, \dots, \langle A_K, \beta_K \rangle) = \tilde{D}.$$

Ainsi l'équation (3.78) pourra être réécrite comme

$$G_3PD + G_4\Delta QD + G_5\Delta \tilde{D}D + AG_3\Delta D + G_3\Delta DA^T + G_3\Delta \tilde{D} = \bar{0} \otimes \bar{0}. \quad (3.79)$$

Ainsi, pour tout $W \in \mathbb{R}^d$,

$$AG_3\Delta D + G_3\Delta DA^T$$

est une matrice diagonale. Comme les éléments de $G_3\Delta$ sont non nuls, ceci prouve, sous les hypothèses (H1) et (H3), que A est une matrice diagonale. Dans ce cas,

$$\tilde{D} = A\tilde{B} \text{ avec } \tilde{B} = Diag(\|\beta_1\|^2, \dots, \|\beta_K\|^2)$$

et l'équation (3.79) pourra être réécrite comme

$$G_3PD + G_4\Delta QD + G_5\Delta A\tilde{B}D + AG_3\Delta D + G_3\Delta DA^T + G_3\Delta A\tilde{B} = \bar{0} \otimes \bar{0}. \quad (3.80)$$

Mais en prenant $W \in \mathbb{R}^d$ tel que, pour $k = 1, \dots, K$ $\beta_k^T W_k = 0$, on a $D = 0$. Dans ce cas l'équation (3.80) est donnée par

$$G_3\Delta A\tilde{B} = \bar{0} \otimes \bar{0}.$$

En utilisant le fait que les entrées de G_3 , Δ et \tilde{B} sont toutes non nulles, on obtient que $A = 0$.

Comme $A = 0$, ceci implique que $U(\beta_k) = 0$, $k = 1, 2, \dots, K$. Ainsi, en utilisant le fait que B est de rang plein, on a à partir des équations (3.75) et (3.76), que

$$G_1P + G_2\Delta Q = \bar{0} \otimes \bar{0}, \quad (3.81)$$

et

$$G_2P + G_3\Delta Q = \bar{0} \otimes \bar{0}. \quad (3.82)$$

En multipliant l'équation (3.81) par G_3 et l'équation (3.82) par G_2 , on a

$$G_1G_3P + G_2G_3\Delta Q = \bar{0} \otimes \bar{0}, \quad (3.83)$$

et

$$G_2^2P + G_2G_3\Delta Q = \bar{0} \otimes \bar{0}. \quad (3.84)$$

En prenant la différence les équations (3.83) et (3.84), on a

$$(G_1G_3 - G_2^2)P = \bar{0} \otimes \bar{0}$$

Comme les éléments de $G_1G_3 - G_2^2$ sont non nuls, cela veut dire que $P = 0$. Et comme les éléments de $G_3\Delta$ sont aussi non nuls, on a $Q = 0$. Ainsi, sous les hypothèses (H1), (H3) et (H4), la matrice V est inversible.

3.5 Applications

Nous présenterons dans cette partie le package **R** préparé dans le but d'estimer les paramètres du modèle présenté à la section 3.2.2. Nous ferons ensuite une étude de simulations pour illustrer notre étude. Nous présenterons à la fin une méthode de sélection de variables en utilisant l'algorithme spectral.

3.5.1 Package **R**

Le package **R** relié à cette étude est appelé **morpheus**. Il s'articule autour de quatre fonctions principales, qui sont les suivantes :

1. **generateSampleIO** permet de simuler des données de mélange de modèles linéaires généralisés. C'est-à-dire un échantillon $(X_n, Y_n) = ((x_1, y_1), \dots, (x_n, y_n))$ i.i.d de même loi que (X, Y) avec $X \sim N(0, I_d)$ et $Y|X$ un mélange de modèles linéaires généralisés. Elle prend en entrée les arguments suivants :
 - **n** : la taille de l'échantillon,
 - **p** : le vecteur de poids du mélange,
 - **β** : la matrice de vecteurs de paramètres du mélange. elle est de taille $d \times K$, avec d la taille de X et K le nombre de composants (populations) du mélange.
 - **b** : le vecteur des intercepts. il est de taille K .
 - **link** : la fonction lien utilisée (**logit** ou **probit**).
Elle renvoie à la sortie à une liste contenant les données (X_n, Y_n) .
2. **computeMoments** permet d'estimer les moments croisés M_1 , M_2 et M_3 définis à la section 3.3.1. Elle prend comme arguments :
 - **Xn** la matrice des covariables, de taille $n \times d$.

- Y_n le vecteur réponse de taille n .
- Elle renvoie en sortie une liste contenant les moments M_1 , M_2 et M_3 .
3. `computeMu` permet d'estimer les vecteurs normalisés de la matrice β . Elle prend comme arguments :
- la matrice des données X_n ,
 - le vecteur de réponses Y_n ,
- et une liste contenant des arguments optionnels, tels que :
- `jd_method` la méthode de diagonalisation jointe utilisée à partir du package `jointDiag` (`uwedge` (par défaut) ou `jedi`).
 - `jd_nvects` nombre de vecteurs à fournir pour la diagonalisation jointe (ou 0 pour prendre les d vecteurs de la base canonique de \mathbb{R}^d).
 - `M` la liste contenant les moments d'ordre 1, 2 et 3. Les moments seront estimés s'ils ne sont pas fournis.
 - `K` le nombre de populations (estimé avec le rang de M_2 s'il n'est pas fourni)
4. `optimParams` permet d'estimer tous les paramètres du modèle. Elle prend comme arguments :
- `K` le nombre de populations.
 - `link` le type de lien à utiliser (`logit` ou `probit`).
- Mais aussi une liste contenant des arguments optionnels :
- `M` la liste contenant les moments d'ordre 1, 2 et 3. Les moments seront estimés s'ils ne sont pas fournis.
 - `Xn, Yn` : la matrice des données X_n et le vecteur de réponses Y_n , à fournir si les moments ne sont pas donnés.
- Elle fournit à la sortie un objet `op` de classe `OptimParams`. `op$run(x0)` renvoie une liste contenant les estimateurs de :
- `p` le vecteur de proportions ou poids du mélange de taille K
 - β la matrice des vecteurs de paramètres de régression de taille $d \times K$.
 - `b` le vecteur des intercepts de taille K .
- x_0 est un vecteur initial (point initial de l'algorithme) contenant respectivement les $K - 1$ premiers éléments de p , β par colonne et le vecteur b . x_0 est alors de taille $K(d + 2) - 1$.

Il existe aussi d'autres fonctions secondaires dans le package `morpheus`, à savoir `multiRun` pour des estimations des statistiques par Monte-Carlo ou Bootstrap; `alignMatrices` pour gérer le problème de "label Switching"; `plotHist` pour afficher les statistiques calculées.

3.5.2 Simulations

Dans cette partie, nous essayerons de comparer les différents algorithmes de diagonalisation jointe utilisés pour estimer les directions. Nous présenterons ensuite des résultats sur l'estimation des paramètres du modèle présenté à la section 3.2. Nous présenterons aussi des résultats de normalité asymptotique et comparerons la méthode spectrale et la méthode du maximum de vraisemblance. Nous présenterons à la fin des résultats de sélection de variables.

Algorithmes de diagonalisation jointe

Beaucoup d’algorithmes de diagonalisation jointe pour un ensemble fini de matrices carrés sont implémentés dans la littérature. Nous comparons ici deux algorithmes implémentés par Cédric Gouy-Paillier ([78]) pendant sa thèse : `jedi` et `uwedge` qui sont basés sur les travaux de A. Souloumiac ([39]) et de Tichavsky et al. ([79]) respectivement. Ces algorithmes seront utilisés dans la suite pour estimer les directions des vecteurs de paramètres du modèle.

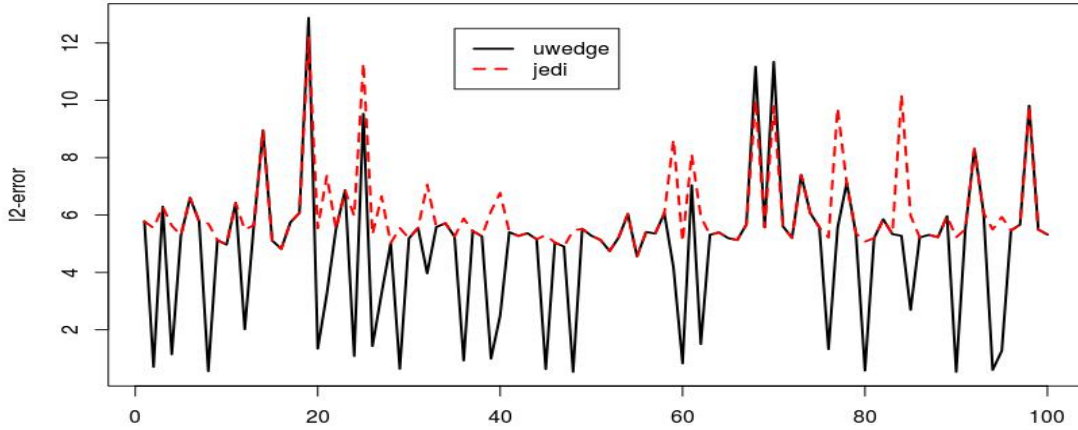


FIGURE 3.1 – Expérience 1 : L’erreur L2 des estimateurs entre les algorithmes `uwedge` et `jedi`. sur 100 échantillons. Lien [logit](#).

La figure 3.1 présente les erreurs L2 calculées sur $\hat{\theta}$ obtenue avec les algorithmes “`uwedge`” et “`jedi`”. Sur 100 échantillons différents, l’erreur suivante a été calculée

$$err = \left(\sum_{j=1}^q (\hat{\theta}_j - \theta_j^*)^2 \right)^{1/2}$$

pour chaque méthode. Avec θ^* le vecteur de dimension q contenant les vraies valeurs ω^* , β^* et b^* . Pour chaque méthode, $\hat{\theta}$ est le vecteur de dimension q contenant les estimateurs $\hat{\omega}$, $\hat{\beta}$ et \hat{b} . On peut voir à partir de cette figure, que sur les 100 points estimés, l’erreur calculée pour “`uwedge`” est très souvent inférieure à celle calculée pour “`jedi`”. En plus de ça, la méthode “`uwedge`” est plus stable si d et K augmentent. Dans toute la suite, on utilisera la méthode “`uwedge`” pour faire l’estimation. Cette méthode est utilisée dans le package `Morpheus`.

Estimation

Dans cette partie, on commencera d’abord par estimer les directions des vecteurs de paramètres. C’est-à-dire, estimer

$$\beta_k / \|\beta_k\|, \quad k = 1, \dots, K.$$

Ensuite nous estimerons tous les paramètres du modèle (ω , β et b) en prenant comme point initial les directions estimées. On comparera à la fin l'estimateur du maximum de vraisemblance et l'estimateur par la méthode des moments.

(a). **Estimation des directions**

Dans la suite de cette partie, nous utiliserons les trois expériences résumées dans le tableau 3.1. Dans ce tableau, n représente la taille de l'échantillon, d la taille des vecteurs de paramètres, K le nombre de composantes ou populations ($K \geq 2$), ω le vecteur des poids du mélange, b le vecteur des intercepts et β la matrice des vecteurs de paramètres.

<i>Paramètre</i> \ <i>Expérience</i>	Expérience 1	Expérience 2	Expérience 3
n	1e4	1e5	1e6
d	2	4	5
K	2	2	3
ω	$\begin{pmatrix} 0.5 & 0.5 \end{pmatrix}$	$\begin{pmatrix} 0.6 & 0.4 \end{pmatrix}$	$\begin{pmatrix} 0.3 & 0.4 & 0.3 \end{pmatrix}$
b	$\begin{pmatrix} -0.5 & 0.5 \end{pmatrix}$	$\begin{pmatrix} 0.3 & 0.1 \end{pmatrix}$	$\begin{pmatrix} 1 & 0.5 & -0.4 \end{pmatrix}$
β	$\begin{pmatrix} 1 & 3 \\ -2 & 1 \end{pmatrix}$	$\begin{pmatrix} -1 & 1.8 \\ 2 & -2 \\ -3 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 1 & -1 & -2 \\ 2 & 1 & 1.6 \\ -3 & 1.5 & -2.5 \\ 2.3 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}$

TABLE 3.1 – Les différentes expériences de simulations utilisées.

Pour les trois expériences décrites dans le tableau 3.1, les matrices des vecteurs normalisés associées sont

$$\mu = \begin{pmatrix} 0.45 & 0.95 \\ -0.89 & 0.32 \end{pmatrix} \text{ pour l'expérience 1,}$$

$$\mu = \begin{pmatrix} -0.26 & 0.63 \\ 0.52 & -0.69 \\ -0.77 & 0.35 \\ 0.26 & 0.00 \end{pmatrix} \text{ pour l'expérience 2}$$

et

$$\mu = \begin{pmatrix} 0.23 & -0.33 & -0.52 \\ 0.45 & 0.33 & 0.41 \\ -0.68 & 0.49 & -0.65 \\ 0.52 & 0.66 & -0.26 \\ 0.00 & -0.33 & 0.26 \end{pmatrix} \text{ pour l'expérience 3.}$$

En utilisant notre algorithme d'estimation des directions $\mu_k = \beta_k / \|\beta_k\|$, $k = 1, \dots, K$, dans les trois expériences du tableau 3.1, pour des liens logit et probit, on arrive à estimer les directions. Les tableaux 3.2, 3.3 et 3.4 présentent les résultats de ces différentes expériences. Pour chaque élément de la matrice de vecteurs normalisés μ , l'estimateur est donné, de même que la racine carrée de l'erreur quadratique moyen (RMSE), l'erreur absolue moyen [MAE] et l'écart type.

Le tableau 3.2 montre que pour des dimensions petites, l'estimation des directions est très bonne et meilleure dans le cas *probit*. Pour des dimensions assez grandes, l'estimation reste assez instable et similaire dans les deux cas.

		Directions			
<i>Lien</i>		0.45	-0.89	0.95	0.32
<i>logit</i>		0.441	-0.890	0.944	0.320
		(0.101)	(0.052)	(0.026)	(0.074)
		[0.080]	[0.041]	[0.020]	[0.060]
		{0.100}	{0.052}	{0.026}	{0.074}
<i>probit</i>		0.442	-0.894	0.945	0.320
		(0.067)	(0.033)	(0.020)	(0.060)
		[0.053]	[0.027]	[0.016]	[0.047]
		{0.066}	{0.033}	{0.020}	{0.060}

TABLE 3.2 – Estimation des directions. Résultats de simulations pour l'expérience 1 avec $N = 1000$ réplifications. Pour chaque lien, l'estimateur est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l'écart type associé.

		Directions							
<i>Lien</i>		-0.26	0.52	-0.77	0.26	0.63	-0.69	0.35	0.00
<i>logit</i>		-0.278	0.492	-0.611	0.182	0.360	-0.332	0.013	0.111
		(0.282)	(0.234)	(0.306)	(0.279)	(0.475)	(0.598)	(0.568)	(0.420)
		[0.209]	[0.157]	[0.201]	[0.194]	[0.337]	[0.412]	[0.444]	[0.329]
		{0.281}	{0.232}	{0.262}	{0.268}	{0.391}	{0.471}	{0.458}	{0.405}
<i>probit</i>		-0.309	0.504	-0.604	0.192	0.392	-0.343	0.030	0.087
		(0.278)	(0.220)	(0.305)	(0.254)	(0.449)	(0.587)	(0.555)	(0.399)
		[0.208]	[0.154]	[0.204]	[0.180]	[0.307]	[0.397]	[0.424]	[0.306]
		{0.274}	{0.220}	{0.256}	{0.245}	{0.382}	{0.467}	{0.453}	{0.389}

TABLE 3.3 – Estimation des directions. Résultats de simulations pour l'expérience 2 avec $N = 1000$ réplifications. Pour chaque lien, l'estimateur est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l'écart type associé.

Directions

<i>Lien</i>	0.23	0.45	-0.68	0.52	0.00	-0.33	0.33	0.49	0.66	-0.33	-0.52	0.41	-0.65	-0.26	0.26
	0.223	0.387	-0.576	0.453	-0.005	-0.296	0.286	0.467	0.601	-0.303	-0.465	0.361	-0.571	-0.241	0.240
<i>logit</i>	(0.185)	(0.235)	(0.290)	(0.252)	(0.215)	(0.172)	(0.189)	(0.165)	(0.201)	(0.186)	(0.202)	(0.225)	(0.253)	(0.191)	(0.182)
	[0.100]	[0.114]	[0.132]	[0.121]	[0.110]	[0.071]	[0.075]	[0.069]	[0.076]	[0.072]	[0.095]	[0.103]	[0.106]	[0.099]	[0.082]
	{0.185}	{0.224}	{0.271}	{0.243}	{0.215}	{0.168}	{0.184}	{0.163}	{0.192}	{0.184}	{0.194}	{0.217}	{0.240}	{0.190}	{0.181}
	0.219	0.374	-0.575	0.454	-0.015	-0.297	0.295	0.479	0.610	-0.310	-0.476	0.360	-0.565	-0.239	0.223
<i>probit</i>	(0.184)	(0.249)	(0.300)	(0.247)	(0.227)	(0.160)	(0.179)	(0.146)	(0.181)	(0.143)	(0.175)	(0.227)	(0.268)	(0.181)	(0.214)
	[0.101]	[0.121]	[0.135]	[0.119]	[0.116]	[0.059]	[0.067]	[0.057]	[0.063]	[0.049]	[0.086]	[0.102]	[0.111]	[0.093]	[0.095]
	{0.184}	{0.234}	{0.282}	{0.238}	{0.226}	{0.157}	{0.175}	{0.146}	{0.174}	{0.142}	{0.169}	{0.219}	{0.254}	{0.180}	{0.211}

TABLE 3.4 – Estimation des directions. Résultats de simulations pour l'expérience 3 avec $N = 1000$ réplifications. Pour chaque lien, l'estimateur est donné, de même que (\cdot) pour le RMSE, $[\cdot]$ pour le MAE et $\{\cdot\}$ pour l'écart type associé.

(b). **Estimation de tous les paramètres**

Dans cette partie, on essayera d'estimer les paramètres du modèle, à savoir ω , β et b . Rappelons que le point initial considéré ici est l'estimateur des vecteurs normalisés pour β et 0 (le vecteur rempli de 0 de taille K , avec K est le nombre de composantes du mélange) pour le vecteur b . Nous utiliserons dans cette partie les trois expériences suivants :

— Expérience 4 (dimension 2) :

$$\begin{aligned}K &= 2 \\ \omega &= (0.5, 0.5) \\ b &= (-0.2, 0.5) \\ \beta &= \begin{pmatrix} 1 & 3 \\ -2 & 1 \end{pmatrix}\end{aligned}$$

— Expérience 5 (dimension 5) :

$$\begin{aligned}K &= 2 \\ \omega &= (0.5, 0.5) \\ b &= (-0.2, 0.5) \\ \beta &= \begin{pmatrix} 1 & 2 \\ 2 & -3 \\ -1 & 0 \\ 0 & 1 \\ 3 & 0 \end{pmatrix}\end{aligned}$$

— Expérience 6 (dimension 10) :

$$\begin{aligned}K &= 3 \\ \omega &= (0.3, 0.3, 0.4) \\ b &= (-0.2, 0, 0.5) \\ \beta &= \begin{pmatrix} 1 & 2 & -1 \\ 2 & -3 & 1 \\ -1 & 0 & 3 \\ 0 & 1 & -1 \\ 3 & 0 & 0 \\ 4 & -1 & 0 \\ -1 & -4 & 2 \\ -3 & 3 & 0 \\ 0 & 2 & 1 \\ 2 & 0 & -2 \end{pmatrix}\end{aligned}$$

Pour différentes valeurs de n ($(n = 5.10^3, 10^4, 10^5, 5.10^5, 10^6)$) et pour $N = 1000$ réplifications, les vecteurs de régressions ($\beta_k, k = 1, \dots, K$), le vecteur des poids du mélange (ω) et le vecteur des intercepts (b) sont estimés. Pour chaque vecteur, la distance entre son estimateur et le vrai paramètre est calculée (l'erreur quadratique moyen). Les figures 3.2, 3.3 et 3.4 présentent les résultats des expériences 4, 5 et

6 respectivement. On peut voir à partir de ces figures la qualité de l'estimateur dépend de la taille de l'échantillon et qu'on a besoin de beaucoup de données pour bien estimer les paramètres. Les figures 3.2 et 3.3 montrent qu'à partir de $n = 10^5$, on commence à bien estimer les paramètres pour des dimensions inférieures ou égales à 5. Pour des dimensions assez grandes, la figure 3.4 montre que qu'il faut un échantillon de taille supérieure à 10^6 .

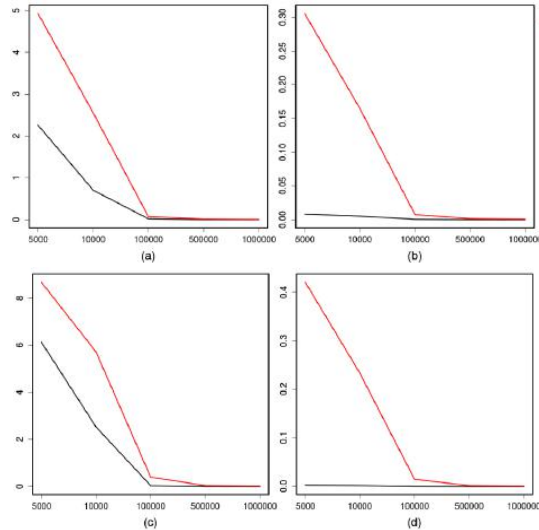


FIGURE 3.2 – Experiment 4 : erreur quadratique moyen (MSE) pour les deux liens (logit en haut et probit en bas). Lien logit : (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$. Lien probit : (c) $MSE(\hat{\beta})$, (d) $MSE(\hat{b}, \hat{p})$.

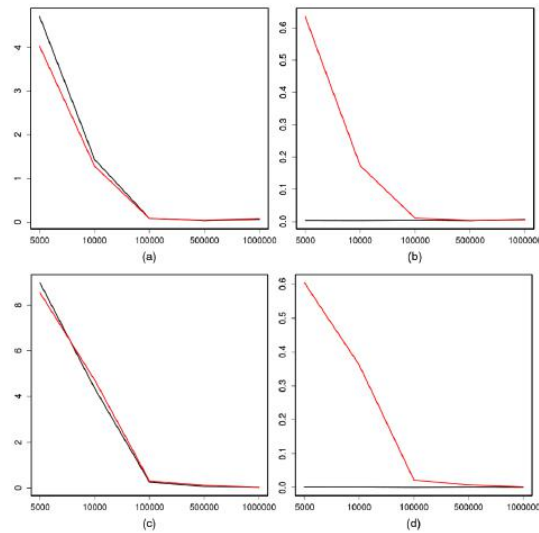


FIGURE 3.3 – Experiment 5 : erreur quadratique moyen (MSE) pour les deux liens (logit en haut et probit en bas). Lien logit : (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$. Lien probit : (c) $MSE(\hat{\beta})$, (d) $MSE(\hat{b}, \hat{p})$.

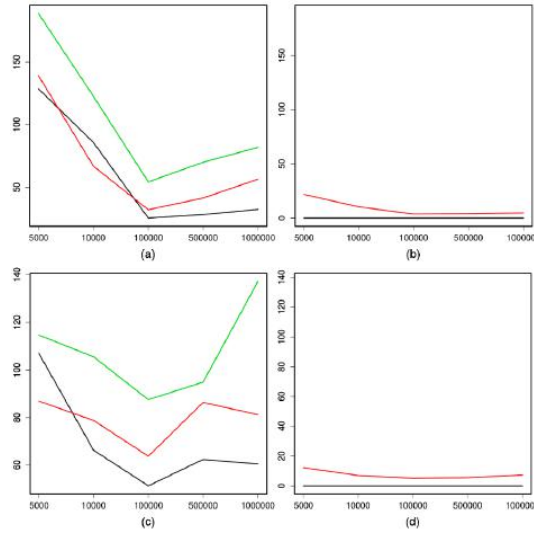


FIGURE 3.4 – Experiment 6 : erreur quadratique moyen (MSE) pour les deux liens (logit en haut et probit en bas). Lien logit : (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$. Lien probit : (c) $MSE(\hat{\beta})$, (d) $MSE(\hat{b}, \hat{p})$.

Méthode des moments versus maximum de vraisemblance

Dans cette partie, on essayera de comparer notre algorithme d'estimation à l'estimation par maximum de vraisemblance. D'abord en temps de calcul, ensuite en terme de biais, de variance et de stabilité.

(a). Temps de calcul

La figure 3.5 présente le temps de calcul en fonction de $\log(n)$ avec des dimensions différentes ($d = 2, 5$ et 10). Les lignes pleines représentent le temps de calcul de notre algorithme d'estimation et les lignes pointillées représentent le temps de calcul de l'estimateur du maximum de vraisemblance.

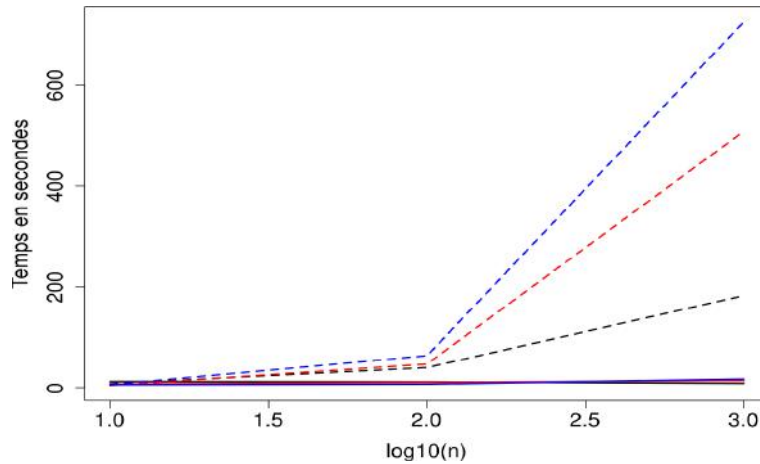


FIGURE 3.5 – EMV vs MM : temps de calcul ; $d = 2$ (noir), $d = 5$ (rouge) et $d = 10$ (bleu)

On peut constater que le temps de calcul reste assez constant pour notre algo-

rithme même si la taille de l'échantillon et la dimension augmentent. Par contre, le temps de calcul du maximum de vraisemblance augmente en fonction de la dimension (d) et de la taille de l'échantillon (n).

(b). **Biais, variance et stabilité**

On présente ici l'estimation des paramètres sur $N = 1000$ échantillons. Pour chaque échantillon, on estime les paramètres avec les deux différentes méthodes (MM et EMV). On calcule ensuite pour chaque méthode, l'estimateur (ligne pointillée rouge) et l'écart type (ligne pointillées rouges) sur les 1000 échantillons. Les figures 3.6 et 3.7 présentent les résultats pour les liens logit et probit respectivement. Dans chaque figure, la ligne pleine noire représente le vrai paramètre (ω^*, β^*, b^*), la ligne pointillée rouge représente l'estimateur ($\hat{\omega}, \hat{\beta}, \hat{b}$) et les lignes pointillées noires représentent l'estimateur plus ou moins l'écart type.

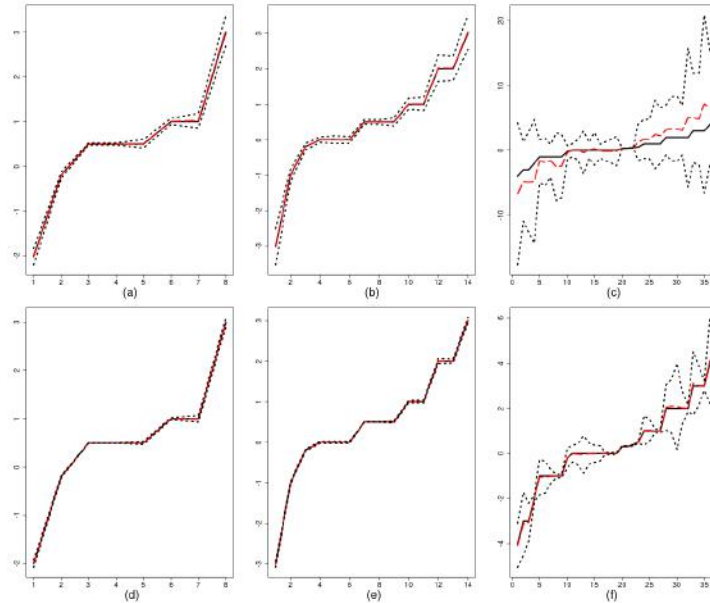


FIGURE 3.6 – Lien *logit*. En haut : la méthode des moments, en bas : le maximum de vraisemblance. De gauche à droite : expérience 4, 5 et 6 respectivement.

Pour des dimensions petites, les deux estimateurs sont assez similaires même si l'EMV reste légèrement meilleur en terme de variance avec le lien logit. Pour des dimensions assez grandes, l'EMV est assez mauvais en terme de biais dans le cas d'un lien probit alors que l'estimateur par la méthode des moments reste assez stable.

Malgré le fait que l'estimateur du maximum de vraisemblance soit meilleur en terme de variance dans certains cas, l'estimateur par la méthode des moments reste assez meilleur en temps de calcul. En grande dimension, il faut avoir plus de données ou explorer plusieurs points initiaux pour bien estimer les paramètres.

Normalité asymptotique

Pour l'expérience 4, on estime les paramètres du modèle pour $N = 1000$ réplifications. Comme prouvé à la section 3.4.3, la loi asymptotique des estimateurs est normale quelle

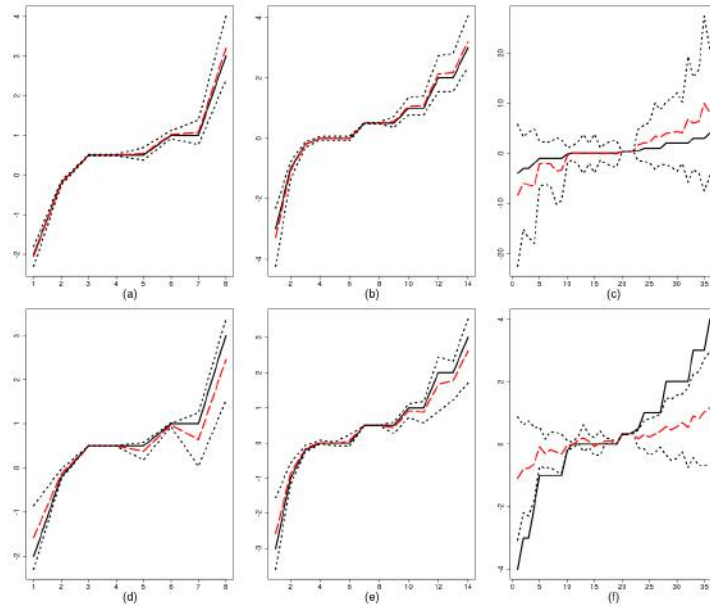


FIGURE 3.7 – Lien *probit*. En haut : la méthode des moments, en bas : le maximum de vraisemblance. De gauche à droite : expérience 4, 5 et 6 respectivement.

que soit la fonction lien g . Les figures 3.8 et 3.9 présentent la loi asymptotique des estimateurs pour des lien *logit* et *probit* respectivement.

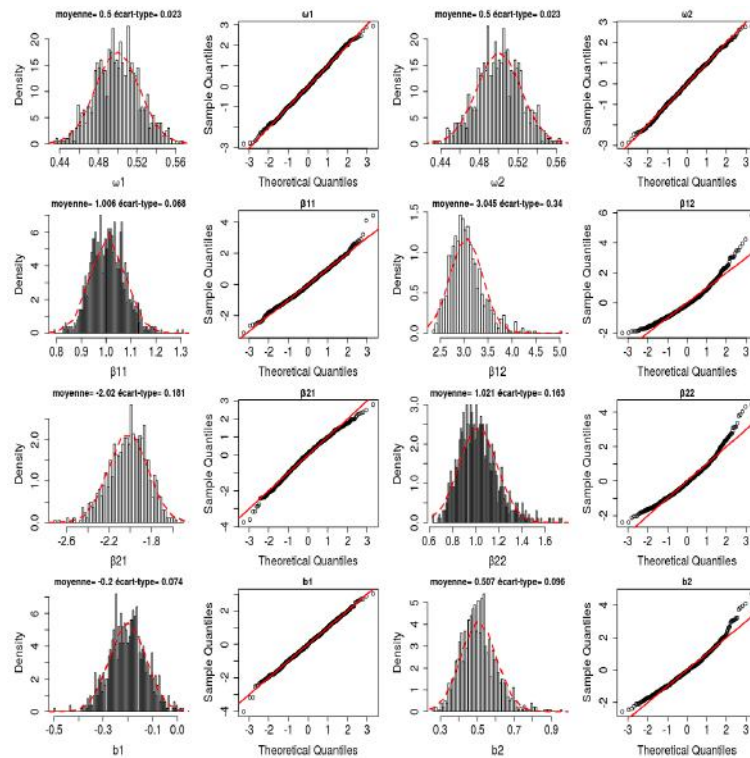


FIGURE 3.8 – Expérience 4 : Normalité Asymptotique des estimateurs pour $N = 1000$ réplifications. Lien *logit*.

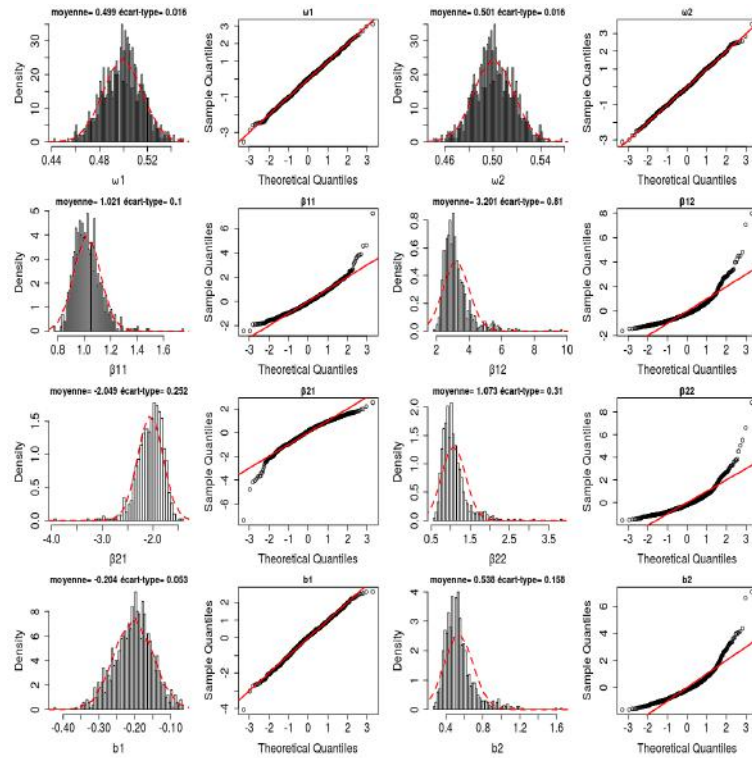


FIGURE 3.9 – Expérience 4 : Normalité Asymptotique des estimateurs pour $N = 1000$ réplifications. Lien *probit*.

On peut constater à partir des figures (3.8) et (3.9) que le lien *logit* présente de meilleurs résultats pour l'approximation de la loi. On peut aussi remarquer qu'en terme de variance, le lien *logit* reste meilleur dans la mesure où il présente, pour chaque paramètre (hormis le vecteur de poids et le vecteur d'intercepts), une variance plus faible que dans le cas *probit*.

Globalement, on obtient la normalité asymptotique des estimateurs quelle que soit la fonction lien utilisée ou la dimension. Quelle que soit la dimension, le lien *logit* est meilleur en terme d'approximation de la loi, de biais et de variance. Mais en grande dimension les deux liens restent assez similaires.

3.5.3 Sélection de variables

Pour sélectionner les variables importantes, il suffit d'effectuer un test de nullité des paramètres associés. Par exemple, pour sélectionner les variables importantes dans la population (ou composante) k , il faut tester pour $j = 1, \dots, d$, l'hypothèse

$$H_0 : \beta_{jk} = 0 \text{ vs } \beta_{jk} \neq 0.$$

En utilisant la méthode du bootstrap, on peut effectuer ce test de deux manières différentes : (1) en effectuant un test purement bootstrap ou (2) en estimant par bootstrap la variance obtenue à partir de la normalité asymptotique (théorème 3.4.3.1, section 3.4.3). Dans ce cas, on a le corollaire suivant :

Corollaire 3.5.3.1 *Sous les hypothèses du théorème 3.4.3.1, on peut estimer Σ de manière consistante par $\hat{\Sigma}$. Dans ce cas, on peut tester l'hypothèse*

$$H_0 : \rho_i = 0 \text{ vs } H_1 : \rho_i \neq 0, \quad i \in \{1, \dots, d \times K\},$$

où ρ est le vecteur contenant les éléments des K vecteurs de paramètres. On a alors, sous H_0 , que

$$W = \frac{\hat{\rho}_i^2}{\hat{\Sigma}_{ii}} \text{ converge en loi vers } \chi^2(1), \quad i = 1, \dots, d \times K.$$

Dans la suite, nous utiliserons les expériences résumées dans le tableau 3.5 pour illustrer la sélection de variables. Pour sélectionner les variables par la méthode bootstrap, on n'a pas nécessairement besoin d'estimer tous les paramètres du modèle. Il suffit d'utiliser l'estimation des directions présentée à la section 3.3.1.

<i>Paramètre</i> \ <i>Expérience</i>	Expérience 7	Expérience 8
n	1e5	1e6
d	4	5
K	2	3
ω	$\begin{pmatrix} 0.6 & 0.4 \end{pmatrix}$	$\begin{pmatrix} 0.3 & 0.4 & 0.3 \end{pmatrix}$
b	$\begin{pmatrix} 0.3 & 0.1 \end{pmatrix}$	$\begin{pmatrix} 1 & 0.5 & -0.4 \end{pmatrix}$
β	$\begin{pmatrix} 0 & 1.8 \\ 2 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 1.6 \\ 0 & 1.5 & 0 \\ 2.3 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$
μ	$\begin{pmatrix} 0.00 & 0.60 \\ 0.62 & 0.00 \\ 0.00 & 0.80 \\ 0.78 & 0.00 \end{pmatrix}$	$\begin{pmatrix} 0.56 & 0.00 & 0.00 \\ 0.47 & 0.37 & 0.51 \\ 0.00 & 0.74 & 0.00 \\ 0.68 & 0.56 & 0.00 \\ 0.00 & 0.00 & 0.86 \end{pmatrix}$

TABLE 3.5 – Les différentes expériences de simulations utilisées pour la sélection de variables.

- **Expérience 7** : Le tableau 3.6 présente les résultats de sélection de variables pour l'expérience 7. Pour chaque fonction lien, l'estimateur des directions $\hat{\mu}$, l'estimateur des vecteurs de paramètres $\hat{\beta}$, les p-values bootstrap (pv-b) et les p-values test (pv-w) sont donnés. A partir de ce tableau, on constate que, hormis la première composante de la dernière ligne (pv-w) pour le lien logit, on sélectionne bien les variables pour les deux méthodes utilisées.
- **Expérience 8** : Le tableau 3.7 présente les résultats de sélection de variables pour l'expérience 8. Comme dans le cas de l'expérience 7, l'estimateur des directions $\hat{\mu}$, l'estimateur des vecteurs de paramètres $\hat{\beta}$, les p-values bootstrap (pv-b) et les p-values test (pv-w) sont donnés pour chaque fonction lien. On peut constater ici qu'on sélectionne bien les variables pour les deux méthodes utilisées.

<i>Estimateur</i> \ <i>Lien</i>	<i>Logit</i>	<i>Probit</i>
$\hat{\mu}$	$\begin{pmatrix} \mathbf{0.015} & 0.604 \\ 0.609 & \mathbf{-0.062} \\ \mathbf{0.253} & 0.793 \\ 0.751 & \mathbf{-0.038} \end{pmatrix}$	$\begin{pmatrix} \mathbf{0.019} & 0.600 \\ 0.64 & \mathbf{0.032} \\ \mathbf{-0.005} & 0.800 \\ 0.765 & \mathbf{0.005} \end{pmatrix}$
$\hat{\beta}$	$\begin{pmatrix} \mathbf{0.053} & 2.563 \\ 3.868 & \mathbf{-0.005} \\ \mathbf{-0.086} & 3.426 \\ 4.750 & \mathbf{-0.003} \end{pmatrix}$	$\begin{pmatrix} \mathbf{0.103} & 2.474 \\ 3.569 & \mathbf{0.051} \\ \mathbf{-0.030} & 3.454 \\ 4.3506 & \mathbf{0.064} \end{pmatrix}$
pv-b	$\begin{pmatrix} \mathbf{0.819} & 0.095 \\ 0.054 & \mathbf{0.737} \\ \mathbf{0.683} & 0.051 \\ 0.045 & \mathbf{0.822} \end{pmatrix}$	$\begin{pmatrix} \mathbf{0.769} & 0.096 \\ 0.083 & \mathbf{0.744} \\ \mathbf{0.934} & 0.067 \\ 0.074 & \mathbf{0.960} \end{pmatrix}$
pv-w	$\begin{pmatrix} \mathbf{3.0e-1} & 2.2e-05 \\ 1.7e-03 & \mathbf{7.9e-01} \\ \mathbf{5.9e-02} & 2.1e-05 \\ 1.7e-03 & \mathbf{9.1e-01} \end{pmatrix}$	$\begin{pmatrix} \mathbf{0.063} & 0.079 \\ 0.049 & \mathbf{0.11} \\ \mathbf{0.132} & 0.079 \\ 0.049 & \mathbf{0.09} \end{pmatrix}$

TABLE 3.6 – Expérience 7 : résultats de sélection de variables. pv-b, pour les p-values obtenues par sélection purement bootstrap et pv-w pour les p-values obtenues en utilisant le test présenté dans le corollaire 3.5.3.1.

Dans l'ensemble, la sélection de variables reste stable quelle que soit la fonction lien utilisée ou la dimension. Mais les p-values semblent meilleures en grande dimension. Ce qui peut être justifié par la conclusion de la section 3.5.2. C'est-à-dire la méthode spectrale est meilleure en grande dimension. En plus de cette stabilité, elle peut être meilleure en temps de calcul, du fait que la méthode bootstrap n'utilise que l'estimation des directions.

<i>Estimateur</i> \ <i>Lien</i>	<i>Logit</i>	<i>Probit</i>
$\hat{\mu}$	$\begin{pmatrix} 0.562 & \mathbf{-0.013} & \mathbf{-0.012} \\ 0.503 & 0.369 & 0.517 \\ \mathbf{0.007} & 0.749 & \mathbf{0.006} \\ 0.666 & 0.549 & \mathbf{0.001} \\ \mathbf{0.016} & \mathbf{-0.011} & 0.856 \end{pmatrix}$	$\begin{pmatrix} 0.588 & \mathbf{0.008} & \mathbf{0.044} \\ 0.466 & 0.388 & 0.476 \\ \mathbf{0.003} & 0.740 & \mathbf{0.017} \\ 0.660 & 0.548 & \mathbf{-0.011} \\ \mathbf{0.024} & \mathbf{-0.006} & 0.877 \end{pmatrix}$
$\hat{\beta}$	$\begin{pmatrix} 2.397 & \mathbf{-0.381} & \mathbf{-0.004} \\ 2.080 & 2.220 & 3.265 \\ \mathbf{0.001} & 4.436 & \mathbf{0.025} \\ 2.844 & 3.267 & \mathbf{0.010} \\ \mathbf{0.082} & \mathbf{-0.066} & 5.469 \end{pmatrix}$	$\begin{pmatrix} 2.848 & \mathbf{-0.057} & \mathbf{0.049} \\ 2.336 & 1.735 & 2.812 \\ \mathbf{0.093} & 3.308 & \mathbf{-0.037} \\ 3.388 & 2.480 & \mathbf{-0.008} \\ \mathbf{0.047} & \mathbf{0.029} & 4.601 \end{pmatrix}$
pv-b	$\begin{pmatrix} 0.05 & \mathbf{0.63} & \mathbf{0.69} \\ 0.076 & 0.09 & 0.04 \\ \mathbf{0.87} & 0.02 & \mathbf{0.79} \\ 0.06 & 0.04 & \mathbf{0.96} \\ \mathbf{0.66} & \mathbf{0.57} & 0.02 \end{pmatrix}$	$\begin{pmatrix} 0.056 & \mathbf{0.840} & \mathbf{0.491} \\ 0.097 & 0.086 & 0.060 \\ \mathbf{0.949} & 0.035 & \mathbf{0.721} \\ 0.092 & 0.051 & \mathbf{0.710} \\ \mathbf{0.517} & \mathbf{0.824} & 0.028 \end{pmatrix}$
pv-w	$\begin{pmatrix} 6.0e-07 & \mathbf{7.2e-01} & \mathbf{0.9} \\ 2.1e-07 & 1.9e-13 & 0.0 \\ \mathbf{9.9e-01} & 3.4e-11 & \mathbf{0.5} \\ 6.3e-07 & 7.9e-13 & \mathbf{0.7} \\ \mathbf{7.1e-01} & \mathbf{8.9e-02} & 0.0 \end{pmatrix}$	$\begin{pmatrix} 2.8e-05 & \mathbf{5.6e-01} & \mathbf{5.6e-02} \\ 5.5e-05 & 3.8e-07 & 3.8e-15 \\ \mathbf{3.9e-01} & 2.8e-06 & \mathbf{1.2e-01} \\ 1.1e-05 & 1.8e-07 & \mathbf{6.4e-01} \\ \mathbf{4.0e-01} & \mathbf{1.8e-01} & 1.4e-14 \end{pmatrix}$

TABLE 3.7 – Expérience 8 : résultats de sélection de variables. pv-b, pour les p-values obtenues par sélection purement bootstrap et pv-w pour les p-values obtenues en utilisant le test présenté dans le corollaire 3.5.3.1.

Chapitre 4

Extensions des mélanges de modèles linéaires généralisés

Sommaire

4.1	Introduction	100
4.2	Mélanges de modèles linéaires généralisés	100
4.2.1	Covariables continues	100
4.2.2	Covariables continues et catégorielles	102
4.3	Données longitudinales	103

4.1 Introduction

Comme précisé dans le chapitre 3, le modèle étudié jusqu'à présent est un cas particulier des mélanges de modèles linéaires généralisés. En effet, nous nous sommes intéressés aux modèles de type régression logistique où g est la fonction lien est connue (par exemple la fonction logistique ou la fonction probit).

Dans ce chapitre, nous nous intéresserons à l'extension de ce modèle en supposant que g et K sont inconnus. Nous montrons dans ce cas que le modèle reste identifiable pour des covariables continues mais aussi dans le cas où le vecteur des covariables est constitué d'une partie continue et d'une partie catégorielle. On montre aussi que le modèle reste identifiable dans le cas où les données sont longitudinales (répétitions indépendantes).

Ce chapitre est organisé comme suit. La section 4.2 présentera le modèle et les résultats d'identifiabilité dans le cas où le vecteur des covariables est continue mais aussi dans le cas où le vecteur des covariables est constitué d'une partie continue et d'une partie catégorielle. On présentera à la section 4.3 les résultats d'identifiabilité pour des données longitudinales.

4.2 Mélanges de modèles linéaires généralisés

Soit (X, Y) un vecteur aléatoire de loi \mathbb{P}_θ telle que $X \in \mathbb{R}^d$, $Y \in \mathbb{R}$. On suppose que la loi de Y sachant X est un mélange donné par

$$E(Y|X) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + b_k),$$

avec $\omega = (\omega_1, \dots, \omega_K)$, $\beta = [|\beta_1|, \dots, |\beta_K|] \in \mathbb{R}^{d \times K}$, et $b = (b_1, \dots, b_K) \in \mathbb{R}^K$. On suppose que pour tout k , $\omega_k \geq 0$, que $\sum_{k=1}^K \omega_k = 1$, et que g est une fonction à valeurs dans $(0, 1)$.

Si $Y \in \{0, 1\}$ et $\theta = (\omega, \beta, b)$, la loi de Y sachant X est un mélange de lois de Bernoulli. Ce qui revient aux cas étudiés au chapitre 3.

Si $\theta = (K, g, \omega, \beta, b)$, on cherche à retrouver les paramètres inconnus K , g , ω , β et b . Il est évident qu'il est nécessaire de fixer l'origine et l'échelle, ce qui peut se faire en fixant par exemple $g(0)$ et $g(1)$ (sans perte de généralité). On notera $\mu_k = \beta_k / \|\beta_k\|$ et $\lambda_k = \|\beta_k\|$, $k = 1, \dots, K$, de sorte que $\beta_k = \lambda_k \mu_k$.

4.2.1 Covariables continues

On considère que le vecteur des covariables X est continu. Dans ce cas, on peut retrouver le nombre de composantes K et les vecteurs normalisés μ_k , sous les hypothèses suivantes :

- (H1) Le support de X est \mathbb{R}^d .
- (H2) Pour tous $j \neq k$, on a $\mu_j \neq \mu_k$ et $\mu_j \neq -\mu_k$.
- (H3) La fonction $g : \mathbb{R} \rightarrow]0, 1[$ est strictement croissante, de limite 0 en $-\infty$, de limite 1 en $+\infty$, et elle est continûment dérivable et sa dérivée a pour limite 0 en $-\infty$ et en $+\infty$. Par ailleurs, $g(0) < g(1)$ sont fixés (par exemple à $1/2$ et $3/4$).

Remarque : on ne fait aucune hypothèse sur K .

Théorème 4.2.1.1 *Sous les hypothèses (H1), (H2) et (H3), la connaissance de $\mathbb{P}_{g,\omega,\beta,b}$ permet de retrouver K et μ_1, \dots, μ_K .*

Preuve du Théorème 4.2.1.1.

Si on connaît la loi de (Y, X) alors on connaît la fonction

$$x \mapsto H(x) = \sum_{k=1}^K \omega_k g(\lambda_k \langle \mu_k, x \rangle + b_k)$$

sur le support de X , donc sur \mathbb{R}^d . On connaît aussi la fonction

$$DH(x) = \sum_{k=1}^K \omega_k g'(\lambda_k \langle \mu_k, x \rangle + b_k) \mu_k$$

et on a que si $V \in \mathbb{R}^d$, $\lim_{t \rightarrow +\infty} \|DH(tV)\|$ vaut 0 sauf si V est orthogonal à au moins un des μ_k . L'ensemble des $V \in \mathbb{R}^d$ tels que $\lim_{t \rightarrow +\infty} \|DH(tV)\| \neq 0$ est donc $\cup_{k=1}^K \langle \mu_k \rangle^\perp$, réunion d'espaces vectoriels disjoints de dimension $d-1$, dont la connaissance permet de connaître les espaces $\langle \mu_k \rangle^\perp$, et donc K et toutes les droites $\langle \mu_k \rangle$. Comme pour tout k , $\omega_k g'(b_k) > 0$, cela permet de retrouver tous les μ_k .

Comme dans le chapitre 3, le but est de retrouver tous les paramètres K, g, ω, β et b . Si on remplace l'hypothèse H2 par l'hypothèse suivante :

— (H2bis) Les vecteurs μ_1, \dots, μ_K sont linéairement indépendants,

On peut retrouver les paramètres K, g, ω, β et b .

Remarque : (H2bis) implique que $K \leq d$.

Théorème 4.2.1.2 *Sous les hypothèses (H1), (H2bis) et (H3), le modèle est identifiable : la connaissance de $\mathbb{P}_{g,\omega,\beta,b}$ permet de retrouver K, g, ω, β et b .*

Preuve du Théorème 4.2.1.2.

Par le Théorème 4.2.1.1, on connaît les K et les μ_k , on veut retrouver les λ_k , les b_k , les ω_k et g .

En considérant les espaces orthogonaux à tous les μ_k sauf un, comme les μ_k sont linéairement indépendants, on voit que l'on connaît les fonctions h_1, \dots, h_K données par, pour $j = 1, \dots, K$:

$$t \mapsto h_j(t) = \omega_j g(\lambda_j t + b_j) + \sum_{k=1, k \neq j}^K \omega_k g(b_k).$$

On a ensuite :

$$\begin{aligned} h_j(0) &= \sum_{k=1}^K \omega_k g(b_k), \\ \lim_{t \rightarrow +\infty} h_j(t) &= \omega_j + \sum_{k=1, k \neq j}^K \omega_k g(b_k), \\ \lim_{t \rightarrow -\infty} h_j(t) &= \sum_{k=1, k \neq j}^K \omega_k g(b_k). \end{aligned}$$

On peut donc retrouver les ω_j et les $g(b_j)$. Et donc on connaît les fonctions

$$t \mapsto \ell_j(t) = g(\lambda_j t + b_j).$$

Comme $g(0) = \ell_j(-b_j/\lambda_j)$ et $g(1) = \ell_j((1 - b_j)/\lambda_j)$ sont fixés, on peut retrouver les λ_j et les b_j , et donc pour finir la fonction g .

4.2.2 Covariables continues et catégorielles

On considère maintenant qu'une partie $X \in \mathbb{R}^d$ des covariables est continue et une autre partie Z , de dimension d' , est constituée de variables catégorielles à valeurs dans $\{z_1, \dots, z_m\} \subset \mathbb{R}^{d'}$, et que

$$E(Y|X, Z) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + \langle \gamma_k, Z \rangle + b_k).$$

On note $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ la loi de (X, Y) , avec $\omega = (\omega_1, \dots, \omega_K)$, $\beta = [\beta_1, \dots, \beta_K] \in \mathbb{R}^{d \times K}$, $\gamma = [\gamma_1, \dots, \gamma_K] \in \mathbb{R}^{d' \times K}$, et $b = (b_1, \dots, b_K) \in \mathbb{R}^K$.

Sous ce modèle, on a l'identifiabilité en ajoutant l'hypothèse suivante

— (H4) La matrice $\begin{pmatrix} 1 & z_1^T \\ 1 & z_2^T \\ \vdots & \vdots \\ 1 & z_m^T \end{pmatrix}$ est de rang plein.

Remarques :

Remarquons que (H4) implique $d' + 1 \leq m$ et que les covariables continues qui permettent d'identifier g .

Théorème 4.2.2.1 *Sous les hypothèses (H1), (H2bis), (H3) et (H4), le modèle est identifiable : la connaissance de $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ permet de retrouver K , g , ω , β , γ et b .*

Preuve du Théorème 4.2.2.1.

Par le Théorème 4.2.1.1, appliqué aux lois de Y sachant X et $Z = z$ pour tout $z \in \{z_1, \dots, z_m\}$, la connaissance de $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ permet de retrouver K , g , ω , β , et $A_k = (a_{k,i})_{1 \leq i \leq m}$, $k = 1, \dots, K$, avec

$$a_{k,i} = b_k + \langle \gamma_k, z_i \rangle.$$

On connaît alors pour tout k

$$A_k = \begin{pmatrix} 1 & z_1^T \\ 1 & z_2^T \\ \vdots & \vdots \\ 1 & z_m^T \end{pmatrix} \begin{pmatrix} b_k \\ \gamma_k \end{pmatrix}$$

ce qui permet de retrouver les b_k et γ_k si (H4) est vérifiée.

4.3 Données longitudinales

On suppose ici qu'on est dans le cas des données longitudinales. C'est-à-dire, on a des répétitions indépendantes (conditionnellement à la population) avec des covariables différentes. L'observation Y est m -dimensionnelle (si on a m répétitions) et on a m covariables. Le modèle est donné dans ce cas par

$$E(Y|X, Z) = \sum_{k=1}^K \omega_k \otimes_{j=1}^m g(\langle \beta_k, X_j \rangle + \langle \gamma_k, Z_j \rangle + b_k).$$

On a alors l'identifiabilité sous l'hypothèse faible (H2) dès que l'on a au moins 3 répétitions.

Théorème 4.3.0.1 *On suppose $m \geq 3$. Sous les hypothèses (H1), (H2), (H3) et (H4), le modèle est identifiable : la connaissance de $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ permet de retrouver K , g , ω , β , γ et b .*

Preuve du Théorème 4.3.0.1.

Si l'on connaît la loi de Y , alors, pour tout z , on connaît la fonction H de $(\mathbb{R}^d)^m$ dans $(0, 1)^m$ donnée par

$$H(x_1, \dots, x_m) = \sum_{k=1}^K \omega_k \otimes_{j=1}^m g(\langle \beta_k, x_j \rangle + \tilde{b}_k(z))$$

avec $\tilde{b}_k(z) = b_k + \langle \gamma_k, z_i \rangle$. Montrons tout d'abord que pour tout z , les fonctions $g(\langle \beta_k, \cdot \rangle + \tilde{b}_k(z))$ sont linéairement indépendantes. En effet, si $\alpha_1, \dots, \alpha_K$ sont des réels tels que pour tout $x \in \mathbb{R}^d$,

$$\sum_{k=1}^K \alpha_k g(\langle \beta_k, x \rangle + \tilde{b}_k(z)) = 0,$$

alors par dérivation pour tout $x \in \mathbb{R}^d$,

$$\sum_{k=1}^K \alpha_k g'(\langle \beta_k, x \rangle + \tilde{b}_k(z)) \beta_k = 0.$$

En prenant $V \in \langle \beta_k \rangle^\perp$ tel que $V \notin \langle \beta_j \rangle^\perp$, $j \neq k$, ce qui est possible sous (H2), puis $x = tV$ et t tend vers l'infini, on obtient que $\alpha_k g'(\tilde{b}_k(z)) \beta_k = 0$, et donc $\alpha_k = 0$.

Par la méthode spectrale analogue à celle permettant de montrer l'identifiabilité des mélanges multidimensionnels non paramétriques, on voit que la connaissance de H permet de retrouver K , et, pour tout z , les ω_k et fonctions $g(\langle \beta_k, \cdot \rangle + \tilde{b}_k(z))$.

Si l'on connaît la fonction $x \mapsto g(\lambda_k \langle \mu_k, x \rangle + \tilde{b}_k(z))$ on retrouve μ_k par dérivation, puis g , puis les $\tilde{b}_k(z)$ comme pour le Théorème 4.2.1.2 puis les γ_k et les b_k comme pour le Théorème 4.2.2.1.

Chapitre 5

Mélange de valeurs extrêmes en présence de censure

Sommaire

5.1	Introduction	108
5.2	Modèle de mélange et valeurs extrêmes	110
5.3	Estimation des paramètres	111
5.3.1	Vraisemblance en dessous du seuil u	111
5.3.2	Vraisemblance au dessus du seuil u	112
5.3.3	Estimation	113
5.4	Estimation des quantiles extrêmes	115
5.4.1	Par la fonction de répartition du modèle extrême	115
5.4.2	Par la méthode de reparamétrisation	116
5.5	Étude de simulation	116
5.5.1	Pour un seuil u fixé	117
5.5.2	Pour un seuil u inconnu	119
5.5.3	Conclusion	121
5.6	Discussion & conclusion	122

Résumé

La théorie des valeurs extrêmes appelée “Extreme value theory” EVT en anglais, est une vaste théorie dont le but est d’étudier les événements rares. C’est-à-dire, les événements dont la probabilité d’apparition est faible. Par exemple les intempéries, les inondations, les catastrophes naturelles, . . . Il est donc important de pouvoir déterminer un seuil suffisamment grand au dessus duquel les données sont considérées extrêmes. Dans ce chapitre, on observe un mélange à deux composantes : (1) une composante en dessous du seuil, appelé “bulk” modèle et (2) une composante au dessus du seuil appelé “tail” modèle. Dans ce chapitre, nous considérons qu’en dessous du seuil, on a un modèle paramétrique (exemple Weibull) et une GPD (Generalized Pareto Distribution) au dessus du seuil. De plus, on suppose que les données au dessus du seuil ne sont pas complètement observées. Elles sont censurées aléatoirement à droite par une variable de loi extrême.

Cependant, il faut noter que la vraisemblance de ce modèle ne peut pas pleinement tirer profit de l’algorithme EM, qui est couramment utilisé dans l’étude des modèles de mélange. En effet, le seuil u est un paramètre commun aux deux composantes et poids du mélange. De ce fait des méthodes bayésiennes paramétriques et non-paramétriques sont souvent utilisées pour estimer les paramètres des mélanges de valeurs extrêmes.

Pour stabiliser l’estimation par maximum de vraisemblance, nous proposons ici une méthode d’estimation en deux étapes : (1) on estime d’abord par maximum de vraisemblance les paramètres du modèle en fixant la valeur du seuil u . Ensuite (2) on répète la procédure (1) sur une grille de valeurs de u pour en sélectionner celle qui correspond à la plus grande vraisemblance. Avec cette méthode d’estimation, nous montrons par simulation que l’augmentation de la censure diminue la qualité de l’estimation des paramètres au dessus du seuil. Ainsi, en cas de forte censure, il faut suffisamment de données pour bien estimer les paramètres par le maximum de vraisemblance.

Mots clés : Censure aléatoire, Maximum de vraisemblance, Mélange de valeurs extrêmes, Théorie des valeurs extrêmes.

Abstract

Extreme Value Theory (EVT) is used to develop models for studying rare events, that are events with low probability of occurrence. For example : bad weather, floods natural disasters, So, it is important to be able to determine a threshold sufficiently large above of which the data are considered extreme. In this case, we observe a mixture with two components : (1) a component below the threshold, called “bulk” model and (2) a component above the threshold called “tail” model. In this chapter, we considered that below the threshold we have a parametric model (example Weibull) and a GPD (Generalized Pareto Distribution) beyond the threshold. Moreover, we assume that the data above the threshold are not completely observed. They are randomly censored on the right by a variable with an extreme law.

However, it should be noted that the likelihood of this model can not fully benefit from the EM algorithm, which is commonly used in the study of mixing models. Indeed, the threshold u is a parameter common to both components and weights of the mixture. Therefore, parametric and non-parametric Bayesian methods are often used to study this type of model.

To stabilize the maximum likelihood estimation, we propose here a two-step estimation method : (1) we first estimate the model parameters using maximum likelihood by setting the threshold value u . Next (2) repeat the procedure (1) for several values of u to select the value with corresponds to the highest likelihood. With this estimation method, we show by simulation that the increase of the censoring decreases the quality of the parameters estimation above the threshold. Thus, in case of strong censorship, more data is needed to properly estimate the parameters.

Keywords : Extreme value mixture, Extreme value theory, Maximum likelihood, Random censoring.

5.1 Introduction

La théorie des valeurs extrêmes appelée “Extreme value theory” EVT en anglais, est une vaste théorie dont le but est d’étudier les événements rares ([41]). C’est-à-dire, les événements dont la probabilité d’apparition est faible. Par exemple les intempéries, les inondations, les catastrophes naturelles, . . .

L’un des défis pour les modèles des valeurs extrêmes est de déterminer un seuil suffisamment élevé, au dessus duquel les données peuvent être considérées comme extrêmes. Dans ce cas, le modèle est vu comme un mélange composé de deux sous-modèles : (1) un sous-modèle en dessous du seuil, appelé “bulk” modèle et (2) un autre au dessus du seuil appelé “tail” modèle. Le modèle résultant est appelé modèle de mélange de valeurs extrêmes (“Extreme value mixture model” en anglais, voir Figure 5.1). Dans le passé, des choix de seuil étaient faits en utilisant des outils graphiques. Aujourd’hui le challenge est de considérer le seuil comme un paramètre du modèle à estimer.

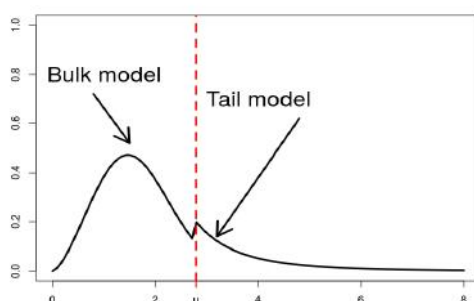


FIGURE 5.1 – Exemple de modèle de mélange de valeurs extrêmes. La ligne pointillée représente le seuil u .

Plusieurs auteurs sont récemment allés dans ce sens en utilisant différentes méthodes d’estimation. En 2002, Frigessi et al.([80]) ont utilisé un modèle dynamique pondéré en combinant une loi de Weibull pour le bulk modèle avec une loi de Pareto généralisée (GPD) pour le tail modèle. Ils ont considéré les poids du mélange comme fonction de la fonction de répartition d’une loi de Cauchy, ce qui augmente le nombre de paramètre à estimer. Hu ([51]) montre dans sa thèse que si le paramètre d’échelle de la loi de Cauchy est proche de 0, la qualité de l’estimation diminue. D’autres auteurs ont essayé d’utiliser la distribution en dessous du seuil pour définir les poids du mélange. C’est le cas de Behrens et al. ([81]) et de Mendes et Lopes ([82]) (qui utilisent deux bulk modèles en même temps, c’est-à-dire deux seuils u_1 et u_2 pour définir deux modèles extrêmes). En 2006, Trancedi et al.([83]) ont essayé de construire le modèle en considérant que les poids du mélange ne dépendent que de la probabilité de dépasser le seuil. Mais ils ont utilisé une méthode non-paramétrique en dessous du seuil pour estimer les paramètres. Par la suite plusieurs auteurs ont essayé d’aller dans la même direction en utilisant une méthode bayésienne pour estimer les paramètres (MacDonald et al. ([57]), Zhao et al.([84]), Zhao et al. ([85])). En 2012, Lee et al. ([86]) utilisent un modèle exponentiel en dessous du seuil avec une GPD pour le tail modèle en utilisant une méthode dite “Peaks-over threshold” (voir la thèse de Toulemonde [44] pour plus de détails) pour estimer les pa-

ramètres extrêmes. Un peu plus tôt, en 2011, Nascimento et al. ([87]) ont combiné une méthode semi-paramétrique et la méthode bayésienne mais en utilisant en dessous du seuil un mélange de distributions gamma. L'un des problèmes du modèle de mélange extrême est la régularité au point u . Carreau et Bengio ([88, 89]) ont commencé par ajouter des contraintes de continuité de la densité et de ses dérivées au point u . Comme précisé dans MacDonald et al. ([57]), la forme paramétrique est la plus simple des modèles de mélanges extrêmes (Frigessi et al. (2002, [80]), Behrens et al. (2004, [81]) et Zhao et al. (2010, [84])). Cependant, il faut noter que la vraisemblance de ce modèle ne peut pas pleinement tirer profit de l'algorithme EM, qui est couramment utilisé dans l'étude des modèles de mélange. En effet, le seuil u est un paramètre commun aux deux composantes et poids du mélange. Ce qui fait que la méthode bayésienne est souvent utilisée dans le cadre de l'estimation des paramètres extrêmes (Coles et al. [90, 91]).

Dans ce chapitre, nous considérons qu'en dessous du seuil, on a un modèle paramétrique (exemple Weibull) et une GPD au dessus du seuil. De plus, on suppose que les données au dessus du seuil ne sont pas complètement observées. Elles sont censurées aléatoirement à droite par une variable de loi extrême.

Nous proposons dans ce chapitre une méthode d'estimation en deux étapes : (1) on estime d'abord par maximum de vraisemblance les paramètres du modèle en fixant la valeur du seuil u . Ensuite (2) on répète la procédure (1) sur une grille de valeurs de u pour sélectionner celle qui correspond à la plus grande vraisemblance.

Ce chapitre est organisé comme suit. La section 5.2 présente le modèle étudié en détail. Les sections 5.3 et 5.4 présentent la méthode d'estimation de paramètres et des quantiles extrêmes respectivement. La section 5.5 présente des résultats de simulation qui illustrent l'estimation des paramètres. Une discussion est donnée à la section 5.6.

5.2 Modèle de mélange et valeurs extrêmes

Soit $X_1 \dots, X_n$ un échantillon indépendant et identiquement distribué de même loi que X . La fonction de répartition F_X , de X est définie, comme dans les travaux de MacDonald et al. [57], par :

$$F_X(x) = \begin{cases} (1 - \phi_u) \frac{H(x|\beta, \lambda)}{H(u|\beta, \lambda)} & \text{si } x \leq u \\ (1 - \phi_u) + \phi_u G(x|u, \sigma_u, \xi) & \text{si } x > u \end{cases} \quad (5.1)$$

avec $\phi_u = \mathbb{P}(X > u)$ et $G(\cdot|u, \sigma_u, \xi)$ la fonction de répartition de la loi de Pareto généralisée (GPD), définie par :

$$G(z|u, \sigma_u, \xi) = \begin{cases} 1 - \left[1 + \xi \left(\frac{z-u}{\sigma_u}\right)_+\right]^{-\frac{1}{\xi}} & \text{si } \xi \neq 0 \\ 1 - \exp\left[-\left(\frac{z-u}{\sigma_u}\right)_+\right] & \text{si } \xi = 0. \end{cases} \quad (5.2)$$

On suppose qu'en dessous du seuil u , la variable X est modélisée par une loi de Weibull et au dessus du seuil par une loi de Pareto généralisée. C'est-à-dire que la fonction de répartition est définie en dessous du seuil u par la fonction de répartition de la loi de Weibull $H(\cdot|\beta, \lambda)$ qui dépend de β (paramètre de forme) et de λ (paramètre d'échelle). Elle est définie par

$$H(z|\beta, \lambda) = 1 - e^{-\left(\frac{z}{\lambda}\right)^\beta}, \quad \beta, \lambda > 0, \quad x \in \mathbb{R}_+. \quad (5.3)$$

La densité est définie dans ce cas par

$$h(x|\beta, \lambda) = \frac{\beta}{\lambda} \left(\frac{x}{\lambda}\right)^{\beta-1} e^{-\left(\frac{x}{\lambda}\right)^\beta}, \quad \beta, \lambda > 0, \quad x \in \mathbb{R}_+. \quad (5.4)$$

Au dessus du seuil u , la fonction de répartition de X est définie par la fonction de répartition de la loi de Pareto généralisée (GPD) $G(\cdot|u, \sigma_u, \xi)$ donnée par l'équation (5.2). En plus du seuil u , elle dépend des paramètres extrêmes (σ_u, ξ) .

Comme présentée dans [57], la densité de X peut être vue comme une densité d'un mélange à deux composantes :

$$f_X(x) = (1 - \phi_u) f_1(x) + \phi_u f_2(x) \quad (5.5)$$

avec

$$\begin{aligned} f_1(x) &= \frac{h(x|\beta, \lambda)}{H(u|\beta, \lambda)} \mathbf{1}_{\{x \leq u\}} \\ f_2(x) &= g(x|u, \sigma_u, \xi) \mathbf{1}_{\{x > u\}} \end{aligned}$$

où $g(\cdot|u, \sigma_u, \xi)$ est la densité de la GPD définie par

$$g(z|u, \sigma_u, \xi) = \begin{cases} \frac{1}{\sigma_u} \left[1 + \xi \left(\frac{z-u}{\sigma_u} \right) \right]_+^{-\frac{\xi+1}{\xi}} & \text{si } \xi \neq 0 \\ \frac{1}{\sigma_u} \exp \left[- \left(\frac{z-u}{\sigma_u} \right)_+ \right] & \text{si } \xi = 0 \end{cases} \quad (5.6)$$

et

$$\mathbb{1}_{\{x>u\}} = \begin{cases} 1 & \text{si } x > u \\ 0 & \text{si } x \leq 0. \end{cases}$$

On suppose ici que la variable d'intérêt X n'est pas complètement observée au dessus du seuil u . Un moyen de modéliser cette situation est d'introduire une variable C , indépendante de X , telle que

$$Z = X \wedge C \text{ et } \delta = \mathbb{1}_{\{X \leq C\}}$$

sont observées. Supposons que la variable C est modélisée par une loi de Pareto généralisée (GPD(u, σ'_u, ξ')) qui dépend du seuil u et des paramètres extrêmes (σ'_u, ξ'). La fonction de répartition de C est donnée dans ce cas par :

$$F_C(x) = \begin{cases} 0 & \text{si } x \leq u \\ G(x|u, \sigma'_u, \xi') & \text{si } x > u \end{cases} \quad (5.7)$$

5.3 Estimation des paramètres

Soit $(\mathbf{Z}, \delta) = ((Z_1, \delta_1), \dots, (Z_n, \delta_n))$ un échantillon indépendant et identiquement distribué de même loi que (Z, δ) .

5.3.1 Vraisemblance en dessous du seuil u

En dessous du seuil u , les données sont complètement observées avec une densité donnée par

$$\frac{(1 - \phi_u)}{H(u|\beta, \lambda)} h(x|\beta, \lambda) \mathbb{1}_{\{x \leq u\}},$$

où h est définie par l'équation (5.4). Définissons $A(u) = \{j : 0 \leq z_j \leq u\}$, C'est-à-dire, l'ensemble des indices des observations en dessous du seuil u . La vraisemblance est donnée dans ce cas par

$$\mathcal{L}_W(\mathbf{Z}|\beta, \lambda, u) = \left\{ \frac{(1 - \phi_u)}{H(u|\beta, \lambda)} \right\}^{|A(u)|} \prod_{i \in A(u)} h(z_i|\beta, \lambda), \quad (5.8)$$

où h est définie par l'équation (5.4).

5.3.2 Vraisemblance au dessus du seuil u

Rappelons que la variable d'intérêt X n'est pas complètement observée au dessus du seuil. Soit $B(u) = \{j : z_j > u\}$, l'ensemble des indices des observations au dessus du seuil u . On observe alors $(\mathbf{Z}, \delta) = (Z_j, \delta_j)_{j \in B(u)}$ avec $Z_j = X_j \wedge C_j$ et $\delta_j = \mathbf{1}_{\{X_j \leq C_j\}}$, $j \in B(u)$. Rappelons aussi que les variables X et C sont indépendantes avec les densités respectives f_X et f_C . Les fonctions de répartition associées sont F_X pour X et F_C pour C . Ici f_X est égale à une constante près, à la densité d'une GPD qui dépend de (u, σ_u, ξ) et f_C est égale à la densité d'une GPD qui dépend de (u, σ'_u, ξ') . Rappelons que $\theta = (\beta, \lambda, u, \sigma_u, \xi)$. Ainsi, la contribution d'un individu i à la vraisemblance est proportionnelle à

$$f_{Z_i, \delta_i}(z_i, \delta_i | \theta) dz = \mathbb{P}\left(Z_i \in [z_i, z_i + dz], \delta_i = 1 | \theta\right)^{\delta_i} \times \mathbb{P}\left(Z_i \in [z_i, z_i + dz], \delta_i = 0 | \theta\right)^{1-\delta_i} \quad (5.9)$$

Posons

$$L_i = \mathbb{P}\left(Z_i \in [z_i, z_i + dz], \delta_i = 1 | \theta\right)^{\delta_i} \times \mathbb{P}\left(Z_i \in [z_i, z_i + dz], \delta_i = 0 | \theta\right)^{1-\delta_i}.$$

Comme

$$\delta_i = \begin{cases} 1 & \text{si } X_i \leq C_i \text{ (on observe } X_i) \\ 0 & \text{si } X_i > C_i \text{ (} X_i \text{ est censurée)} \end{cases}$$

on a

$$\begin{aligned} L_i &= \mathbb{P}\left(X_i \in [z_i, z_i + dz], X_i \leq C_i | \theta\right)^{\delta_i} \times \mathbb{P}\left(C_i \in [z_i, z_i + dz], X_i > C_i | \theta\right)^{1-\delta_i} \\ &= \mathbb{P}\left(X_i \in [z_i, z_i + dz], z_i \leq C_i | \theta\right)^{\delta_i} \times \mathbb{P}\left(C_i \in [z_i, z_i + dz], X_i > z_i | \theta\right)^{1-\delta_i}. \end{aligned}$$

En utilisant l'indépendance entre X et C , on a

$$\begin{aligned} L_i &= \left(\mathbb{P}(X_i \in [z_i, z_i + dz] | \theta) \times \mathbb{P}(z_i \leq C_i)\right)^{\delta_i} \times \left(\mathbb{P}(C_i \in [z_i, z_i + dz]) \times \mathbb{P}(X_i > z_i | \theta)\right)^{1-\delta_i} \\ &= \left(\mathbb{P}(X_i \leq z_i + dz | \theta) - \mathbb{P}(X_i \leq z_i | \theta)\right)^{\delta_i} \times \mathbb{P}(z_i \leq C_i)^{\delta_i} \\ &\times \left(\mathbb{P}(C_i \leq z_i + dz) - \mathbb{P}(C_i \leq z_i)\right)^{1-\delta_i} \times \mathbb{P}(X_i > z_i | \theta)^{1-\delta_i} \\ &= \left(f_X(z_i | \theta) dz \cdot S'(z_i)\right)^{\delta_i} \times \left(f_C(z_i) dz \cdot S(z_i | \theta)\right)^{1-\delta_i} \end{aligned}$$

avec S' et S les fonctions de survie de C et X respectivement. En remplaçant dans l'équation (5.9), on a

$$\begin{aligned} f_{Z, \delta}(z_i, \delta_i | \theta) &= \left(f_X(z_i | \theta) \cdot S'(z_i)\right)^{\delta_i} \times \left(f_C(z_i) \cdot S(z_i | \theta)\right)^{1-\delta_i} \\ &\propto f_X(z_i | \theta)^{\delta_i} \times S(z_i | \theta)^{1-\delta_i} \\ &\propto g(z_i | u, \sigma_u, \xi)^{\delta_i} \times S(z_i | u, \sigma_u, \xi)^{1-\delta_i}. \end{aligned} \quad (5.10)$$

La vraisemblance est donnée, dans ce cas, par

$$\mathcal{L}_{G\hat{P}D}(z|u, \sigma_u, \xi) \propto \prod_{j \in B} g(z_j|u, \sigma_u, \xi)^{\delta_j} \times S(z_j|u, \sigma_u, \xi)^{1-\delta_j} \quad (5.11)$$

où g et S sont la densité et la fonction de survie d'une GPD respectivement.

5.3.3 Estimation

Dans cette partie, nous allons estimer les paramètres du modèle défini à la section 5.2. On commencera par supposer que le seuil u , au dessus duquel les données sont considérées extrêmes, est connu. Ensuite, on étudiera le cas où u n'est pas connu et qu'il est considéré comme un paramètre inconnu à estimer.

La méthode habituelle du maximum de vraisemblance revient à maximiser sous $\theta = (\beta, \lambda, u, \sigma_u, \xi)$ le produit des équations (5.8) et (5.11). Mais comme précisé dans [57], la maximisation de cette vraisemblance ne marche pas très bien. En effet, le support de la densité du modèle dépend du paramètre u à estimer. Dans ce cas, nous proposons de faire une estimation par maximum de vraisemblance en deux étapes :

Estimation avec un seuil u connu

Comme le seuil u est connu, le vecteur de paramètres à estimer est composé des paramètres de la loi de Weibull (β, λ) et des paramètres extrêmes (σ_u, ξ) . C'est-à-dire le vecteur de paramètres à estimer est $\theta = (\beta, \lambda, \sigma_u, \xi)$. Dans ce cas, la vraisemblance du modèle est donnée en dessous du seuil u par

$$\mathcal{L}_W(\mathbf{Z}|\beta, \lambda, u) = \left\{ \frac{(1 - \phi_u)}{H(u|\beta, \lambda)} \right\}^{|A(u)|} \prod_{i \in A(u)} \left[\frac{\beta}{\lambda} \left(\frac{z_i}{\lambda} \right)^{\beta-1} e^{-\left(\frac{z_i}{\lambda} \right)^\beta} \right], \quad (5.12)$$

et au dessus du seuil, par

$$\mathcal{L}_{G\hat{P}D}(z|u, \sigma_u, \xi) \propto \prod_{j \in B(u)} g(z_j|u, \sigma_u, \xi)^{\delta_j} \times S(z_j|u, \sigma_u, \xi)^{1-\delta_j}. \quad (5.13)$$

En utilisant l'équation (5.6) et

$$S(z_j|u, \sigma_u, \xi) = 1 - G(z_j|u, \sigma_u, \xi) = \begin{cases} \left[1 + \xi \left(\frac{z_j - u}{\sigma_u} \right) \right]_+^{-\frac{1}{\xi}} & \text{si } \xi \neq 0 \\ \exp \left[- \left(\frac{z_j - u}{\sigma_u} \right) \right]_+ & \text{si } \xi = 0 \end{cases} \quad (5.14)$$

la vraisemblance au dessus du seuil est donnée par

$$\begin{aligned} \mathcal{L}_{G\hat{P}D}(z|u, \sigma_u, \xi) &\propto \phi_u^{|B|} \prod_{j \in B} \left(\frac{1}{\sigma_u} \right)^{\delta_j} \left\{ \left[1 + \xi \left(\frac{z_j - u}{\sigma_u} \right) \right]_+^{-\frac{(1+\xi\delta_j)}{\xi}} \mathbf{1}_{\{\xi \neq 0\}} \right. \\ &\quad \left. + \exp \left[- \left(\frac{z_j - u}{\sigma_u} \right) \right]_+ \mathbf{1}_{\{\xi = 0\}} \right\}. \end{aligned} \quad (5.15)$$

La vraisemblance complète du modèle est donnée par le produit des deux équations (5.15) et (5.12). C'est-à-dire

$$\begin{aligned}
\mathcal{L}(\mathbf{Z}|\theta) &\propto \left\{ \frac{(1 - \phi_u)}{1 - e^{-\left(\frac{u}{\lambda}\right)^\beta}} \right\}^{|A|} \prod_{i \in A} \frac{\beta}{\lambda} \left(\frac{z_i}{\lambda}\right)^{\beta-1} e^{-\left(\frac{z_i}{\lambda}\right)^\beta} \\
&\times \phi_u^{|B|} \prod_{j \in B} \left(\frac{1}{\sigma_u}\right)^{\delta_j} \left\{ \left[1 + \xi \left(\frac{z_j - u}{\sigma_u}\right) \right]_+^{-\frac{(1+\xi\delta_j)}{\xi}} \mathbf{1}_{\{\xi \neq 0\}} \right. \\
&\left. + \exp \left[-\left(\frac{z_j - u}{\sigma_u}\right) \right]_+ \mathbf{1}_{\{\xi = 0\}} \right\}. \tag{5.16}
\end{aligned}$$

Pour estimer les paramètres, on utilise la méthode du maximum de vraisemblance. Pour un échantillon (\mathbf{Z}, δ) de taille n , on cherche l'estimateur $\hat{\theta}_n$ qui maximise l'équation (5.16) ou son logarithme. C'est-à-dire

$$\hat{\theta}_n = \underset{\theta}{\operatorname{argmax}} \log(\mathcal{L}(\mathbf{Z}|\theta)) \tag{5.17}$$

Avec un seuil u fixé, on peut bien maximiser cette vraisemblance de manière stable. Comme le seuil u est un paramètre du modèle à estimer, nous proposons une deuxième étape permettant de prendre en compte l'estimation de u en se basant sur la vraisemblance donnée par l'équation (5.16).

Estimation avec le seuil u comme paramètre

On suppose dans cette partie que le seuil u n'est pas connu et qu'il fait partie des paramètres du modèle à estimer. Dans ce cas, le vecteur de paramètres θ est composé du seuil u , des paramètres de la loi de Weibull (β, λ) et des paramètres extrêmes (σ_u, ξ) . C'est-à-dire le vecteur de paramètres à estimer est donné par $\theta = (\beta, \lambda, u, \sigma_u, \xi)$.

Comme évoqué à la section 5.1, la vraisemblance du modèle présente dans ce cas des problèmes de régularité au point u . Ce qui fait qu'on ne peut pas maximiser directement la vraisemblance donnée par l'équation (5.16). Rappelons aussi que pour chaque u_l , $l = 1, \dots, L$, $A(u_l)$ et $B(u_l)$ sont données par :

$$A(u_l) = \{j : 0 \leq z_j \leq u_l\}$$

et

$$B(u_l) = \{j : z_j > u_l\}.$$

En pratique, pour chaque valeur u_l fixée, on peut estimer les $A(u_l)$, $B(u_l)$ associés et leurs cardinales.

Pour estimer les paramètres dans ce cas, on se donne une grille de L valeurs de u pour ensuite chercher l'estimateur du maximum de vraisemblance associé à chacune de ces valeurs. Soient u_1, u_2, \dots, u_L , les L différentes valeurs du seuil choisi (en pratique, on peut choisir deux quantiles extrêmes empiriques $(q_1$ et $q_2)$ puis discrétiser l'intervalle $[q_1, q_2]$). Supposons pour tout $l = 1, \dots, L$, que $\hat{\theta}_n^l$ est l'estimateur du maximum de vraisemblance associé à u_l . C'est-à-dire $\hat{\theta}_n^l$ maximise le logarithme de l'équation (5.16) avec $u = u_l$.

En effet, on a

$$\hat{\theta}_n^l = (\hat{\beta}_l, \hat{\lambda}_l, \hat{\sigma}_{ul}, \hat{\xi}_l) = \operatorname{argmax}_{\beta, \lambda, \sigma_u, \xi} \log (\mathcal{L}(\mathbf{Z}|\beta, \lambda, \sigma_u, \xi)) \quad (5.18)$$

avec $\mathcal{L}(\mathbf{Z}|\beta, \lambda, \sigma_u, \xi)$ donné par l'équation (5.16).

Notons par $\mathcal{LL}(\hat{\theta}_n^l)$ le logarithme de la vraisemblance (5.16) au point $\hat{\theta}_n^l$ avec $u = u_l$, $l = 1, \dots, L$. Ainsi, pour chaque u_l , on a l'estimateur $\hat{\theta}_n^l = (\hat{\beta}_l, \hat{\lambda}_l, \hat{\sigma}_{ul}, \hat{\xi}_l)$ et le logarithme de la vraisemblance associée $\mathcal{LL}(\hat{\theta}_n^l)$. On obtient l'estimateur $\hat{\theta}_n$ du modèle en maximisant $\mathcal{LL}(\hat{\theta}_n^l)$, $l = 1, \dots, L$. C'est-à-dire $\hat{\theta}_n$ est donné par

$$\hat{\theta}_n = (\hat{\beta}, \hat{\lambda}, \hat{u}, \hat{\sigma}_u, \hat{\xi}) = \operatorname{argmax}_{1, \dots, L} \mathcal{LL}(\hat{\theta}_n^l). \quad (5.19)$$

5.4 Estimation des quantiles extrêmes

Un quantile extrême x_p d'ordre $1 - p$ est tel que

$$F(x_p) = 1 - p,$$

où p est une valeur petite telle que $0 < p < 1$. En utilisant cette définition et l'équation (5.1), on peut estimer le quantile x_p de deux manières différentes : (1) par la méthode classique avec l'inverse de la fonction de répartition ou (2) par la reparamétrisation de la vraisemblance.

5.4.1 Par la fonction de répartition du modèle extrême

En utilisant l'inverse de la fonction de répartition du modèle extrême, le quantile x_p donné par

$$x_p = \begin{cases} u + \frac{\sigma_u}{\xi} [(\phi_u^{-1}p)^{-\xi} - 1] & \text{si } \xi \neq 0 \\ u - \sigma_u \log [\phi_u^{-1}p] & \text{si } \xi = 0 \end{cases} \quad (5.20)$$

— Si u est connu, on a $\hat{\theta}_n = (\hat{\beta}, \hat{\lambda}, \hat{\sigma}_u, \hat{\xi})$ et l'estimateur de x_p est donné dans ce cas

$$\hat{x}_p = \begin{cases} u + \frac{\hat{\sigma}_u}{\hat{\xi}} [(\hat{\phi}_u^{-1}p)^{-\hat{\xi}} - 1] & \text{si } \hat{\xi} \neq 0 \\ u - \hat{\sigma}_u \log [\hat{\phi}_u^{-1}p] & \text{si } \hat{\xi} = 0 \end{cases} \quad (5.21)$$

— Si u est inconnu, on a $\hat{\theta}_n = (\hat{\beta}, \hat{\lambda}, \hat{u}, \hat{\sigma}_u, \hat{\xi})$ et l'estimateur de x_p est donné dans ce cas

$$\hat{x}_p = \begin{cases} \hat{u} + \frac{\hat{\sigma}_u}{\hat{\xi}} [(\hat{\phi}_u^{-1}p)^{-\hat{\xi}} - 1] & \text{si } \hat{\xi} \neq 0 \\ \hat{u} - \hat{\sigma}_u \log [\hat{\phi}_u^{-1}p] & \text{si } \hat{\xi} = 0. \end{cases} \quad (5.22)$$

Remarquons que la qualité de l'estimation des quantiles par cette méthode dépend fortement de la qualité de l'estimation initiale des paramètres du modèle. Dans le but d'estimer simultanément les paramètres du modèle et les quantiles extrêmes, nous proposons d'utiliser la méthode de reparamétrisation.

5.4.2 Par la méthode de reparamétrisation

En utilisant la forme du quantile donné par l'équation (5.20), on peut reparamétriser la vraisemblance donnée par l'équation (5.16). Ainsi le seuil u peut être donné par l'équation (5.20) sous la forme suivante :

$$u = \begin{cases} x_p - \frac{\sigma_u}{\xi} [(\phi_u^{-1}p)^{-\xi} - 1] & \text{si } \xi \neq 0 \\ x_p + \sigma_u \log [\phi_u^{-1}p] & \text{si } \xi = 0. \end{cases} \quad (5.23)$$

On peut ainsi réécrire la vraisemblance (5.16) comme

$$\begin{aligned} \mathcal{L}(\mathbf{Z}|\beta, \lambda, x_p, \sigma_u, \xi) &\propto (1 - \phi_u)^{|A|} \phi_u^{|B|} \prod_{i \in A} \frac{\beta}{\lambda} \left(\frac{z_i}{\lambda}\right)^{\beta-1} e^{-(\frac{z_i}{\lambda})^\beta} \times \\ &\left[\left\{ 1 - \exp \left[-\frac{1}{\lambda^\beta} \left(x_p - \frac{\sigma_u}{\xi} [(\phi_u^{-1}p) - 1] \right)^\beta \right] \right\}^{-|A|} \right. \\ &\times \prod_{j \in B} \left(\frac{1}{\sigma_u} \right)^{\delta_j} \left[1 + \frac{\xi}{\sigma_u} \left(z_j - x_p + \frac{\sigma_u}{\xi} [(\phi_u^{-1}p) - 1] \right) \right]_+^{-\frac{(1+\xi\delta_j)}{\xi}} \mathbf{1}_{\{\xi \neq 0\}} \\ &+ \left\{ 1 - \exp \left[-\frac{1}{\lambda^\beta} \left(x_p + \sigma_u \log [\phi_u^{-1}p] \right)^\beta \right] \right\}^{-|A|} \\ &\times \left. \prod_{j \in B} \left(\frac{1}{\sigma_u} \right)^{\delta_j} \exp \left[-\frac{1}{\sigma_u} (z_j - x_p - \sigma_u \log [\phi_u^{-1}p]) \right]_+ \mathbf{1}_{\{\xi=0\}} \right]. \quad (5.24) \end{aligned}$$

Ainsi, par maximum de vraisemblance, on peut estimer les quantiles x_p

$$\hat{\theta}'_n = (\hat{\beta}, \hat{\lambda}, \hat{x}_p, \hat{\sigma}_u, \hat{\xi}) = \underset{\beta, \lambda, x_p, \sigma_u, \xi}{\operatorname{argmax}} \log (\mathcal{L}(\mathbf{Z}|\beta, \lambda, x_p, \sigma_u, \xi)). \quad (5.25)$$

5.5 Étude de simulation

Dans cette partie, nous présentons les résultats d'estimation basés sur des simulations. Nous commencerons par présenter les résultats dans le cas où u est fixé et connu pour ensuite présenter le cas où u est inconnu et considéré comme paramètre à estimer. Dans toute cette partie, les simulations sont faites avec les valeurs de θ suivantes :

$$\begin{cases} \beta = 2.5 \\ \lambda = 1.2 \\ u = 1.451613 \\ \sigma_u = 1 \\ \xi = 0.3 \end{cases}$$

où la valeur de u est choisie comme le quantile 0.8 de la loi Weibull de paramètre (β, λ) .

5.5.1 Pour un seuil u fixé

Comme précisé à la section 5.3.3, on suppose dans cette partie que u est fixé et connu. C'est-à-dire, on connaît le seuil au dessus duquel les données sont extrêmes. Dans ce cas, le vecteur de paramètre θ à estimer est donné par $\theta = (\beta, \lambda, \sigma_u, \xi)$.

Estimation des paramètres du modèle

Soient n la taille de l'échantillon utilisée et N le nombre de réplifications (ici fixé à $N = 1000$.) Pour différentes valeurs de n ($n = 100, 500, 1000, 1500$) et différents taux de censure donnés, on peut estimer les paramètres du modèle pour mesurer l'effet de la censure sur la qualité de nos estimateurs. Les tableaux 5.1 à 5.4 présentent les résultats pour des taux de censure de 0%, 10%, 30%, 50%, 70%, et 90%. Dans chaque tableau, l'estimateur est donné, de même que l'erreur quadratique moyenne (*RMSE*), l'erreur absolue moyenne [*MAE*] et l'écart type {*SD*}.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	2.559	2.541	2.540	2.564	2.568	2.529
	(0.390)	(0.365)	(0.374)	(0.399)	(0.387)	(0.382)
	[0.308]	[0.288]	[0.295]	[0.308]	[0.303]	[0.306]
	{0.386}	{0.362}	{0.372}	{0.394}	{0.381}	{0.381}
500	2.507	2.506	2.509	2.511	2.520	2.515
	(0.164)	(0.169)	(0.159)	(0.162)	(0.168)	(0.164)
	[0.130]	[0.135]	[0.126]	[0.129]	[0.136]	[0.129]
	{0.164}	{0.169}	{0.159}	{0.161}	{0.167}	{0.163}
1000	2.505	2.504	2.509	2.507	2.508	2.507
	(0.120)	(0.115)	(0.113)	(0.114)	(0.113)	(0.115)
	[0.096]	[0.093]	[0.090]	[0.091]	[0.090]	[0.091]
	{0.120}	{0.115}	{0.113}	{0.114}	{0.113}	{0.114}
1500	2.500	2.507	2.504	2.504	2.502	2.509
	(0.095)	(0.094)	(0.095)	(0.090)	(0.095)	(0.095)
	[0.077]	[0.075]	[0.075]	[0.072]	[0.076]	[0.076]
	{0.095}	{0.094}	{0.094}	{0.090}	{0.095}	{0.095}

TABLE 5.1 – Résultats de simulations pour l'estimateur de β avec $N = 1000$ réplifications. Pour chaque taux de censure l'estimateur de β est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l'écart type associé.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	1.595	1.45	1.55	1.392	1.432	1.933
	(2.875)	(1.955)	(2.762)	(1.781)	(2.441)	(6.493)
	[0.496]	[0.352]	[0.456]	[0.298]	[0.339]	[0.834]
	{2.849}	{1.940}	{2.741}	{1.771}	{2.432}	{6.455}
500	1.2103	1.212	1.210	1.209	1.212	1.210
	(0.089)	(0.092)	(0.083)	(0.087)	(0.089)	(0.091)
	[0.066]	[0.066]	[0.063]	[0.064]	[0.066]	[0.067]
	{0.088}	{0.091}	{0.083}	{0.086}	{0.089}	{0.090}
1000	1.204	1.207	1.203	1.204	1.204	1.204
	(0.059)	(0.056)	(0.054)	(0.058)	(0.056)	(0.056)
	[0.047]	[0.044]	[0.042]	[0.045]	[0.043]	[0.044]
	{0.059}	{0.056}	{0.054}	{0.058}	{0.056}	{0.056}
1500	1.204	1.202	1.203	1.203	1.203	1.200
	(0.047)	(0.046)	(0.046)	(0.045)	(0.047)	(0.045)
	[0.037]	[0.036]	[0.036]	[0.036]	[0.036]	[0.035]
	{0.047}	{0.046}	{0.046}	{0.045}	{0.047}	{0.045}

TABLE 5.2 – Résultats de simulations pour l’estimateur de λ avec $N = 1000$ réplifications. Pour chaque taux de censure l’estimateur de λ est donné, de même que $(.)$ pour le RMSE, $[.]$ pour le MAE et $\{.\}$ pour l’écart type associé.

Les tableaux 5.1 et 5.2 présentent les résultats pour l’estimation des paramètres du modèle en dessous du seuil u . C’est-à-dire les estimateurs de β et de λ . On peut voir à partir de ces tableaux que la censure n’a pas d’effet spécifique sur l’estimation. En effet, en augmentant ou en diminuant le taux de censure, on ne constate pas d’augmentation ni diminution du biais ou de l’écart type des estimateurs. Ce qui est normal dans le sens où la censure ne concerne que le modèle au dessus du seuil à savoir les paramètres σ_u et ξ . Par contre, l’augmentation de la taille des échantillons n améliore bien les résultats d’estimation au sens de l’erreur quadratique mais aussi en termes d’erreur absolue moyenne et d’écart type.

Les tableaux 5.3 et 5.4 présentent les résultats d’estimation des paramètres au dessus du seuil u (σ_u et ξ). On peut constater à partir de ces tableaux que plus la censure augmente, plus les estimateurs se dégradent. L’augmentation de la taille de l’échantillon n permet d’améliorer ces résultats. Ce qui veut dire qu’il faut nécessairement avoir beaucoup de données pour bien estimer les paramètres par maximum de vraisemblance en cas de forte censure.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	1.190	1.258	1.315	1.381	1.507	21.088
	(0.549)	(0.719)	(0.812)	(1.013)	(7.451)	(53.881)
	[0.383]	[0.457]	[0.527]	[0.656]	[0.933]	[20.597]
	{0.515}	{0.671}	{0.749}	{0.939}	{7.437}	{50.021}
500	1.035	1.029	1.046	1.042	1.031	1.243
	(0.179)	(0.181)	(0.204)	(0.246)	(0.350)	(1.049)
	[0.138]	[0.142]	[0.158]	[0.188]	[0.262]	[0.618]
	{0.175}	{0.179}	{0.199}	{0.243}	{0.349}	{1.021}
1000	1.009	1.018	1.017	1.010	1.017	1.106
	(0.116)	(0.122)	(0.134)	(0.151)	(0.231)	(0.545)
	[0.092]	[0.097]	[0.105]	[0.120]	[0.179]	[0.379]
	{0.116}	{0.121}	{0.133}	{0.151}	{0.231}	{0.535}
1500	1.011	1.006	1.012	1.011	1.022	1.069
	(0.096)	(0.099)	(0.108)	(0.122)	(0.195)	(0.408)
	[0.077]	[0.080]	[0.086]	[0.098]	[0.152]	[0.299]
	{0.095}	{0.099}	{0.107}	{0.122}	{0.193}	{0.403}

TABLE 5.3 – Résultats de simulations pour l’estimateur de σ_u avec $N = 1000$ réplifications. Pour chaque taux de censure l’estimateur de σ_u est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l’écart type associé.

5.5.2 Pour un seuil u inconnu

On suppose dans cette partie que le seuil u est inconnu et qu’il est considéré comme un paramètre à estimer. Dans ce cas, le vecteur de paramètre θ à estimer est donné par $\theta = (\beta, \lambda, u, \sigma_u, \xi)$.

Estimation des paramètres du modèle

Comme présenté à la section 5.5.1, on note n la taille de l’échantillon utilisée et N le nombre de réplifications fixé à $N = 1000$. Pour $n = 100, 500, 1000$ et 1500 , on estime les paramètres du modèle en faisant varier le taux de censure de 0%, 10%, 30%, 50%, 70%, et 90%. Les tableaux de 5.5 à 5.9 présentent les résultats d’estimation, de l’erreur quadratique moyenne (*RMSE*), de l’erreur absolue moyenne [*MAE*] et de l’écart type {*SD*}.

Les tableaux 5.5 et 5.6 présentent les résultats pour l’estimation des paramètres du modèle en dessous du seuil u (β et λ). On peut voir à partir de ces tableaux qu’il n’y a pas de comportement spécifique des estimateurs en fonction du taux de censure. En effet,

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	0.153	0.117	0.046	0.078	1.086	-1.887
	(0.417)	(0.477)	(0.681)	(1.296)	(6.053)	(44.428)
	[0.317]	[0.361]	[0.503]	[0.964]	[3.045]	[25.213]
	{0.391}	{0.441}	{0.632}	{1.277}	{6.005}	{44.397}
500	0.265	0.278	0.260	0.267	0.414	0.919
	(0.147)	(0.152)	(0.202)	(0.370)	(1.401)	(15.035)
	[0.116]	[0.121]	[0.158]	[0.283]	[1.031]	[9.776]
	{0.142}	{0.151}	{0.198}	{0.369}	{1.397}	{15.030}
1000	0.287	0.287	0.277	0.298	0.388	0.498
	(0.097)	(0.100)	(0.134)	(0.233)	(0.949)	(8.257)
	[0.076]	[0.078]	[0.105]	[0.183]	[0.736]	[6.050]
	{0.096}	{0.099}	{0.132}	{0.233}	{0.945}	{8.258}
1500	0.289	0.292	0.287	0.293	0.311	0.237
	(0.077)	(0.084)	(0.106)	(0.186)	(0.754)	(5.893)
	[0.061]	[0.067]	[0.083]	[0.148]	[0.596]	[4.459]
	{0.076}	{0.084}	{0.105}	{0.186}	{0.755}	{5.896}

TABLE 5.4 – Résultats de simulations pour l'estimateur de ξ avec $N = 1000$ répliques. Pour chaque taux de censure l'estimateur de ξ est donné, de même que (.) pour le RMSE, [.] pour le MAE et { . } pour l'écart type associé.

en augmentant ou en diminuant le taux de censure, on ne constate pas une diminution ou une augmentation de l'erreur absolue moyenne, de l'écart type ou de l'erreur quadratique moyenne des estimateurs. Mais l'augmentation de la taille des échantillons n améliore aussi fortement les résultats d'estimation au sens de l'erreur quadratique et de l'écart type.

Le tableau 5.7 présente les résultats pour l'estimation du seuil u . On peut voir à partir de ce tableau que l'estimateur de u s'améliore en fonction du taux de censure. En effet, en augmentant le taux de censure, on constate une diminution de l'écart type et de l'erreur quadratique moyenne des estimateurs. De même l'augmentation de la taille des échantillons n améliore les résultats d'estimation au sens de l'erreur quadratique mais aussi en terme d'écart type.

Les tableaux 5.8 et 5.9 présentent les résultats d'estimation des paramètres au dessus du seuil u (σ_u et ξ). On peut constater que les résultats se dégradent en fonction de la censure. En effet, une augmentation de la censure détériore les résultats en termes

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	2.863	2.794	2.720	2.716	2.730	2.664
	(0.568)	(0.543)	(0.482)	(0.492)	(0.481)	(0.386)
	[0.451]	[0.420]	[0.375]	[0.375]	[0.368]	[0.305]
	{0.437}	{0.457}	{0.430}	{0.442}	{0.423}	{0.305}
500	2.868	2.719	2.612	2.602	2.561	2.666
	(0.485)	(0.405)	(0.318)	(0.327)	(0.295)	(0.311)
	[0.404]	[0.312]	[0.243]	[0.249]	[0.221]	[0.247]
	{0.316}	{0.341}	{0.298}	{0.310}	{0.289}	{0.263}
1000	2.842	2.716	2.603	2.587	2.558	2.670
	(0.430)	(0.339)	(0.239)	(0.231)	(0.188)	(0.281)
	[0.358]	[0.254]	[0.163]	[0.151]	[0.128]	[0.206]
	{0.260}	{0.261}	{0.215}	{0.214}	{0.179}	{0.224}
1500	2.866	2.755	2.669	2.672	2.638	2.705
	(0.431)	(0.335)	(0.231)	(0.235)	(0.190)	(0.290)
	[0.373]	[0.266]	[0.178]	[0.176]	[0.144]	[0.219]
	{0.228}	{0.218}	{0.157}	{0.161}	{0.130}	{0.205}

TABLE 5.5 – Estimation avec u inconnu. Résultats de simulations pour l'estimateur de β avec $N = 1000$ réplifications. Pour chaque taux de censure l'estimateur de β est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l'écart type associé.

d'erreur quadratique moyenne, d'erreur absolue moyenne et d'écart type. L'augmentation de la taille de l'échantillon n permet d'améliorer ces résultats. Ce qui veut dire qu'il faut nécessairement avoir beaucoup de données pour bien estimer les paramètres par maximum de vraisemblance en cas de forte censure.

5.5.3 Conclusion

Globalement, on constate que l'estimation est meilleure dans le cas où on connaît le seuil u que dans le cas où u est inconnu. De plus la censure n'a pas d'effet spécifique sur l'estimation des paramètres en dessous du seuil. L'effet de la censure est mis en évidence dans l'estimation des paramètres au dessus du seuil. Ce qui est normal dans la mesure où la censure n'est observée qu'au dessus du seuil. Notons aussi que la taille de l'échantillon semble plus influente que la censure.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	9.493	11.384	10.249	8.371	11.291	11.772
	(16.573)	(18.457)	(16.839)	(14.804)	(18.264)	(21.023)
	[8.539]	[10.395]	[9.236]	[7.363]	[10.282]	[10.711]
	{14.356}	{15.401}	{14.208}	{12.958}	{15.231}	{18.181}
500	1.255	1.575	1.788	1.827	1.890	1.516
	(0.577)	(0.866)	(1.015)	(1.045)	(1.084)	(0.755)
	[0.377]	[0.619]	[0.778]	[0.806]	[0.852]	[0.480]
	{0.575}	{0.782}	{0.828}	{0.837}	{0.836}	{0.686}
1000	1.119	1.267	1.423	1.465	1.491	1.267
	(0.262)	(0.330)	(0.389)	(0.404)	(0.413)	(0.268)
	[0.242]	[0.300]	[0.361]	[0.378]	[0.387]	[0.211]
	{0.249}	{0.323}	{0.319}	{0.305}	{0.293}	{0.259}
1500	1.071	1.155	1.267	1.293	1.311	1.184
	(0.195)	(0.190)	(0.187)	(0.185)	(0.185)	(0.144)
	[0.194]	[0.189]	[0.186]	[0.184]	[0.184]	[0.134]
	{0.146}	{0.184}	{0.174}	{0.160}	{0.148}	{0.143}

TABLE 5.6 – Estimation avec u inconnu. Résultats de simulations pour l’estimateur de λ avec $N = 1000$ répliquions. Pour chaque taux de censure l’estimateur de λ est donné, de même que $(.)$ pour le RMSE, $[.]$ pour le MAE et $\{.\}$ pour l’écart type associé.

5.6 Discussion & conclusion

Pour des raisons évoquées à la section 5.1, la méthode habituelle du maximum de vraisemblance n’est pas très robuste pour estimer les paramètres dans le cas d’un modèle de mélange de valeurs extrêmes où le seuil u est inconnu.

L’objectif de ce chapitre était de montrer que même si les données extrêmes sont censurées, on arrive à bien estimer les paramètres et les quantiles. Pour ce faire, nous avons proposé, dans cette étude, une méthode basée sur le maximum de vraisemblance, mais en deux étapes. Avec cette méthode, on arrive à estimer les paramètres du modèle dans le cas où les données sont censurées mais aussi dans le cas où il n’y a pas de censure. Si le seuil u est connu, la méthode reste efficace mais présente un problème de biais dans le cas où la taille de l’échantillon est petite. Si u est inconnu, on arrive à bien estimer les paramètres mais la qualité de l’estimation reste meilleure dans le cas où u est connu. Dans les deux cas, l’estimation pourra être améliorée en augmentant la taille de l’échantillon même si la censure est forte. Rappelons que la méthode reste assez coûteuse en temps de calcul du fait de la recherche de la valeur optimale de u dans la deuxième étape.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	1.647	1.568	1.500	1.487	1.482	1.439
	(0.371)	(0.311)	(0.242)	(0.212)	(0.188)	(0.149)
	[0.326]	[0.276]	[0.221]	[0.196]	[0.175]	[0.130]
	{0.316}	{0.288}	{0.237}	{0.209}	{0.185}	{0.149}
500	1.694	1.563	1.477	1.461	1.439	1.471
	(0.319)	(0.248)	(0.197)	(0.182)	(0.172)	(0.125)
	[0.303]	[0.236]	[0.191]	[0.178]	[0.169]	[0.118]
	{0.228}	{0.222}	{0.195}	{0.182}	{0.171}	{0.124}
1000	1.684	1.541	1.423	1.396	1.381	1.468
	(0.308)	(0.234)	(0.184)	(0.173)	(0.167)	(0.124)
	[0.297]	[0.226]	[0.181]	[0.172]	[0.166]	[0.118]
	{0.223}	{0.216}	{0.181}	{0.165}	{0.151}	{0.123}
1500	1.678	1.540	1.406	1.376	1.359	1.468
	(0.299)	(0.230)	(0.179)	(0.168)	(0.165)	(0.122)
	[0.290]	[0.223]	[0.177]	[0.167]	[0.164]	[0.117]
	{0.196}	{0.182}	{0.173}	{0.151}	{0.136}	{0.121}

TABLE 5.7 – Estimation avec u inconnu. Résultats de simulations pour l'estimateur de u avec $N = 1000$ réplifications. Pour chaque taux de censure l'estimateur de u est donné, de même que $(.)$ pour le RMSE, $[.]$ pour le MAE et $\{.\}$ pour l'écart type associé.

Dans le cas où la censure est forte et que la taille de l'échantillon est petite, on peut penser à utiliser une méthode bayésienne comme le suggère MacDonald et al. (cas sans censure, [57]). A court terme, on essaiera d'utiliser la méthode bayésienne pour améliorer les résultats dans ce cas précis mais aussi à présenter des résultats numériques sur l'estimation des quantiles extrêmes.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	0.752	0.988	1.319	1.487	2.028	25.074
	(1.172)	(1.534)	(1.968)	(2.368)	(11.745)	(71.488)
	[0.948]	[1.083]	[1.342]	[1.555]	[2.217]	[25.342]
	{1.147}	{1.535}	{1.942}	{2.318}	{11.706}	{67.346}
500	0.386	0.61	0.809	0.934	1.223	1.101
	(0.727)	(0.736)	(0.757)	(0.850)	(1.160)	(2.155)
	[0.675]	[0.704]	[0.722]	[0.823]	[1.081]	[1.548]
	{0.460}	{0.636}	{0.733}	{0.848}	{1.139}	{2.153}
1000	0.41	0.634	0.947	1.146	1.509	1.024
	(0.674)	(0.682)	(0.688)	(0.727)	(1.086)	(1.848)
	[0.719]	[0.734]	[0.792]	[0.799]	[1.016]	[1.442]
	{0.447}	{0.575}	{0.636}	{0.713}	{0.959}	{1.848}
1500	0.409	0.627	0.978	1.186	1.645	0.98
	(0.614)	(0.625)	(0.637)	(0.680)	(1.006)	(1.683)
	[0.611]	[0.626]	[0.633]	[0.647]	[0.928]	[1.387]
	{0.336}	{0.533}	{0.587}	{0.631}	{0.898}	{1.684}

TABLE 5.8 – Estimation avec u inconnu. Résultats de simulations pour l'estimateur de σ_u avec $N = 1000$ réplifications. Pour chaque taux de censure l'estimateur de σ_u est donné, de même que $(.)$ pour le RMSE, $[.]$ pour le MAE et $\{.\}$ pour l'écart type associé.

n	Taux de censure					
	0%	10%	30%	50%	70%	90%
100	0.328	0.246	0.072	0.099	-0.263	-6.785
	(0.627)	(0.739)	(1.104)	(1.535)	(4.770)	(105.191)
	[0.479]	[0.548]	[0.816]	[1.129]	[2.059]	[21.036]
	{0.627}	{0.738}	{1.080}	{1.523}	{4.739}	{105.005}
500	0.578	0.495	0.451	0.359	-0.440	-1.170
	(0.394)	(0.399)	(0.482)	(0.667)	(1.622)	(5.028)
	[0.348]	[0.354]	[0.409]	[0.587]	[1.263]	[1.695]
	{0.279}	{0.348}	{0.458}	{0.664}	{1.444}	{4.810}
1000	0.578	0.490	0.377	0.204	-0.619	-0.689
	(0.363)	(0.369)	(0.401)	(0.529)	(1.534)	(4.211)
	[0.315]	[0.338]	[0.339]	[0.480]	[1.213]	[1.232]
	{0.243}	{0.309}	{0.393}	{0.520}	{1.229}	{4.095}
1500	0.585	0.499	0.358	0.210	-0.750	-0.665
	(0.354)	(0.361)	(0.384)	(0.497)	(1.657)	(3.667)
	[0.309]	[0.316]	[0.322]	[0.444]	[1.019]	[1.049]
	{0.237}	{0.305}	{0.380}	{0.469}	{1.182}	{3.540}

TABLE 5.9 – Résultats de simulations pour l'estimateur de ξ avec $N = 1000$ réplifications. Pour chaque taux de censure l'estimateur de ξ est donné, de même que (.) pour le RMSE, [.] pour le MAE et {.} pour l'écart type associé.

Chapitre 6

Conclusion et perspectives

Le premier objectif de cette thèse était de fournir un outil de diagnostic de la co-infection et de proposer une recommandation de traitement. En utilisant les données de Kédougou, nous avons proposé une méthodologie qui pourra aider le médecin à élaborer un bon diagnostic de la co-infection. Pour la recommandation de traitement, il n'existe pas de données permettant de quantifier la part de chaque maladie (paludisme et arbovirose) en cas de co-infection. Mais notre étude donne des indications sur la maladie à traiter, en cas de co-infection, en se basant sur les symptômes du patient.

Notre analyse est basée sur deux jeux de données construites à partir d'un jeu de données réelles : le jeu de données *IgM-data* (qui contient les vrais positifs aux arbovirus mais très peu, 39 cas sur 12288) et le jeu de données *IgM/IgG-data* (qui est bien équilibré mais qui présente des cas d'arbovirus qui ne sont pas forcément positifs). Les deux jeux de données ont été utilisés dans l'analyse et l'identification des facteurs influents de chaque maladie. Dans l'analyse prédictive, seul le jeu de données *IgM/IgG-data* a été utilisé. Ce qui nous donne une classification avec une précision globale de 65%.

Dans le but d'améliorer les résultats de classification et de recommandation de traitement, nous envisageons dans le future une récolte de données de co-infection où on suivra les patients qu'à l'après traitement. Nous envisageons aussi d'appliquer la méthodologie aux données de co-infection entre le paludisme et d'autres agents pathogènes plus facilement détectables, au début de l'infection, que les arbovirus.

Dans la partie méthodologie, nous avons d'abord considéré que les données provenaient d'un mélange de population. Et comme les réponses étaient binaires, il était judicieux d'étudier des mélanges de modèles linéaires généralisés. Dans ce cas, nous n'avons pas utilisé les méthodes habituelles (méthodes variationnelles de Bayes, heuristiques basées sur le maximum de vraisemblance (EM), ...) qui peuvent converger vers des optimums locaux et peuvent présenter des temps de calcul assez longs. Nous avons testé des méthodes spectrales qui sont récemment utilisées dans l'apprentissage des mélanges mais elles restent instables en pratique. Dans le but de stabiliser l'estimation des paramètres, nous avons utilisé une méthode des moments basée sur l'algorithme spectral et les moindres carrés. Dans ce cas, on arrive à montrer qu'on peut bien retrouver les paramètres du modèle à partir des trois premiers moments (identifiabilité). On arrive aussi à montrer qu'on a des estimateurs consistants et asymptotiquement gaussiens. A partir d'une étude de simulation, nous avons montré qu'on estime bien les paramètres dans le cas des liens logit et probit. Nous montrons aussi que le temps de calcul reste assez faible comparé

à la méthode du maximum de vraisemblance. Mais que pour des dimensions petites, le maximum de vraisemblance reste meilleur en terme de biais et de variance.

Dans le jeu de données de co-infection, nous avons des covariables continues mais aussi des covariables discrètes. De plus, l'étude de la co-infection est faite dans le chapitre 2 avec une réponse multinomiale. Dans le but d'apporter plus de réponses à la question initiale en prenant compte de la nature des données, nous avons travaillé sur l'extension des mélanges de modèles linéaires généralisés dans le chapitre 3. Dans cette partie, nous avons supposé que la fonction lien était inconnue puis montré des résultats d'identifiabilité dans les cas suivants : (1) le vecteur de covariable est continu, (2) le vecteur de covariables est composé d'une partie continue et d'une partie catégorielle. Un résultat d'identifiabilité a aussi été montré dans le cas où les données seraient longitudinales.

Au terme de ces résultats d'identifiabilité, il serait intéressant de regarder, à court terme,

- l'estimation en s'inspirant de la Slice Inverse Regression (voir l'article de Babichev et Bach [92]).
- l'estimation non paramétrique de tous les paramètres par méthode de sélection de modèles.

Quelques travaux ayant un lien avec les questions posées : [93, 94, 95, 96, 97, 98, 99].

La troisième partie de cette thèse porte sur les mélanges des valeurs extrêmes en présence de censure. La méthode du maximum de vraisemblance habituelle ne marche pas très bien pour ces types de modèles mais nous avons montré qu'en pratique, on peut bien estimer les paramètres en utilisant une méthode à deux étapes : (1) une première étape basée sur le maximum de vraisemblance à un seuil fixé et (2) une deuxième étape basée sur la maximisation de la vraisemblance sur une grille de seuil. Dans ce cas, les quantiles extrêmes peuvent être estimés par reparamétrisation ou par la méthode classique.

A court terme, nous envisageons d'utiliser une méthode bayésienne pour améliorer les résultats dans le cas où la censure est forte et que la taille de l'échantillon est petite. Nous envisageons aussi d'appliquer la méthode d'estimation sur tous les modèles de mélanges d'extrêmes présentés dans la thèse de Hu ([51]) mais en supposant qu'on a de la censure sur les données extrêmes.

Annexe A

Multinomial logistic model for coinfection diagnosis between arbovirus and malaria in Kedougou

Abstract

In tropical regions, populations continue to suffer morbidity and mortality from malaria and arboviral diseases. In Kedougou (Senegal), these illnesses are all endemic due to the climate and its geographical position. The co-circulation of malaria parasites and arboviruses can explain the observation of coinfecting cases. Indeed there is strong resemblance in symptoms between these diseases making problematic targeted medical care of coinfecting cases. This is due to the fact that the origin of illness is not obviously known. Some cases could be immunized against one or the other of the pathogens, immunity typically acquired with factors like age and exposure as usual for endemic area. Then, coinfection needs to be better diagnosed. Using data collected from patients in Kedougou region, from 2009 to 2013, we adjusted a multinomial logistic model and selected relevant variables in explaining coinfection status. We observed specific sets of variables explaining each of the diseases exclusively and the coinfection. We tested the independence between arboviral and malaria infections and derived coinfection probabilities from the model fitting. In case of a coinfection probability greater than a threshold value to be calibrated on the data, duration of illness above 3 days and age above 10 years-old are mostly indicative of arboviral disease while body temperature higher than 40°C and presence of nausea or vomiting symptoms during the rainy season are mostly indicative of malaria disease.

keyword Arbovirus, coinfection, malaria, multinomial logistic regression, random forest classification, variable selection.

A.1 Introduction

Concurrent infections are often observed among vector borne diseases such as malaria and arthropod-borne viral diseases (arbovirus) in tropical regions ([1, 2]). It is the case for malaria and dengue in American, African and Asian tropical regions where their endemic

areas overlap largely ([3, 4, 5, 6, 7, 8, 9]). Malaria can be easily ascribed to other febrile illnesses because its clinical symptoms are often indistinguishable from those initially seen in dengue or chikungunya for instance ([10]). Since the introduction of the Rapid Diagnostic Test (RDT) in 2007 in Senegal, malaria has been better diagnosed and an important decrease had been noticed in the prevalence of malaria. Thus we may think that malaria has been overestimated for some time at the expense of other febrile diseases such as arbovirus or bacteria ([11, 12]). Presumptive treatment of fever with antimalarial is widely practiced to reduce malaria attributable mortality. This practice means that ill patients may be inappropriately treated, particularly where rapid diagnosis test kits are not readily available, or if the opportunity to test for arboviral infections is missed. Thus, misdiagnosis of arbovirus coinfections as malaria infections may be a cause for underestimating emerging arbovirus infections. In 2009, surveillance of acute febrile illness (AFI) was implemented in Kedougou for early detection of arbovirus outbreaks and malaria and in order to accurately measure disease morbidity and mortality in this geographical region. Due to co-circulation of malaria parasites and arbovirus, that were mainly dengue (DEN), chikungunya (CHIK), Zika (ZIK), yellow fever (YF) and Rift Valley fever viruses (RVFV) in this region (neglecting the prevalence of other arboviral infections), concurrent infections were observed and posed a challenge for medical diagnosis ([13]). Here we compare clinical profiles of coinfecting patients to clinical profiles of mono-infected patients through the statistical analysis of a data set collected from febrile patients in the Kedougou region, Senegal from 2009 to 2013. Our study aims to characterize the risk factors of coinfection and to provide statistical indicators that improve differential diagnosis of febrile cases for arbovirus.

The data of our study were provided by the Institut Pasteur de Dakar (IPD) at Kedougou (southern-east Senegal). In this region, malaria and arbovirus are endemic due to the climate and the population movements. Data were collected through seven healthcare centers in the region : *Ninefsha rural hospital, Kedougou and Saraya Health Centers, Bandafassi and Khossanto health posts, the Kedougou military health post, and the Catholic Mission mobile team*. Inclusion criteria were (i) being at least one year old at the date of the visit, (ii) having fever (i.e., body temperature $\geq 38^{\circ}C$) and (iii) manifesting at least one clinical sign within a list of symptoms. Patients satisfying inclusion criteria were enrolled once a written informed consent was signed.

In the present paper, we propose a multinomial logistic model to analyse coinfections between arbovirus and malaria. There were four outcomes determining four groups of patients : arbovirus mono-infections (with respect to the 5 tested arbovirus), malaria mono-infections, arbovirus-malaria coinfections and controls defined as patients negative for malaria and for the 5 tested arbovirus. Febrile episodes from this control group were probably due to other circulating pathogens for which all groups were supposed to be equally exposed. We performed a covariable selection using random forests based on the variable importance measure ([58]). Then we fitted a parametric multinomial logistic model including the selected covariables and we proposed a Wald-type test to test the correlation between malaria infection and arboviral infection. As the independence hypothesis was rejected, we were able to predict the probability that a patient be coinfecting given that malaria is observed. From the analysis of the influent factors on the different outcomes, we investigated the following questions : Which factors can explain coinfection ? Which risk factors enable to distinguish between malaria and arbovirus ?

The paper is organized as follows. In Section A.2, we present the working data set. Section A.3 describes the statistical model and the variable selection. In Section A.4, we present the independence test between arbovirus and malaria infections and we propose a predictive analysis. A concluding discussion is given in Section A.5.

A.2 Data description

We based our analysis on the data from the Institut Pasteur de Dakar (IPD) at Kedougou. The initial data set included 15 523 patients and collected various features : patients' data (like sex, age, occupation, location, . . .), clinical symptoms, climate indicators and three binary infections status variables indicating (*i*) the presence or absence of malaria parasites in blood, (*ii*) detection of virus or IgM antibodies against virus and (*iii*) detection of IgG antibodies against virus. Malaria diagnosis relied on the identification of haematozoa using the thick blood smear (TBS) method. Arboviral infections were investigated by the detection of specific anti-arbovirus IgM and/or IgG antibodies using ELISA (enzyme-linked immunosorbent assay). We considered an *arboviral* case as any individual tested positive to the infection with at least one of the five arbovirus (DEN, CHIK, ZIK, YF and RVF).

Based on these data we created a new categorical response variable built from four possible combinations of the three infection status variables as follows :

$$Y = \begin{cases} 0 & \text{“Other febrile illnesses (O)”} \\ 1 & \text{“Arboviral monoinfection (A)”} \\ 2 & \text{“Malaria monoinfection (M)”} \\ 3 & \text{“Coinfection (C)”} \end{cases}$$

Category 0 corresponds to individuals that are negative for both malaria and the tested arboviral infections ; their symptoms could be due to other unknown febrile illnesses. Category 1 corresponds to individuals positive for at least one of the five tested arbovirus and negative for malaria. Category 2 corresponds to individuals negative for tested arbovirus and positive for malaria. Category 3 represents individuals simultaneously positive for malaria and for at least one of the tested arbovirus. The subjects of category 3 are said “coinfected” with malaria and arbovirus.

Our aim is to differentiate febrile syndroms that could be due to arbovirus from febrile syndroms that could be due to malaria. As coinfection in a single patient may change the spectrum of clinical symptoms, we want to identify those features that predict arboviral infection to improve medical and treatment diagnosis in the primary care setting.

A.2.1 Data set

In this study, arboviral cases are diagnosed by the detection of IgM or IgG antibodies. We can have two different ways of defining an arboviral case : (1) by considering only the detection of IgM antibodies or (2) by considering the detection of both IgM and IgG antibodies. Biologically, IgM detection among patients means that they have a recent arboviral infection. So we considered that positive IgM cases are positive arboviral cases. Ignoring individuals with missing data (974 missing data on Malaria response and 803

missing data on the covariates values), we obtained a data set of size $n = 12\,288$, called the *IgM* data set. We noticed that the distributions of the different variables with and without missing data remain similar. A summary of the *IgM data* is given in Table A.1. We can see that this data set is very imbalanced (3 arboviral or coinfecting cases per 1 000 patients) and it will require a specific statistical analysis.

Arbovirus \ Malaria	+	-	Total
+	18 (0.15%)	21 (0.16%)	39 (<i>A+</i>)
-	7 069 (57.53%)	5 180 (42.16%)	12 305 (<i>A-</i>)
Total	7 087 (<i>M+</i>)	5 201 (<i>M-</i>)	12 288

TABLE A.1 – *IgM* data. *A+* for the individuals positive to arboviral infection, *A-* for the individuals negative to arboviral infection, *M+* for the individuals positive to Malaria and *M-* for the individuals negative to Malaria.

The diagnosis of arboviral infection at the time of an acute episode is ideally based on the presence in the serum of a patient of detectable IgM. However, to obtain a more balanced data set, we decided to build a separate data set by considering arboviral infected patients as individuals who were tested positive to IgM or IgG. As 13 412 missing values were recorded on the IgG variable, the size of the data set was drastically reduced and we obtained a data set of size $n = 1\,976$ which is called *IgM/IgG* data and summarized in Table A.2. For this data set, we compared the distributions of each covariate with and without missing data on the response IgG. Except for the variable *Nasal Congestion* which is over-represented (60% of positive cases in the sample compared to 40% in the initial data set), the distributions of the other variables are similar. So we considered that ignoring individuals with missing data did not affect the predictive analysis.

Arbovirus \ Malaria	+	-	Total
+	397 (20.10%)	263 (13.31%)	633 (<i>A+</i>)
-	751 (38.00%)	565 (28.59%)	1318 (<i>A-</i>)
Total	1148 (<i>M+</i>)	828 (<i>M-</i>)	1976

TABLE A.2 – *IgM/IgG* data. *A+* for the individuals positive to arboviral infection, *A-* for the individuals negative to arboviral infection, *M+* for the individuals positive to Malaria and *M-* for the individuals negative to Malaria.

Thereafter, we will consider two data sets that are derived from the same original data set using two different encoding : 1. the *IgM/IgG* data set which is suitable to apply our entire methodology ; 2. the *IgM* data set containing the true arboviral status (from a biological point of view) which is strongly imbalanced. We will use in the next section a re-sampling strategy to handle this problem.

A.2.2 Covariates

In this data set, there are four quantitative covariables : the measured body temperature (in Celsius degrees), the number of sick days defined as the number of days between

the date of symptoms onset and the date of consultation, the patient’s age (in year) and the rainfall measure (in millimeters) which is a proxy for the season (rainy or dry). The individual rainfall measure corresponds to the rainfall measure of the patient’s month of consultation. The eleven qualitative covariables are the patient’s gender and ten other binary variable, which record presence or absence of ten symptoms : headache, eye pain, muscle pain, joint pain, cough, nausea or vomiting, chills, diarrhea, nasal congestion and icterus and/or jaundice. All the variables of the data sets are summarized in Figure A.1.

Designation	For categorical variables			quantitative variables			
	# levels	0 (%)	1 (%)	mean	median	min	max
Age				19.5	16.5	1	90
Temperature				38.97	39	38	42
Number of sick days				3.039	3	0	19
Rainfall				147.5	76.1	0	500.2
Sex (F=0 and H=1)	2	42	58				
Cephalalgia	2	6	94				
Nausea/vomiting	2	50	50				
Diarrhea	2	83	17				
Chills	2	45	55				
Cough	2	64	36				
Eye pain	2	95	5				
Joint pain	2	77	23				
Muscl pain	2	71	29				
Nasal congestion	2	54	46				
Ictere/jaudice	2	95	5				
Malaria	2	42	58				
IgM	2	99	1				
IgG	2	95	5				

FIGURE A.1 – List of variables

In our data set, females represented 42% of the population and males represented 58% of the population. In the *IgM* data set, the two categories “Coinfection” and “Arboviral monoinfection” are underrepresented, which results in irrelevant descriptive graphs. A descriptive analysis of the *IgM/IgG* data set shows that the age is positively correlated to arboviral infections whereas the temperature, nausea or vomiting, and rainfall variables are associated with malaria. For example, among the patients having nausea or vomiting symptoms, 45% had malaria monoinfection, 10% had arboviral monoinfection and 23% were coinfecting. Among the patients having a nasal congestion symptom, 31% were positive to malaria monoinfection, 21% were coinfecting and 14% were positive to arboviral monoinfection. Figure A.2 displays the distributions of age, rainfall and number of sick days over the four classes of the *IgM/IgG* data set. Overall, Figure A.2 shows that arboviral-infected patients are older than malaria-infected patients and the duration of illness is longer for many arboviral cases. Higher fevers were observed for malaria and coinfection illnesses. Figure A.2(b) shows that high values of rainfall are observed in the coinfection and malaria groups.

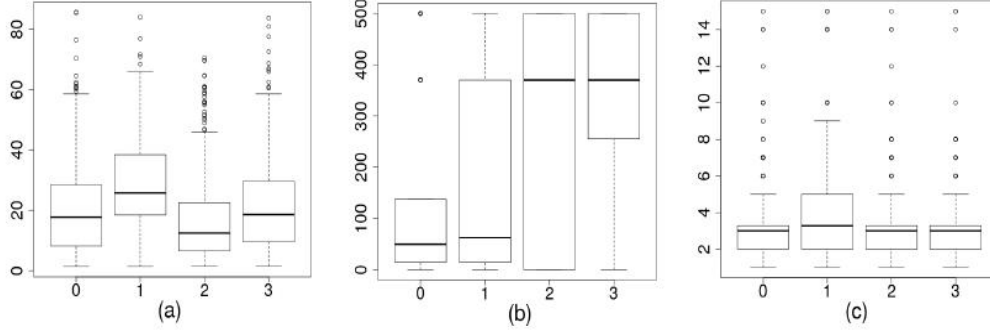


FIGURE A.2 – *IgM/IgG* data set; boxplots of the empirical distributions of the covariates (a) *age*, (b) *rainfall* and (c) *number of sick days* for the four modalities of the response variable Y : 0 (other febrile illnesses), 1 (arboviral monoinfection), 2 (malaria monoinfection) and 3 (coinfection).

A.3 Statistical analysis of the coinfection influential factors

The objective of this section is to propose a methodology that can identify the important symptoms for the arbovirus diagnosis and can help making decision for arbovirus treatment in absence of laboratory confirmation.

Variable selection is appreciable in medical data analysis as the diagnosis of the disease could be done on a minimum number of clinical measures. Reducing the number of relevant covariates may also produce more accurate classification results. In a first step, we select relevant covariates that explain the disease status typically via a multinomial logistic model. The statistical analysis is challenging because of the small number of instances of the arboviral class (39) with respect to the total number of observations (12 288). The cases that are more important for the study are rare and few exist on the available training set. We face what is usually known as a problem of imbalanced data sets. This results in models that poorly represent the rare class examples or simply ignore the observations of the minority class. To handle this problem, we proposed two data pre-processing approaches. The first one is based on biological considerations and extends the arboviral cases from 39 to 633 by merging patients with either blood positive IgM or blood positive IgG. We obtained the balanced *IgM/IgG* data set described in the previous section, which contains 1976 observations. The second approach involves randomly removing observations from the majority class to prevent its signal from dominating the fitting procedure. We applied to the imbalanced *IgM* data set a common undersampling technique to obtain a more balanced data distribution. As the data distribution is changed, it is expected that the fitted models are biased to the goals of the user and are more interpretable in terms of these goals.

In a second step we investigate the robustness of the variable selection using random forests. Introduced by [62], random forests (RF hereafter) are a robust nonparametric method to deal with classification problems. They require only mild conditions on the data generating model. They are also less sensitive to weaknesses in the data, because the randomized tree generation procedure ensures that all covariates are more equally eva-

luated. Moreover, RF decision trees often perform well on imbalanced data sets because ensemble methods offer ways to rebalance the distributions in varied ways. In this study, RF models have the advantage of providing a ranking of covariates using the RF score of variable importance that is a useful and effective tool to find important covariates for interpretation.

In a third step, we quantify the effects of the selected covariates using odds ratios. We compute odds ratios for one disease category relative to an other one and we contrast the effects of the covariates on the disease category, arboviral monoinfection, malaria monoinfection and coinfection.

A.3.1 Multinomial logit model

We recall that Y is the response variable indicating the class of the disease : “Other febrile illnesses” ($Y = 0$), “arboviral monoinfection” ($Y = 1$), “malaria monoinfection” ($Y = 2$) and “coinfection” ($Y = 3$). Let $X = (1, X_1, \dots, X_p)$ be the vector of the p covariates. For an individual with covariates $X = x$, we want to predict the probability of belonging to the class k given x ,

$$\pi_k(x) = P(Y = k|X = x), \quad k = 0, 1, 2, 3.$$

The multinomial logit model assumes the existence of $\beta_1, \beta_2, \beta_3 \in \mathbb{R}^{p+1}$ such that, for each $k = 1, 2, 3$ and each vector of covariates x ,

$$\log \frac{P(Y = k|X = x)}{P(Y = 0|X = x)} = \langle x, \beta_k \rangle \quad (\text{A.1})$$

where

$$\langle x, \beta_k \rangle = \sum_{j=0}^p x_j \beta_{kj}$$

and $x_0 = 1$ to include the intercept parameters β_{k0} , $k = 1, 2, 3$. The reference modality is class 0.

Consequently, for each $k = 1, 2, 3$ and each vector of covariates x ,

$$P(Y = k|X = x) = \frac{\exp(\langle x, \beta_k \rangle)}{1 + \sum_{l=1}^3 \exp(\langle x, \beta_l \rangle)}$$

and

$$P(Y = 0|X = x) = \frac{1}{1 + \sum_{l=1}^3 \exp(\langle x, \beta_l \rangle)}.$$

From the computation of the maximum likelihood estimates $\hat{\beta}_k$, we derive for $k = 1, 2, 3$,

$$\hat{\pi}_k(x) = \frac{e^{\langle x, \hat{\beta}_k \rangle}}{1 + \sum_{l=1}^3 e^{\langle x, \hat{\beta}_l \rangle}}. \quad (\text{A.2})$$

Application to the *IgM/IgG* data

We first give the results for the *IgM/IgG* data set since they are based on a standard logit analysis. The multinomial model was fitted to the *IgM/IgG* data by using either the `multinom` function or the `vglm` function of the `nnet` and the `VGAM` R packages. A stepwise procedure based on the AIC criterion selected eight significant covariates : *age*, *temperature*, *number of sick days*, *rainfall*, *nausea or vomiting*, *cough*, *nasal congestion and joint pain*. Likelihood-ratio tests of the sub-models obtained by removing one covariate at a time from the final model confirmed that each selected covariate was significant, with p-values less than 10^{-9} except for the variable *joint pain* that displayed a p-value of $7.44 \cdot 10^{-3}$.

Fitting strategy for handling imbalanced *IgM* data

The *IgM* data set contains 18 arboviral monoinfection cases, 21 coinfection cases, 5 180 other febrile illness cases and 7 069 malaria monoinfection cases. Trained on the original *IgM* data set, the fitted logit model only predicted classes 0 and 2, which means it ignores the two minority classes 1 and 3 in favour of the majority classes.

Applying resampling strategies to obtain a more balanced data sample is an effective solution to the imbalance problem (see [61] for a survey of existing methods). Two of the most simple resampling approaches are undersampling and oversampling. Since the *IgM* is highly imbalanced with a large number of observations in the two majority classes, we used a random undersampling strategy that removes observations and reduces the sample size. We sampled without replacement 50 cases from each of the two majority classes to create a balanced sub-sample of size $18 + 21 + 50 + 50 = 139$. Trained on a sub-sample, the model predicted four classes.

Undersampling results in loss of information and the risk of removing relevant observations is present. To overcome this problem, we repeated the sampling step a thousand times and worked with 1 000 balanced sub-samples of the *IgM* data set. The multinomial model was fitted to each sub-sample and a stepwise covariate selection was performed. The observed variability of the 1 000 covariate selections raised robustness questions. To answer this point, we conducted a nonparametric analysis based on the RF algorithm. In recent years, several methods involving the combination of resampling and ensemble learning have appeared in the imbalanced distributions literature ([61]). We found that the importance score based on random forests yielded a convenient way to summarize the information obtained from the 1 000 sub-samples.

A.3.2 Variable selection using random forests

A random forest is an ensemble of unpruned trees, induced from bootstrap samples of the training data, that uses random covariate selection in the tree construction process. Prediction is made by aggregating the predictions of the ensemble, using the majority vote rule.

One of the most widely used RF score of importance of a given variable is the Mean Decrease of Accuracy (*MDA*) in predictions. It is based on the out-of-bag (OOB) error. For each tree t of the forest, consider the associated OOB_t sample (data not included in the bootstrap sample used to construct t). Denote by $errOOB_t$ the misclassification rate

of tree t computed on this OOB_t sample. Then, randomly permute the observed values of covariate X_j in OOB_t to get a perturbed sample and compute $errOOB_t^j$, the error of t on the perturbed sample. Variable importance of X_j is then given by

$$MDA(X_j) = \frac{1}{ntree} \sum_{t=1}^{ntree} (errOOB_t^j - errOOB_t),$$

where $ntree$ denotes the number of trees of the RF. The higher the MDA , the more important the variable is. Several variable selection procedures using RF are based on this quantification of variable importance.

Using R packages, we made the following implementation choices : `randomForest` for RF fitting and MDA calculation, `VSURF` for selecting the important variables. The main parameters of `randomForest` were calibrated and set to their default values, `ntree`=500 and `mtry`= \sqrt{p} =3 (number of variables tried at each split of a tree of the RF). The variable selection strategy of `VSURF` is based on a two-stage procedure ([66]) : 1. the covariates are ranked by sorting their variable importance measures in descending order and the covariates whose importance is less than a threshold (the minimum value of the standard deviations of the importance measures) are eliminated ; 2. a sequence of nested models starting from the one with only the most important variable and ending with the one involving all important variables kept previously is considered ; the variables of the model leading to the smallest OOB error are selected. An advantage of using `VSURF` is that this procedure does not require the choice of tuning parameters.

Application to the *IgM/IgG* data

A graphical representation of the variable importance of the 15 covariates is shown in Figure A.3. The variable with the largest MDA is *rainfall*, which is indicative of the rainy season. This is expected since the development of malaria parasites is observed mostly during the rainy season. A second group of less important individual covariates are the disease symptoms : *nasal congestion*, *age* and *number of sick days*. The other covariates are ranked from the most to the least important. The `VSURF` procedure led to select the model with seven covariates : *rainfall*, *nasal congestion*, *age*, *number of sick days*, *nausea or vomiting*, *cough* and *temperature*. This result is in agreement with the logit selection variable that selected the same seven covariates and added *joint pain*.

Application to the *IgM* data

Figure A.4 ranks the variable importances (MDA) of the 15 covariates across the 1 000 sub-samples. First, *rainfall* is the most important covariate ; a second group of less important covariates is formed by *cough*, *age* and *joint pain* ; then comes a group of five covariates : *number of sick days*, *temperature*, *nausea or vomiting*, *eye pain* and *nasal congestion* ; finally, six unimportant covariates are displayed : *muscle pain*, *chills*, *cephalalgia*, *jaundice*, *diarrhea* and *sex* . The boundary between the two last groups is not clear and we used the `VSURF` procedure to separate the important covariates from the other ones. We can notice on the plot that both MDA level and variability are larger for relevant variables ; as explained by [58], this is expected and the `VSURF` threshold value is based on MDA standard deviation estimation. Figure A.5 summarizes the results of

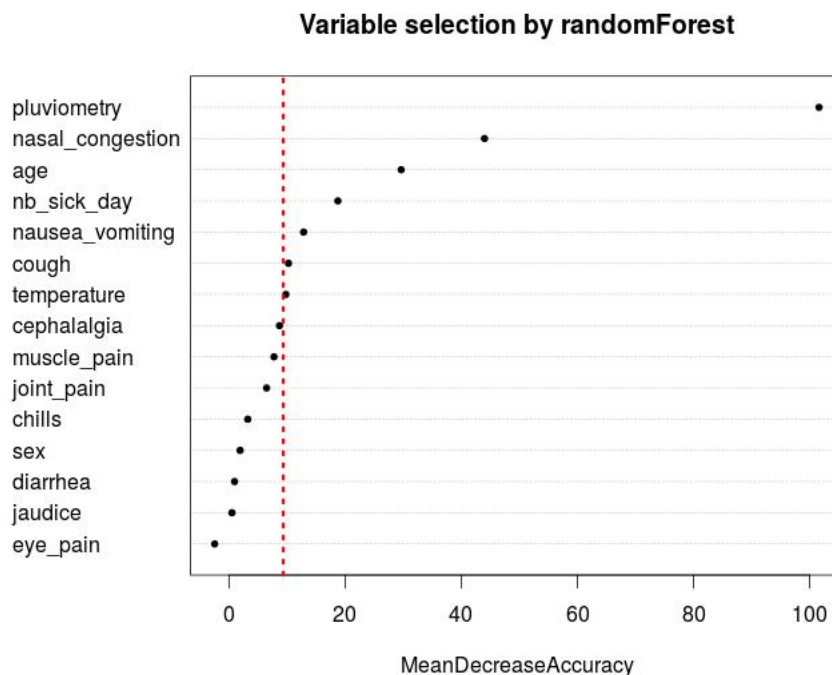


FIGURE A.3 – A variable importance plot for the *IgM/IgG* data set : mean decrease of accuracy (MDA) of the covariates, by increasing order. The variables whose *MDA* is to the right of the dotted line are selected by the VSURF procedure.

the VSURF selection procedure based on the 1 000 sub-samples. The covariate *rainfall* (95.2%) is almost always selected. Next, the more often selected variables are *cough* (29.1%), *age* (28.3%), *joint pain* (19.8%), *nausea or vomiting* (16.4%), *number of sick days* (16.1%), *temperature* (16.1%) and *nasal congestion* (11%), in decreasing order. The other covariates are selected in less than 10% of the samples.

We set different random seeds and we found that, for our purpose of selecting significant covariates, aggregation of 1 000 RF classifiers learned from 1 000 randomly balanced sub-samples yielded stable selected variable sets.

A.3.3 Influence of selected covariates on disease status

In the previous sections, we carried out a comparison between RF and multinomial logit covariate selections on the *IgM/IgG* data set and the conclusion is that the results are in agreement. The RF variable importance results on the *IgM* sub-samples produced a robust ranking of the covariates. The same group of seven important variables was selected by RF algorithm (see Figures A.4 and A.5); an eighth supplementary variable, *joint pain*, was added in the stepwise selection of the *IgM/IgG* data set. In conclusion, we decided to fit the same multinomial model with eight covariates to the data sets of our analysis and to further quantify the effects of the covariates in this model.

Within the multinomial logit model, we can quantify the effect of a variable in terms of an odds ratio or its logarithm. The odds that $Y = k$ occurs for an individual with covariates $X = x$ is the ratio of $P(Y = k|X = x)$ divided by $P(Y = 0|X = x)$, $k = 1, 2, 3$.

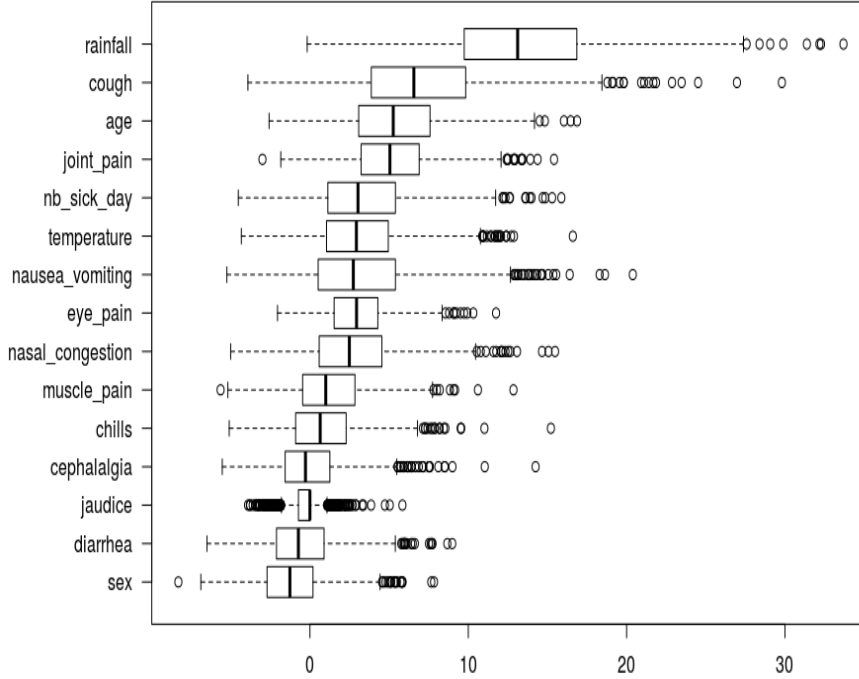


FIGURE A.4 – A variable importance plot for the *IgM* data set. Each boxplot summarizes the distribution of the variable importance among 1000 *IgM* sub-samples.

Then, the log odds of category k is given by Equation (A.1) :

$$\log \text{odds}(Y = k | X = x) = \langle x, \beta_k \rangle.$$

Thus the multinomial logit model is a linear regression model in the log odds. The parameter component β_{kj} can be interpreted as the change in the log odds per unit change in the continuous covariate X_j , if all other covariates are held constant. The odds ratio (OR) of category k for a d units increase of X_j , all other covariates remaining constant, is defined as

$$OR_k(d) = \frac{P(Y = k | X_j + d) / P(Y = 0 | X_j + d)}{P(Y = k | X_j) / P(Y = 0 | X_j)} = \exp(\beta_{kj}d).$$

Once β is estimated, one can estimate any odds or odds ratios. An OR equal to one means that a change in covariate X_j has no effect on the odds of category k ; if $OR_k(d) > 1$ ($OR_k(d) < 1$), the effect of an increase of X_j is to increase (decrease) the odds of category k . An odds ratio is a popular description of an effect in a probability model since it can be constant. On the contrary, the risk ratio $P(Y = k | X_j + d) / P(Y = k | X_j)$, which could be more interpretable in terms of predicted probabilities instead of odds, depends on the values of all other covariates. ORs are similar to risk ratios if the risk is small, otherwise ORs overestimate risk ratios.

For each covariate, we computed the odds ratios OR_k , $k = 1, 2, 3$ and their confidence intervals for each disease. Table A.3 for the *IgM/IgG* data set and Figure A.7 for the

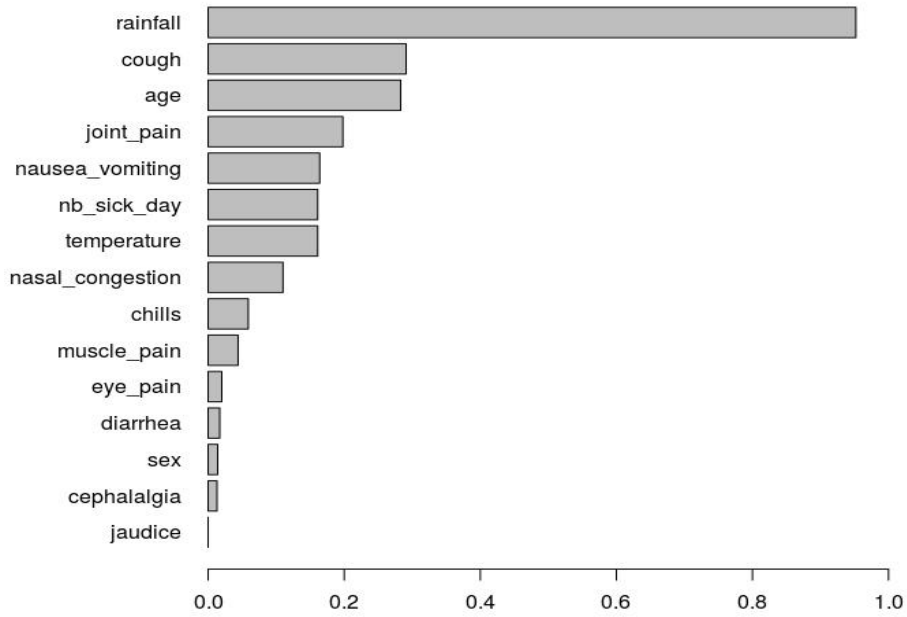


FIGURE A.5 – Ranking by VSURF : for each variable, the length of the bar corresponds to the empirical probability to be selected by VSURF among 1000 *IgM* sub-samples

IgM data set display the OR by which the odds increases for a certain change in a covariate, holding all other covariates constant. The ORs associated with binary variables (*Nausea/vomiting*, *Cough*, *Nasal congestion* and *Joint pain*) were computed by comparing the two modalities : 0 for absence and 1 for presence of the symptom. We computed the ORs resulting from increasing *Temperature* from 38 to 40 degrees Celsius ($d = 2$) and from increasing *Number of sick days* from 2 to 6 days ($d = 4$). The outer quartiles of *Age* are 8 and 28 years ($d = 20$), so we computed the half-sample OR for age. Similarly, we computed the half-sample OR for a *rainfall* of 14 mm compared to a *rainfall* of 370 mm ($d = 356$).

The ORs defined previously are relative to the reference category $Y = 0$. We also computed the ORs between two diseases $Y = k$ and $Y = l$ in order to differentiate the effect of each covariable between the three clinical groups, arbovirus vs malaria, coinfection vs arbovirus and coinfection vs malaria :

$$OR_{k|l}(d) = \frac{P(Y = k|X_j + d)/P(Y = l|X_j + d)}{P(Y = k|X_j)/P(Y = l|X_j)} = \exp((\beta_{kj} - \beta_{lj})d).$$

These results are displayed in Figure A.6 (*IgM/IgG* data set) and Figure A.8 (*IgM* data set). The confidence intervals are derived from the fitted multinomial logit model and their accuracy is based on the parametric assumption that the true data generating distribution does fall in the model.

Diseases Variables	Arbovirus	Coinfection	Malaria
<i>Age</i>	1.71 [1.42; 2.07]	1.12 [0.92; 1.36]	0.61 [0.50; 0.73]
<i>Temperature</i>	1.02 [0.69; 1.49]	2.16 [1.52; 3.07]	2.47 [1.82; 3.35]
<i>Number-of-sick-days</i>	2.54 [1.91; 3.37]	1.43 [1.04; 1.96]	1.04 [0.77; 1.39]
<i>Rainfall</i>	2.19 [1.53; 3.14]	17.0 [12.0; 24.0]	9.81 [7.18; 13.4]
<i>Nausea /vomiting</i>	0.83 [0.60; 1.13]	2.07 [1.55; 2.78]	2.15 [1.67; 2.77]
<i>Cough</i>	0.79 [0.58; 1.10]	0.46 [0.33; 0.63]	0.57 [0.44; 0.74]
<i>Nasal congestion</i>	0.52 [0.35; 0.75]	0.13 [0.09; 0.20]	0.10 [0.07; 0.13]
<i>Joint pain</i>	1.52 [0.99; 2.32]	1.90 [1.26; 2.83]	1.74 [1.21; 2.50]

TABLE A.3 – *IgM/IgG* data : odds ratios with respect to the reference modality and 95% confidence intervals.

Results for the *IgM/IgG* data

From Table A.3, we can say that the effect of increasing temperature from 38 to 40 is to double the odds of coinfection or to increase the odds of malaria by a factor of 2.5. The odds of arboviral monoinfection is multiplied by 1.71 for an adult compared to a child, whereas the odds of malaria decreases by a factor of 0.61. An increase of the number of sick days from 2 to 6 increases the odds of arboviral monoinfection by a factor of 2.54. The presence of nausea or vomiting symptoms increases the odds of malaria or the odds of coinfection by a factor of 2.07 and 2.15 respectively.

To summarize these results, we can say that a high temperature and presence of nausea or vomiting symptoms are risk factors for malaria and coinfection; a number of sick days greater than 2 and age above eight-years old are risk factors for arbovirus and coinfection.

Figure A.6(a) displays the odds ratios between malaria monoinfection and arboviral monoinfection. We can say that *Nasal congestion*, *Number of sick days* and *Age* are correlated to arbovirus; *Temperature*, *Rainfall* and *Nausea or vomiting* are correlated to malaria monoinfections. The variables *joint pain* and *cough* are not significant in distinguishing malaria and arboviral monoinfections. Figure A.6(b) suggests that vomiting symptoms and a high fever are indicative of coinfection among patients exhibiting arboviral monoinfection. But these covariates are not significant to differentiate coinfection from malaria monoinfection (Figure A.6(c)). Figure A.6(c) suggests that *Age*, *Number of sick days* and *Nasal congestion* are significantly correlated with coinfecting patients compared to patients with single malaria disease.

Results for the *IgM* data

Figure A.7 and Figure A.8 display the sampling distribution of ORs based on the fitting of the 1000 sub-samples of the *IgM* data set.

From Figure A.7 and Figure A.8, we can say that temperature, rainfall and vomiting symptoms are significantly correlated with malaria monoinfections whereas joint pain, age and number of sick days are correlated with arboviral monoinfections. The odds of coinfection increases with high fever and high rainfall values, and the presence of vomiting

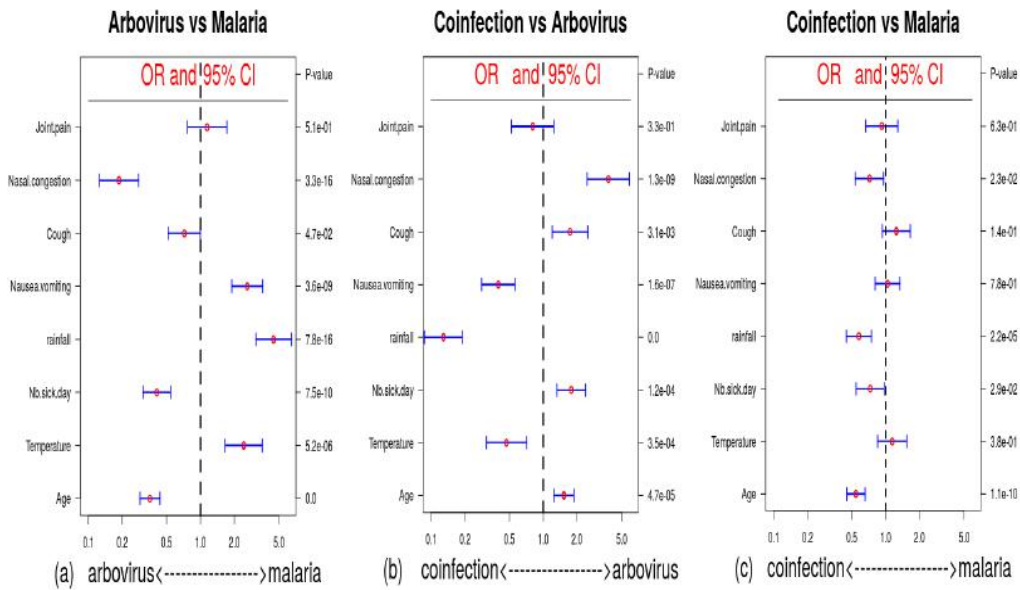


FIGURE A.6 – *IgM/IgG* data : odds ratios between two diseases and 95% confidence intervals ; (a) Arbovirus vs Malaria (b) Coinfection vs Arbovirus (c) Coinfection vs Malaria.

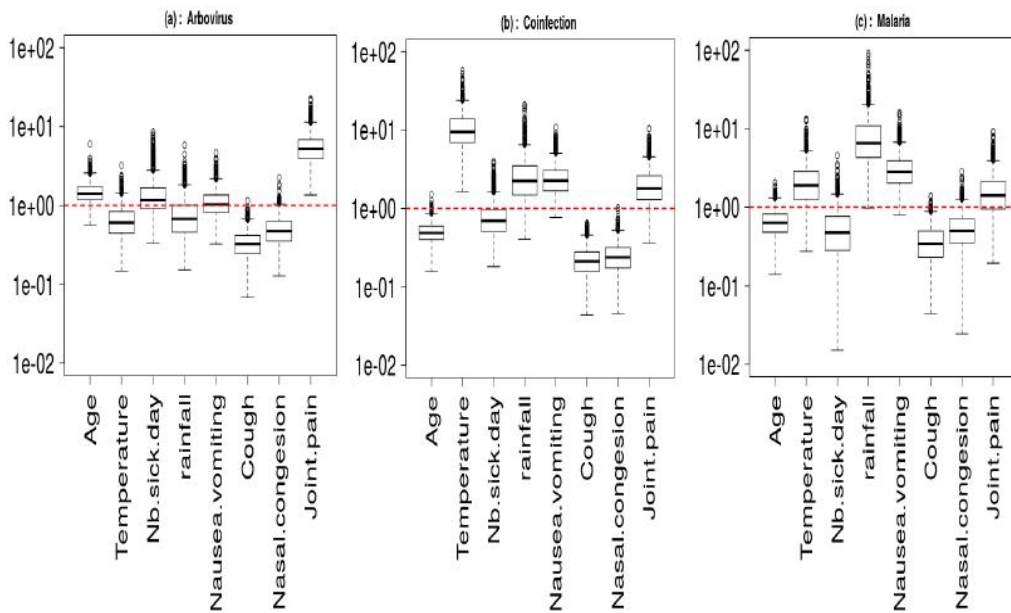


FIGURE A.7 – *IgM* data : boxplots of 1000 odds ratios with respect to the reference category ; (a) Arbovirus (b) Coinfection (c) Malaria.

and joint pain symptoms.

Conclusion

The results based on both data sets show that a high temperature and the presence of nausea or vomiting symptoms are mostly indicative to malaria parasite infections whereas

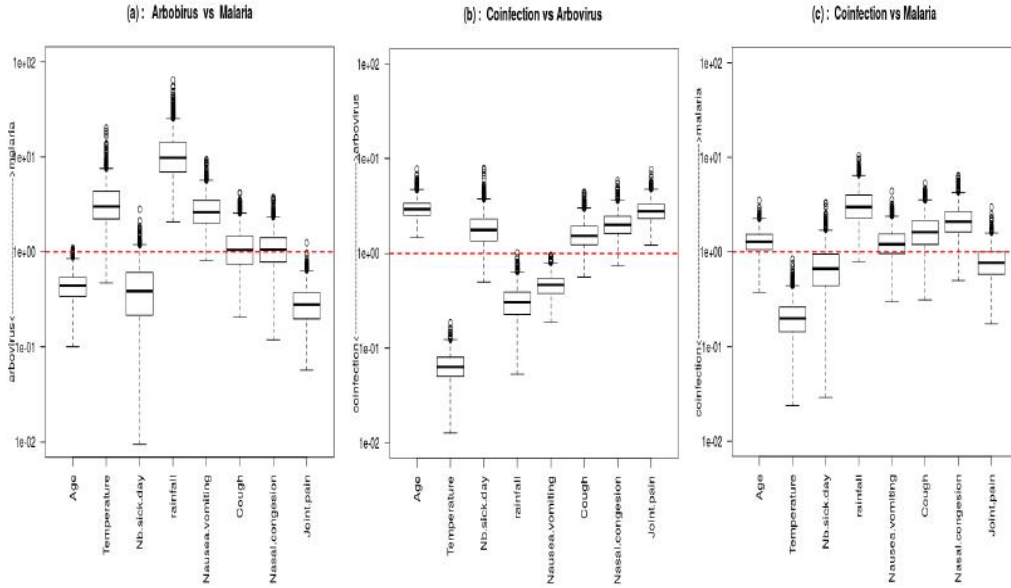


FIGURE A.8 – *IgM* data : boxplots of 1000 odds ratios between two categories; (a) Arbovirus *vs* Malaria (b) Coinfection *vs* Arbovirus (c) Coinfection *vs* Malaria.

an increase of the number of sick days and the age are indicative to arboviral infections. The effects of the nasal congestion and joint pain symptoms on the disease status are not clear enough to be interpreted. The main question of the study was to identify risk factors that can help doctors to diagnose a concurrent malaria and arbovirus infection. From these results, *Temperature* is the only risk factor that differentiates between coinfection and single infections.

A.4 Predictive analysis

In this section we aim to propose a methodology that can help make timely decisions for targeted treatment in pathogens coinfection cases. We show that we can derive a predictive analysis to discriminate arbovirus positive and arbovirus negative cases among coinfecting patients.

A.4.1 Testing independence between arbovirus and malaria

In the multinomial model given by Equation (A.1) in Section A.3.1, we can test the independence between arboviral and malaria infections.

The joint statistical distribution of arboviral infection (A^+) and malaria infection (M^+) is given in Table A.4. As in Table A.1 and Table A.2, A^+ corresponds to an individual belonging to categories 1 or 3 of the response Y , and M^+ corresponds to an individual belonging to categories 2 or 3 of the response Y .

Independence between arboviral and malaria infections means that for all (l_1, l_2) in $\{0, 1\}$,

$$P(M^+ = l_1, A^+ = l_2) = P(M^+ = l_1) \times P(A^+ = l_2).$$

	$A = 0$	$A = 1$	law of M^+
$M = 0$	π_0	π_1	$P(M^+ = 0) = \pi_0 + \pi_1$
$M = 1$	π_2	π_3	$P(M^+ = 1) = \pi_2 + \pi_3$
law of A^+	$P(A = 0) = \pi_0 + \pi_2$	$P(A = 1) = \pi_1 + \pi_3$	1

TABLE A.4 – Joint distribution of arboviral infection and malaria infection

The independence hypothesis can be written in terms of parameters as :

$$H_0 : \quad “\beta_3 = \beta_1 + \beta_2”.$$

The Wald statistic to test H_0 against its two-sided alternative is computed as

$$W = h(\hat{\beta})^T \Sigma^{-1} h(\hat{\beta}),$$

with $h(\hat{\beta}) = \hat{\beta}_3 - \hat{\beta}_1 - \hat{\beta}_2$ and $\Sigma = DV D^T$ where $D = (-Id_{p+1}, -Id_{p+1}, Id_{p+1})$; Id_p is the $p \times p$ identity matrix and V is an estimator of the variance of $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3)^T$. Under H_0 , W is asymptotically distributed as a chi-square variable with $(p + 1)$ degrees of freedom. Under H_1 , W converges to infinity as the sample size goes to infinity.

We fitted model (A.1) including the eight covariates selected in Section A.3.3 and we computed the independence test. Based on *IgM/IgG* data, the independence hypothesis was rejected with a p-value equal to $1.46.10^{-6}$. We studied the robustness of the test decision with respect to the variable selection. Whatever the selected number of variables, we obtained p-values with order less than or equal to 10^{-3} . Thus, we can consider that arbovirus and malaria are correlated.

Applying the test on the *IgM* data, we computed the 1000 p-values corresponding to the 1000 sub-samples and obtained that 42.5% of them were less than 0.05. So we can not reject the independence hypothesis in a majority of sub-samples. It can be explained by the fact that the size of the sub-samples is small (139) and the asymptotic approximation of the law of the test statistic is not accurate. Moreover, the *IgM* data set may not contain enough information to explain coinfection. Which means that the independence test lacks of power.

In the following, we will only consider the *IgM/IgG* data set to propose a predictive analysis.

A.4.2 Diagnosis of arboviral disease

In this section, we present a methodology to help doctors to diagnose the arboviral infected patients whose symptoms are masked by malaria symptoms. We propose to base the diagnosis on the conditional probability $P(C|M)$ to be coinfecting given that malaria infection is observed. This probability is the quantity of interest because arboviral infections are considered by healthcare workers only if malaria tests are negative. In absence of rapid arbovirus detection tests, the aim is to provide a decision support tool to determine if an arbovirus could be responsible for the clinical symptoms of the patient coinfection.

Based on the previous results of the *IgM/IgG* data set, the independence test of Section A.4.1 displays an association between malaria and arboviral infections. Then the probability $P(C|M)$ can be computed in function of the π_k probabilities estimated from the multinomial logit model. For an individual with covariate x ,

$$P(C|M) = \frac{\widehat{\pi}_3(x)}{\widehat{\pi}_3(x) + \widehat{\pi}_2(x)} = \frac{e^{\langle x, \widehat{\beta}_3 \rangle}}{e^{\langle x, \widehat{\beta}_3 \rangle} + e^{\langle x, \widehat{\beta}_2 \rangle}}.$$

This probability can be used to differentiate whether the illness to be treated should be arbovirus or malaria. We propose a binary classification rule and we predict an arbovirus illness if the estimated coinfection probability is greater than a threshold value γ :

$$\begin{cases} \text{If } P(C|M) \geq \gamma : & \text{arbovirus positive case,} \\ \text{If } P(C|M) < \gamma : & \text{arbovirus negative case.} \end{cases}$$

The evaluation of the classification is based on the confusion matrix and the overall classification accuracy. The confusion matrix is used to compute true arbovirus positives (TP), false arbovirus positives (FP), true negatives (TN) and false negatives (FN). A global performance measure is the miss-classification rate (MCR) defined as :

$$\text{MCR} = \frac{FP + FN}{N},$$

with $N = TP + FP + TN + FN$.

The analysis performed in this section is based on 1148 instances of the *IgM/IgG* data set corresponding to the patients infected with malaria parasites. The multinomial logit model was trained on 66.7% of the data, namely 1317 instances and tested on the remaining 377 individuals positive to malaria. To choose the classification threshold value γ , standard practice is to minimize the miss-classification rate. We computed the five-fold cross-validation estimator of the MCR. We can see on Figure A.9 that the optimal threshold is around $\gamma = 0.5$. Five-fold cross-validation was run different times, each with a different split of the data and the optimal value of γ was found to be quite stable. Then, a classification with $\gamma = 0.5$ was used to predict the type of illness that has affected the patient based on his clinical symptoms. Predicted and actual arbovirus cases were compared using the test set, as presented in Table A.5. The rows of the matrix are actual classes and the columns are the predicted classes. We observe that the corresponding MCR is 38%, and the number of FN is quite high. In applications such as disease diagnosis, it is desirable to have a classifier that reduces the number of FN, since a false negative could be more dangerous to the care of a patient, who then may not be treated, whereas with a false positive, the patient would most likely undergo more testing before treatment. Different strategies can be adopted. One possibility is to reduce the number of FN by minimizing a weighted version of the MCR,

$$\text{WMCR} = \frac{FP + cFN}{N}, \quad c > 1.$$

A weight coefficient c higher than one increases the cost of classification mistakes on the FN. We tried empirical values of $c = 2, 3, 4$ and found that they resulted in a decrease

<i>True</i>	<i>Predicted</i>	
	0	1
0	211	29
1	114	23

TABLE A.5 – Confusion table with $\gamma = 0.5$.

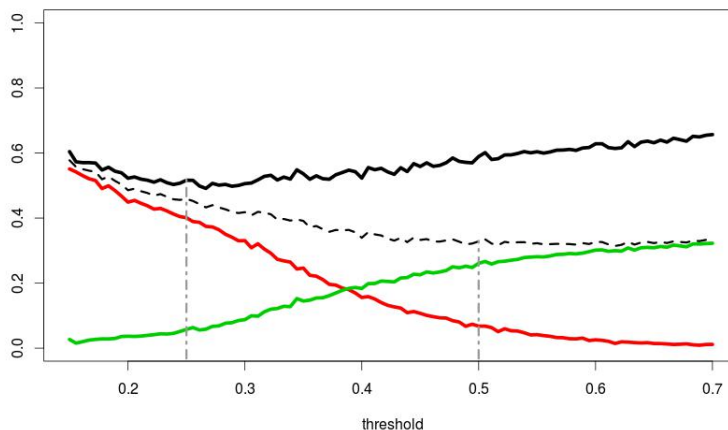


FIGURE A.9 – *IgM/IgG* data : estimated cross-validation miss-classification rate. The WMCR is shown in black as a full line. The MCR is shown as a black dotted line. Increasing γ increases the number of FN (green full line) and decreases the FP (red full line).

<i>True</i>	<i>Predicted</i>	
	0	1
0	88	152
1	24	113

TABLE A.6 – Confusion table with $\gamma = 0.25$.

of the FN rate at the cost of an increase of the *WMCR*. With a choice of $c = 2$, the threshold value that minimizes the *WMCR* is 0.25. With this γ choice, we observe on Table A.6 that the number of FN is reduced but the MCR remains too high (46.7%).

In a next step, we proposed to select, among the positive predicted patients, those individuals with age greater than 10 and number of sick days greater than 3. Indeed we concluded in Section A.3.3 that these two variables are mostly indicative of arboviral disease. The threshold values were again chosen to minimize the *WMCR* using cross-validation. Table A.7 gives the corresponding results : the MCR is decreased to 36% while the number of FN is smaller than the number of FN of Table A.5 and the number of TP is doubled.

The objective of these predictions was to assign patients to either a “Malaria” group or a “Arbovirus” group and to handle mystifying cases due to the similarity of the initial

<i>True</i>	<i>Predicted</i>	0	1
	0		190
1		85	52

TABLE A.7 – Confusion table with $\gamma = 0.25$, $Age = 10$ and Number of sick days = 3.

symptoms in both diseases. The classification procedure is based on the computation of the conditional probability $P(C|M)$. The threshold parameter γ is calibrated to minimize the weighted miss-classification rate. To improve the accuracy of the classification, we propose to take advantage of the covariates that were selected in Section A.3.3 as arbovirus specific covariates. Based on these specific covariates, we filtered the positive predicted patients and obtained better results.

The performance of a classification procedure is greatly affected by the quality of data source. We based our analysis on the *IgM/IgG* data collected from patients who were tested positive to *IgM* or *IgG*. As *IgG* antibodies appear later in time in blood than *IgM* antibodies, they may not reflect a current infection, thus minimizing the possibility of finding a true correlation with the current recorded symptoms. This drawback may reduce the prediction capacity of the classification procedure. Also, false positives and false negatives from biological tests could impact these results. However, the diagnostic tests used in the study displayed high sensitivity and specificity parameters; then their impact on the results should be negligible.

A.5 Discussion

Misdiagnosis of arbovirus coinfections as malaria infections may increase the spread of arbovirus diseases in areas where fast diagnostic assays are not available. This study proposes an appropriate statistical methodology that can assist doctors in the elaboration of the differential diagnosis of febrile cases for arboviruses.

Our analysis is based on a real-life medical data set. In the original *IgM* data set, arbovirus positive individuals are identified as individuals likely to be in the early stages of arbovirus illness. It is the relevant data set for the classification problem. However, the positive cases constitute only a very small minority class of the data (39 positive cases over 12288 individuals). Several sampling strategies have been developed to learn from imbalanced data sets [67] and to correct classification of the rare class. [61] proposed a categorisation of these approaches into two main categories : data pre-processing and modifications on the learning algorithms. Algorithm level strategies including random forest solutions have been discussed to deal with the imbalanced data classification problem ([69], [68]). They require a good insight into the modified learning algorithm and a precise identification of reasons for its failure in mining imbalanced data sets. As we were interested in the statistical methodology that could be applied to a more relevant data set, we took solutions that pre-process the given imbalanced data set. Since the data distribution is changed to make standard methods focus on the cases that are more relevant for our problem, the results should be interpreted carefully.

To analyze coinfection data we propose a methodology with three steps : 1. a variable

selection with random forests ; 2. an analysis of the influent factors through multinomial model fitting and odd ratios computation ; 3. a predictive analysis based on coinfection probabilities. From our experiments, we can say that the random forests algorithm is a robust method to select the important variables for the different diseases. The analysis of the odd ratios allows to identify the risk factors that characterize each disease. We observed that higher values of number of sick days and of age are mostly indicative of arboviral disease while higher values of temperature and presence of nausea or vomiting symptoms during the rainy season are mostly indicative of malaria disease. The results also pointed out that a high-grade fever could be considered as a differential diagnostic for malaria and arbovirus coinfection, which is in agreement with the study of [13]. The classification rule based on coinfection probability, age and number of sick days identifies coinfecting patients to be treated for arbovirus with global accuracy of 65%. The results could be improved on a more suitable data set. A future study will apply this methodology to coinfection data between malaria and other pathogens more easily detectable in the early stages of infection than arboviruses.

Acknowledgements : The authors thank two referees for detailed and helpful comments that improved the manuscript.

Annexe B

Mixture of generalized linear models : identifiability and applications

Abstract

We consider finite mixtures of generalized linear models with binary output. We prove that cross moments till order 3 are sufficient to identify all parameters of the model. We propose a least squares estimation method and we prove the consistency and the Gaussian asymptotic behavior of the estimator. An R-package is developed to apply our method, we give numerical experiments to compare with likelihood methods. We then provide new identifiability results for several finite mixtures of generalized linear models with binary output and unknown link function including both continuous and categorical covariates, and possibly longitudinal data.

B.1 Introduction

Logistic models, or more generally multinomial regression models that fit covariates to discrete responses through a link function, are very popular for use in various application fields. When the data under study come from several groups that have different characteristics, using mixture models is also a very popular way to handle heterogeneity. Thus, many algorithms were developed to deal with various mixture models, see for instance the book [100]. Most of them use likelihood methods or Bayesian methods that are likelihood dependent. Indeed, the now well known expectation-maximization (EM) methodology or its randomized versions makes it often easy to build algorithms. However one problem of such methods is that they can converge to local spurious maxima so that it is necessary to explore many enough initial points. Recently, spectral methods were developed to bypass EM algorithms and they were proved able to recover the directions of the regression parameter in models with known link function and random covariates, see [37].

One aim of this paper is to extend such moment methods using least squares to get estimators of the whole parameters, and to provide theoretical guarantees of this estimation

method. The setting is that of regression models with binary outputs, random covariates and known link function, detailed in Section 2. We first prove that cross moments up to order 3 between the output and the regression variables are enough to recover all the parameters of the model, see Theorem B.3.1.1 for the probit link function and Theorem B.3.1.2 for general link functions. We then obtain consistency and asymptotic normality of our least squares estimators as usual, see Theorem B.3.2.1. The algorithm is described at the end of Section 3, and to apply it, we developed the R-package `morpheus` available on the CRAN ([101]). We then compare experimentally our method to the maximum likelihood estimator computed using the R-package `flexmix` ([102]). We show that our estimator may be better for the probit link function with finite samples when the dimension increases, though keeping very small computation times when that of `flexmix` increases with dimension. The experiments are presented in Section 4.

Another aim of this paper is to investigate identifiability in various mixture of non linear regression models with binary outputs. Indeed, identifiability results for such models are still few and not enough to give theoretical guarantees of available algorithms. Let us review what is known up to our knowledge. In [103], the identifiability is proved for finite mixtures of logistic regression models where only the intercept varies with the population [104]. In [105], finite mixtures of multinomial logit models with varying and fixed effects are investigated, the proofs of identifiability results use the explicit form of the logit function. In [104], further non parametric identifiability of the link function is proved, but only for models where the base exponential models are identifiable for mixtures, which does not apply to binary data (Bernoulli models).

We provide in Section 5 several identifiability results, that for example are useful to get theoretical guarantees in applications such as the one in [106]. We prove that with known smooth enough link function, the directions of the covariates may be recovered under the only assumption that they are distinct, see Theorem B.5.1.1. Then, under the strengthened assumption that they are linearly independent, we prove that the link function may be non parametrically recovered, see Theorem B.5.1.2. We then study the simultaneous use of continuous and categorical covariates and further give assumptions under which parameters and link function may be recovered, see Theorem B.5.2.1. We finally prove that, with longitudinal data having at least 3 repetitions for each individual, the whole model is identifiable under the weakest assumption that the regression directions are distinct, see Theorem B.5.3.1.

B.2 Model and notations

Let us denote $[n]$ the set $\{1, 2, \dots, n\}$ and $e_i \in \mathbb{R}^d$, the i -th canonical basis vector of \mathbb{R}^d . Denote also $I_d \in \mathbb{R}^{d \times d}$ the identity matrix in \mathbb{R}^d . The tensor product of p euclidean spaces \mathbb{R}^{d_i} , $i \in [p]$ is noted $\bigotimes_{i=1}^p \mathbb{R}^{d_i}$. T is called a real p -th order tensor if $T \in \bigotimes_{i=1}^p \mathbb{R}^{d_i}$. For $p = 1$, T is a vector in \mathbb{R}^d and for $p = 2$, T is a $d \times d$ real matrix. The (i_1, i_2, \dots, i_p) -th coordinate of T with respect the canonical basis is denoted $T[i_1, i_2, \dots, i_p]$, $i_1, i_2, \dots, i_p \in [d]$.

Let $X \in \mathbb{R}^d$ be the vector of covariates and $Y \in \{0, 1\}$ be the binary output.

A binary regression model assumes that for some link function g , the probability that $Y = 1$ conditionally to $X = x$ is given by $g(\langle \beta, x \rangle + b)$, where $\beta \in \mathbb{R}^d$ is the vector of regression coefficients and $b \in \mathbb{R}$ is the intercept. Popular examples of link functions are the logit link function where for any real z , $g(z) = e^z / (1 + e^z)$ and the probit link function where $g(z) = \Phi(z)$, with Φ the cumulative distribution function of the standard normal $\mathcal{N}(0, 1)$.

If now we want to modelise heterogeneous populations, let K be the number of populations and $\omega = (\omega_1, \dots, \omega_K)$ their weights such that $\omega_j \geq 0$, $j = 1, \dots, K$ and $\sum_{j=1}^K \omega_j = 1$. Define, for $j = 1, \dots, K$, the regression coefficients in the j -th population by $\beta_j \in \mathbb{R}^d$ and the intercept in the j -th population by $b_j \in \mathbb{R}$. Let $\omega = (\omega_1, \dots, \omega_K)$, $b = (b_1, \dots, b_K)$, $\beta = [\beta_1 | \dots | \beta_K]$ the $d \times K$ matrix of regression coefficients and denote $\theta = (\omega, \beta, b)$. The model of population mixture of binary regressions is given by :

$$\mathbb{P}_\theta(Y = 1 | X = x) = \sum_{k=1}^K \omega_k g(\langle \beta_k, x \rangle + b_k). \quad (\text{B.1})$$

We assume that the random variable X has a Gaussian distribution. We now focus on the situation where $X \sim \mathcal{N}(0, I_d)$, I_d being the identity $d \times d$ matrix. All results may be easily extended to the situation where $X \sim \mathcal{N}(m, \Sigma)$, $m \in \mathbb{R}^d$, Σ a positive and symmetric $d \times d$ matrix.

Define the cross moments between the response Y and the covariable X , up to order 3 :

- $M_1(\theta) := \mathbb{E}_\theta[Y.X]$, first-order moment,
- $M_2(\theta) := \mathbb{E}_\theta \left[Y.(X \otimes X - \sum_{j \in [d]} Y.e_j \otimes e_j) \right]$, second-order moment and
- $M_3(\theta) := \mathbb{E}_\theta \left[Y(X \otimes X \otimes X - \sum_{j \in [d]} [X \otimes e_j \otimes e_j + e_j \otimes X \otimes e_j + e_j \otimes e_j \otimes X]) \right]$ third-order moment.

Let, for $k = 1, \dots, K$, $\lambda_k = \|\beta_k\|$ and $\mu_k = \beta_k / \|\beta_k\|$. Using Stein's identity, Anandkumar et al. ([37]) prove the following lemma :

Lemme B.2.0.1 ([37]) *Under enough smoothness and integrability of the link function (which hold for the logit and probit link functions, or under our assumption (H3) below) the moments can be rewritten :*

$$\begin{aligned} M_1(\theta) &= \sum_{k=1}^K \omega_k \lambda_k \mathbb{E}[g'(\lambda_k \langle X, \mu_k \rangle + b_k)] \mu_k, \\ M_2(\theta) &= \sum_{k=1}^K \omega_k \lambda_k^2 \mathbb{E}[g''(\lambda_k \langle X, \mu_k \rangle + b_k)] \mu_k \otimes \mu_k, \\ M_3(\theta) &= \sum_{k=1}^K \omega_k \lambda_k^3 \mathbb{E}[g^{(3)}(\lambda_k \langle X, \mu_k \rangle + b_k)] \mu_k \otimes \mu_k \otimes \mu_k. \end{aligned}$$

It is proved in [37] that the knowledge of $M_3(\theta)$ leads to the knowledge of μ_1, \dots, μ_K up to their sign as soon as they are linearly independent. In the next section, we prove that the knowledge of all cross moments till order 3 allows to recover all parameters for the

probit link function under the same assumption on the regression coefficients. We also prove that for a general link function satisfying some weak assumption, the knowledge of all cross moments till order 3 allows to recover all parameters provided they are not too far from 0.

B.3 Moment identifiability and estimation

To prove our moment identifiability result, we shall use the following assumptions :

- (H1) The vectors β_1, \dots, β_K are linearly independent and the weights are positive : $\omega_k > 0, k = 1, \dots, K$.
- (H2) The link function g is strictly increasing from 0 at $-\infty$ to 1 at $+\infty$, it has continuous derivatives till order 4, decreasing first derivative on $[0, +\infty[$, and it satisfies

$$\forall z \in \mathbb{R}, g(z) + g(-z) = 1.$$

- (H3) There exists a neighborhood \mathcal{O} of $(0, 0)$ in $\mathbb{R}_+^* \times \mathbb{R}$ and any functions $L_s, s = 1, 2, 3$, such that $\forall z \in \mathbb{R}, \forall (\lambda, b) \in \mathcal{O}$, we have

$$(|z| + 1) \left| \frac{\partial g^{(s+1)}}{\partial \lambda}(\lambda z + b) \right| \leq L_s(z)$$

and further for $s = 1, 2, 3$,

$$\int_{\mathbb{R}} L_s(z) e^{-z^2/2} dz < +\infty.$$

Notice that (H1) implies that $d \geq K$, and that (H2) and (H3) hold in particular for the logistic link function and the probit link function.

From (H2), one gets that

- (P1) The function g' is positive and satisfies $g'(x) = g'(-x)$ for all $x \in \mathbb{R}$,
- (P2) The function g'' satisfies $g''(x) = -g''(-x)$ for all $x \in \mathbb{R}$ and $g''(x) < 0$ for $x > 0$.

To prove the limiting Gaussian distribution of our moment estimator, we shall need more assumptions. For $j = 1, \dots, 5$, let G_j be the $K \times K$ diagonal matrix having the $\mathbb{E}[g^{(j)}(\langle \beta_k, X \rangle + b_k)]$'s on the diagonal.

- (H4) All diagonal coefficients of G_3 are non zero.
- (H5) All diagonal coefficients of $G_1 G_3 - G_2^2$ are non zero.

B.3.1 Identifiability results

In the whole section we assume that (H1) holds. Under (H1), we see by Lemma B.2.0.1 that K is the rank of $M_2(\theta)$.

It is proved in [37] we can recover the μ_k 's up to sign from the knowledge of $M_2(\theta)$ and $M_3(\theta)$, but since under (H1) $M_1(\theta)$ is a linear combination of the μ_k 's with positive coefficients, the knowledge of $M_1(\theta)$ allows to recover the signs. It is then seen that using $M_1(\theta)$, $M_2(\theta)$ and $M_3(\theta)$, one may recover the 3-uples

$$(\omega_k E[g'(\langle \beta_k, X \rangle + b_k)] \lambda_k; \omega_k E[g''(\langle \beta_k, X \rangle + b_k)] \lambda_k^2; \omega_k E[g^{(3)}(\langle \beta_k, X \rangle + b_k)] \lambda_k^3),$$

$k = 1, \dots, K$. Thus, one gets identifiability as soon as the function from $]0, +\infty[\times]0, +\infty[\times \mathbb{R}$ to its image that associates (ω, λ, b) to

$$\left(\omega \lambda \int g'(\lambda z + b) e^{-z^2/2} dz; \omega \lambda^2 \int g''(\lambda z + b) e^{-z^2/2} dz; \omega \lambda^3 \int g^{(3)}(\lambda z + b) e^{-z^2/2} dz \right)$$

is one-to-one. Using integration by parts this is equivalent to the fact that the function from $]0, +\infty[\times]0, +\infty[\times \mathbb{R}$ to its image that associates (ω, λ, b) to

$$\lambda \left(\omega \int g'(\lambda z + b) e^{-z^2/2} dz; \omega \int z g'(\lambda z + b) e^{-z^2/2} dz; \omega \int z^2 g'(\lambda z + b) e^{-z^2/2} dz \right)$$

is one-to-one. This is again equivalent to the fact that the function from $]0, +\infty[\times \mathbb{R}$ to its image that associates (λ, b) to

$$\left(\frac{\int z g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz}; \frac{\int z^2 g'(\lambda z + b) e^{-z^2/2} dz}{\int g'(\lambda z + b) e^{-z^2/2} dz} \right)$$

is one-to-one. For any $(b, \lambda) \in \mathbb{R} \times]0, +\infty[$, define

$$dQ_{(b,\lambda)}(z) = \frac{g'(\lambda z + b) e^{-z^2/2}}{\int g'(\lambda z + b) e^{-z^2/2} dz} dz. \quad (\text{B.2})$$

Then it is equivalent to prove that the knowledge of

$$(E_{(b,\lambda)}(Z); E_{(b,\lambda)}(Z^2)) := \left(\int z dQ_{(b,\lambda)}(z); \int z^2 dQ_{(b,\lambda)}(z) \right) \quad (\text{B.3})$$

implies the knowledge of (b, λ) . When g is the probit link function, $Q_{(b,\lambda)}$ is a Gaussian distribution and computations detailed in Section B.6.1 leads to the following identifiability result.

Théorème B.3.1.1 (Probit identifiability) *If (H1) holds and if g is the probit link function, one may recover K and $\theta = (\omega, \beta, b)$ from the knowledge of $M_1(\theta)$, $M_2(\theta)$ and $M_3(\theta)$.*

In the general situation, identifiability holds at least in an open set. To prove it, one just has to prove that for some $B > 0$ and $L > 0$, if (H2) holds, then, the function that associates $(b, \lambda) \in]-B, B[\times]0, L[$ to $(E_{(b,\lambda)}(Z); E_{(b,\lambda)}(Z^2))$ is one-to-one on its image. This leads to the following identifiability result whose proof is postponed to Section B.6.2.

Théorème B.3.1.2 (General identifiability) *If (H1), (H2), (H3) hold and $g^{(3)}(0) \neq 0$, there exist $L > 0$ and $B > 0$ such that as soon as $\|\beta_k\| < L$ and $|b_k| < B$ for all $k = 1, \dots, K$, then one may recover K and $\theta = (\omega, \beta, b)$ from the knowledge of $M_1(\theta)$, $M_2(\theta)$ and $M_3(\theta)$.*

Since the proof uses Taylor expansions, it only proves the existence of *small enough* positive L and B such that the result holds. However, numerical study of the function $(b, \lambda) \mapsto (E_{(b,\lambda)}(Z); E_{(b,\lambda)}(Z^2))$ when the link function g is the logit function shows that identifiability seems to hold at least with $L = 8$ and $B = 8$.

B.3.2 The least squares moment estimator

In the previous section we showed that the parameters can be recovered by matching the cross-moments till order 3. Those moments are unknown, so that we estimate them empirically using :

$$\begin{aligned}\widehat{M}_1 &= \frac{1}{n} \sum_{i=1}^n Y_i X_i \\ \widehat{M}_2 &= \frac{1}{n} \sum_{i=1}^n \left[Y_i (X_i \otimes X_i - \sum_{j \in [d]} e_j \otimes e_j) \right] \\ \widehat{M}_3 &= \frac{1}{n} \sum_{i=1}^n \left[Y_i (X_i \otimes X_i \otimes X_i - \sum_{j \in [d]} [X_i \otimes e_j \otimes e_j + e_j \otimes X_i \otimes e_j + e_j \otimes e_j \otimes X_i]) \right].\end{aligned}$$

It is not possible to match the empirical moments exactly, so that we use a least-squares estimator. Define for all θ :

$$Q_n(\theta) = \sum_{j \in [d]} \left\{ \widehat{M}_1[j] - M_1(\theta)[j] \right\}^2 + \sum_{j, k \in [d]} \left\{ \widehat{M}_2[j, k] - M_2(\theta)[j, k] \right\}^2 + \sum_{j, k, l \in [d]} \left\{ \widehat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right\}^2$$

and the estimator

$$\widehat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmin}} Q_n(\theta). \quad (\text{B.4})$$

Théorème B.3.2.1 *Assume that (H1), (H2), (H3) hold, and that Θ is compact and included in the set of identifiable parameters. Then $\widehat{\theta}_n$ is consistent.*

If moreover (H4) and (H5) hold, then $\sqrt{n} (\widehat{\theta}_n - \theta^)$ converges in distribution under \mathbb{P}_{θ^*} to a centered Gaussian distribution.*

The proof of Theorem B.3.2.1 is detailed in Section B.6.3 and follows the usual analysis of the asymptotic behavior of Z -estimators, the more delicate part of the proof being to prove that the Hessian of $Q_n(\theta)$ has an invertible limiting value.

B.3.3 Algorithm

The estimator $\widehat{\theta}_n$ is computed by using the representation of the regression vectors β_k through their direction μ_k and norm λ_k , $k = 1, \dots, K$. In a first step, we compute a preliminary estimate of $[\mu_1, \dots, \mu_K]$ using a spectral method. In a second step, we search the minimizer of Q_n using usual optimization methods. The preliminary estimator for the directions is used as initial point for the directions in the optimization procedure.

The preliminary estimation of the directions is based on the spectral method. For any vector $z \in \mathbb{R}^p$, define $B(z)$ the $d \times d$ matrix such that

$$B(z)[i, j] := \sum_{s=1}^d M_3(\theta)[i, j, s] z_s,$$

Algorithm *M3LS* : Estimation of all parameters

input : X, Y, K, g

1 : Estimate the directions μ_1, \dots, μ_K using Algorithm *InitDir*

2 : Optimize $Q_n(\theta)$ using the estimators of 1. as initial directions

Output : The estimated parameter $\hat{\theta}$

so that, using Lemma B.2.0.1, we get

$$B(z) = \sum_{k=1}^K r^{(3)} \omega_k \lambda_k^3 \mathbb{E}[g^{(3)}(\lambda_k \langle X, \mu_k \rangle + b_k)] \langle \mu_k, z \rangle \mu_k^{\otimes 2}.$$

It is proved in [40] that it is possible to recover the directions by joint diagonalisation of $B(z_1), \dots, B(z_P)$ for distinct vectors z_1, \dots, z_P , $P \geq 2$. Joint diagonalisation of $B(z_1), \dots, B(z_P)$ means finding a matrix V such that the matrices $VB(z_p)V^T$ are the most diagonal possible. The normalized vectors μ_1, \dots, μ_K are obtained up to sign and label switching by taking the first K vectors of V^{-1} . Let us denote U the matrix of these K vectors. Let $O = U_*^{-1} M_1(\theta) \in \mathbb{R}^K$, with U_*^{-1} the general inverse of U . The real numbers $\omega_k \lambda_k \mathbb{E}[g'(\lambda_k \langle X, \mu_k \rangle + b_k)]$, $k = 1, \dots, K$, are given up to sign by the elements of O . Since they are positive, the sign of the μ_k 's are obtained by multiplying -1 all the vectors associated to the negative values of O .

In practice, the vectors μ_1, \dots, μ_K are estimated using the joint diagonalisation method applied to the matrices $\hat{B}(z_p)$, $p = 1, \dots, P$, computed using \hat{M}_3 .

Algorithm *InitDir* : Joint diagonalisation algorithm to estimate the directions

input : X, Y, K

1 : Estimate the cross moments \hat{M}_1, \hat{M}_2 and \hat{M}_3 as explained in section B.3.2

2 : Choose vectors $\{z_1, z_2, \dots, z_P\} \subseteq \mathbb{R}^d$ (for instance : the canonical basis e_1, e_2, \dots, e_P of \mathbb{R}^d)

3 : Compute $\hat{B}(z_p)$ for all $p \in \{1, 2, \dots, P\}$

4 : Joint diagonalisation : compute V such that $V\hat{B}(z_p)V^T$ are the most diagonal possible

5 : Compute $U = V^{-1}[1 : K]$ the K -first vectors of V^{-1} (by ordering the diagonal values in decreasing absolute value)

6 : Compute $O = \text{ginv}(U)\hat{M}_1$

7 : Multiply by -1 all the vectors of U corresponding to the negative values of O $U[, O < 0] = -U[, O < 0]$

Output : The preliminar estimators of μ_1, \dots, μ_K

B.4 Simulations

B.4.1 R package

The developed R-package is called `morpheus` [101] and divided into two main parts :

1. the computation of the directions matrix μ , based on the empirical cross-moments as described in the previous sections ;
2. the optimization of all parameters (including μ), using the initially estimated directions as a starting point.

The former is a straightforward translation of the mathematical formulas (file `R/computeMu.R`), while the latter calls `R constrOptim()` method on the objective function expression and its derivative (file `R/optimParams.R`). For usage examples, please refer to the package help.

B.4.2 Experiments

In this section, we evaluate our algorithm in a first step using mean squared error (MSE). In a second step, we compare experimentally our moments method (morpheus package [101]) and the likelihood method (with `felxmix` package [102]). We arbitrarily choose the parameters for the simulations, which should be discovered by the algorithms (ours, and the EM algorithm).

Experiment 1 (dimension 2) :

$$\begin{aligned}
 K &= 2 \\
 p &= (0.5, 0.5) \\
 b &= (-0.2, 0.5) \\
 \beta &= \begin{pmatrix} 1 & 3 \\ -2 & 1 \end{pmatrix}
 \end{aligned}$$

Experiment 2 (dimension 5) :

$$\begin{aligned}
 K &= 2 \\
 p &= (0.5, 0.5) \\
 b &= (-0.2, 0.5) \\
 \beta &= \begin{pmatrix} 1 & 2 \\ 2 & -3 \\ -1 & 0 \\ 0 & 1 \\ 3 & 0 \end{pmatrix}
 \end{aligned}$$

Experiment 3 (dimension 10) :

$$K = 3$$

$$p = (0.3, 0.3, 0.4)$$

$$b = (-0.2, 0, 0.5)$$

$$\beta = \begin{pmatrix} 1 & 2 & -1 \\ 2 & -3 & 1 \\ -1 & 0 & 3 \\ 0 & 1 & -1 \\ 3 & 0 & 0 \\ 4 & -1 & 0 \\ -1 & -4 & 2 \\ -3 & 3 & 0 \\ 0 & 2 & 1 \\ 2 & 0 & -2 \end{pmatrix}$$

For all three experiments we use both logit and probit links. Computations are always run on the same data both for our package and flexmix – which is a reference for this kind of estimation, using an iterative algorithm to maximize the log-likelihood. Results are aggregated over $N = 1000$ Monte-Carlo runs.

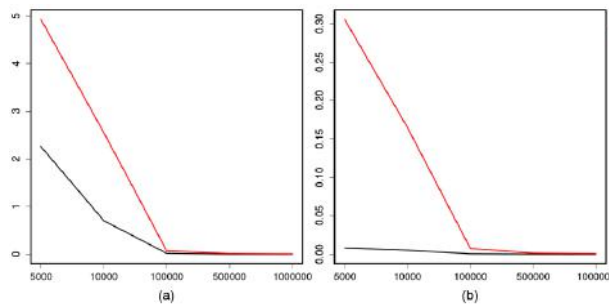


FIGURE B.1 – Experiment 1 : *logit* link function for our algorithm; (a) $\text{MSE}(\hat{\beta})$, (b) $\text{MSE}(\hat{b}, \hat{p})$.

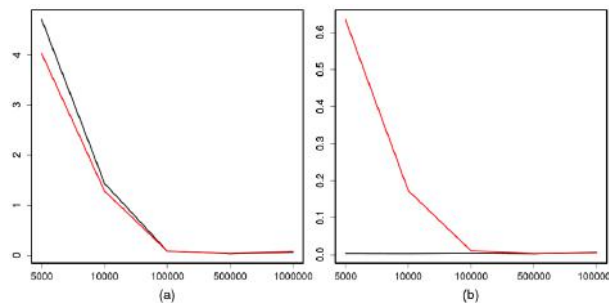


FIGURE B.2 – Experiment 2 : *logit* link function for our algorithm; (a) $\text{MSE}(\hat{\beta})$, (b) $\text{MSE}(\hat{b}, \hat{p})$.

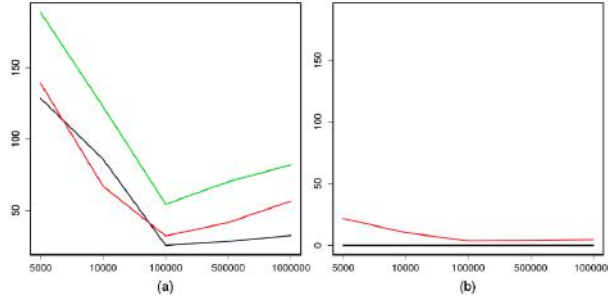


FIGURE B.3 – Experiment 3 : *logit* link function for our algorithm; (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$.

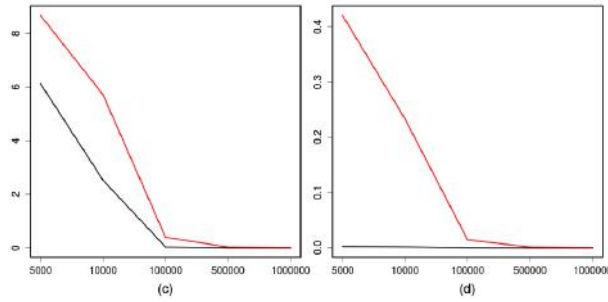


FIGURE B.4 – Experiment 1 : *probit* link function for our algorithm; (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$.

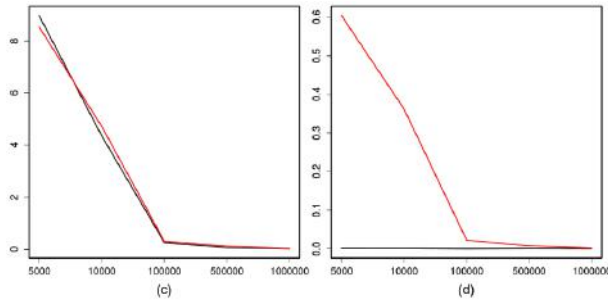


FIGURE B.5 – Experiment 2 : *probit* link function for our algorithm; (a) $MSE(\hat{\beta})$, (b) $MSE(\hat{b}, \hat{p})$.

Mean squared error (MSE). Graphical representations of the MSE are given in figures B.1 to B.6 (figures B.1 to B.3 for logit and figures B.4 to B.6 for probit) versus the sample size (n). In each figure, we represent the MSE associated with each parameter vector versus the sample size. We can see that the goodness of the estimation depends on the sample size and that enough observations is needed to properly estimate the parameters. We have from figure B.1, figure B.2, figure B.4 and figure B.5 that for our moment method, with dimension less or equal to 5, the necessary sample size is around of 10^5 . For large dimension, figures B.3 and B.6 show that 10^6 is not enough to estimate the parameters vectors β .

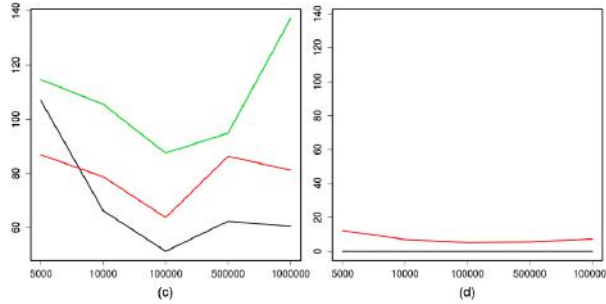


FIGURE B.6 – Experiment 3 : *probit* link function for our algorithm ; (a) $\text{MSE}(\hat{\beta})$, (b) $\text{MSE}(\hat{b}, \hat{p})$.

Algorithms performance. To evaluate algorithms performance, the total number of sample points is fixed to $n = 10^5$. This value is still enough to observe correct performances, yet small enough to remain realistic.

On the figures B.7 and B.8, all (true) parameters (p, b, β) are re-ordered in a real vector, of size $K \times (d + 2) - 1$. This vector is plotted as a line to improve visualization experience, but it must be noted that this does not represent any curve data. Dotted lines corresponds to the computed values plus or minus one standard deviation, and computed values themselves are represented with long dashed lines. The leftmost column corresponds to experiment 1 ($d = 2$), the middle one to experiment 2 ($d = 5$), and the rightmost column corresponds to experiment 3 ($d = 10$).

Figure B.7 : logit link. While most of the times flexmix finds a better solution than our proposed algorithm (smaller variance), both methods are good on average for $d \leq 5$. The case $d = 10$ is not handled well neither by the EM algorithm nor by our package (the latter showing even poorer accuracy). Indeed in this relatively high dimension the number of observations should be much higher.

Figure B.8 : probit link. In this case our algorithm performs slightly better than its flexmix counterpart for $d \leq 5$. However, again, the variance in the case $d = 10$ is way too high – in fact even the average value is generally wrong, when coefficients are non-zero. We can increase n by a factor 100 to obtain more accurate results.

Considering both links with $d \leq 5$, the algorithms performances are comparable with a global small advantage to our method. The case $d = 10$ is handled better by the flexmix algorithm, although clearly not well. Finally, concerning our method we observe a tendency to overestimate large parameters while underestimating small ones. This observation is inverted for flexmix.

It must be noted that our algorithm timing does not depend on n , since it operates on matrices of size $O(d \times K)$. Thus it can be more suitable for very large datasets, where the variability is clearly reduced. Figures B.9 and B.10 (time by seconds versus $\log_{10} n$) illustrate this fact : the timings clearly favor our package – because increasing n has almost no impact on the running time. However, flexmix timings are not that high : just a few minutes for the longest run, on average, on one million sample points. To obtain the data shown on the figure we averaged 1000 runs on random parameters.

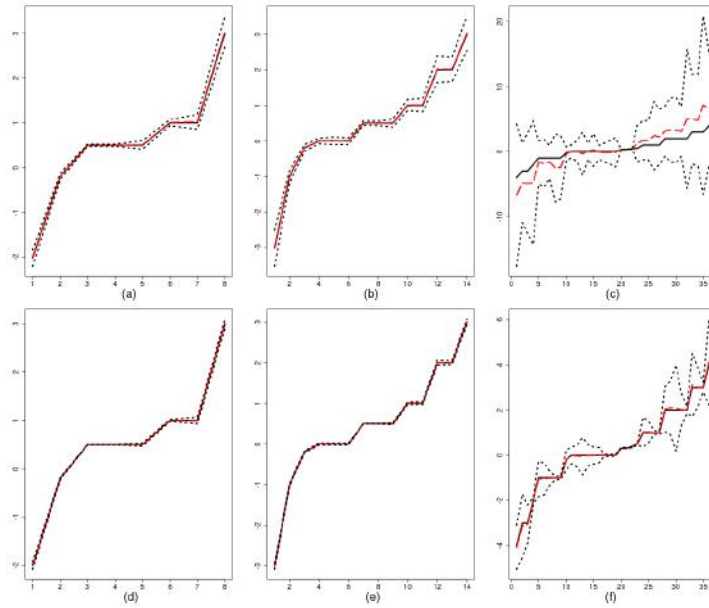


FIGURE B.7 – *Logit* link function. Top : our package, bottom : flexmix. From left to right : experiment 1, 2 and 3 respectively

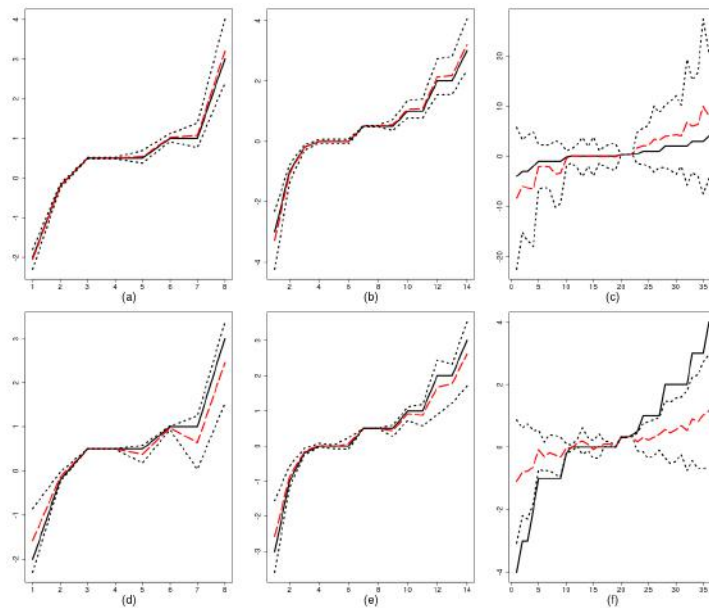


FIGURE B.8 – *Probit* link function. Top : our package, bottom : flexmix. From left to right : experiment 1, 2 and 3 respectively

B.5 Some other identifiability results

In this section, we provide several further identifiability results for mixtures of generalized linear models (GLMs) under various assumptions.

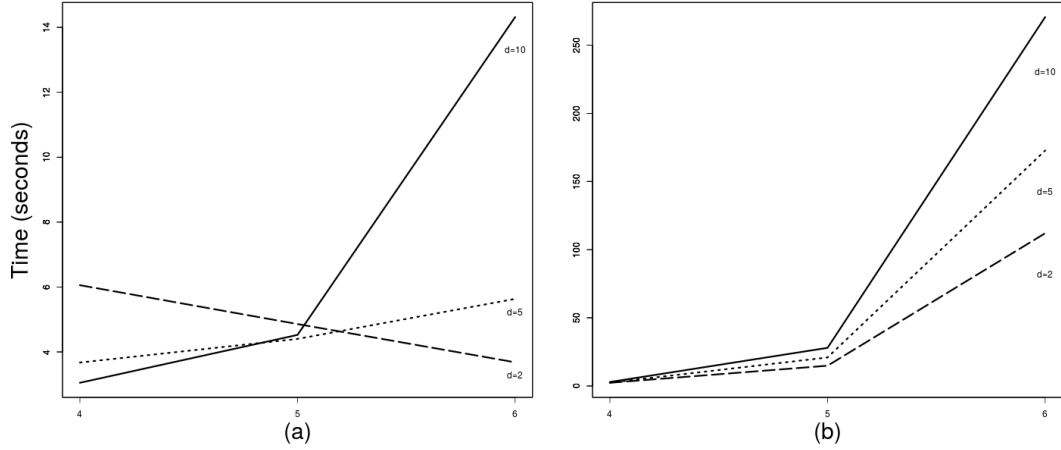


FIGURE B.9 – *Logit* link function. Timings versus sample size. Our algorithm on the left, flexmix on the right. $d = 2$ in long-dashed line, $d = 5$ in dotted line, $d = 10$ in solid line.

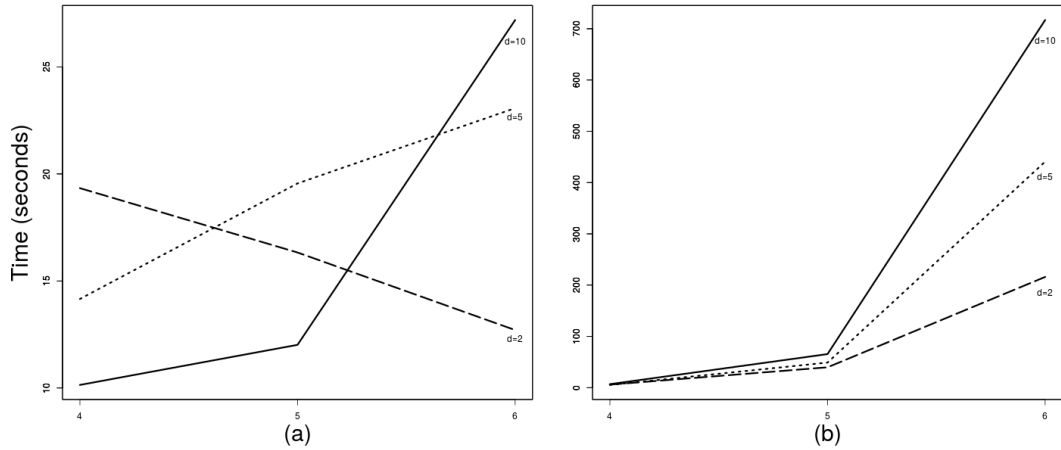


FIGURE B.10 – *Probit* link function. Timings versus sample size. Our algorithm on the left, flexmix on the right. $d = 2$ in long-dashed line, $d = 5$ in dotted line, $d = 10$ in solid line.

B.5.1 Continuous covariates

We first consider the setting where the random vector (X, Y) , $X \in \mathbb{R}^d$, $Y \in \mathbb{R}$ is such that

$$E(Y|X) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + b_k).$$

We assume that for all k , $\omega_k \geq 0$, that $\sum_{k=1}^K \omega_k = 1$, and that g takes value in $(0, 1)$. In case Y takes binary values, this is exactly the model we considered in the previous sections. We show below that the directions of the regression vectors may be recovered as soon as they are distinct, even if the link function is unknown.

Denote $\mathbb{P}_{g, \omega, \beta, b}$ the probability distribution of (X, Y) , with $\omega = (\omega_1, \dots, \omega_K)$, $\beta = [\beta_1, \dots, \beta_K] \in \mathbb{R}^{d \times K}$, and $b = (b_1, \dots, b_K) \in \mathbb{R}^K$. When g is unknown obviously it is needed to fix origin and scale, we choose to fix $g(0)$ and $g(1)$ (with no loss of genera-

lity). Denot $\mu_k = \beta_k / \|\beta_k\|$ and $\lambda_k = \|\beta_k\|$, $k = 1, \dots, K$, so that $\beta_k = \lambda_k \mu_k$.

We introduce the assumptions :

- (S1) The support of the law of X is \mathbb{R}^d .
- (S2) For all $j \neq k$, $\mu_j \neq \mu_k$ and $\mu_j \neq -\mu_k$.
- (S3) The function $g : \mathbb{R} \rightarrow]0, 1[$ is increasing, has limit 0 in $-\infty$, limit 1 in $+\infty$, and it is continuously derivable with derivative having limit 0 in $-\infty$ and in $+\infty$. Also, $g(0) < g(1)$ are fixed.

Remark : There is no assumption on K with respect to d .

Théorème B.5.1.1 *Under assumptions (S1), (S2) and (S3), knowledge of $\mathbb{P}_{g,\omega,\beta,b}$ allows to recover K and μ_1, \dots, μ_K .*

Proof.

If one knows the law of (Y, X) then the function

$$x \mapsto H(x) = \sum_{k=1}^K \omega_k g(\lambda_k \langle \mu_k, x \rangle + b_k)$$

is known on the support of X , thus on \mathbb{R}^d . Then the function

$$DH(x) = \sum_{k=1}^K \omega_k g'(\lambda_k \langle \mu_k, x \rangle + b_k) \mu_k$$

is known, and if $V \in \mathbb{R}^d$, $\lim_{t \rightarrow +\infty} \|DH(tV)\| = 0$ except in case V is orthogonal to at least one of the μ_k 's. The set of $V \in \mathbb{R}^d$ such that $\lim_{t \rightarrow +\infty} \|DH(tV)\| \neq 0$ is then $\cup_{k=1}^K \langle \mu_k \rangle^\perp$, union of disjoint vectorial spaces of dimension $d - 1$, which allows to recover the orthogonal space of $\langle \mu_k \rangle^\perp$ for all k , thus to recover K and all one dimensional spaces $\langle \mu_k \rangle$. Since for all k , $\omega_k g'(b_k) > 0$, this allows to recover the μ_k 's.

Under the more stringent assumption that the regression vectors are linearly independent, it is possible to recover all parameters and the link function.

- (S2bis) The vectors μ_1, \dots, μ_K are linearly independent.

Remark : (H2bis) implies that $K \leq d$.

Théorème B.5.1.2 *Under assumptions (S1), (S2bis) and (S3), the mixture model is identifiable : the knowledge of $\mathbb{P}_{g,\omega,\beta,b}$ allows to recover K , g , ω , β and b .*

Proof.

Using Theorem B.5.1.1, one knows K and the μ_k 's. Since the μ_k 's are linearly independent, by considering the spaces that are orthogonal to all U_k 's except one, we see that the following functions are known : h_1, \dots, h_K given for $j = 1, \dots, K$ by :

$$t \mapsto h_j(t) = \omega_j g(\lambda_j t + b_j) + \sum_{k=1, k \neq j}^K \omega_k g(b_k).$$

Then :

$$\begin{aligned}
h_j(0) &= \sum_{k=1}^K \omega_k g(b_k), \\
\lim_{t \rightarrow +\infty} h_j(t) &= \omega_j + \sum_{k=1, k \neq j}^K \omega_k g(b_k), \\
\lim_{t \rightarrow -\infty} h_j(t) &= \sum_{k=1, k \neq j}^K \omega_k g(b_k).
\end{aligned}$$

This allows to recover ω_j and $g(b_j)$ for $j = 1, \dots, K$. Thus the functions

$$t \mapsto \ell_j(t) = g(\lambda_j t + b_j)$$

are known. Since $g(0) = \ell_j(-b_j/\lambda_j)$ and $g(1) = \ell_j((1 - b_j)/\lambda_j)$ are fixed, one can find λ_j and b_j , and then the function g .

B.5.2 Continuous and categorical covariates

We now consider the situation where part of the covariates are categorical, we denote them Z , and $\{z_1, \dots, z_m\} \subset \mathbb{R}^d$ their possible values. We still denote $X \in \mathbb{R}^d$ the continuous covariates. Now

$$E(Y|X, Z) = \sum_{k=1}^K \omega_k g(\langle \beta_k, X \rangle + \langle \gamma_k, Z \rangle + b_k),$$

and we denote $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ the probability distribution of (X, Y) , with $\omega = (\omega_1, \dots, \omega_K)$, $\beta = [\beta_1 | \dots | \beta_K] \in \mathbb{R}^{d \times K}$, $\gamma = [\gamma_1 | \dots | \gamma_K] \in \mathbb{R}^{d' \times K}$, and $b = (b_1, \dots, b_K) \in \mathbb{R}^K$. We introduce

– (S4) The matrix $\begin{pmatrix} 1 & z_1^T \\ 1 & z_2^T \\ \vdots & \vdots \\ 1 & z_m^T \end{pmatrix}$ is full rank.

Remark : (S4) implies that $d' + 1 \leq m$.

It is the continuous covariates that allow to identify g .

Théorème B.5.2.1 *Under assumptions (S1), (S2bis), (S3) and (S4), the model is identifiable : the knowledge of $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ allows to recover K , g , ω , β , γ et b .*

Proof.

Using Theorem B.5.1.1 applied to the distributions of Y conditional to X and $Z = z$ for all $z \in \{z_1, \dots, z_m\}$, the knowledge of $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ allows to recover K , g , ω , β , and $A_k = (a_{k,i})_{1 \leq i \leq m}$, $k = 1, \dots, K$, with

$$a_{k,i} = b_k + \langle \gamma_k, z_i \rangle.$$

We then know for all k

$$A_k = \begin{pmatrix} 1 & z_1^T \\ 1 & z_2^T \\ \vdots & \vdots \\ 1 & z_m^T \end{pmatrix} \begin{pmatrix} b_k \\ \gamma_k \end{pmatrix}$$

which allows to recover the b_k 's and γ_k 's when (S4) holds.

B.5.3 Longitudinal observations

We now consider the situation where for each individual Y , conditional to the membership of a population, we have p independent experiments with several covariates X_1, \dots, X_p . Thus the random variable Y has dimension m , and

$$E(Y|X, Z) = \sum_{k=1}^K \omega_k (g(\langle \beta_k, X_j \rangle + \langle \gamma_k, Z_j \rangle + b_k))_{1 \leq j \leq p}.$$

As soon as the number of experiments is at least 3, we do not need the linear independence of the regression vectors to get identifiability.

Théorème B.5.3.1 *Assume that $p \geq 3$. If (S1), (S2), (S3) and (S4) hold, then the model is identifiable : the knowledge of $\mathbb{P}_{g, \omega, \beta, \gamma, b}$ allows to recover $K, g, \omega, \beta, \gamma$ and b .*

Proof.

If one knows the law of Y , then, for all fixed $z \in \{z_1, \dots, z_m\}$, one knows the function $H : (\mathbb{R}^d)^p \rightarrow (0, 1)^p$ given by

$$H(x_1, \dots, x_p) = \sum_{k=1}^K \omega_k \left(g(\langle \beta_k, x_j \rangle + \tilde{b}_k(z)) \right)_{1 \leq j \leq p}$$

with $\tilde{b}_k(z) = b_k + \langle \gamma_k, z_i \rangle$. Let us first prove that for all z , the functions $g(\langle \beta_k, \cdot \rangle + \tilde{b}_k(z))$ are linearly independent. Indeed, if $\alpha_1, \dots, \alpha_K$ are such that for all $x \in \mathbb{R}^d$,

$$\sum_{k=1}^K \alpha_k g(\langle \beta_k, x \rangle + \tilde{b}_k(z)) = 0,$$

then by taking the derivative, for all $x \in \mathbb{R}^d$,

$$\sum_{k=1}^K \alpha_k g'(\langle \beta_k, x \rangle + \tilde{b}_k(z)) \beta_k = 0.$$

Since (S2) holds, there exists $V \in \langle \beta_k \rangle^\perp$ such that $V \notin \langle \beta_j \rangle^\perp, j \neq k$. Then taking $x = tV$ and t tending to infinity, we get that $\alpha_k g'(\tilde{b}_k(z)) \beta_k = 0$, and then $\alpha_k = 0$.

Now, following the spectral method of proof developed in [36] to prove that multidimensional mixtures are identifiable, we see that the knowledge of H allows to recover K , the ω_k 's and, for all z , the functions $g(\langle \beta_k, \cdot \rangle + \tilde{b}_k(z))$.

Then, if one knows the function $x \mapsto g(\lambda_k \langle \mu_k, x \rangle + \tilde{b}_k(z))$ one can recover μ_k by taking the derivative, then g and the $\tilde{b}_k(z)$'s as in the proof of Theorem B.5.1.2 then the γ_k 's and the b_k 's as in the proof of Theorem B.5.2.1.

B.5.4 Some perspectives

Identifiability of a model is a first step to obtain theoretical guarantees for practical estimation procedures. In this paper, we proposed one moment method as an estimation strategy in the particular case of binary outcomes and gaussian covariates, for which we proved the asymptotic Gaussian behaviour. As soon as the identifiability of a model is known, any reasonable estimation strategy leads to consistent estimators. Considering the non parametric estimation of the link function, model selection methods should lead to well behaved estimators. Our identifiability results open the way to build estimators for which theoretical guarantees could be obtained. In particular, for parametric maximum likelihood estimators in mixture models for which algorithms already exist, consistency is a consequence of our identifiability theorems by applying the usual theory.

B.6 Proofs

B.6.1 Proof of Theorem B.3.1.1

When the link function g is probit, then $g'(z) = \frac{1}{\sqrt{2\pi}}e^{-z^2/2}$. Replacing in equation (B.2), we have

$$dQ_{(b,\lambda)}(z) = \frac{e^{-\frac{1}{2}((\lambda z+b)^2+z^2)}}{\int e^{-\frac{1}{2}((\lambda z+b)^2+z^2)} dz} dz,$$

which after some computations leads to

$$Q_{(b,\lambda)} = \mathcal{N}\left(-\frac{\lambda b}{\lambda^2 + 1}; \frac{1}{\lambda^2 + 1}\right).$$

Its first two moments are then given by

$$(\alpha_1, \alpha_2) := (\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) = \left(-\frac{\lambda b}{\lambda^2 + 1}; \frac{\lambda^2 b^2 + \lambda^2 + 1}{(\lambda^2 + 1)^2}\right).$$

We can then recover b and λ by

$$b = -\alpha_1 \frac{(\lambda^2 + 1)}{\lambda}$$

and

$$\lambda = \sqrt{(\alpha_2 - \alpha_1^2)^{-1} - 1}.$$

B.6.2 Proof of Theorem B.3.1.2

Using (H3) and integration by parts, we get that for (λ, b) in a neighborhood of $(0, 0)$

$$(\mathbb{E}_{(b,\lambda)}(Z); \mathbb{E}_{(b,\lambda)}(Z^2)) = \left(\frac{\lambda \int g''(\lambda z + b)e^{-z^2/2} dz}{\int g'(\lambda z + b)e^{-z^2/2} dz}; 1 + \frac{\lambda^2 \int g^{(3)}(\lambda z + b)e^{-z^2/2} dz}{\int g'(\lambda z + b)e^{-z^2/2} dz}\right). \quad (\text{B.5})$$

Using (H2),

- (P3) $g'(0) > 0$
- (P4) $g''(0) = g^{(4)}(0) = 0$

Let us define the functions K_s , $s = 1, 2, 3$ such that

$$K_s : \mathbb{R}_+^* \times \mathbb{R} \rightarrow \mathbb{R}$$

$$(\lambda, b) \mapsto K_s(\lambda, b) = \int g^{(s)}(\lambda z + b)e^{-z^2/2} dz$$

Using (H3), the functions K_s , $s = 1, 2, 3$, are differentiable in a neighborhood of $(0, 0)$ and Taylor expansion writes :

$$K_s(\lambda, b) = K_s(0, 0) + \langle \nabla K_s(0, 0), (\lambda, b) \rangle + o(\lambda^2 + b^2). \quad (\text{B.6})$$

Now

$$\frac{\partial K_s}{\partial \lambda}(0, 0) = \int z g^{(s+1)}(0) e^{-z^2/2} dz \quad (\text{B.7})$$

and

$$\frac{\partial K_s}{\partial b}(0, 0) = \int g^{(s+1)}(0) e^{-z^2/2} dz \quad (\text{B.8})$$

so that

$$K_s(\lambda, b) = g^{(s)}(0) \int e^{-z^2/2} dz + g^{(s+1)}(0) \int (\lambda z + b) e^{-z^2/2} dz + o(\lambda^2 + b^2). \quad (\text{B.9})$$

Using (P4) and (B.9), we have

$$\int g'(\lambda z + b) e^{-z^2/2} dz = \sqrt{2\pi} g'(0) + o(\lambda^2 + b^2), \quad (\text{B.10})$$

$$\int g''(\lambda z + b) e^{-z^2/2} dz = \sqrt{2\pi} g^{(3)}(0) b + o(\lambda^2 + b^2), \quad (\text{B.11})$$

and

$$\int g^{(3)}(\lambda z + b) e^{-z^2/2} dz = \sqrt{2\pi} g^{(3)}(0) + o(\lambda^2 + b^2). \quad (\text{B.12})$$

Therefore, replacing (B.10) to (B.12) in (B.5), we get

$$E_{(b,\lambda)}(Z) = \frac{g^{(3)}(0)}{g'(0)} \lambda b + o(\lambda^2 + b^2)$$

and

$$E_{(b,\lambda)}(Z^2) = 1 + \frac{g^{(3)}(0)}{g'(0)} \lambda^2 + o(\lambda^2 + b^2),$$

which easily leads to

$$\lambda^2 = \frac{g'(0)}{g^{(3)}(0)} (E_{(b,\lambda)}(Z)^2 - 1) + o(|E_{(b,\lambda)}(Z)^2 - 1| + |E_{(b,\lambda)}(Z)|)$$

and

$$\lambda b = \frac{g'(0)}{g^{(3)}(0)} E_{(b,\lambda)}(Z) + o(|E_{(b,\lambda)}(Z)^2 - 1| + |E_{(b,\lambda)}(Z)|).$$

This proves that the function $(\lambda, b) \mapsto (E_{(b,\lambda)}(Z), E_{(b,\lambda)}(Z^2))$ is invertible in a neighborhood of $(0, 0)$.

B.6.3 Proof of Theorem B.3.2.1

Let θ^* be the true value of the parameter. For each θ , by the law of large numbers, $Q_n(\theta)$ converges to

$$Q(\theta) := \sum_{j \in [d]} \left\{ M_1(\theta^*)[j] - M_1(\theta)[j] \right\}^2 + \sum_{j, k \in [d]} \left\{ M_2(\theta^*)[j, k] - M_2(\theta)[j, k] \right\}^2 + \sum_{j, k, l \in [d]} \left\{ M_3(\theta^*)[j, k, l] - M_3(\theta)[j, k, l] \right\}^2$$

Define

$$S = \sup_{\theta \in \Theta} \left| Q_n(\theta) - Q(\theta) \right|.$$

Since $Q(\theta)$ has θ^* as unique minimum (up to label switching), to prove the consistency of $\hat{\theta}_n$, it is enough to prove that S converges to 0 in probability, see Theorem 5.7 in [77]. We easily get

$$\begin{aligned} S &\leq \sum_{j \in [d]} \left(\left| \hat{M}_1[j] - M_1(\theta^*)[j] \right| \right) \left(\left| \hat{M}_1[j] \right| + \left| M_1(\theta^*)[j] \right| + 2 \sup_{\theta \in \Theta} \left| M_1(\theta)[j] \right| \right) \\ &+ \sum_{j, k \in [d]} \left(\left| \hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right| \right) \left(\left| \hat{M}_2[j, k] \right| + \left| M_2(\theta^*)[j, k] \right| + 2 \sup_{\theta \in \Theta} \left| M_2(\theta)[j, k] \right| \right) \\ &+ \sum_{j, k, l \in [d]} \left(\left| \hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right| \right) \left(\left| \hat{M}_3[j, k, l] \right| + \left| M_3(\theta^*)[j, k, l] \right| + 2 \sup_{\theta \in \Theta} \left| M_3(\theta)[j, k, l] \right| \right), \end{aligned}$$

and since the functions $\theta \mapsto M_r(\theta)$, $r = 1, 2, 3$ are continuous and Θ is compact, then there exist c_1 , c_2 and c_3 such that

$$\begin{aligned} S &\leq \sum_{j \in [d]} \left(c_1 + \left| \hat{M}_1[j] \right| \right) \left(\left| \hat{M}_1[j] - M_1(\theta^*)[j] \right| \right) \\ &+ \sum_{j, k \in [d]} \left(c_2 + \left| \hat{M}_2[j, k] \right| \right) \left(\left| \hat{M}_2[j, k] - M_2(\theta^*)[j, k] \right| \right) \\ &+ \sum_{j, k, l \in [d]} \left(c_3 + \left| \hat{M}_3[j, k, l] \right| \right) \left(\left| \hat{M}_3[j, k, l] - M_3(\theta^*)[j, k, l] \right| \right) \end{aligned}$$

which converges to 0 by the law of large numbers, which ends the proof of the consistency of $\hat{\theta}_n$.

Let us define Z_n as $Z_n(\theta) = \nabla_{\theta} Q_n(\theta)$. The r -th coordinate of $Z_n(\theta)$ can be obtained by

$$\begin{aligned} \frac{\partial Q_n(\theta)}{\partial \theta_r} &= -2 \left\{ \sum_{j \in [d]} \frac{\partial M_1(\theta)[j]}{\partial \theta_r} \left[\hat{M}_1[j] - M_1(\theta)[j] \right] \right. \\ &+ \sum_{j, k \in [d]} \frac{\partial M_2(\theta)[j, k]}{\partial \theta_r} \left[\hat{M}_2[j, k] - M_2(\theta)[j, k] \right] \\ &\left. + \sum_{j, k, l \in [d]} \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_r} \left[\hat{M}_3[j, k, l] - M_3(\theta)[j, k, l] \right] \right\} \end{aligned}$$

Using Taylor expansion, we get

$$Z_n(\hat{\theta}_n) = Z_n(\theta^*) + \int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] (\hat{\theta}_n - \theta^*) dt \quad (\text{B.13})$$

where $D_1 Z_n$ is the first derivative matrix of Z_n . Since $Z_n(\hat{\theta}_n) = 0$, we have

$$-\sqrt{n}Z_n(\theta^*) = \left[\int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] dt \right] \sqrt{n}(\hat{\theta}_n - \theta^*) \quad (\text{B.14})$$

Let us set

$$\hat{M} = \left(\hat{M}_1[j], \hat{M}_2[j, k], \hat{M}_3[j, k, l] \right)_{1 \leq j, k, l \leq d}$$

and

$$M(\theta^*) = \left(M_1(\theta^*)[j], M_2(\theta^*)[j, k], M_3(\theta^*)[j, k, l] \right)_{1 \leq j, k, l \leq d}$$

Applying the central limit theorem and the delta method we get that $\sqrt{n}Z_n(\theta^*)$ is asymptotically Gaussian.

The (r_1, r_2) -th coordinate of $D_1 Z_n(\theta) = \nabla_{\theta}^2 Q_n(\theta)$ are given by

$$\begin{aligned} \frac{\partial^2 Q_n(\theta)}{\partial \theta_{r_1} \partial \theta_{r_2}} &= -2 \sum_{j \in [d]} \frac{\partial^2 M_1(\theta)[j]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_1[j] - M_1(\theta)[j]] - 2 \sum_{j, k \in [d]} \frac{\partial^2 M_2(\theta)[j, k]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_2[j, k] - M_2(\theta)[j, k]] \\ &\quad - 2 \sum_{j, k, l \in [d]} \frac{\partial^2 M_3(\theta)[j, k, l]}{\partial \theta_{r_1} \partial \theta_{r_2}} \times [\hat{M}_3[j, k, l] - M_3(\theta)[j, k, l]] + V_{r_1 r_2}(\theta) \end{aligned}$$

with

$$\begin{aligned} V_{r_1 r_2}(\theta) &= 2 \sum_{j \in [d]} \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_1}} \times \frac{\partial M_1(\theta)[j]}{\partial \theta_{r_2}} + 2 \sum_{j, k \in [d]} \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_1}} \times \frac{\partial M_2(\theta)[j, k]}{\partial \theta_{r_2}} \\ &\quad + 2 \sum_{j, k, l \in [d]} \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_1}} \times \frac{\partial M_3(\theta)[j, k, l]}{\partial \theta_{r_2}}. \end{aligned}$$

It is not difficult to prove that $\int_0^1 D_1 Z_n[\theta^* + t(\hat{\theta}_n - \theta^*)] dt$ converges in probability to $V(\theta^*)$ so that the proof is completed by showing that the matrix $V = V(\theta^*)$ is invertible.

V is a $q \times q$ matrix with $q = K(2+d) - 1$. Let $U \in \mathbb{R}^q$. We shall denote the coordinates of U according to the parameters. Using the form of V we get that $U^T V U = 0$ if and only if :

$$U^T D M_1(\theta)[j] = 0, \quad j = 1, \dots, q, \quad (\text{B.15})$$

and

$$U^T D M_2(\theta)[j, l] = 0, \quad j, l = 1, \dots, q, \quad (\text{B.16})$$

and

$$U^T D M_3(\theta)[j, l, m] = 0, \quad j, l, m = 1, \dots, q. \quad (\text{B.17})$$

Here, $DM[\cdot]$ is the gradient vector of the involved coordinate of M . Denote $U(\beta_k)$ the d -dimensional vector involving the coordinates of U according to parameter β_k . Denote $\bar{0}$ the d -dimensional zero vector, $\bar{0} \otimes \bar{0}$ the $d \times d$ -dimensional zero matrix and $\bar{0} \otimes \bar{0} \otimes \bar{0}$ the $d \times d \times d$ -dimensional zero third order tensor. Then, the equation (B.15) can be rewritten as :

$$U^T DM_1(\theta)[j] = \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_1(\theta)[j]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_1(\theta)[j]}{\partial b_k} + \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_1(\theta)[j]}{\partial \beta_{mk}}, \quad (\text{B.18})$$

the equation (B.16) can be rewritten as

$$U^T DM_2(\theta)[j, l] = \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_2(\theta)[j, l]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_2(\theta)[j, l]}{\partial b_k} + \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_2(\theta)[j, l]}{\partial \beta_{mk}}, \quad (\text{B.19})$$

and the equation (B.17) can be rewritten as

$$U^T DM_3(\theta)[j, l, m] = \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial \omega_k} + \sum_{k=1}^K U(b_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial b_k} + \sum_{k=1}^K \sum_{s=1}^d U(\beta_{sk}) \frac{\partial M_3(\theta)[j, l, m]}{\partial \beta_{sk}} \quad (\text{B.20})$$

Using the fact that $\sum_{k=1}^d \omega_k = 1$, the first terms of the equations (B.18) to (B.20) are rewritten as :

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_1(\theta)[j]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \right. \\ &\quad \left. - \mathbb{E} [g'(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \right\}, \end{aligned} \quad (\text{B.21})$$

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_2(\theta)[j, l]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g''(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \right. \\ &\quad \left. - \mathbb{E} [g''(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \beta_K(l) \right\} \end{aligned} \quad (\text{B.22})$$

and

$$\begin{aligned} \sum_{k=1}^{K-1} U(\omega_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial \omega_k} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \beta_k(m) \right. \\ &\quad \left. - \mathbb{E} [g^{(3)}(\langle x, \beta_K \rangle + b_K)] \cdot \beta_K(j) \beta_K(l) \beta_K(m) \right\} \end{aligned} \quad (\text{B.23})$$

respectively. Likewise the seconds terms of equations (B.18) to (B.20) are rewritten as :

$$\sum_{k=1}^K U(b_k) \frac{\partial M_1(\theta)[j]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g'(\langle x, \beta_k \rangle + b_k)] \cdot \beta(j), \quad (\text{B.24})$$

$$\sum_{k=1}^K U(b_k) \frac{\partial M_2(\theta)[j, l]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g^{(3)}(\langle x, \beta_k \rangle + b_k)] \cdot \beta(j) \beta(l) \quad (\text{B.25})$$

and

$$\sum_{k=1}^K U(b_k) \frac{\partial M_3(\theta)[j, l, m]}{\partial b_k} = \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g^{(4)} (\langle x, \beta_k \rangle + b_k)] \cdot \beta_k(j) \beta_k(l) \beta_k(m) \quad (\text{B.26})$$

respectively. Derivating with respect to the β_k 's coordinates and using Stein's identity, the last terms of equations (B.18) to (B.20) are rewritten as :

$$\begin{aligned} \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_1(\theta)[j]}{\partial \beta_{mk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(3)} (\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g' (\langle x, \beta_k \rangle + b_k)] U(\beta_k(j)), \end{aligned} \quad (\text{B.27})$$

$$\begin{aligned} \sum_{k=1}^K \sum_{m=1}^d U(\beta_{mk}) \frac{\partial M_2(\theta)[j, l]}{\partial \beta_{mk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(4)} (\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \beta_k(l) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g'' (\langle x, \beta_k \rangle + b_k)] \left\{ \beta_k(j) U(\beta_k(l)) \right. \\ &+ \left. \beta_k(l) U(\beta_k(j)) \right\}, \end{aligned} \quad (\text{B.28})$$

and

$$\begin{aligned} \sum_{k=1}^K \sum_{s=1}^d U(\beta_{sk}) \frac{\partial M_3(\theta)[j, l, m]}{\partial \beta_{sk}} &= \sum_{k=1}^K \omega_k \mathbb{E} [g^{(5)} (\langle x, \beta_k \rangle + b_k)] \langle \beta_k, U(b_k) \rangle \beta_k(j) \beta_k(l) \beta_k(s) \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g^{(3)} (\langle x, \beta_k \rangle + b_k)] \left\{ \beta_k(j) \beta_k(l) U(\beta_k(s)) \right. \\ &+ \left. \beta_k(j) U(\beta_k(l)) \beta_k(s) + U(\beta_k(j)) \beta_k(l) \beta_k(s) \right\} \end{aligned} \quad (\text{B.29})$$

respectively. Then using equations (B.18) to (B.29), we can rewrite equation (B.15) as :

$$\begin{aligned} \bar{0} &= \sum_{k=1}^{K-1} U(\omega_k) \left\{ \mathbb{E} [g' (\langle x, \beta_k \rangle + b_k)] \beta_k - \mathbb{E} [g' (\langle x, \beta_K \rangle + b_K)] \beta_K \right\} \\ &+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E} [g'' (\langle x, \beta_k \rangle + b_k)] \beta_k + \sum_{k=1}^K \omega_k \mathbb{E} [g^{(3)} (\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k \\ &+ \sum_{k=1}^K \omega_k \mathbb{E} [g' (\langle x, \beta_k \rangle + b_k)] U(\beta_k), \end{aligned} \quad (\text{B.30})$$

rewrite equation (B.16) as :

$$\begin{aligned}
\bar{0} \otimes \bar{0} &= \sum_{k=1}^{K-1} U(\omega_k) \left[\mathbb{E}[g''(\langle x, \beta_k \rangle + b_k)] \beta_k \otimes \beta_k - \mathbb{E}[g''(\langle x, \beta_K \rangle + b_K)] \beta_K \otimes \beta_K \right] \\
&+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \beta_k \otimes \beta_k \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g^{(4)}(\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k \otimes \beta_k \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g''(\langle x, \beta_k \rangle + b_k)] \left(U(\beta_k) \otimes \beta_k + \beta_k \otimes U(\beta_k) \right), \tag{B.31}
\end{aligned}$$

and rewrite equation (B.17) as :

$$\begin{aligned}
\bar{0}^{\otimes 3} &= \sum_{k=1}^{K-1} U(\omega_k) \left[\mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \beta_k^{\otimes 3} - \mathbb{E}[g^{(3)}(\langle x, \beta_K \rangle + b_K)] \beta_K^{\otimes 3} \right] \tag{B.32} \\
&+ \sum_{k=1}^K \omega_k U(b_k) \mathbb{E}[g^{(4)}(\langle x, \beta_k \rangle + b_k)] \beta_k^{\otimes 3} + \sum_{k=1}^K \omega_k \mathbb{E}[g^{(5)}(\langle x, \beta_k \rangle + b_k)] \langle U(\beta_k), \beta_k \rangle \beta_k^{\otimes 3} \\
&+ \sum_{k=1}^K \omega_k \mathbb{E}[g^{(3)}(\langle x, \beta_k \rangle + b_k)] \left(U(\beta_k) \otimes \beta_k \otimes \beta_k + \beta_k \otimes U(\beta_k) \otimes \beta_k + \beta_k \otimes \beta_k \otimes U(\beta_k) \right).
\end{aligned}$$

We shall first prove that the vectors $U(\beta_1), \dots, U(\beta_K)$ all belong to the linear space spanned by β_1, \dots, β_K .

Let W be any vector that is orthogonal to this linear space. By multiplying (B.32) on the right by W , and by using the fact that β_1, \dots, β_K are linearly independent by (H1), we get that

$$\forall k = 1, \dots, K, \omega_k (G_3)_k \langle U(\beta_k), W \rangle = 0.$$

Using (H1) we have $\omega_k > 0$, $k = 1, \dots, K$ so that we get

$$\forall k = 1, \dots, K, (G_3)_k \langle U(\beta_k), W \rangle = 0.$$

Then, if (H4) holds, we get that for any k and any W , $\langle U(\beta_k), W \rangle = 0$, which proves that the vectors $U(\beta_1), \dots, U(\beta_K)$ all belong to the linear space spanned by β_1, \dots, β_K . Let B be the $d \times K$ matrix having β_1, \dots, β_K as column vectors. Let $U(\beta)$ be the $d \times K$ matrix having $U(\beta_1), \dots, U(\beta_K)$ as column vectors. We thus have that there exists a $K \times K$ matrix $A = (A_1, \dots, A_K)$ such that $UU(\beta) = BA$.

Set

$$\begin{aligned}
U(\omega) &= \left(U(\omega_1), \dots, U(\omega_{K-1}), - \sum_{k=1}^{K-1} U(\omega_k) \right), \\
U(b) &= (U(b_1), \dots, U(b_K))
\end{aligned}$$

and recall that

$$\omega = \left(\omega_1, \dots, \omega_{K-1}, 1 - \sum_{k=1}^{K-1} \omega_k \right).$$

Whenever R is a K -dimensional vector, denote $Diag(R)$ the $K \times K$ diagonal matrix having the R_k 's on the diagonal.

Let P , Q and Δ be diagonal matrices such that $P = Diag(U(\omega))$, $Q = Diag(U(b))$ and $\Delta = Diag(\omega)$. For $W \in \mathbb{R}^d$, set, $D = Diag(\langle \beta_1, W \rangle, \dots, \langle \beta_K, W \rangle)$. Then using the fact that B is full rank, (B.32) gives that

$$G_3PD + G_4\Delta QD + G_5 + AG_3\Delta D + G_3\Delta DA^T + \Delta Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle), BA_K)D + G_3Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle) = \bar{0} \otimes \bar{0}. \quad (\text{B.33})$$

Since $U(\beta) = BA$, then $U(\beta_k) = \sum_{r=1}^K \beta_r A_{rk} = BA_k$. This implies that

$$Diag(\langle U(\beta_1), \beta_1 \rangle, \dots, \langle U(\beta_K), \beta_K \rangle) = Diag(\langle BA_1, \beta_1 \rangle, \dots, \langle A_K, \beta_K \rangle) = \tilde{D}.$$

So (B.33) can be rewritten as

$$G_3PD + G_4\Delta QD + G_5\Delta \tilde{D}D + AG_3\Delta D + G_3\Delta DA^T + G_3\Delta \tilde{D} = \bar{0} \otimes \bar{0}, \quad (\text{B.34})$$

So that for all $W \in \mathbb{R}^d$, $W \in \mathbb{R}^d$, $AG_3\Delta D + G_3\Delta DA^T$ is a diagonal matrix. Since $G_3\Delta$ has no zero entries, this proves that, under (H1) and (H3), A is a diagonal matrix. In such a case,

$$\tilde{D} = A\tilde{B} \text{ avec } \tilde{B} = Diag(\|\beta_1\|^2, \dots, \|\beta_K\|^2)$$

and (B.34) can be rewritten as

$$G_3PD + G_4\Delta QD + G_5\Delta A\tilde{B}D + AG_3\Delta D + G_3\Delta DA^T + G_3\Delta A\tilde{B} = \bar{0} \otimes \bar{0}. \quad (\text{B.35})$$

But by taking, for $k = 1, \dots, K$, W_k such that $\beta_k^T W_k = 0$, we have $D = 0$. In this case, (B.35) is given by

$$G_3\Delta A\tilde{B} = \bar{0} \otimes \bar{0},$$

and using the fact that we get that G_3 , Δ and \tilde{B} have no zero entries we get that $A = 0$. This implies that $U(\beta_k) = 0$, $k = 1, 2, \dots, K$. Then using the fact that B is full rank, we conclude from (B.30) and (B.31) that

$$G_1P + G_2\Delta Q = \bar{0} \otimes \bar{0}, \quad (\text{B.36})$$

and

$$G_2P + G_3\Delta Q = \bar{0} \otimes \bar{0}. \quad (\text{B.37})$$

Multiplying (B.36) by G_3 and (B.37) by G_2 , we have

$$G_1G_3P + G_2G_3\Delta Q = \bar{0} \otimes \bar{0}, \quad (\text{B.38})$$

and

$$G_2^2P + G_2G_3\Delta Q = \bar{0} \otimes \bar{0}. \quad (\text{B.39})$$

Taking the difference (B.38)-(B.39), we get

$$(G_1G_3 - G_2^2)P = \bar{0} \otimes \bar{0},$$

and since $G_1G_3 - G_2^2$ has no zero entries, this leads to $P = 0$. Moreover, since $G_3\Delta$ has no zero entries, this leads also $Q = 0$. Thus, under (H1), (H4) and (H5), the matrix V is full rank.

Annexe C

Morpheus-package : Estimate Parameters of Mixtures of Logistic Regressions

Sommaire

C.1 R topics documented :	174
C.1.1 morpheus-package	174
C.1.2 alignMatrices	175
C.1.3 computeMoments	175
C.1.4 computeMu	176
C.1.5 generateSampleIO	177
C.1.6 multiRun	177
C.1.7 normalize	179
C.1.8 optimParams	180
C.1.9 plotBox	181
C.1.10 plotCoefs	181
C.1.11 plotHist	182
C.1.12 plotQn	182

Description

Mixture of logistic regressions parameters (H)estimation with (U)spectral methods. The main methods take d -dimensional inputs and a vector of binary outputs, and return parameters according to the GLMs mixture model (General Linear Model). For more details see chapter 3 in the PhD thesis of Mor-Absa Loum : <http://www.theses.fr/s156435>, available here <https://www.math.u-psud.fr/~loum/IMG/pdf/these.compressed.pdf>.

Version 0.2 – 0

Author Benjamin Auder “Benjamin.Auder@u-psud.fr” [aut,cre],

Mor-Absa Loum “Mor-Absa.Loum@u-psud.fr” [aut]

Maintainer Benjamin Auder “Benjamin.Auder@u-psud.fr”

Depends $R(\geq 3.0.0)$,

Imports MASS, jointDiag, methods, pracma

Suggests devtools, flexmix, parallel, testthat, roxygen2, tensor, nloptr

License MIT + file LICENSE

RoxygenNote 5.0.1

Collate 'utils.R' 'A_NAMESPACE.R' 'computeMu.R' 'multiRun.R' 'optimParams.R' 'plot.R' 'sampleIO.R'

NeedsCompilation yes

Repository CRAN

Date/Publication 2018-06-14 18 :48 :14 UTC

C.1 R topics documented :

C.1.1 morpheus-package

Description

Mixture of logistic regressions parameters (H)estimation with (U)spectral methods. The main methods take d -dimensional inputs and a vector of binary outputs, and return parameters according to the GLMs mixture model (General Linear Model). For more details see chapter 3 in the PhD thesis of Mor-Absa Loum : <http://www.theses.fr/s156435>, available here <https://www.math.u-psud.fr/~loum/IMG/pdf/these.compressed.pdf>.

Details

The package devtools should be useful in development stage, since we rely on testthat for unit tests, and roxygen2 for documentation. knitr is used to generate the package vignette. Concerning the other suggested packages :

- **tensor** is used for comparing to some reference functions initially coded in R ; it should not be required in further package versions ;
- **jointDiag** allows to solve a joint diagonalization problem, providing a more robust solution compared to a single diagonalization ;

- `parallel` (generally) permits to run the bootstrap method faster.

The three main functions are located in `R/main.R` :

- `getParamsDirs_ref` : reference method to estimate parameters directions ;
- `getParamsDirs` : method of choice to estimate parameters directions, using a spectral decomposition of inputs/outputs ;
- `getBootstrapParams` : run `getParamsDirs` on B bootstrap replicates.

Author(s)

Benjamin Auder “Benjamin.Auder@u-psud.fr” [aut,cre],
 Mor-Absa Loum “Mor-Absa.Loum@u-psud.fr” [aut]
 Maintainer : Benjamin Auder “Benjamin.Auder@u-psud.fr”

C.1.2 alignMatrices

Description

Align a set of parameters matrices, with potential permutations.

Usage

```
alignMatrices(Ms, ref, ls_mode)
```

Arguments

<code>Ms</code>	A list of matrices, all of same size $d \times K$
<code>ref</code>	Either a reference matrix or "mean" to align on empirical mean
<code>ls_mode</code>	How to compute the labels assignment : “exact” for exact algorithm (default, but might be time-consuming, complexity is $O(K^3)$), or “approx1”, or “approx2” to apply a greedy matching algorithm (heuristic) which for each column in reference (resp. in current row) compare to all unassigned columns in current row (resp. in reference)

Value

The aligned list (of matrices), of same size as `Ms`

C.1.3 computeMoments

Description

Compute cross-moments of order 1, 2, 3 from `X`, `Y`

Usage

```
computeMoments(X, Y)
```

Arguments

X	Matrix of input data (size $n \times d$)
Y	Vector of binary outputs (size n)

Value

A list L where L[[i]] is the i-th cross-moment

C.1.4 computeMu

Description

Estimate the normalized columns mu of the beta matrix parameter in a mixture of logistic regressions models, with a spectral method described in the package vignette.

Usage

```
computeMu(X, Y, optargs = list())
```

Arguments

X	Matrix of input data (size $n \times d$)
Y	Vector of binary outputs (size n)
optargs	List of optional argument : <ul style="list-style-type: none">• 'jd_method', joint diagonalization method from the package <code>jointDiag</code>: 'uwedge'(default) or 'jedi'.• 'jd_nvects', number of random vectors for joint-diagonalization (or 0 for $p = d$, canonical basis by default)• 'M', moments of order 1, 2, 3 : will be computed if not provided• 'K', number of populations (estimated with ranks of M2 if not given)

Value

The estimated normalized parameters as columns of a matrix mu of size $d \times K$.

See Also

`multiRun` to estimate statistics based on mu, and `generateSampleIO` for I/O random generation.

Examples

```
io = generateSampleIO(10000, 1/2, matrix(c(1,0,0,1),ncol=2), c(0,0), "probit")
mu = computeMu(io$X, io$Y, list(K=2)) #or just X and Y for estimated K
```

C.1.5 generateSampleIO

Description

Generate input matrix X of size $n \times d$ and binary output of size n , where Y is subdivided into K groups of proportions p . Inside one group, the probability law $\mathbb{P}(Y = 1)$ is described by the corresponding column parameter in the matrix $\text{beta} + \text{intercept } b$.

Usage

```
generateSampleIO(n, p, beta, b, link)
```

Arguments

<code>n</code>	Number of individuals
<code>p</code>	Vector of $K - 1$ populations relative proportions (sum ≤ 1)
<code>beta</code>	Vectors of model parameters for each population, of size $d \times K$
<code>b</code>	Vector of intercept values (use <code>rep(0, K)</code> for no intercept)
<code>link</code>	Link type; ‘‘logit’’ or ‘‘probit’’

Value

A list with

- `X`: the input matrix (size $n \times d$)
- `Y`: the output vector (size n)
- `index`: the population index (in $1 : K$) for each row in X

C.1.6 multiRun

Description

Estimate N times some parameters, outputs of some list of functions. This method is thus very generic, allowing typically bootstrap or Monte-Carlo estimations of matrices μ or beta . Passing a list of functions opens the possibility to compare them on a fair basis (exact same inputs). It’s even possible to compare methods on some deterministic design of experiments.

Usage

```
multiRun(fargs, estimParams, prepareArgs = function(x, i) x, N = 10,  
         ncores = 3, agg = lapply, verbose = FALSE)
```

Arguments

<code>fargs</code>	List of arguments for the estimation functions
<code>estimParams</code>	List of <code>nf</code> function(s) to apply on <code>fargs</code> - shared signature
<code>prepareArgs</code>	Prepare arguments for the functions inside <code>estimParams</code>
<code>N</code>	Number of runs
<code>ncores</code>	Number of cores for parallel runs (≤ 1 : sequential)
<code>agg</code>	Aggregation method (default : <code>lapply</code>)
<code>verbose</code>	TRUE to indicate runs + methods numbers

Value

A list of `nf` aggregates of `N` results (matrices).

Examples

```
beta <- matrix(c(1,-2,3,1),ncol=2)

# Bootstrap + computeMu, morpheus VS flexmix ; assumes fargs first 3 elts X,Y,K
io <- generateSampleIO(n=1000, p=1/2, beta=beta, b=c(0,0), "logit")
mu <- normalize(beta)
res <- multiRun(list(X=io$X,Y=io$Y,optargs=list(K=2,jd_nvects=0)), list(
# morpheus
function(fargs) {
  library(morpheus)
  ind <- fargs$ind
  computeMu(fargs$X[ind,],fargs$Y[ind],fargs$optargs)
},
# flexmix
function(fargs) {
  library(flexmix)
  ind <- fargs$ind
  K <- fargs$optargs$K
  dat = as.data.frame( cbind(fargs$Y[ind],fargs$X[ind,]) )
  out = refit( flexmix( cbind(V1, 1 - V1) ~ 0+., data=dat, k=K,
  model=FLXMRglm(family="binomial") ) )
  normalize( matrix(out@coef[1:(ncol(fargs$X)*K)], ncol=K) )
} ),
prepareArgs = function(fargs,index) {
if (index == 1)
  fargs$ind <- 1:nrow(fargs$X)
else
  fargs$ind <- sample(1:nrow(fargs$X),replace=TRUE)
fargs
}, N=10, ncores=3)
for (i in 1:2)
  res[[i]] <- alignMatrices(res[[i]], ref=mu, ls_mode="exact")
```

```

# Monte-Carlo + optimParams from X,Y, morpheus VS flexmix ; first args n,p,beta,
res <- multiRun(list(n=1000,p=1/2,beta=beta,b=c(0,0),optargs=list(link="logit"))
# morpheus
function(fargs) {
  library(morpheus)
  K <- fargs$optargs$K
  mu <- computeMu(fargs$X, fargs$Y, fargs$optargs)
  V <- list( p=rep(1/K,K-1), beta=mu, b=c(0,0) )
  optimParams(V,fargs$optargs)$beta
},
# flexmix
function(fargs) {
  library(flexmix)
  K <- fargs$optargs$K
  dat <- as.data.frame( cbind(fargs$Y,fargs$X) )
  out <- refit( flexmix( cbind(V1, 1 - V1) ~ 0+., data=dat, k=K,
  model=FLXMRglm(family="binomial") ) )
  sapply( seq_len(K), function(i) as.double( out@components[[1]][[i]][,1] ) )
} ),
prepareArgs = function(fargs,index) {
  library(morpheus)
  io = generateSampleIO(fargs$n, fargs$p, fargs$beta, fargs$b, fargs$optargs$link)
  fargs$X = io$X
  fargs$Y = io$Y
  fargs$optargs$K = ncol(fargs$beta)
  fargs$optargs$M = computeMoments(io$X,io$Y)
  fargs
}, N=10, ncores=3)
for (i in 1:2)
  res[[i]] <- alignMatrices(res[[i]], ref=beta, ls_mode="exact")

```

C.1.7 normalize

Description

Normalize a vector or a matrix (by columns), using euclidian norm

Usage

```
normalize(X)
```

Arguments

X Vector or matrix to be normalized

Value

The normalized matrix (1 column if **X** is a vector)

C.1.8 optimParams

Description

Optimize the parameters of a mixture of logistic regressions model, possibly using `mu <- computeMu(...)` as a partial starting point.

Usage

```
optimParams(K, link = c("logit", "probit"), optargs = list())
```

Arguments

- | | |
|----------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| K | Number of populations. |
| link | The link type, 'logit' or 'probit'. |
| optargs | a list with optional arguments : <ul style="list-style-type: none">• 'M' : list of moments of order 1,2,3 : will be computed if not provided.• 'X, Y' : input/output, mandatory if moments not given• 'exact' : use exact formulas when available? |

Value

An object 'op' of class OptimParams, initialized so that `op$run(x0)` outputs the list of optimized parameters

- **p** : proportions, size K
- **beta** : regression matrix, size $d \times K$
- **b** : intercepts, size K

x0 is a vector containing respectively the $K - 1$ first elements of **p**, then **beta** by columns, and finally **b** : `x0 = c(p[1 : (K - 1)], as.double(beta), b)`.

See Also

`multiRun` to estimate statistics based on `mu`, and `generateSampleIO` for I/O random generation.

Examples

```
# Optimize parameters from estimated mu
io = generateSampleIO(10000, 1/2, matrix(c(1,-2,3,1),ncol=2), c(0,0), "logit")
mu = computeMu(io$X, io$Y, list(K=2))
M <- computeMoments(io$X, io$Y)
o <- optimParams(2, "logit", list(M=M))
x0 <- c(1/2, as.double(mu), c(0,0))
```

```
par0 <- o$run(x0)
# Compare with another starting point
x1 <- c(1/2, 2*as.double(mu), c(0,0))
par1 <- o$run(x1)
o$f( o$linArgs(par0) )
o$f( o$linArgs(par1) )
```

C.1.9 plotBox

Description

Draw boxplot

Usage

```
plotBox(mr, x, y)
```

Arguments

mr	Output of <code>multiRun()</code> , list of lists of functions results
x	Row index of the element inside the aggregated parameter
y	Column index of the element inside the aggregated parameter

Examples

```
#See example in ?plotHist
```

C.1.10 plotCoefs

Description

Draw coefs estimations + standard deviations

Usage

```
plotCoefs(mr, params)
```

Arguments

mr	Output of <code>multiRun()</code> , list of lists of functions results
params	True value of parameters matrix

Examples

```
#See example in ?plotHist
```

C.1.11 plotHist

Description

Plot histogram

Usage

```
plotHist(mr, x, y)
```

Arguments

mr	Output of <code>multiRun()</code> , list of lists of functions results
x	Row index of the element inside the aggregated parameter
y	Column index of the element inside the aggregated parameter

Examples

```
beta <- matrix(c(1,-2,3,1),ncol=2)
mr <- multiRun(...) #see bootstrap example in ?multiRun : return lists of mu_hat
mu <- normalize(beta)
for (i in 1:2)
  mr[[i]] <- alignMatrices(res[[i]], ref=mu, ls_mode="exact")
plotHist(mr, 2, 1) #second row, first column
```

C.1.12 plotQn

Description

Draw 3D map of objective function values

Usage

```
plotQn(N, n, p, beta, b, link)
```

Arguments

N	Number of starting points
n	Number of points in sample
p	Vector of proportions
beta	Regression matrix (target)
b	Vector of biases
link	Link function (logit or probit)

Examples

```
#See example in ?plotQn
```


Bibliographie

- [1] Manoj Kumar Mohapatra, P Patra, and Ritesh Ku Agrawala. Manifestation and outcome of concurrent malaria and dengue infection. *Journal of Vector Borne Diseases*, 49(4) :262–265, 12 2012.
- [2] M.B. Mushtaq, Mehmood Qadri, and A Rashid. Concurrent infection with dengue and malaria : An unusual presentation. *Hindawi Publishing cooperation Case report in Medecine*, 2013 :2, 3 2013.
- [3] Bernard Carme, Severine Matheus, Gerd Donutil, Olivia Raulin, Mathieu Nacher, and Jacques Morvan. Concurrent dengue and malaria in cayenne hospital, french guiana. *Emerging Infections Diseases*, 15(4) :668–671, 4 2009.
- [4] C. Subhash Arya, Lalit K. Mehta, Nirmala Agarwal, K. Agarwal Bharat, George Mathai, and Arun Moondhara. Episodes of concurrent dengue and malaria. *Dengue Bulletin*, 29, 01 2005.
- [5] Stan Deresinski. Concurrent plasmodium vivax malaria and dengue. *Emerging Infectious Diseases*, 12(11) :1802, 11 2006.
- [6] Nadir Ali, Asif Nadeem, Masood Anwar, Waheed Uz Zaman Tariq, and Rashid A Chotani. Dengue fever in malaria endemic areas. *Journal of theCollege of Physicians and Surgeons–Pakistan : JCPSP*, 16(5) :340–342, 6 2006.
- [7] Nicolas Senn, Dagwin Suarkia, Doris Manong, Peter Max Siba, and William John Hannan Mcbride. Contribution of dengue fever to the burden of acute febrile illnesses in papua new guinea : An age-specific prospective study. *The American Journal of Tropical Medecine and Hygiene*, 85(1) :132–137, 7 2011.
- [8] Remi N. Charrel, Philippe Brouqui, Cedric Foucault, and Xavier De Lamballerie. Concurrent dengue and malaria. *Emerging Infections Diseases*, 11(7) :1153–1154, 7 2005.
- [9] Vitor R Mendonça, Bruno B Andrade, Belisa M L Souza, Ligia C L adn Magalhães, Maria P G Mourão, Marcus V G Lacerda, and Manoel Barral-Netto. Unravelling the patterns of host immune responses in plasmodium vivax malaria and dengue co-infection. 14 :315, 2015.
- [10] Marycelin Baba, Christopher Hugh Logue, Bamidele Oderinde, Hauwa Abdulmaleek, Joshua Williams, James Lewis, Thomas R Laws, Roger Hewson, Alessandro Marcello, and Pierlanfranco D’ Agaro. Evidence of arbovirus co-infection in suspected febrile malaria and typhoid patients in nigeria. *The Journal of Infection in Developing Countries*, 7(1) :51–59, 1 2013.
- [11] Sylla Thiam, Moussa Thior, Babacar Faye, Médoune N diop, Mamadou Lamine Diouf, Mame Birame Diouf, Ibrahima Diallo, Fatou Ba Fall, Jean Louis Ndiaye,

- Audrey Albertini, Evan Lee, Pernille Jorgensen, Oumar Gaye, and David Bell. Major reduction in anti-malarial drug consumption in senegal after nation-wide introduction of malaria rapid diagnostic tests. *PloS One*, 6(4) :1–7, 4 2011.
- [12] ANSD. *programme national de lutte contre le paludisme au Senegal.*, 2009.
- [13] Abdourahmane Sow, Cheikh Loucoubar, Diawo Diallo, Oumar Faye, Youssoupha Ndiaye, Cheikh Saadibou Senghor, Anta Tal Dia, Ousmane Faye, Scott C. Weaver, Mawlouth Diallo, Denis Malvy, and Amadou Alpha Sall. Concurrent malaria and arbovirus infections in kedougou, southeastern senegal. *Malaria Journal*, 15 :47, 01 2016.
- [14] J. A. Nelder and R. W. M. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)*, 135(3) :370–384, 1972.
- [15] P. McCullagh and Nelder J. A. *Generalized linear models, second edition*. Chapman and Hall, 1989.
- [16] A. Agresti. *Categorical Data Analysis, 3rd edition*. Wiley, 2013.
- [17] Anestis Antoniadis, Jacques Berruyer, and Rene Carmona. Régression non linéaire et applications. 01 1992.
- [18] David W. Hosmer and Standley Lemeshow. *Applied Logistic Regression, 2nd edition*. Wiley, 2001.
- [19] Simon Newcomb. A generalized theory of the combination of observations so as to obtain the best result. *American Journal of Mathematics*, 8(4) :343–366, 08 1886.
- [20] Karl Pearson. Iii. contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London A : Mathematical, Physical and Engineering Sciences*, 185 :71–110, 1894.
- [21] Sylvia Frühwirth-Schnatter. *Finite mixture and Markov switching models*. Springer.
- [22] Bettina Grün and Friedrich Leisch. Finite mixtures of generalized linear regression models. In *Recent Advances in Linear Models and Related Areas : Essays in Honour of Helge Toutenburg*, pages 205–230, Heidelberg, 2008. Physica-Verlag HD.
- [23] Murray Aitkin. A general maximum likelihood analysis of overdispersion in generalized linear models. *Statistics and Computing*, 6(3) :251–262, Sep 1996.
- [24] Dean A. Follmann and Diane Lambert. Generalizing logistic regression by nonparametric mixing. *Journal of the American Statistical Association*, 84(405) :295–300, 1989.
- [25] Peiming Wang, Martin L. Puterman, Iain Cockburn, and Nhu Le. Mixed poisson regression models with covariate dependent rates. *Biometrics*, 52(2) :381–400, 1996.
- [26] Peiming Wang and Martin L. Puterman. Mixed logistic regression models. *Journal of Agricultural, Biological, and Environmental Statistics*, 3(2) :175–200, 1998.
- [27] G. Boiteau, M. Singh, R. P. Singh, G. C. C. Tai, and T. R. Turner. Rate of spread of pvyn by *alatemyzus persicae* (sulzer) from infected to healthy plants under laboratory conditions. *Potato Research*, 41(4) :335–344, Dec 1998.
- [28] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1) :1–38, 1977.

- [29] Christopher M. Bishop and Markus Svensén. Bayesian hierarchical mixtures of experts. *CoRR*, abs/1212.2447, 2012.
- [30] Michael I. Jordan and Robert A. Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural Comput.*, 6(2) :181–214, March 1994.
- [31] Lei Xu, Michael I. Jordan, and Geoffrey E Hinton. An alternative model for mixtures of experts. In G. Tesauro, D. S. Touretzky, and T. K. Leen, editors, *Advances in Neural Information Processing Systems 7*, pages 633–640. MIT Press, 1995.
- [32] Bettina Grün. Fitting finite mixtures of linear mixed models with the EM algorithm. In Paula Brito, editor, *Compstat 2008—Proceedings in Computational Statistics*, volume II, pages 165–173. Physica Verlag, Heidelberg, Germany, 2008.
- [33] Peter McCullagh. *Tensor methods in statistics*. Monographs on statistics and applied probability (Series). Chapman and Hall, london ; new york edition, 1987.
- [34] Hanie Sedghi, Majid Janzamin, and Anima Anandkumar. Provable tensor methods for learning mixtures of generalized linear models. In Arthur Gretton and Christian C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 1223–1231, Cadiz, Spain, 09–11 May 2016. PMLR.
- [35] A. Anandkumar, D. Hsu, and S. M. Kakade. A method of moments for mixture models and hidden markov models. 2012.
- [36] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variables models. *Journal of machine learning*, 2014.
- [37] Anima Anandkumar, Majid Janzamin, and Hanie Sedghi. Provable tensor methods for learning mixtures of generalized linear model. 2015.
- [38] Anastasia Podosinnikova, Francis Bach, and Simon Lacoste-Julien. Rethinking LDA : Moment Matching for Discrete ICA. In *NIPS 2015 - Advances in Neural Information Processing Systems 28*, Montreal, Canada, December 2015.
- [39] A. Souloumiac. Nonorthogonal joint diagonalization by combining givens and hyperbolic rotations. *IEEE Transactions on Signal Processing*, 57(6) :2222–2231, June 2009.
- [40] Bijan Afsari. Sensitivity analysis for the problem of matrix joint diagonalization. *SIAM Journal on Matrix Analysis and Applications*, 30(3) :1148–1171, 09 2008.
- [41] Stuart Coles. *An introduction to Statistical Modeling of Extreme Values*. Springer, 2001.
- [42] M.R. Leadbetter, G. Lindgren, and H. Rootzen. *Extremes and Related Properties of Random Sequences and Processes*. 3Island Press, 1983.
- [43] James Pickands. Statistical inference using extreme order statistics. *Annals of Statistics*, 3(1) :119–131, 1975.
- [44] Gwladys Toulemonde. *Estimation et tests en théorie des valeurs extrêmes*. PhD thesis, Université Paris VI- Pierre et Marie Curie, October 2008.
- [45] Pathé Ndao. *Modélisation de valeurs extrêmes conditionnelles en présence de censure*. PhD thesis, Université Gaston Berger de Saint-Louis, August 2015.

- [46] Myriam Garrido. *Modélisation des évènements rares et estimation des quantiles extrêmes, méthodes de sélection de modèles pour les queues de distribution*. PhD thesis, Université Joseph-Fourier - Grenoble I, June 2002.
- [47] Alexandre Lekina. *Estimation non-paramétrique des quantiles extrêmes conditionnels*. PhD thesis, Université Joseph-Fourier - Grenoble I, October 2010.
- [48] Gilles Stupfler. *Un modèle de Markov caché en assurance et Estimation de frontière et de point terminal*. PhD thesis, Université de Strasbourg, November 2011.
- [49] Laurent Gardes. *Estimation d'une fonction quantile extrême*. PhD thesis, Université Montpellier II - Sciences et Techniques du Languedoc, October 2003.
- [50] Jonathan EL Methni. *Contributions à l'estimation de quantiles extrêmes. Applications à des données environnementales*. PhD thesis, Université de Grenoble, October 2013.
- [51] Yang Hu. *Extreme Value Mixture Modelling with Simulation Study and Applications in Finance and Insurance*. PhD thesis, University of Canterbury, New Zealand, July 2013.
- [52] Anna Elizabeth MacDonald. *Extreme Value Mixture Modelling with Medical and Industrial Applications*. PhD thesis, University of Canterbury, 2012.
- [53] Mor Absa Loum, Marie Anne Poursat, Abdourahmane Sow, Amadou Alpha Sall, Cheikh Loucoubar, and Elisabeth Gassiat. Multinomial logistic model for coinfection diagnosis between arbovirus and paludisme in kedougou. *The International Journal of Biostatistics*, 2017.
- [54] Mor Absa Loum, Benjamin Auder, and Elisabeth Gassiat. Mixture of generalized linear models : identifiability and applications. *preprint*, 2018.
- [55] Mor Absa Loum and Benjamin Auder. Package r : Morpheus, mixture of generalized linear models. *preprint*, 2018.
- [56] Mor Absa Loum, Gilles Celeux, and Aliou Diop. Extreme value mixture under random censoring. *preprint*, 2018.
- [57] A. MacDonald, C.J. Scarrott, D. Lee, B. Darlow, M. Reale, and G. Russell. A flexible extreme value mixture model. *Computational Statistics and Data Analysis*, 55(6) :2137 – 2157, 2011.
- [58] Robin Genuer, Jean-Michel Poggi, and Christine Tuleau-Malot. Variable selection using random forests. *Pattern Recognition Letters*, 31(14) :2225–2236, 2010.
- [59] Hirotugu Akaike. A new look at the statistical model identification. *IEEE transactions on automatic control*, 19(6) :716–722, 12 1974.
- [60] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning : With Applications in R*. Springer Publishing Company, Incorporated, 2014.
- [61] Paula Branco, Luis Torgo, and Rita P. Ribeiro. A survey of predictive modeling on imbalanced domains. *ACM Computing Surveys (CSUR)*, 49(31) :31 :1–31 :50, 11 2016.
- [62] Leo Breiman. Random forest. *Machine Learning*, 45 :5–32, 2001.

- [63] Leo Breiman, J.H Friedman, R. A Olshen, and C. J Stone. *Classification And Regression Trees*. Chapman and Hall, 1984.
- [64] Robin Genuer. *Forêts aléatoires : aspects théoriques, sélection de variables et applications*. PhD thesis, Université Paris-Sud 11, Novembre 2010.
- [65] Robin Genuer and Jean-Michel Poggi. Arbres cart et forêts aléatoires, importance et sélection de variables. January 2017. working paper or preprint.
- [66] Robin Genuer, Jean-Michel Poggi, and Christine Tuleau-Malot. Vsurf : An r package for variable selection using random forests. *The R Journal*, 7(2) :19–33, 11 2015.
- [67] Chao Chen, Andy Liaw, and Leo Breiman. Using random forests to learn imbalanced data. Technical report, University of Berkeley, 2006.
- [68] Nitesh V. Chawla. Data mining for imbalanced datasets : An overview. In : *O. Maimon, L. Rokach (eds.) Data Mining and Knowledge Discovery Hand-book, SPRINGER*, pages 875–886, 2010.
- [69] Bartosz Krawczyk. Learning from imbalanced data : open challenges and future directions. *Progress in Artificial Intelligence*, 5(4) :221–232, 2016.
- [70] Sham Kakade, Adam Tauman Kalai, Varun Kanade, and Ohad Shamir. Efficient learning of generalized linear and single index models with isotonic regression. *CoRR*, abs/1104.2018, 2011.
- [71] Ariadna Quattoni, Michael Collins, and Trevor Darrell. Conditional random fields for object recognition. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 1097–1104. MIT Press, 2005.
- [72] Y. Wang and G. Mori. Max-margin hidden conditional random fields for human action recognition. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 872–879, June 2009.
- [73] Slav Petrov and Dan Klein. Discriminative log-linear grammars with latent variables. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 1153–1160. Curran Associates, Inc., 2008.
- [74] Percy Liang, Alexandre Bouchard-Côté, Dan Klein, and Ben Taskar. An end-to-end discriminative approach to machine translation. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 761–768. Association for Computational Linguistics, 2006.
- [75] Hanie Sedghi, Majid Janzamin, and Anima Anandkumar. Provable tensor methods for learning mixtures of generalized linear models. In Arthur Gretton and Christian C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 1223–1231, Cadiz, Spain, 09–11 May 2016. PMLR.
- [76] Charles Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2 : Probability Theory*, pages 583–602, Berkeley, Calif., 1972. University of California Press.

- [77] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, cambridge edition, 1998.
- [78] Cédric Gouy-Pailler. *Vers une modélisation dynamique de l'activité cérébrale pour la conception d'interfaces cerveau-machines asynchrones*. PhD thesis, Université Joseph Fourier (Grenoble), 2009.
- [79] P. Tichavsky and A. Yeredor. Fast approximate joint diagonalization incorporating weight matrices. *IEEE Transactions on Signal Processing*, 57(3) :878–891, March 2009.
- [80] Arnaldo Frigessi, Ola Haug, and Håvard Rue. A dynamic mixture model for unsupervised tail estimation without threshold selection. *Extremes*, 5(3) :219–235, September 2002.
- [81] Cible N. Behrens, Hedibert F Lopes, and Dani Gamerman. Bayesian analysis of extreme events with threshold estimation. *Statistical Modelling*, 4 :227–244, 2004.
- [82] Beatriz Vaz de Melo Mendes and Hedibert Freitas Lopes. Data driven estimates for mixtures. *Computational Statistics and Data Analysis*, 47(3) :583 – 598, 2004.
- [83] Andrea Tancredi, Clive Anderson, and Anthony O'Hagan. Accounting for threshold uncertainty in extreme value estimation. *Extremes*, 9(2) :87, Aug 2006.
- [84] Xin Zhao, Carl Scarrott, Les Oxley, and Marco Reale. Extreme value modelling for forecasting market crisis impacts. *Applied Financial Economics*, 20(1-2) :63–72, 2010.
- [85] Xin Zhao, Carl John Scarrott, Les Oxley, and Marco Reale. Garch dependence in extreme value models with bayesian inference. *Mathematics and Computers in Simulation*, 81(7) :1430 – 1440, 2011. Selected Papers of the Combined IMACS World Congress and MSSANZ 18th Biennial Conference on Modelling and Simulation, Cairns, Australia, 13-17 July, 2009.
- [86] David Lee, Wai Keung Li, and Tony Siu Tung Wong. Modeling insurance claims via a mixture exponential model combined with peaks-over-threshold approach. *Insurance : Mathematics and Economics*, 51(3) :538 – 550, 2012.
- [87] Fernando Ferraz do Nascimento, Dani Gamerman, and Hedibert Freitas Lopes. A semiparametric bayesian approach to extreme value estimation. *Statistics and Computing*, 22(2) :661–675, Mar 2012.
- [88] Julie Carreau and Yoshua Bengio. A hybrid pareto mixture for conditional asymmetric fat-tailed distributions. *IEEE Transactions on Neural Networks*, 20(7) :1087–1101, July 2009.
- [89] Julie Carreau and Yoshua Bengio. A hybrid pareto model for asymmetric fat-tailed data : the univariate case. *Extremes*, 12(1) :53–76, Mar 2009.
- [90] Stuart G. Coles and Jonathan A. Tawn. A bayesian analysis of extreme rainfall data. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 45(4) :463–478, 1996.
- [91] Stuart G. Coles and Elwyn A. Powell. Bayesian methods in extreme value modelling : A review and new developments. *International Statistical Review / Revue Internationale de Statistique*, 64(1) :119–136, 1996.

- [92] Dmitry Babichev and Francis Bach. Slice inverse regression with score functions. working paper or preprint, October 2016.
- [93] Shaoli Wang, Weixin Yao, and Mian Huang. A note on the identifiability of nonparametric and semiparametric mixtures of glms. *Statistics and Probability Letters*, 93 :41–45, 2014.
- [94] Bettina Grün and Friedrich Leisch. Identifiability of finite mixtures of multinomial logit models with varying and fixed effects. *Journal of Classification*, 25(2) :225–247, Nov 2008.
- [95] Dean A. Follmann and Diane Lambert. Identifiability of finite mixtures of logistic regression models. *Journal of Statistical Planning and Inference*, 27(3) :375 – 381, 1991.
- [96] Peirong Xu, Heng Peng, and Tao Huang. Unsupervised learning of mixture regression models for longitudinal data. *Computational Statistics and Data Analysis*, 125 :44 – 56, 2018.
- [97] Arnost Komarek and Lenka Komarkova. Clustering for multivariate continuous and discrete longitudinal data. *The Annals of Applied Statistics*, 7(1) :177–200, 2013.
- [98] Hiroyuki Kasahara and Katsumi Shimotsu. Nonparametric identification of finite mixture models of dynamic discrete choices. *Econometrica*, 77(1) :135–175, 2009.
- [99] Marco Mondelli and Andrea Montanari. On the connection between learning two-layers neural networks and tensor decomposition. *CoRR*, abs/1802.07301, 2018.
- [100] Sylvia Frühwirth-Schnatter. *Finite mixture and Markov switching models*. Springer Series in Statistics. Springer, New York, 2006.
- [101] Benjamin Auder and Mor Absa Loum. Morpheus : An R package to estimate parameters of logistic regressions mixtures. *CRAN*, June 2018.
- [102] Bettina Grün and Friedrich Leisch. Flexmix : An R package for finite mixture modelling. *R News*, 7(1) :8–13, April 2007.
- [103] Dean A. Follmann and Diane Lambert. Identifiability of finite mixtures of logistic regression models. *J. Statist. Plann. Inference*, 27(3) :375–381, 1991.
- [104] Shaoli Wang, Weixin Yao, and Mian Huang. A note on the identifiability of nonparametric and semiparametric mixtures of GLMs. *Statist. Probab. Lett.*, 93 :41–45, 2014.
- [105] Bettina Grün and Friedrich Leisch. Identifiability of finite mixtures of multinomial logit models with varying and fixed effects. *J. Classification*, 25(2) :225–247, 2008.
- [106] Arnost Komárek and Lenka Komárková. Clustering for multivariate continuous and discrete longitudinal data. *Ann. Appl. Stat.*, 7(1) :177–200, 2013.

Titre : Modèle de mélange et modèles linéaires généralisés, application aux données de co-infection (arbovirus & paludisme)

Mots Clefs : Co-infection, Méthode des moments, Méthode spectrale, Modèle de mélange, Modèle linéaire généralisé, Théorie de valeurs extrêmes.

Résumé :

Nous nous intéressons, dans cette thèse, à l'étude des modèles de mélange et des modèles linéaires généralisés, avec une application aux données de co-infection entre les arbovirus et les parasites du paludisme.

Après une première partie consacrée à l'étude de la co-infection par un modèle logistique multinomial, nous proposons dans une deuxième partie l'étude des mélanges de modèles linéaires généralisés. La méthode proposée pour estimer les paramètres du mélange est une combinaison d'une méthode des moments et d'une méthode spectrale. Nous proposons à la fin une dernière partie consacrée aux mélanges de valeurs extrêmes en présence de censure. La méthode d'estimation proposée dans cette partie se fait en deux étapes basées sur la maximisation d'une vraisemblance.

Title : Mixture model and generalized linear model, application to co-infection data (arbovirus & malaria)

Keys words : Co-infection, Extreme value theory, Generalized linear model, Mixture model, Moments method, Spectral method.

Abstract :

We are interested, in this thesis, to the study of mixture models and generalized linear models, with an application to co-infection data between arboviruses and malaria parasites.

After a first part dedicated to the study of co-infection using a multinomial logistic model, we propose in a second part to study the mixtures of generalized linear models. The proposed method to estimate the parameters of the mixture is a combination of a moment method and a spectral method. Finally, we propose a final section for studying extreme value mixtures under random censoring. The estimation method proposed in this section is done in two steps based on the maximization of a likelihood.

