



**HAL**  
open science

# Modèles statistiques pour des données de survie corrélées

Tristan Lorino

► **To cite this version:**

Tristan Lorino. Modèles statistiques pour des données de survie corrélées. Sciences du Vivant [q-bio]. Institut national agronomique paris-grignon - INA P-G, 2002. Français. NNT: . tel-00003672

**HAL Id: tel-00003672**

**<https://theses.hal.science/tel-00003672>**

Submitted on 3 Nov 2003

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

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

Institut National Agronomique Paris-Grignon

# THÈSE

pour l'obtention du

DOCTORAT DE L'INSTITUT NATIONAL AGRONOMIQUE PARIS-GRIGNON

Discipline : Statistique

présentée et soutenue publiquement le 21 mai 2002

par

**M. Tristan LORINO**

titulaire du DEA de statistique et modèles aléatoires en Économie et Finance de l'Université  
Denis Diderot - Paris 7

## MODÈLES STATISTIQUES POUR DES DONNÉES DE SURVIE CORRÉLÉES

**Directeurs de thèse**

MM. Jean-Jacques DAUDIN et Moez SANAA

**Composition du jury**

*Rapporteurs :* Monsieur le Professeur J.-M. AZAÏS  
Monsieur le Docteur D. COMMENGES

*Examineurs :* Monsieur le Professeur J.-J. DAUDIN  
Monsieur le Docteur C. DUCROT  
Monsieur le Docteur S. ROBIN  
Monsieur le Docteur M. SANAA



« JE NE PEUX PAS me tromper au sujet de  $12 \times 12 = 144$ . Et on ne peut pas opposer la sûreté de la *mathématique* au relatif manque de sûreté de propositions empiriques. En effet la proposition mathématique a été obtenue par une série d'actions qui ne se différencient d'aucune façon du reste des actions de la vie et qui sont tout aussi sujettes à l'oubli, l'inadvertence et l'illusion. »

L. WITTGENSTEIN, *De la certitude*.

« ON APPELLE ÇA, un peu obscurément, la loi des grands nombres. Par quoi l'on peut dire à peu près que, si un homme se tue pour telle raison et un autre pour telle autre, dès qu'on a affaire à un très grand nombre, le caractère arbitraire et personnel de ces motifs disparaît, et il ne demeure... précisément, qu'est-ce qui demeure? Voilà ce que j'aimerais vous entendre dire. Ce qui reste, en effet, vous le voyez vous-même, c'est ce que nous autres profanes appelons tout bonnement la moyenne, c'est-à-dire quelque chose dont on ne sait absolument pas ce que c'est. Permettez-moi d'ajouter que l'on a tenté d'expliquer logiquement cette loi des grands nombres en la considérant comme une sorte d'évidence. On a prétendu, au contraire, que cette régularité dans des phénomènes qu'aucune causalité ne régit ne pouvait s'expliquer dans le cadre de la pensée traditionnelle; sans parler de mainte autre analyse, on a aussi défendu l'idée qu'il ne s'agissait pas seulement d'événements isolés, mais de lois, encore inconnues, régissant la totalité. Je ne veux pas vous ennuyer avec les détails, d'autant que je ne les ai plus présents à l'esprit, mais personnellement, il m'importerait beaucoup de savoir s'il faut chercher là-dedans quelque mystérieuse loi de la totalité ou si tout simplement, par une ironie de la Nature, l'exceptionnel provient de ce qu'il ne se produit rien d'exceptionnel, et si le sens ultime du monde peut être découvert en faisant la moyenne de tout ce qui n'a pas de sens! L'une ou l'autre de ces deux conceptions ne devrait-elle pas avoir une influence décisive sur notre sentiment de la vie? Quoi qu'il en soit, en effet, la possibilité d'une vie ordonnée repose toute entière sur cette loi des grands nombres; si cette loi de compensation n'existait pas, il y aurait des années où il ne se produirait rien, et d'autres où plus rien ne serait sûr; les famines alterneraient avec l'abondance, les enfants seraient en défaut ou en excès et l'humanité voletterait de côté et d'autre entre ses possibilités célestes et ses possibilités infernales comme les petits oiseaux quand on s'approche de leur cage. »

R. MUSIL, *L'homme sans qualités*.



Je désire témoigner ici de ma reconnaissance envers tous ceux qui ont suivi, soutenu et guidé ce travail de recherche de quatre années. Qu'ils en soient très chaleureusement remerciés.

Je désire adresser une marque tout particulière de ma gratitude

à M. Moez Sanaa, instigateur de ce travail, pour sa présence et son aide toutes deux considérables, et à qui le fruit de notre longue collaboration revient pour grande partie,

à M. Jean-Jacques Daudin, pour son acceptation du titre de directeur de cette thèse, pour la constance de son soutien et la justesse des orientations qu'il a bien voulu me suggérer,

à M. Stéphane Robin, pour ses conseils sans cesse pertinents et ses levées de doute,

à M<sup>elle</sup> Sylvie Escolano et M. Michel Chavance, pour leur participation essentielle au suivi de ce travail,

à MM. Jean-Marc Azais et Daniel Commenges, pour l'attention assidue qu'ils ont bien voulu prêter à la lecture, puis à la correction de ce manuscrit,

à tous les membres du Laboratoire d'Épidémiologie et d'Analyse des Risques (ENVA) et du Département Organisation et Modélisation de l'Information et des Processus (INA P-G), présents à mes côtés tout au long de mon travail,

ainsi qu'aux utilisateurs aguerris de (L<sup>A</sup>)T<sub>E</sub>X et S-Plus, pour leur aide si précieuse apportée au travers des groupes de discussion (`fr.comp.text.tex`) et listes de diffusion (`gut@ens.fr` et `s-news@lists.biostat.wustl.edu`).

Enfin, de ma famille dont j'ai vu le visage se modifier grandement durant ces quatre années, je remercie infiniment chaque membre – pour son soutien, sa patience.

Que le dernier de mes remerciements, tout particulier, soit pour Elsa.



*À la mémoire de ma mère.*



# Sommaire

Table des figures et tableaux . . . . .	4
<b>I MODÉLISATION DE LA SURVIE POUR DES DONNÉES INDÉPENDANTES</b>	<b>5</b>
Introduction . . . . .	7
Notations . . . . .	11
<b>1 DÉFINITIONS, NOTATIONS ET PROPRIÉTÉS</b>	<b>13</b>
1.1 Fonctions de survie et de risque . . . . .	14
1.2 Censure . . . . .	15
1.3 Constitution des observations et processus de comptage . . . . .	17
1.4 Estimation du risque cumulé . . . . .	19
1.5 Estimation de la fonction de survie . . . . .	21
1.6 Comparaison de la survie de plusieurs groupes . . . . .	25
<b>2 MODÈLE DE COX</b>	<b>29</b>
2.1 Définitions et notations . . . . .	30
2.2 Estimation . . . . .	30
2.3 Tests . . . . .	34
2.4 Critères d'adéquation au modèle de Cox . . . . .	35
2.5 Considération des ex-æquo . . . . .	42
2.6 Extensions du modèle . . . . .	43
<b>II MODÉLISATION DE LA SURVIE POUR DES DONNÉES CORRÉLÉES</b>	<b>45</b>
Introduction . . . . .	47
<b>3 MODÈLE MARGINAL DE COX</b>	<b>51</b>
3.1 Estimation et tests dans le cadre général . . . . .	52
3.2 Estimations et tests dans un cadre restreint . . . . .	55
3.3 Critères d'adéquation . . . . .	55
3.4 Coefficient de dispersion ( <i>design effect</i> ) . . . . .	59
<b>4 MODÈLE DE FRAGILITÉ DE COX</b>	<b>61</b>
4.1 Écriture du modèle . . . . .	62
4.2 Concernant la fragilité . . . . .	63
4.3 Cas où la fragilité suit une loi gamma . . . . .	66

4.4	Tests portant sur la variance de la fragilité . . . . .	69
4.5	Compléments . . . . .	71
4.6	Prédiction de la fragilité . . . . .	73
<b>III</b>	<b>DEUX ÉTUDES POUR UNE COMPARAISON EMPIRIQUE DES DIFFÉRENTS MODÈLES</b>	<b>81</b>
<b>5</b>	<b>ÉTUDE PAR SIMULATIONS</b>	<b>83</b>
5.1	Présentation des procédures S-PLUS . . . . .	84
5.2	Présentation des simulations . . . . .	92
5.3	Paramétrage des simulations et choix des sorties . . . . .	95
5.4	Résultats . . . . .	97
5.5	Conclusions et perspectives . . . . .	106
<b>6</b>	<b>ÉTUDE ÉPIDÉMIOLOGIQUE DES GENN</b>	<b>109</b>
6.1	Introduction . . . . .	110
6.2	Protocole . . . . .	111
6.3	Statistiques descriptives . . . . .	117
6.4	Statistiques analytiques . . . . .	119
<b>IV</b>	<b>CONCLUSIONS</b>	<b>135</b>
<b>V</b>	<b>BIBLIOGRAPHIE</b>	<b>139</b>
<b>VI</b>	<b>ANNEXES</b>	<b>149</b>
<b>A</b>	<b>THÉORIE STATISTIQUE</b>	<b>151</b>
A.1	Processus aléatoires et intégrales stochastiques . . . . .	151
A.2	Processus de comptage . . . . .	152
A.3	Théorème de la limite centrale . . . . .	154
A.4	Produit infini (ou intégral) . . . . .	155
A.5	Vraisemblances complète et partielle . . . . .	156
A.6	Résultats complémentaires . . . . .	160
<b>B</b>	<b>DÉMONSTRATIONS COMPLÉMENTAIRES</b>	<b>161</b>
B.1	Test non-paramétrique d'égalité des fonctions de risque . . . . .	161
B.2	Concernant le modèle marginal de Cox . . . . .	164
<b>C</b>	<b>COMPLÉMENT AU MODÈLE DE FRAGILITÉ GAMMA DE COX</b>	<b>171</b>
<b>D</b>	<b>PROCÉDURE DES SIMULATIONS SOUS S-PLUS</b>	<b>173</b>
<b>E</b>	<b>RÉSULTATS COMPLETS DES SIMULATIONS</b>	<b>175</b>
<b>F</b>	<b>COMPLÉMENTS À L'ÉTUDE ÉPIDÉMIOLOGIQUE DES GENN</b>	<b>193</b>
<b>G</b>	<b>GLOSSAIRE ÉPIDÉMIOLOGIQUE</b>	<b>197</b>
	<b>Index</b>	<b>201</b>

# Table des figures et tableaux

## Figures

2.1	Les vraisemblances successives. . . . .	32
2.2	Schéma représentant la corrélation. . . . .	47
4.1	Modèle bayésien. . . . .	72
5.1	Mécanisme des simulations pas à pas. . . . .	93
5.2	Étude de la puissance ( $\beta_0 = 0$ et taux de censure de 60 %) . . . . .	102
5.3	Étude de la puissance ( $\beta_0 = 0,7$ et taux de censure de 60 %) . . . . .	104
6.1	Répartition des deux types d'observations . . . . .	117
6.2	Estimation de Nelson-Aalen de la courbe de survie . . . . .	118
6.3	Valeurs des estimations de la fragilité $\exp(\hat{\Theta})$ . . . . .	125
6.4	Résidus de la déviance pour le modèle de fragilité. . . . .	127

## Tableaux

5.1	Récapitulatif des configurations choisies pour les simulations . . . . .	96
5.2	Biais relatif de l'estimation du coefficient de régression . . . . .	98
5.3	Coefficient de dispersion et critère d'Akaike ( $\beta_0 = 0$ ) . . . . .	100
5.4	Coefficient de dispersion et critère d'Akaike ( $\beta_0 = 0,7$ ) . . . . .	101
5.5	Étude de la puissance (taux de censure de 60 %) . . . . .	103
5.6	Étude du taux de recouvrement à 95 % de $\sigma^2$ . . . . .	105
5.7	Tableau récapitulatif des résultats pour chaque paramètre . . . . .	107
6.1a	Variables explicatives étudiées dans le cadre des GENN . . . . .	113
6.1b	Variables explicatives étudiées dans le cadre des GENN (suite) . . . . .	114
6.1c	Variables explicatives étudiées dans le cadre des GENN (fin) . . . . .	115
6.2	Tableau récapitulatif des types d'observations . . . . .	117
6.3	Étude de la fonction de survie . . . . .	118
6.4a	Comparaison des ajustements naïf et marginal . . . . .	121
6.4b	Comparaison des ajustements avec « effet élevage » . . . . .	122
6.5	Ajustements ne prenant en compte que les variables individuelles. . . . .	124
6.6	Étude des élevages les plus « fragiles » . . . . .	126
6.7	Tests de la proportionnalité des risques . . . . .	127
6.8	Valeur de la log-vraisemblance pour les différents ajustements. . . . .	132
6.9	Procédures algorithmiques pour l'ajustement du modèle de fragilité . . . . .	133
E.1	Estimations (indépendance, premier regroupement et $\beta_0 = 0$ ) . . . . .	176
E.2	Estimations (indépendance, second regroupement et $\beta_0 = 0$ ) . . . . .	177
E.3	Estimations (indépendance, premier regroupement et $\beta_0 = 0,7$ ) . . . . .	178
E.4	Estimations (indépendance, second regroupement et $\beta_0 = 0,7$ ) . . . . .	179
E.5	Estimations (loi gamma de variance 1, premier regroupement et $\beta_0 = 0$ ) . . . . .	180

TABLE DES FIGURES ET TABLEAUX

---

E.6	Estimations (loi gamma de variance 1, second regroupement et $\beta_0 = 0$ ) . . .	181
E.7	AIC et TR (loi gamma et $\beta_0 = 0$ ) . . . . .	182
E.8	Estimations (loi gamma de variance 4, premier regroupement et $\beta_0 = 0$ ) . .	183
E.9	Estimations (loi gamma de variance 4, second regroupement et $\beta_0 = 0$ ) . .	184
E.10	Estimations (loi gamma de variance 1, premier regroupement et $\beta_0 = 0, 7$ )	185
E.11	Estimations (loi gamma de variance 1, second regroupement et $\beta_0 = 0, 7$ ) .	186
E.12	AIC et TR (loi gamma et $\beta_0 = 0, 7$ ) . . . . .	187
E.13	Estimations (loi gamma de variance 4, premier regroupement et $\beta_0 = 0, 7$ )	188
E.14	Estimations (loi gamma de variance 4, second regroupement et $\beta_0 = 0, 7$ ) .	189
E.15	Estimations (loi stable positive, second regroupement et $\beta_0 = 0$ ) . . . . .	190
E.16	Estimations (loi stable positive, second regroupement et $\beta_0 = 0, 7$ ) . . . . .	191
E.17	AIC et TR (loi stable positive et second regroupement) . . . . .	192
F.1a	Ajustement naïf pour décembre et janvier . . . . .	194
F.1b	Ajustement naïf pour février et mars . . . . .	194
F.2	Comparaison des ajustements naïf et marginal (interactions incluses) . . .	195
F.3	Ajustement mixte (interactions incluses) . . . . .	196

Première partie

**MODÉLISATION DE LA SURVIE  
POUR DES DONNÉES  
INDÉPENDANTES**



# Introduction

## Historique

L'analyse des données de survie voit le jour au XVII<sup>e</sup> siècle, dans le domaine de la démographie. L'objectif des analystes de ce siècle est l'estimation, à partir des registres de décès, de diverses caractéristiques de la population – son effectif, sa longévité, etc. Ces analyses, très générales, ne sont affinées qu'à partir du XIX<sup>e</sup> siècle, avec l'apparition de catégorisations suivant des « variables exogènes » (sexe, nationalité, catégories socio-professionnelles...). Durant ce siècle, apparaissent également les premières modélisations concernant la probabilité de mourir à un certain âge, probabilité qui sera par la suite désignée sous le terme de « fonction de risque ». Enfin, l'analyse des données de survie commence de déborder le cadre stricte de la démographie pour investir, au XX<sup>e</sup> siècle, toutes les disciplines susceptibles d'avoir recours à de tels types de données : l'actuariat, la physique (avec l'apparition de la théorie de la fiabilité), l'industrie (pharmaceutique, biomédicale)...

Jusqu'en 1950, la communauté des statisticiens s'intéresse peu à l'analyse des données de survie, la principale contribution étant celle de Greenwood (1926), qui propose une formule pour l'erreur standard d'une table de survie.

En 1951, Weibull conçoit un modèle paramétrique dans le domaine de la fiabilité ; à cet effet, il fournit une nouvelle distribution de probabilité qui sera par la suite fréquemment utilisée en analyse de la survie : la « loi de Weibull ».

En 1958, Kaplan et Meier présentent d'importants résultats concernant l'estimation non-paramétrique de la fonction de survie ; de l'estimateur résultant, ils étudient l'espérance, la variance et les propriétés asymptotiques.

L'année 1972 se révèle être une date fondamentale : en effet, un modèle statistique semi-paramétrique voit le jour, grâce aux travaux de Cox. Ce modèle comporte des variables exogènes qui sont introduites, dans la fonction de risque, au moyen d'une composante de régression paramétrique – le reste de cette fonction de risque, non-paramétrique, demeurant indéterminée.

De ce modèle, sans nul doute le plus utilisé en analyse des données de survie, seront tirées quantités de variantes, et notamment des formulations permettant de stratifier l'effet des covariables, d'introduire une dépendance vis-à-vis du temps, ou encore de prendre en compte une possible interdépendance des durées de vie observées.

## Modélisation

L'analyse des données de survie a pour première particularité de ne concerner que des **variables aléatoires positives** (modélisant les durées de vie). Une conséquence de cette par-

ticularité est que la loi normale ne sera plus ici la référence en matière de distribution. Le plus souvent, toute autre loi issue de la famille exponentielle, et à support dans  $\mathbb{R}^+$ , lui sera préférée.

Une deuxième particularité de cette analyse est l'**incomplétude des données** – différente de la troncature, qui équivaut à une perte d'information. Analysant la survenue d'un certain type d'événement, nous qualifierons de « donnée complète » un temps correspondant à l'observation de la survenue de l'événement, et de « donnée incomplète » un temps correspondant à l'absence d'observation de cet événement. Nous emploierons respectivement, par la suite, les termes de **donnée non censurée** et de **donnée censurée**.

Tout temps d'observation, qu'il corresponde ou non à une date d'apparition de l'événement étudié – qu'il soit ou non *censuré* – contribue à l'information de départ : le modèle statistique pour données de survie exploite donc toute l'information qui peut être recueillie. Par conséquent, et comparativement à tout autre modèle statistique pour des données quantitatives, le modèle pour données de survie accroît la qualité de l'inférence (la contrepartie de cette amélioration étant la complexification des processus statistiques sous-tendant le modèle).

Enfin, la troisième particularité de l'analyse des données de survie est la **terminologie** s'y rattachant. Outre les termes conçus à l'origine de cette analyse (tels que *fonction de survie*, *fonction de risque*), cette terminologie est essentiellement due au contexte épidémiologique. Les principaux termes sont définis dans l'annexe G.

Concernant les modèles statistiques proprement dits, trois approches sont possibles : paramétrique, non-paramétrique et semi-paramétrique.

L'**approche paramétrique** stipule l'appartenance de la loi de probabilité *réelle* des observations à une classe particulière de lois, qui dépendent d'un certain nombre (fini) de paramètres.

L'avantage de cette approche est la facilitation attendue de la phase d'estimation des paramètres, ainsi que de l'obtention d'intervalles de confiance et de la construction de tests.

L'inconvénient de la méthode paramétrique est l'inadéquation pouvant exister entre le phénomène étudié et le modèle retenu.

L'**approche non-paramétrique** ne nécessite aucune hypothèse quant à la loi de probabilité *réelle* des observations – et c'est là son principal avantage. Il s'agit dès lors d'un problème d'*estimation fonctionnelle*, avec les ambiguïtés que cela implique – par exemple, la fonction de survie, qui est continue, sera estimée par une fonction discontinue.

L'inconvénient d'une telle approche est la nécessité de disposer d'un nombre important d'observations, le problème de l'estimation d'un paramètre fonctionnel étant délicat puisqu'il appartient à un espace de dimension infinie.

L'**approche semi-paramétrique** est une sorte de *compromis* entre les deux approches précédentes. La loi de probabilité *réelle* des observations est supposée appartenir à une classe de lois pour partie dépendant de paramètres, et pour partie s'écrivant sous forme de fonction(s) non-paramétrique(s). Relativement récente – elle est apparue au cours des années soixante-dix –, cette approche est très répandue en analyse de la survie, notamment au travers du **modèle de régression de Cox** (1972).

## La corrélation

La statistique fait de l'indépendance des observations une hypothèse forte. Or cette hypothèse s'avère quelquefois irrecevable. Ainsi, un échantillon peut présenter, du fait même de sa constitution, une structure de corrélation<sup>1</sup>. Des « modèles-types » d'échantillons laissent immédiatement percevoir cette corrélation.

Parmi eux, citons les échantillons avec **répétition de mesures** : ces échantillons, qui correspondent à un suivi longitudinal, sont constitués de mesures effectuées à plusieurs reprises sur les mêmes sujets. Il est évident que les mesures propres à un sujet ne satisfont pas l'hypothèse d'indépendance.

Citons également les échantillons avec **emboîtement de données** : ces échantillons sont constitués de données pouvant être regroupées selon certains critères. Ainsi, des mesures effectuées sur des sujets suivis en milieu hospitalier peuvent présenter une structure de corrélation : les sujets appartenant à un même établissement hospitalier sont susceptibles de fournir des données plus homogènes que celles observées chez des sujets issus d'établissements distincts.

Lorsque l'hypothèse d'indépendance n'est plus vérifiée, l'analyse statistique doit prendre en compte l'**hétérogénéité** des données, afin de corriger les estimations produites : en effet, l'ignorance de cette hétérogénéité entraîne un biais des estimateurs et une mésestimation de leur variance.

À l'origine – c'est-à-dire dans le modèle de régression linéaire –, cette prise en compte consiste en l'inclusion d'une variable explicative supplémentaire. Cette variable, modélisant l'hétérogénéité, peut être fixe ou bien aléatoire : le choix s'effectue en fonction des objectifs que vise l'analyse statistique. Dans le cas d'une variable fixe, il n'est question que d'affiner l'analyse en tenant compte de l'hétérogénéité (l'inférence demeure relative) : dans l'exemple des sujets suivis en milieu hospitalier, l'analyse ne sera valable que pour les hôpitaux qui auront été retenus. Dans le cas d'une variable aléatoire, l'objectif est de valider pleinement l'inférence qui en résultera : le caractère aléatoire de la variable « hôpital » autorise l'interprétation des résultats pour d'autres hôpitaux que ceux qui ont été retenus pour l'analyse.

Par la suite, concernant les modèles non linéaires, deux grandes approches ont vu le jour : l'**approche marginale** et l'**approche mixte**.

La première se caractérise par une modélisation de la réponse d'un sujet conditionnellement aux covariables, et non aux réponses des autres sujets. Il s'agit donc de la modélisation d'une « réponse moyenne » sur l'ensemble des groupes (hôpitaux).

La seconde, quant à elle, modélise une réponse individuelle (mesure sur un sujet dans un hôpital donné) conditionnellement aux covariables et aux réponses des autres sujets de ce groupe (c'est-à-dire des autres sujets de ce même hôpital). Il s'agit donc ici de la modélisation d'une réponse spécifique à un groupe (hôpital).

Signalons qu'une comparaison, d'un point de vue statistique, des avantages liés à l'introduction d'un effet aléatoire plutôt que fixe a été entreprise en ce qui concerne le modèle de Cox (Andersen *et al.*, 1999).

---

1. Le terme de « corrélation » est employé (abusivement) pour signifier « non-indépendance ».

## Constitution et perspectives de la thèse

La thèse se compose de six parties principales.

Les deux premières sont constituées par le recensement bibliographique réalisé sur les données de survie. Les différentes notations qui se présentaient dans les ouvrages de référence et les articles ont été harmonisées ; certains résultats, qui paraissaient essentiels à la compréhension de la théorie des données de survie, ont été repris *in extenso* dans ce travail (leur énoncé étant suivi d'une démonstration). Plus précisément, la première partie expose les définitions, les notations et les résultats concernant la théorie des données de survie et le modèle de Cox ; la deuxième partie développe les extensions du modèle de Cox au cas de données de survie corrélées.

La troisième partie présente deux études comparatives des différents modèles de Cox en présence de corrélation. La première, menée au travers de simulations, permet de souligner le rôle joué par différents paramètres (taux de censure, taille des groupes constituant l'échantillon étudié, variabilité de l'« effet groupe ») et de répondre à la question qui se pose couramment à l'épidémiologiste : quel type de modèle pour quelle(s) perspective(s) ? La seconde consiste en l'étude épidémiologique des gastro-entérites néonatales (GENN) chez le veau. Cette étude fait suite aux travaux doctoraux entrepris en 1999, au sein de l'École Nationale Vétérinaire d'Alfort, par M. Bendali.

La quatrième partie vient conclure le travail entrepris dans les trois précédentes.

La cinquième partie recense les références bibliographiques rencontrées tout au long de la thèse.

La sixième partie, quant à elle, est constituée des différentes annexes : la théorie statistique des processus de comptage, des résultats complémentaires concernant certaines assertions rencontrées au cours de la thèse, les tableaux des résultats des simulations et le programme de ces simulations en langage S (langage du logiciel S-Plus) sont présentés dans cette partie.

Enfin, un index est à la disposition du lecteur à la toute fin du document.

Notons que l'apprentissage de la théorie mathématique des données de survie peut se faire suivant deux voies distinctes. La première consiste en une approche qui peut être qualifiée d'« heuristique » : le modèle de Cox, tel qu'il fût exposé en 1972, ne reposait pas encore sur une théorie mathématique pouvant en justifier la conception.

La seconde confère au modèle de Cox une pleine justification mathématique, basée sur la théorie des processus de comptage : cette approche, qui présente l'avantage d'une grande rigueur mathématique, peut en revanche apparaître – du moins au néophyte en matière de processus aléatoires – absconse. La complexité des notions mathématiques auquel le modèle de Cox fait appel peut cependant être déjouée. Aussi ne devrait-elle pas rebuter un épidémiologiste dont les connaissances en matière de processus aléatoires pourraient n'être que parcellaires.

C'est pourquoi nous avons opté pour la seconde approche : en annexe sont présentées les bases relatives à la théorie mathématique des processus de comptage. Cette annexe, alliée à la première partie – qui expose de façon précise les notions utilisées en analyse des données de survie – devrait fournir tous les éléments de compréhension de la statistique appliquée à la survie. Nous espérons avoir ainsi rendu ce travail accessible, tant au mathématicien qui, retrouvant nombre d'outils statistiques, découvrira les perspectives qu'ouvrent ces outils dans le domaine de la santé, qu'à l'épidémiologiste, qui complètera ainsi ses connaissances en statistique, après avoir acquis quelques notions mathématiques certes complexes, mais indispensables à la bonne compréhension de l'analyse de la survie.

## Notations

Nous donnons ci-dessous les conventions d'écriture adoptées :

– Abréviations et symboles :

AIC : Critère d'Information d'Akaike (*Akaike Information Criterion*)

cadlag : continu à droite avec une limite à gauche

v.a. : variable aléatoire

p.s. : presque sûrement

$\mathbb{P}$  : probabilité

$\mathbb{E}$  : espérance

$\mathbb{V}$  : variance

Cov : covariance

$\mathbb{E}(X | Y)$  : espérance conditionnelle de  $X$  sachant  $Y$

$\mathcal{F}_t$  : filtration à l'instant  $t$

$D([a, b])$  : espace des fonctions cadlag sur  $[a, b]$

$x \wedge y$  :  $\inf(x, y)$

$\mathbb{1}_{\{A\}}$  : fonction indicatrice de l'événement  $A$

$\mathbf{x}_+ = \sum_{i=1}^n x_i$  pour tout vecteur  $(x_1, \dots, x_n)$

$\mathbf{J}(\hat{\beta})$  : matrice d'information de Fisher relative au paramètre  $\beta$

$\mathcal{L}_X(s) = \mathbb{E}[\exp(-sX)]$  : transformée de Laplace de  $X$

$\mathcal{P}_{]s, t]}$  ou  $\overset{t}{\mathcal{P}}_s$  : produit-intégrale (ou produit infini) sur  $]s, t]$

$(x_1, \dots, x_n)^t$  : vecteur-colonne transposé

$\otimes$  : produit direct

$\mathbf{a}^{\otimes 2} = \mathbf{a}\mathbf{a}^t$  pour tout vecteur colonne  $\mathbf{a}$

$\xrightarrow{\mathbb{P}}$  : convergence en probabilité

$\xrightarrow{\mathcal{L}}$  : convergence en loi

$X \rightsquigarrow \chi_n^2$  :  $X$  suit une loi du chi-deux à  $n$  degrés de libertés

i.i.d. : identiquement et indépendamment distribuées

– Typographie des mathématiques :

– caractères italiques réguliers : scalaires,

– caractères italiques gras : vecteurs et matrices,

– caractères italiques minuscules : variables non aléatoires,

– caractères italiques majuscules : variables aléatoires,

– caractères grecs : paramètres.

– Numérotation des équations : elle comprend le numéro du chapitre courant, suivi du numéro de l'équation au sein de ce chapitre.



# Chapitre 1

## DÉFINITIONS, NOTATIONS ET PROPRIÉTÉS

### Contenu

---

<b>1.1 Fonctions de survie et de risque</b> . . . . .	<b>14</b>
<b>1.2 Censure</b> . . . . .	<b>15</b>
1.2.1 Définition . . . . .	15
1.2.2 Caractéristiques . . . . .	16
<b>1.3 Constitution des observations et processus de comptage</b> . . . . .	<b>17</b>
<b>1.4 Estimation du risque cumulé</b> . . . . .	<b>19</b>
1.4.1 Estimateur . . . . .	19
1.4.2 Loi asymptotique . . . . .	20
<b>1.5 Estimation de la fonction de survie</b> . . . . .	<b>21</b>
1.5.1 Estimateurs . . . . .	21
1.5.2 Loi asymptotique . . . . .	23
<b>1.6 Comparaison de la survie de plusieurs groupes</b> . . . . .	<b>25</b>

---

## 1.1 Fonctions de survie et de risque

Dans ce qui suit, nous considérons un individu susceptible de subir une fois et une seule un certain type d'événement. L'observation de la survenue – *donnée non censurée* – ou non – *donnée censurée* – de cet événement chez l'individu constitue la donnée basale pour une modélisation de la survie.

Nous donnons ci-dessous les définitions des principaux outils utilisés en analyse de la survie. Pour chacun d'eux, nous précisons sa signification statistique d'une part, son interprétation en épidémiologie d'autre part.

**Définition I.1** — La « *durée de vie* » d'un individu est une variable aléatoire (v.a.)  $X$  positive et continue. Sa fonction de répartition

$$F(x) = \mathbb{P}(X \leq x)$$

est la probabilité que l'événement se produise entre 0 et  $x$ .  
Par la suite,  $F$  sera supposée dérivable.

**Définition I.2** — La *fonction de survie* est définie par

$$\begin{aligned} S(x) &= 1 - F(x) \\ &= \mathbb{P}(X > x) . \end{aligned}$$

Remarquons que la théorie de la survie ayant son origine dans l'observation et le décompte de décès, le vocabulaire est resté marqué par les termes de « durée de vie », « décès », « exclu vivant »... Cependant, cette théorie s'applique à divers types d'observations : la « durée de vie » peut ainsi être l'âge d'apparition d'une maladie, un délai de séroconversion, le temps de sortie du chômage, etc.

**Définition I.3** — La *fonction de risque instantané* est la fonction

$$\alpha(x) = \frac{f(x)}{S(x)}$$

où  $f$  est la densité de probabilité de  $X$ .

Notons que cette fonction n'est pas une densité de probabilité.

Le risque instantané est la probabilité que l'événement se produise en  $x$  (sachant qu'il ne s'est pas produit auparavant) :

$$\begin{aligned} \alpha(x) dx &= \frac{\mathbb{P}(x < X \leq x + dx)}{\mathbb{P}(X > x)} \\ &= \mathbb{P}(X \in ]x, x + dx] \mid X > x) . \end{aligned} \tag{1.1}$$

**Définition I.4** — La *fonction de risque cumulé* est donnée par

$$A(x) = \int_0^x \alpha(u) du .$$

Nous pouvons maintenant énoncer les deux résultats suivants :

**Proposition I.1** — *La définition de la distribution de probabilité de  $X$  repose sur l'une des quatre données suivantes, qui sont équivalentes :  $S(x)$ ,  $f(x)$ ,  $\alpha(x)$  et  $A(x)$ .*

**Proposition I.2** — *Nous avons*

$$\begin{aligned} S(x) &= \mathbb{P}(X > x) \\ &= \mathcal{P}_{]0,x]} [1 - \alpha(s) ds] \\ &= \exp \left( - \int_0^x \alpha(s) ds \right) \\ &= \exp \left( - A(x) \right) . \end{aligned}$$

$\mathcal{P}_{]0,x]}$  désigne le **produit infini** (ou **produit-intégrale**). Cette notion est définie en page 155.

## 1.2 Censure

### 1.2.1 Définition

**Définition I.5** — *La **variable de censure**  $C$  est définie par la possible non-observation de l'événement. Si l'on observe  $C$ , et non  $X$ , et que l'on sait que  $X > C$  (respectivement  $X < C$ ,  $C_1 < X < C_2$ ), on dit qu'il y a **censure à droite** (respectivement **censure à gauche**, **censure par intervalle**).*

Si l'événement se produit,  $X$  est « réalisée ». S'il ne se produit pas (l'individu étant perdu de vue, ou bien exclu vivant), c'est  $C$  qui est « réalisée ».

$X$  peut être considérée comme la durée séparant un événement initial A d'un événement terminal B, ou comme la durée pendant laquelle un sujet reste dans un état donné (auquel cas A désigne l'entrée dans cet état et B la sortie de cet état – par exemple le chômage). La censure à droite, dont il sera essentiellement question par la suite, est due à la non-observation de B, dont on sait seulement qu'il sera postérieur à la dernière date d'observation du sujet.

Par ailleurs, la censure se distingue de la **troncature** : on dit qu'il y a troncature à droite (respectivement à gauche) lorsque la variable d'intérêt  $X_i$  (durée de vie du  $i^e$  individu) n'est pas observable quand elle est supérieure (respectivement inférieure) à un seuil  $c > 0$  fixé.

Dans le cas de la censure, on sait que la variable  $X$  non observée est supérieure ou inférieure à une valeur  $C$  qui, elle, a été observée. La troncature, quant à elle, élimine de l'étude une partie des  $X_i$ , ce qui a pour conséquence de faire porter l'analyse uniquement sur la loi de  $X$  conditionnellement à l'événement  $\{X < c\}$  (respectivement  $\{X > c\}$ ).

Enfin, le mécanisme de censure est habituellement supposé être indépendant de l'événement étudié : on parle de censure **non-informative** (*ignorable*). En pratique, cela veut dire que les individus ne doivent pas être censurés parce qu'ils ont un risque de décès particulièrement

élevé (ou faible). En d'autres termes, les individus exclus-vivants ou perdus de vue à une date  $t$  doivent être représentatifs des individus encore à risque à cet instant  $t$ .

Si la censure est **informative**, alors l'expression classique de la vraisemblance ne correspond plus à une vraisemblance complète, mais à une vraisemblance partielle qui peut être utilisée pour des inférences, bien qu'il y ait une perte d'efficacité des estimateurs produits (car toute l'information n'est pas utilisée). Ainsi, la censure informative est à l'origine d'un biais lors de l'analyse standard basée sur la vraisemblance (Kalbfleisch et Prentice, 1980 ; Schluchter, 1992).

### 1.2.2 Caractéristiques

**Définition I.6** — *La censure est dite **non-aléatoire de type I** si, étant donné un nombre positif fixé  $c$  et un  $n$ -échantillon  $X_1, \dots, X_n$ , les observations consistent en  $(T_i, \delta_i)$ , où*

$$\begin{cases} T_i &= X_i \wedge c \\ \delta_i &= \mathbb{1}_{\{X_i \leq c\}}. \end{cases}$$

Exemples : test de l'efficacité d'une molécule sur un lot de souris, les souris survivantes étant sacrifiées au bout d'un temps déterminé  $c$  ; observation de la durée de fonctionnement de  $n$  machines au cours d'une expérience de durée  $c$ .

**Remarque** — Bien que similaires dans l'écriture de leur définition, la censure non-aléatoire de type I à droite et la troncature à droite doivent être distinguées : en effet, l'inférence statistique diffère grandement, selon qu'elle s'applique à l'un ou l'autre de ces deux types de données de survie. Ainsi, si nous considérons  $n$  observations indépendantes et censurées à droite, la vraisemblance retenue lors de l'étude statistique sera le produit d'un nombre aléatoire (inférieur ou égal à  $n$ ) de facteurs. En revanche, si nous considérons maintenant  $n$  observations indépendantes et tronquées à droite, nous étudierons une vraisemblance qui sera le produit d'un nombre fixe (exactement  $n$ ) de facteurs (cf. Andersen *et al.* (1991), p. 166-167 ou Drosbeke *et al.* (1989) p. 54).

**Définition I.7** — *La censure est dite **aléatoire de type I** si, étant donné un  $n$ -échantillon  $X_1, \dots, X_n$ , il existe une v.a.  $n$ -dimensionnelle  $(C_1, \dots, C_n)$  de  $(\mathbb{R}^+)^n$  telle que les observations consistent en  $(T_i, \delta_i)$ , où*

$$\begin{cases} T_i &= X_i \wedge C_i \\ \delta_i &= \mathbb{1}_{\{X_i \leq C_i\}}. \end{cases}$$

Exemple : lors d'une expérience biologique, on s'intéresse à une cause de décès qui a lieu au bout d'un temps  $X$ , et l'on désire étudier la loi de  $X$  ; cependant, une autre cause de décès peut intervenir auparavant, et donc empêcher l'observation de  $X$  par un mécanisme de censure à droite.

**Définition I.8** — *La censure est dite **de type II** si, étant donné un nombre positif fixé  $r$  et un  $n$ -échantillon  $X_1, \dots, X_n$ , les observations consistent en  $(T_i, \delta_i)$ , où*

$$\begin{cases} T_i &= X_i \wedge X_{(r)} \\ \delta_i &= \mathbb{1}_{\{X_i \leq X_{(r)}\}} \end{cases}$$

où  $X_{(1)} < X_{(2)} < \dots < X_{(n)}$  sont les statistiques d'ordre.

Exemples : test de l'efficacité d'une molécule sur un lot de souris, la durée de l'étude correspondant au temps que mettent  $r$  souris à mourir ; observation de la durée de fonctionnement de  $n$  machines tant que  $r$  d'entre elles ne tombent pas en panne.

Désormais, nous nous placerons dans le cadre d'un **mécanisme de censure à droite, aléatoire de type I et indépendante du mécanisme de survenue de l'événement (censure non-informative)**.

### 1.3 Constitution des observations et processus de comptage

**Définition I.9** — Une observation consiste en  $(T_i, \delta_i)$ , où

$$\begin{cases} T_i &= X_i \wedge C_i \\ \delta_i &= \mathbb{1}_{\{X_i \leq C_i\}}. \end{cases}$$

Nous définissons le **processus de comptage** (de survenue de l'événement) par

$$N_i(t) = \mathbb{1}_{\{T_i \leq t \text{ et } \delta_i=1\}}.$$

Si l'individu  $i$  subit l'événement avant l'instant  $t$ , alors  $N_i(t) = 1$  ; sinon,  $N_i(t) = 0$ .

Nous définissons également la fonction

$$Y_i(t) = \mathbb{1}_{\{T_i \geq t\}}$$

qui indique si l'individu  $i$  est encore « à risque » (*i.e.* n'a pas encore subi l'événement) juste avant l'instant  $t$ .

Nous avons donc (cf. (1.1))

$$\mathbb{P}[dN_i(t) = 1 \mid \mathcal{F}_{t-}] = \alpha_i(t)Y_i(t) dt, \quad (1.2)$$

où  $\mathcal{F}_t$  est la filtration naturelle (toute l'information disponible à l'instant  $t$ ), et où la notation  $dN_i(t)$  renvoie à l'écriture formelle de l'intégrale stochastique

$$\begin{aligned} N_i(t) &= \int_0^t dN_i(s) \\ &= \int_0^t N_i(ds), \end{aligned}$$

écriture rendue possible du fait que  $N_i(t)$  est un processus croissant (cf. VI.1 et VI.5 p. 152).

**Proposition I.3** — Le processus stochastique défini par

$$M_i(t) = N_i(t) - \int_0^t \alpha_i(s)Y_i(s) ds \quad (1.3)$$

est une martingale.

*Démonstration* — Nous avons

$$\begin{aligned} \mathbb{E}[dM_i(t) \mid \mathcal{F}_{t-}] &= \mathbb{E}[dN_i(t) - \alpha_i(t)Y_i(t) dt \mid \mathcal{F}_{t-}] \\ &= \mathbb{P}[dN_i(t) = 1 \mid \mathcal{F}_{t-}] - \alpha_i(t)Y_i(t) dt \\ &= 0 \end{aligned}$$

d'après (1.2). ■

**Définition I.10** — *Les processus*

$$\lambda_i(t) = \alpha_i(t)Y_i(t)$$

et

$$\Lambda_i(t) = \int_0^t \alpha_i(s)Y_i(s) ds$$

sont désignés respectivement sous les termes de **processus d'intensité** et **processus d'intensité cumulée** de  $N_i$  (cf. p. 152).

$\Lambda_i(t)$  est également appelé **compensateur** du processus  $N_i(t)$ .

Notons que  $\lambda_i(t)$  est une variable aléatoire, au contraire de  $\alpha_i(t)$  qui est fixe.

Il est possible d'exprimer  $\lambda_i(t)$  en fonction du processus de comptage  $N_i(t)$  :

$$\lambda_i(t) = \lim_{dt \rightarrow 0} \frac{1}{dt} \left[ \mathbb{P}[N_i(t+dt) - N_i(t) = 1] \mid \mathcal{F}_{t-} \right].$$

**Remarque** — La théorie statistique des processus ponctuels porte sur des processus de comptage simple de la forme

$$N_i^*(t) = \mathbb{1}_{\{T_i \leq t\}}.$$

Le processus que nous considérerons tout au long de ce travail est un processus de comptage « filtré ». Ce processus, que nous avons noté  $N$ , se définit de la sorte :

$$\begin{aligned} N_i(t) &= \int_0^t \mathbb{1}_{\{s \leq C_i\}} dN_i^*(s) \\ &= \int_0^t \mathbb{1}_{\{s \leq C_i\}} d\Lambda_i^*(s) + \int_0^t \mathbb{1}_{\{s \leq C_i\}} dM_i^*(s) \\ &= \Lambda_i(t) + M_i(t), \end{aligned}$$

où  $\Lambda^*$  et  $M^*$  sont respectivement le compensateur et la martingale associés au processus  $N^*$ .

Notons que, relativement à la filtration naturelle, les compensateurs de  $N$  et  $N^*$  coïncident :  $\Lambda^* = \Lambda$ . Ce dernier résultat, ainsi que d'autres éléments mathématiques relatifs aux processus filtrés, sont présentés en annexe page 156.

## 1.4 Estimation du risque cumulé

### 1.4.1 Estimateur

Considérons un échantillon de  $n$  individus. Soit  $\tau$  la date de point (cf. glossaire p. 197). Soient également, pour  $i \in \{1, \dots, n\}$ ,

- $X_i$  la date de survenue de l'événement chez l'individu  $i$  ;
- $C_i$  la date de censure correspondante ;
- $T_i = X_i \wedge C_i$  ;
- $\delta_i = \mathbb{1}_{\{X_i \leq C_i\}}$  ;
- $Y_i(t) = \mathbb{1}_{\{T_i \geq t\}}$ .

Enfin, posons

- $\mathbf{N}(t) = \{N_i(t), 0 \leq t \leq \tau, i = 1, \dots, n\}$  le processus de comptage multivarié pour les  $n$  individus ;
- $\mathbf{\Lambda}(t) = \{\Lambda_i(t), 0 \leq t \leq \tau, i = 1, \dots, n\}$  le compensateur de  $\mathbf{N}(t)$  par rapport à la filtration  $\mathcal{F}_t$ .

Définissons

$$N_+(t) = \sum_{i=1}^n N_i(t)$$

et

$$Y_+(t) = \sum_{i=1}^n Y_i(t)$$

comme étant, respectivement, le nombre total de survenues de l'événement à l'instant  $t$  et le nombre total d'individus encore à risque juste avant l'instant  $t$ .

Pour les démonstrations des résultats de cette section, ainsi que de la suivante, nous renvoyons le lecteur aux ouvrages de référence (Andersen *et al.*, 1991 ; Dacunha-Castelle et Duflo, 1993).

**Définition I.11** — *L'estimateur de Nelson-Aalen de la fonction de risque cumulé est défini par (Nelson, 1972 ; Aalen, 1978)*

$$\widehat{A}(t) = \int_0^t \frac{J(u)}{Y_+(u)} dN_+(u),$$

où  $J(t) = \mathbb{1}_{\{Y_+(t) > 0\}}$ .

L'origine « naturelle » de cet estimateur provient de l'équation (1.3) qui peut être réécrite sous la forme

$$dN_i(t) = \alpha_i(t)Y_i(t) dt + dM_i(t)$$

où  $dM_i(t)$  peut être considéré comme étant un bruit aléatoire.

### 1.4.2 Loi asymptotique

**Théorème I.1** —  $\widehat{A}(t)$  est un estimateur biaisé de  $A(t)$  et, sous l'hypothèse que  $F(t) < 1$  (c.-à-d. que  $A(t) < \infty$ ), nous avons

$$\sqrt{n}(\widehat{A}(t) - A(t)) \xrightarrow{\mathcal{L}} U(t)$$

avec  $U$  martingale gaussienne telle que

$$\begin{cases} U(0) = 0 \\ \mathbb{V}(U(t)) = \int_0^t \frac{\alpha(u)}{y(u)} du, \end{cases}$$

où

$$y(s) = [1 - F(s)][1 - G(s-)]$$

avec  $G$  fonction de répartition de  $C$ .

*Démonstration* — Concernant le biais :

$$\begin{aligned} \mathbb{E}[\widehat{A}(t)] &= \mathbb{E}\left[\int_0^t \frac{J(u)}{Y_+(u)} dN_+(u)\right] \\ &= \mathbb{E}\left[\int_0^t \frac{J(u)}{Y_+(u)} \left\{ dM_+(u) + Y_+(u) \alpha(u) du \right\}\right] \\ &= \mathbb{E}\left[\int_0^t \frac{J(u)}{Y_+(u)} dM_+(u)\right] + \mathbb{E}\left[\int_0^t J(u) \alpha(u) du\right] \\ &= 0 + \int_0^t \mathbb{E}[J(u)] \alpha(u) du \quad (\text{car } M_+(t) \text{ est une martingale}) \\ &= \int_0^t \mathbb{P}(Y_+(u) > 0) \alpha(u) du \\ &= \int_0^t \alpha(u) du - \int_0^t \mathbb{P}(Y_+(u) = 0) \alpha(u) du \\ &= A(t) - \int_0^t \mathbb{P}(Y_+(u) = 0) \alpha(u) du. \end{aligned}$$

Pour expliciter le caractère asymptotique des résultats, nous exprimons

$$\begin{aligned} N^{(n)}(t) &= \sum_{i=1}^n N_i(t) \\ Y^{(n)}(t) &= \sum_{i=1}^n Y_i(t) \end{aligned}$$

et

$$J^{(n)}(t) = \mathbb{1}_{\{Y^{(n)}(t) > 0\}}.$$

Notons aussi  $F$  la fonction de répartition des  $X_i$  et  $G$  celle des  $C_i$  ( $i \in \{1, \dots, n\}$ ). Nous obtenons que la fonction de répartition des  $T_i$  est  $[1 - (1 - F)][1 - G]$ . D'après le théorème de Glivenko-Cantelli (cf. p. 160),

$$\sup_{s \in [0, t]} \left| \frac{Y^{(n)}}{n} - [1 - F(s)][1 - G(s-)] \right| \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty). \quad (1.4)$$

Par ailleurs,

$$J^{(n)}(t) = \mathbb{1}_{\{Y^{(n)}(t) > 0\}}.$$

Nous en déduisons que

$$\begin{aligned} 1 - J^{(n)}(t) &= \mathbb{1}_{\{Y^{(n)}(t) = 0\}} \\ &= \mathbb{1}_{\{B(n, [1 - F(t)][1 - G(t-)]) = 0\}} \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty), \end{aligned}$$

et, par suite, que

$$J^{(n)}(t) \xrightarrow{\mathbb{P}} 1 \quad (n \rightarrow \infty).$$

Par conséquent,

$$\left\langle \sqrt{n} [\widehat{A}^{(n)} - A] \right\rangle (t) = \int_0^t n \frac{J^{(n)}(u)}{Y^{(n)}(u)} \alpha(u) du \xrightarrow{\mathbb{P}} \int_0^t \frac{\alpha(s)}{[1 - F(s)][1 - G(s-)]} ds$$

qui est déterministe.

L'application du théorème de Rebolledo (cf. p. 154) donne le résultat. ■

Remarquons que le biais de cet estimateur est extrêmement faible puisqu'en pratique, la probabilité qu'à un instant  $t$  tous les individus aient, soit subi l'événement, soit été censurés, est proche de zéro.

## 1.5 Estimation de la fonction de survie

### 1.5.1 Estimateurs

**Définition I.12** — *L'estimateur de Kaplan-Meier de la fonction de survie est défini par (Kaplan et Meier, 1958)*

$$\begin{aligned} \widehat{S}(t) &= \mathcal{P}_{s \leq t} (1 - d\widehat{A}(s)) \\ &= \mathcal{P}_{s \leq t} \left( 1 - \frac{J(s)}{Y_+(s)} dN_+(s) \right) \end{aligned}$$

ou encore

$$\begin{aligned}\widehat{S}(t) &= \prod_{s \leq t} (1 - \Delta \widehat{A}(s)) \\ &= \prod_{s \leq t} \left(1 - \frac{J(s) \Delta N_+(s)}{Y_+(s)}\right),\end{aligned}$$

où  $\widehat{A}(t)$  est l'estimateur de Nelson-Aalen et où, pour un processus  $X(t)$  cadlag (continu à droite avec une limite à gauche),

$$\Delta X(t) = X(t) - X(t-).$$

**Remarque** — L'estimateur de Kaplan-Meier repose sur l'idée intuitive qu'être encore en vie après l'instant  $t$ , c'est être en vie juste avant  $t$  et ne pas mourir en  $t$ . Cette idée, traduite en termes probabilistes, mène à la relation

$$\begin{aligned}S(t) &= \mathbb{P}(T \geq t) \\ &= \mathbb{P}(T \geq t \mid T \geq t-1) \mathbb{P}(T \geq t-1) \\ &= \dots \\ &= \mathbb{P}(T \geq t \mid T \geq t-1) \dots \mathbb{P}(T \geq 1 \mid T \geq 0) \mathbb{P}(T \geq 0).\end{aligned}$$

**Proposition I.4** — Un premier estimateur de la variance de  $\widehat{S}(t)/S(t)$  est

$$\widehat{\sigma}^2(t) = \int_0^t \frac{J(s)}{Y_+(s)^2} dN_+(s).$$

*Démonstration* — La variance de  $\widehat{S}(t)/S(t)$  est approchée par celle de  $\widehat{S}(t)/S^*(t)$ .

$$\begin{aligned}\mathbb{V} \left[ \frac{\widehat{S}(t)}{S^*(t)} - 1 \right] &= \mathbb{E} \left[ \left\langle \frac{\widehat{S}}{S^*} - 1 \right\rangle (t) \right] \\ &= \int_0^t \left\{ \frac{\widehat{S}(s-)}{S^*(s)} \right\}^2 \frac{J(s)}{Y_+(s)} \alpha(s) ds,\end{aligned}\tag{1.5}$$

d'où le résultat de la proposition, obtenu en remplaçant d'une part  $\widehat{S}(s-)$  et  $S^*(s)$  par  $\widehat{S}(s)$  (grâce à la continuité de  $S$ ), d'autre part  $dA(s) = \alpha(s) ds$  par  $d\widehat{A}(s)$ . ■

L'estimateur de la variance de  $\widehat{S}(t)$  correspondant est donné par

$$\widehat{V}(\widehat{S}(t)) = [\widehat{S}(t)]^2 \widehat{\sigma}^2(t).$$

**Proposition I.5** — Un second estimateur de la variance de  $\widehat{S}(t)/S(t)$  – l'estimateur de **Greenwood**, défini à l'origine dans le cadre de l'estimation de la fonction de survie par la méthode actuarielle<sup>1</sup> – est donné par (Greenwood, 1926)

$$\widehat{\widehat{\sigma}}^2(t) = \int_0^t \frac{dN_+(s)}{Y_+(s)[Y_+(s) - \Delta N_+(s)]}.$$

---

1. Le principe de cette méthode, due à Böhmer (1912), est similaire à celui de la méthode de Kaplan-Meier, la principale différence venant du fait que les intervalles de temps sur lesquels les probabilités conditionnelles sont calculées, sont fixés *a priori*, et non plus déterminés par les dates des décès observés comme le proposent Kaplan et Meier.

*Démonstration* — Réécrivons (1.5) sous la forme

$$\mathbb{V} \left[ \frac{\widehat{S}(t)}{S^*(t)} - 1 \right] = \int_0^t \left\{ \frac{\widehat{S}(s-)}{S^*(s)} \right\}^2 \frac{J(s)}{Y_+(s)} (1 - \Delta A(s)) dA(s) .$$

En remplaçant  $S^*$  par  $\widehat{S}$  et  $A$  par  $\widehat{A}$ , et en notant que

$$\widehat{S}(t) = \left( 1 - \frac{\Delta N(s)}{Y(s)} \right) \widehat{S}(s-),$$

nous obtenons le résultat. ■

L'estimateur de la variance de  $\widehat{S}(t)$  correspondant est donné par

$$\widehat{\mathbb{V}}[\widehat{S}(t)] = [\widehat{S}(t)]^2 \widehat{\sigma}^2(t) .$$

### 1.5.2 Loi asymptotique

**Théorème I.2** — *L'estimateur  $\widehat{S}(t)$  est biaisé; de plus, si nous supposons que*  
*A. pour tout  $s \in [0, t]$ ,*

$$n \int_0^s \frac{J(u)}{Y(u)} \alpha(u) du \xrightarrow{\mathbb{P}} \sigma^2(s) \quad (n \rightarrow \infty),$$

*B. pour tout  $\epsilon > 0$ ,*

$$n \int_0^t \frac{J(s)}{Y(s)} \alpha(s) \mathbb{1}_{\{\sqrt{n}|J(s)/Y(s)| > \epsilon\}} ds \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty),$$

*C.*

$$\sqrt{n} \int_0^t (1 - J(u)) \alpha(u) du \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty),$$

alors cet estimateur vérifie asymptotiquement

$$\sqrt{n}(\widehat{S}(t) - S(t)) \xrightarrow{\mathcal{L}} -U(t).S(t),$$

où  $U$  est la martingale gaussienne définie en page 20.

*Démonstration* — Concernant le biais, posons

$$\begin{aligned} S^*(t) &= \mathcal{P}_{s \leq t} (1 - dA^*(s)) \\ &= \exp(-A^*(t)) \end{aligned}$$

où

$$A^*(t) = \int_0^t J(u) \alpha(u) du .$$

Nous avons (cf. proposition VI.3 p. 155)

$$\begin{aligned} \frac{\widehat{S}(t)}{S^*(t)} - 1 &= - \int_0^t \frac{\widehat{S}(s-)}{S^*(s)} d(\widehat{A} - A^*)(s) \\ &= - \int_0^\tau \frac{\widehat{S}(s-)J(s)}{S^*(s)Y_+(s)} dM_+(s) \end{aligned} \quad (1.6)$$

pour  $t \in [0, \tau)$ , et avec  $dM_+(t) = dN_+(t) - Y_+(t)\alpha(t) dt$ .

Par suite,  $\widehat{S}/S^* - 1$  étant une martingale localement de carré intégrable sur  $[0, \tau)$ , nous avons

$$\mathbb{E} \left[ \frac{\widehat{S}(t)}{S^*(t)} \right] = 1$$

pour tout  $t \in [0, \tau)$ .

Enfin,

$$\begin{aligned} S(t) &\leq S^* \\ \Rightarrow \frac{\widehat{S}(t)}{S^*(t)} &\leq \frac{\widehat{S}(t)}{S(t)} \\ \Rightarrow \mathbb{E} \left[ \frac{\widehat{S}(t)}{S^*(t)} \right] &\leq \mathbb{E} \left[ \frac{\widehat{S}(t)}{S(t)} \right] \\ \Leftrightarrow 1 &\leq \frac{\mathbb{E}(\widehat{S}(t))}{S(t)} \\ \Leftrightarrow \mathbb{E}(\widehat{S}(t)) &\geq S(t). \end{aligned}$$

Concernant la convergence asymptotique, d'après (1.6),

$$\sqrt{n} \left\{ \frac{\widehat{S}(t)}{S^*(t)} - 1 \right\} = -\sqrt{n} \int_0^t \frac{\widehat{S}(s-)J(s)}{S^*(s)Y(s)} dM(s).$$

D'après les conditions A et B, et d'après le fait que  $\widehat{S}(s)/S^*(s) \leq 1/S(t)$  pour tout  $s \in [0, t]$ , nous déduisons du théorème de Rebolledo que

$$\sqrt{n} \left( \frac{\widehat{S}(t)}{S^*(t)} - 1 \right) \xrightarrow{\mathcal{L}} -U(t)$$

quand  $n \rightarrow \infty$ .

La condition C et le fait que

$$\begin{aligned} \sqrt{n} \left| \frac{S(s)}{S^*(s)} - 1 \right| &= \int_0^s \frac{S(u)}{S^*(u)} d(A - A^*)(u) \\ &\leq \frac{1}{S(t)} (1 - J(u)) \alpha(u) du \end{aligned}$$

entraînent que

$$\sup_{s \in [0, t]} \sqrt{n} \left| \frac{S(s)}{S^*(s)} - 1 \right| \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty).$$

D'où il s'ensuit que

$$\sqrt{n} \frac{\hat{S}(s) - S(s)}{S^*(s)} \xrightarrow{\mathcal{L}} -U \quad (n \rightarrow \infty)$$

et nous obtenons le résultat du théorème. ■

**Remarque** — Dans tout ce qui précède, nous avons supposé **distincts** les temps de survenue de l'événement. Cette hypothèse découle de la définition même d'un **processus de comptage multivarié** – définition énoncée en page 153.

En effet, pour que  $N(t) = \{N_i(t), 0 \leq t \leq \tau, i = 1, \dots, n\}$  soit un processus de comptage multivarié, il est nécessaire d'exclure la possibilité de sauts simultanés pour deux (ou plus) de ses composantes.

## 1.6 Comparaison de la survie de plusieurs groupes

Il existe une grande variété de tests portant sur la fonction de survie, en particulier des tests d'égalité des distributions de survie entre plusieurs groupes.

Considérons un processus de comptage multivarié  $\mathbf{N} = (N_1, \dots, N_k)$  ( $k \geq 2$ ) de processus d'intensité  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_k)$  de la forme  $\lambda_i(t) = Y_i(t) \alpha_i(t)$ .

L'hypothèse nulle à tester est

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_k = \alpha \tag{1.7}$$

La construction de la statistique repose sur la comparaison entre l'estimateur de Nelson-Aalen

$$\hat{A}_h(t) = \int_0^t \frac{J_h(s)}{Y_h(s)} dN_h(s)$$

(où  $J_h/Y_h$  est localement borné pour tout  $h$ ) et l'estimateur de l'hypothétique valeur commune

$$A(t) = \int_0^t \alpha(s) ds .$$

Cette dernière quantité peut être estimée par (cf. notations p. 19)

$$\hat{A}(t) = \int_0^t \frac{J_+(s)}{Y_+(s)} dN_+(s) .$$

Il est suffisant de comparer  $\hat{A}(t)$  et  $\hat{A}_h(t)$  sur l'ensemble des valeurs de  $t$  pour lesquelles  $Y_h(t) > 0$ . Ainsi, définissons

$$\begin{aligned} \tilde{A}_h(t) &= \int_0^t J_h(s) d\hat{A}(s) \\ &= \int_0^t \frac{J_h(s)}{Y_+(s)} dN_+(s) . \end{aligned}$$

Sous l'hypothèse nulle,

$$\hat{A}_h(t) - \tilde{A}_h(t) = \int_0^t \frac{J_h(s)}{Y_h(s)} dM_h(s) - \int_0^t \frac{J_h(s)}{Y_+(s)} dM_+(s) ,$$

qui est une martingale localement de carré intégrable.

Introduisons par ailleurs des processus de pondération prédictibles, positifs et localement bornés  $K_h$ , et définissons les processus

$$\Upsilon_h(t) = \int_0^t K_h(s) d(\hat{A}_h - \tilde{A}_h)(s)$$

pour  $h = 1, \dots, k$ .

Ces processus accumulent les différences (pondérées) entre les incréments de  $\hat{A}_h$  et  $\tilde{A}_h$ , et par suite ils fournissent une base pour la construction d'un test statistique.

Les processus choisis sont habituellement de la forme

$$K_h(t) = Y_h(t)K(t)$$

où  $K$  est un processus

- positif,
- localement borné,
- ne dépendant que du couple de processus  $(N_+, Y_+)$ ,
- nul lorsque  $Y_+$  est nul,
- tel que  $K/Y_+$  est considéré comme nul lorsque  $Y_+$  est nul.

Nous avons alors, pour  $h = 1, \dots, k$ ,

$$\Upsilon_h(t) = \int_0^t K(s) dN_h(s) - \int_0^t K(s) \frac{Y_h(s)}{Y_+(s)} dN_+(s) . \quad (1.8)$$

(Puisque  $J_h Y_h = Y_h$ , le facteur  $J_h$  peut être omis.)

Notons que (1.8) entraîne

$$\sum_{h=1}^k \Upsilon_h(t) = 0 .$$

Par la suite, nous nous restreignons au cas où le processus de pondération est de la forme

$$K_h(t) = Y_h(t)K(t) .$$

Cependant, les arguments avancés s'appliquent également au cas général. Sous l'hypothèse nulle,

$$\begin{aligned} \Upsilon_h(t) &= \int_0^t K(s) dM_h(s) - \int_0^t K(s) \frac{Y_h(s)}{Y_+(s)} dM_+(s) \\ &= \sum_{l=1}^k \int_0^t K(s) \left( \delta_{hl} - \frac{Y_h(s)}{Y_+(s)} \right) dM_l(s) , \end{aligned} \quad (1.9)$$

où  $\delta_{hl}$  est l'indice de Kronecker. Ce résultat découle du fait que les  $\Upsilon_h$  sont des martingales localement de carré intégrable. Leur processus de variation prévisible est, d'après l'équation (A.2) de la proposition VI.7 p. 154,

$$\begin{aligned} \langle \Upsilon_h(t), \Upsilon_j(t) \rangle &= \sum_{l=1}^k \int_0^t K^2(s) \left( \delta_{hl} - \frac{Y_h(s)}{Y_+(s)} \right) \left( \delta_{jl} - \frac{Y_j(s)}{Y_+(s)} \right) \alpha(s) Y_l(s) ds \\ &= \int_0^t K^2(s) \frac{Y_h(s)}{Y_+(s)} \left( \delta_{hj} - \frac{Y_j(s)}{Y_+(s)} \right) \alpha(s) Y_+(s) ds . \end{aligned} \quad (1.10)$$

D'après la proposition VI.6 p. 153, les martingales locales de (1.9) sont de carré intégrable sur  $[0, t]$ , à condition que  $\mathbb{E} \langle \Upsilon_h \rangle (t) < \infty$ .

Une condition suffisante de l'existence des moments du premier et du second ordre de  $\Upsilon_h(t)$  est, d'après (1.10), que

$$\int_0^t \mathbb{E} \{ K^2(s) Y_+(s) \} \alpha(s) ds < \infty .$$

Lorsque cette condition est vérifiée, nous avons, sous l'hypothèse nulle, que  $\mathbb{E} (\Upsilon_h(t)) = 0$  d'une part, et que

$$\text{Cov} (\Upsilon_h(t), \Upsilon_j(t)) = \mathbb{E} \langle Z_h, Z_j \rangle (t)$$

peut être estimé sans biais par

$$\hat{\sigma}_{hj}^2(t) = \int_0^t K^2(s) \frac{Y_h(s)}{Y_+(s)} \left( \delta_{hj} - \frac{Y_j(s)}{Y_+(s)} \right) dN_+(s) . \quad (1.11)$$

Ce résultat provient du fait que la différence entre (1.11) et (1.10) est une martingale localement de carré intégrable (sous  $H_0$ ). On note  $\widehat{\Sigma}(t)$  la matrice  $k \times k$  des éléments de (1.11).

Soit  $\Upsilon(t) = (\Upsilon_1(t), \dots, \Upsilon_k(t))$ . Une statistique de l'hypothèse nulle est la forme quadratique

$$\Upsilon(\tau) \widehat{\sigma}^t(\tau)^- \Upsilon(\tau) , \quad (1.12)$$

où  $\widehat{\sigma}^t(t)^-$  est la notation de l'inverse généralisée.

**Proposition I.6** — *La statistique  $\Upsilon(t) \widehat{\sigma}^t(t)^- \Upsilon(t)$  suit asymptotiquement un  $\chi_{k-1}^2$ .*

La démonstration de ce résultat se trouve en annexe, à la page 161.

**Test de Mantel-Haenszel ou du log-rank** Pour  $K(t) = \mathbb{1}_{\{Y_+(t) > 0\}}$ , nous obtenons le test du log-rank (Mantel, 1966 ; Peto et Peto, 1972 ; Cox, 1972).

**Test de Gehan et de Breslow** Pour  $K(t) = Y_+(t)$ , nous obtenons une généralisation des tests de Wilcoxon et de Kruskal-Wallis, qui sont appelés tests de Gehan (1965) et test de Breslow (1970).

**Tests de Taron et Ware** Tarone et Ware (1977) ont suggéré l'emploi d'une famille de statistiques qui englobe les tests du log-rank, de Gehan et de Breslow. Cette famille est obtenue en choisissant un processus de pondération de la forme  $K(t) = g(Y_+(t))$  pour une fonction  $g$  fixée. En particulier, ils proposent de poser  $g(y) = \sqrt{y}$ .

**Test de Prentice** Prentice (1978) (voir aussi Kalbfleisch et Prentice (1980), chap. 6) suggèrent une autre généralisation des tests de Wilcoxon et Kruskal-Wallis. Si l'on pose

$$\tilde{S} = \prod_{s \leq t} \left( 1 - \frac{\Delta N_+(s)}{Y_+(s) + 1} \right),$$

qui est proche de l'expression de l'estimateur de Kaplan-Meier, alors Prentice suggère d'employer<sup>1</sup>

$$K(t) = \frac{\tilde{S}(t-)Y_+(t)}{Y_+(t) + 1}.$$

Notons que la différence entre les généralisations de Gehan et Breslow d'une part, et celles de Peto et Prentice d'autre part, consiste principalement dans le fait que les premières emploient un processus de pondération ( $K = Y_+$ ) dépendant aussi bien des temps de survenue que des temps de censure, tandis que les secondes utilisent un processus de pondération ( $K \simeq \tilde{S}$ ) qui dépend uniquement de l'expérience de survie dans l'échantillon combiné.

Récemment, d'autres tests ont vu le jour : citons par exemple la famille de tests de Jones et Crowley (1989), qui généralise celle de Tarone et Ware, ou encore un test basé sur les permutations dû à Sun et Sherman (1996).

---

1. La justification d'une telle définition du processus  $K(t)$  – c.-à-d. un processus ne se réduisant pas à  $\tilde{S}(t)$  – réside dans la volonté de rendre prévisible ce processus.

# Chapitre 2

## MODÈLE DE COX

### Contenu

---

<b>2.1</b>	<b>Définitions et notations</b>	<b>30</b>
<b>2.2</b>	<b>Estimation</b>	<b>30</b>
2.2.1	Présentation	30
2.2.2	Résolution numérique	33
<b>2.3</b>	<b>Tests</b>	<b>34</b>
2.3.1	Test du rapport de vraisemblance	34
2.3.2	Test de Wald (ou du maximum de vraisemblance)	34
2.3.3	Test du score	34
<b>2.4</b>	<b>Critères d'adéquation au modèle de Cox</b>	<b>35</b>
2.4.1	Dans sa globalité	35
2.4.2	Concernant la forme fonctionnelle des covariables	35
2.4.3	Concernant la proportionnalité des risques vis-à-vis d'une covariable	36
2.4.4	Justesse du modèle pour chaque sujet	40
2.4.5	Concernant les « observations influentes »	41
<b>2.5</b>	<b>Considération des ex-æquo</b>	<b>42</b>
2.5.1	Vraisemblance partielle de Breslow	43
2.5.2	Vraisemblance partielle d'Efron	43
2.5.3	Vraisemblance partielle exacte	43
<b>2.6</b>	<b>Extensions du modèle</b>	<b>43</b>

---

## 2.1 Définitions et notations

Considérons un échantillon de  $n$  individus.

Soit  $\tau$  la date de point et  $\beta$  le paramètre de régression dont la dimension est égale à  $p$ .

Soient, pour  $i \in \{1, \dots, n\}$ ,

- $X_i$  la date de survenue de l'événement chez l'individu  $i$  ;
- $C_i$  la date de censure correspondant ;
- $T_i = X_i \wedge C_i$  ;
- $\delta_i = \mathbb{1}_{\{X_i \leq C_i\}}$  ;
- $\mathbf{Z}_i = (Z_{1i}, \dots, Z_{pi})$  le vecteur de dimension  $p$  des covariables ;
- $Y_i(t) = \mathbb{1}_{\{T_i \geq t\}}$ .

Soient aussi

- $\mathbf{N}(t) = \{N_i(t), 0 \leq t \leq \tau, i = 1, \dots, n\}$  le processus de comptage multivarié pour les  $n$  individus ;
- $\mathbf{\Lambda}(t) = \{\Lambda_i(t), 0 \leq t \leq \tau, i = 1, \dots, n\}$  le compensateur de  $N(t)$  par rapport à la filtration  $\mathcal{F}_t$ .

Le **modèle de Cox** (1972) spécifie que le risque instantané s'écrit

$$\alpha_i(t) = \alpha_0(t) \exp(\beta^t \mathbf{Z}_i),$$

où  $\alpha_0(t)$  est la **fonction de risque de base**.

Il s'agit d'un **modèle semi-paramétrique à risques proportionnels** :

- semi-paramétrique, du fait de la présence, dans la définition du risque instantané, d'une partie paramétrique (la partie de régression  $\exp(\beta^t \mathbf{Z}_i)$ ) et d'une partie non-paramétrique (le risque de base  $\alpha_0(t)$ ) ;
- à risques proportionnels, car quels que soient  $i$  et  $j$ , le rapport des risques instantanés de deux individus ne varie pas au cours du temps :

$$\frac{\lambda_i(t)}{\lambda_j(t)} = \exp[\beta^t (\mathbf{Z}_i - \mathbf{Z}_j)].$$

## 2.2 Estimation

### 2.2.1 Présentation

Nous présentons, dans l'annexe A, la partie de la théorie mathématique des processus de comptage sur laquelle vont reposer les résultats qui suivent. Ainsi, le théorème VI.6 nous assure que la vraisemblance complète associée à un processus ponctuel  $N^*$  simple – *i.e.* non filtré (cf. remarque p. 18), c.-à-d. non censuré – est de la forme

$$\mathcal{L}^*(\beta) = \prod_{t \leq \tau} \prod_{i=1}^n \left\{ (\alpha_i(t) Y_i(t))^{\Delta N_i^*(t)} \right\} \times \exp \left[ - \sum_i Y_i(\tau) A_i(\tau) \right]. \quad (2.1)$$

Nous passons de la vraisemblance associée au processus simple  $N^*$  – dite **vraisemblance complète** – à celle associée au processus censuré  $N$  – dite **vraisemblance partielle** – en supprimant, dans (2.1), les termes correspondant, pour l'intervalle de temps  $dt$ , à la contribution du processus indicateur prévisible de censure  $C$  (voir l'annexe A pour plus de détails concernant ce passage).

Ainsi, la vraisemblance partielle dans le cadre du modèle de Cox s'écrit

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &= \prod_{t \leq \tau} \prod_{i=1}^n \left\{ (\alpha_i(t) Y_i(t))^{\Delta N_i(t)} \right\} \times \exp \left[ - \sum_i Y_i(\tau) A_i(\tau) \right] \\ &= \prod_t \prod_i \left\{ [\alpha_0(t) Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)]^{\Delta N_i(t)} \right\} \times \exp \left[ - \int_0^\tau S^{(0)}(\boldsymbol{\beta}, u) \alpha_0(u) du \right] \quad (2.2)\end{aligned}$$

avec

$$S^{(0)}(\boldsymbol{\beta}, t) = \sum_j Y_j(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_j).$$

À  $\boldsymbol{\beta}$  fixé, la maximisation de (2.2) suivant  $\Delta A_0(t)$  conduit à

$$\Delta \hat{A}_0(t) = \frac{\Delta N_+(t)}{S^{(0)}(\boldsymbol{\beta}, t)}.$$

Par conséquent, toujours à  $\boldsymbol{\beta}$  fixé, on estime  $A_0(t)$  par

$$\hat{A}_0(t) = \int_0^t \frac{J(u)}{S^{(0)}(\boldsymbol{\beta}, u)} dN_+(u), \quad (2.3)$$

avec  $J(u) = \mathbb{1}_{\{Y_1(u) + \dots + Y_n(u) > 0\}}$ .

$\hat{A}_0(t)$  est appelé **estimateur de Breslow** (Breslow, 1974).

En remplaçant, dans (2.2),  $A_0(t)$  par son estimation obtenue en (2.3), nous obtenons pour expression de la vraisemblance partielle

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &= \prod_t \prod_i \left\{ (d\hat{A}_0(t) Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i))^{\Delta N_i(t)} \right\} \times \exp \left[ - \int_0^\tau S^{(0)}(\boldsymbol{\beta}, u) d\hat{A}_0(u) \right] \\ &= \prod_t \prod_i \left\{ (Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i))^{\Delta N_i(t)} (d\hat{A}_0(t))^{\Delta N_i(t)} \right\} \times \exp \left[ - \int_0^\tau S^{(0)}(\boldsymbol{\beta}, u) \frac{J(u) dN_+(u)}{S^{(0)}(\boldsymbol{\beta}, u)} \right] \\ &= \prod_t \prod_i \left\{ [Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)]^{\Delta N_i(t)} \times \left[ \frac{J(t) dN_+(t)}{S^{(0)}(\boldsymbol{\beta}, t)} \right]^{\Delta N_i(t)} \right\} \times \exp \left[ - \int_0^\tau J(u) dN_+(u) \right] \\ &= \prod_t \prod_i \left\{ \left[ \frac{Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{S^{(0)}(\boldsymbol{\beta}, t)} \right]^{\Delta N_i(t)} \times [J(t) dN_+(t)]^{\Delta N_i(t)} \right\} \times \exp \left[ - \int_0^\tau J(u) dN_+(u) \right] \\ &= L(\boldsymbol{\beta}) \times \prod_t \prod_i \left\{ J(t) dN_+(t) \right\}^{\Delta N_i(t)} \times \exp \left[ - \int_0^\tau J(u) dN_+(u) \right],\end{aligned}$$

avec

$$L(\boldsymbol{\beta}) = \prod_t \prod_i \left( \frac{Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{S^{(0)}(\boldsymbol{\beta}, t)} \right)^{\Delta N_i(t)} \quad (2.4)$$

dépendant de  $\boldsymbol{\beta}$  (le reste de la vraisemblance étant indépendant de  $\boldsymbol{\beta}$ ).

Par définition,  $L(\boldsymbol{\beta})$  est la **vraisemblance partielle de Cox**.

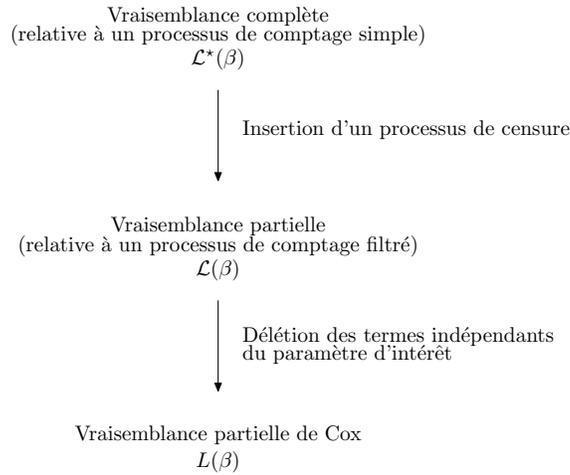


FIG. 2.1 – *Les vraisemblances successives.*

**Remarque** — La présence, dans (2.4), des deux produits (l'un suivant  $t$ , l'autre suivant  $i$ ) découle de la formalité de l'écriture mathématique. Ainsi, ces deux signes peuvent paraître redondant, dans la mesure où l'hypothèse d'absence d'ex-æquo entraîne que chaque processus de comptage  $N_i(t)$  ne vaut 1 que pour une et une seule valeur de  $t$ .

Considérons maintenant la fonction de log-vraisemblance partielle de Cox sur l'intervalle  $[0, t]$  :

$$\log L(\beta, t) = \sum_i \int_0^t \left[ \beta^t \mathbf{Z}_i - \log S^{(0)}(\beta, u) \right] dN_i(u). \quad (2.5)$$

Le vecteur score

$$\mathbf{U}(\beta, t) = \frac{\partial \log L(\beta, t)}{\partial \beta}$$

peut s'écrire

$$\begin{aligned} \mathbf{U}(\beta, t) &= \sum_i \int_0^t \left[ \mathbf{Z}_i - E(\beta, u) \right] dN_i(u) \\ &= \sum_i \int_0^t \left[ \mathbf{Z}_i - E(\beta, u) \right] dM_i(u) \end{aligned} \quad (2.6)$$

où

$$E(\beta, t) = \frac{S^{(1)}(\beta, t)}{S^{(0)}(\beta, t)} \quad (2.7)$$

et

$$S^{(1)}(\beta, t) = \sum_i Y_i(t) \exp(\beta^t \mathbf{Z}_i) \mathbf{Z}_i.$$

$[\mathbf{Z}_i - E(\boldsymbol{\beta}, u)]$  étant un vecteur de processus prévisible,  $\mathbf{U}(\boldsymbol{\beta}, t)$  est une somme de  $n$  martingales vectorielles, et est donc lui-même une martingale.

La suite de martingales  $M^{(n)}(t) = n^{-1/2}\mathbf{U}(\boldsymbol{\beta}, t)$  vérifie les conditions d'application du théorème de Rebolledo. En appliquant ce dernier, nous déduisons de la loi du processus limite  $M^{(\infty)} = M^{(\tau)}$  le résultat suivant :

**Proposition I.7** — Soit  $\hat{\boldsymbol{\beta}}$  l'estimateur du maximum de vraisemblance partielle de Cox, i.e. la quantité vérifiant

$$\mathbf{U}(\hat{\boldsymbol{\beta}}, \tau) = 0. \quad (2.8)$$

Alors

$$\hat{\boldsymbol{\beta}} \xrightarrow{\mathcal{L}} \mathcal{N}(\boldsymbol{\beta}_0, \mathbf{J}^{-1}(\hat{\boldsymbol{\beta}})),$$

où  $\mathbf{J}(\boldsymbol{\beta})$  est la *matrice d'information de Fisher* :

$$\begin{aligned} \mathbf{J}(\boldsymbol{\beta}) &= -\frac{\partial^2 L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^2} \\ &= -\sum_{i=1}^n \int_0^\tau \left\{ \frac{S^{(2)}(\boldsymbol{\beta}, s)}{S^{(0)}(\boldsymbol{\beta}, s)} - E(\boldsymbol{\beta}, s)^{\otimes 2} \right\} dN_i(s) \end{aligned}$$

avec

$$S^{(2)}(\boldsymbol{\beta}, s) = \sum_i Y_i(s) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i) \mathbf{Z}_i^{\otimes 2}.$$

$\mathbf{J}^{-1}(\boldsymbol{\beta})$ , inverse de la matrice d'information de Fisher, fournit une estimation de la variance de  $\hat{\boldsymbol{\beta}}$ .

### 2.2.2 Résolution numérique

Pour résoudre l'équation du score (2.8), l'**algorithme de Newton-Raphson** est habituellement employé.

Partant d'une solution initiale  $\hat{\boldsymbol{\beta}}_0 = 0$ , l'algorithme consiste en la succession d'itérations de la forme

$$\hat{\boldsymbol{\beta}}^{j+1} = \hat{\boldsymbol{\beta}}^j - \left[ \frac{\partial^2 \log L(\hat{\boldsymbol{\beta}}^j, \tau)}{\partial \boldsymbol{\beta}^2} \right]^{-1} \frac{\partial \log L(\hat{\boldsymbol{\beta}}^j, \tau)}{\partial \boldsymbol{\beta}}.$$

Le terme qui suit le signe moins est le pas itératif de l'algorithme de Newton-Raphson.

Si la fonction de vraisemblance évaluée en  $\hat{\boldsymbol{\beta}}^{j+1}$  est inférieure à celle évaluée en  $\hat{\boldsymbol{\beta}}^j$ , alors  $\hat{\boldsymbol{\beta}}^{j+1}$  est recalculé en utilisant, cette fois-ci, la moitié du pas itératif.

Ces étapes se succèdent jusqu'à ce que la convergence soit obtenue, c'est-à-dire jusqu'à ce que  $\hat{\boldsymbol{\beta}}^{m+1}$  soit suffisamment proche de  $\hat{\boldsymbol{\beta}}^m$ . L'estimateur du maximum de vraisemblance de  $\boldsymbol{\beta}$  est alors  $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^{m+1}$ .

## 2.3 Tests

Trois tests de l'hypothèse nulle «  $H_0 : \beta = \beta_0$  » peuvent être déduits du résultat concernant la convergence asymptotique de  $\hat{\beta}$ .

### 2.3.1 Test du rapport de vraisemblance

Ce test, très couramment utilisé en statistique, découle d'un développement de Taylor à l'ordre 2 de  $\log L(\beta)$ , puis de propriétés de convergence en loi (Dacunha-Castelle et Duflo, 1993).

Il s'énonce comme suit :

$$2[\log L(\hat{\beta}) - \log L(\beta_0)] \rightsquigarrow \chi^2(p).$$

Ce test mesure la différence des valeurs prises par le logarithme de la vraisemblance en  $\hat{\beta}$  et  $\beta_0$  ; l'espérance de cette quantité doit être nulle sous  $H_0$ .

### 2.3.2 Test de Wald (ou du maximum de vraisemblance)

D'après le résultat de la page 33, nous avons

$$\sqrt{\mathfrak{J}(\hat{\beta})}(\hat{\beta} - \beta_0) \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1).$$

Or si une v.a.  $X$   $p$ -dimensionnelle suit une loi normale centrée réduite, alors  $X^2$  suit une loi du chi-deux à  $p$  degrés de liberté. Ainsi,

$$(\hat{\beta} - \beta_0)^t \mathfrak{J}(\hat{\beta}) (\hat{\beta} - \beta_0) \rightsquigarrow \chi^2(p).$$

Il mesure l'écart entre  $\hat{\beta}$  et  $\beta_0$ , qui est nul en moyenne sous  $H_0$  car  $\hat{\beta}$  est asymptotiquement sans biais.

### 2.3.3 Test du score

Il est possible de montrer que

$$\left. \frac{\partial \log L(\beta, t)}{\partial \beta} \right|_{\beta=\beta_0} \xrightarrow{\mathcal{L}} \mathcal{N}(0, \mathfrak{J}).$$

En notant  $U(\beta, t) = (\partial \log L(\beta, t)) / \partial \beta$ , nous obtenons par conséquent

$$\left[ U(\beta_0, \tau)^t \mathfrak{J}^{-1}(\beta_0) U(\beta_0, \tau) \right] \rightsquigarrow \chi^2(p).$$

Ce test mesure la pente de la tangente en  $\beta_0$ . Sous  $H_0$ , le maximum de vraisemblance est obtenu pour une valeur  $\hat{\beta}$  proche de  $\beta_0$ . La pente en  $\beta_0$  diffère donc peu de 0, elle est nulle en moyenne sous  $H_0$ .

## 2.4 Critères d'adéquation au modèle de Cox

L'appréciation de la justesse de l'adéquation (*goodness of fit*) passe le plus souvent par une comparaison graphique des résidus des différents modèles que nous désirons comparer ; une alternative, employant certaines statistiques, peut quelquefois se présenter.

### 2.4.1 Dans sa globalité

Définissons les **résidus de Cox-Snell** (Cox et Snell, 1968) comme étant les quantités suivantes :

$$R_i = \widehat{A}_0(T_i) e^{\widehat{\beta}^t \mathbf{Z}_i} .$$

Si le modèle est correct et si les paramètres sont proches de leurs vraies valeurs, les  $R_i$  doivent alors constituer un échantillon censuré de distribution exponentielle unitaire.

Par suite, la représentation graphique de l'estimateur de Nelson-Aalen du risque cumulé en fonction des  $R_i$  doit approcher la première bissectrice.

### 2.4.2 Concernant la forme fonctionnelle des covariables

Il est possible que la forme fonctionnelle des covariables, telle qu'elle est spécifiée par le modèle de Cox – soit  $\exp(\beta^t \mathbf{Z})$  –, ne soit pas exacte. Considérons donc le modèle à deux covariables

$$\begin{aligned} \Lambda(t, \mathbf{Z}, X) &= h(\mathbf{Z}^{(1)}) \exp(\beta^t \mathbf{Z}^{(2)}) \Lambda_0(t) \\ &= \exp(f(\mathbf{Z}^{(1)})) \exp(\beta^t \mathbf{Z}^{(2)}) \Lambda_0(t) , \end{aligned} \quad (2.9)$$

où la forme fonctionnelle pour  $\mathbf{Z}^{(2)}$  est connue, tandis que la fonction positive  $h(\mathbf{Z}^{(1)})$  est, elle, inconnue.

Définissons également les **résidus martingales** par

$$\widehat{M}_i(T_i) = \delta_i - \widehat{A}_0(T_i) \exp(\widehat{\beta}^t \mathbf{Z}_i^{(2)}) ,$$

où  $T_i$  est le temps d'observation concernant le  $i^e$  individu.

À l'instar d'un modèle linéaire classique, les résidus sont la différence entre les valeurs observées et les valeurs attendues.

Leur représentation graphique en fonction de  $\mathbf{Z}^{(1)}$  permet d'obtenir des estimations de  $h$  ou de  $f$ .

Si l'on note  $\widehat{M}$  le résidu martingale lorsque le modèle (2.9) est le modèle adéquat, mais que l'on a ignoré  $\mathbf{Z}^{(1)}$ , et si  $\mathbf{Z}^{(1)}$  et  $\mathbf{Z}^{(2)}$  sont indépendantes, alors (Therneau *et al.*, 1990)

$$\mathbb{E}[\widehat{M}(t) \mid \mathbf{Z}^{(1)}] \approx \left\{ 1 - \frac{\bar{h}}{h(\mathbf{Z}^{(1)})} \right\} \mathbb{E}[N(t) \mid \mathbf{Z}^{(1)}] , \quad (2.10)$$

où

$$\bar{h} = \frac{\mathbb{E}\left\{\exp[f(\mathbf{Z}^{(1)})]\mathbf{Y}(t)\right\}}{\mathbb{E}[\mathbf{Y}(t)]}.$$

Cette dernière équation s'interprète naturellement : le nombre attendu de décès supplémentaires est approximativement égal à 1 moins le taux de risque que multiplie le nombre attendu d'événements.

Puisque  $\widehat{\mathbf{M}}$  et  $\mathbf{N}$  sont connus, on peut inverser (2.10) afin d'obtenir

$$f(\mathbf{Z}^{(1)}) - \bar{f} \approx -\log\left\{1 - \frac{\text{sm}(\widehat{\mathbf{M}}, \mathbf{Z}^{(1)})}{\text{sm}(\mathbf{N}, \mathbf{Z}^{(1)})}\right\}, \quad (2.11)$$

où

- $\bar{f} = \log(\bar{h})$ ,
- $\text{sm}(\widehat{\mathbf{M}}, \mathbf{Z}^{(1)})$  est une estimation lissée (*smoothed*) de  $\mathbb{E}[\widehat{\mathbf{M}}(t) | \mathbf{Z}^{(1)}]$ , qui peut être obtenue en traçant le graphe lissé de  $\widehat{\mathbf{M}}$  en fonction de  $\mathbf{X}$ ,
- $\text{sm}(\widehat{\mathbf{N}}, \mathbf{Z}^{(1)})$  est l'analogue, concernant  $\mathbf{N}$ , de la quantité précédente.

Therneau *et al.* (1990) démontrent que l'équation (2.11) peut, pour  $t = \infty$ , être remplacée par l'approximation suivante<sup>1</sup> :

$$\mathbb{E}[\widehat{\mathbf{M}}(t) | \mathbf{Z}^{(1)}] \approx c \{f(\mathbf{Z}^{(1)}) - \bar{f}\}, \quad (2.12)$$

où  $c$  est le nombre total d'événements divisé par le nombre total de sujets.

Ainsi, un graphe lissé des  $\widehat{M}_i$  suivant une covariable fournira une approximation de la forme fonctionnelle correcte à placer dans l'exponentielle du modèle de Cox.

Enfin, un avantage de (2.12) par rapport à (2.11) réside en son interprétation : l'axe des ordonnées est à l'échelle directe des décès supplémentaires.

Therneau *et al.* ont mené une expérience – limitée – de simulations, qui a montré que l'approximation (2.12) est acceptable lorsque  $\beta \mathbb{E}(\mathbf{N}) < 2$ .

### 2.4.3 Concernant la proportionnalité des risques vis-à-vis d'une covariable

#### Méthodes graphiques

Supposons que nous désirions tester la proportionnalité du modèle vis-à-vis de la covariable  $\mathbf{Z}^{(1)}$ , après ajustement sur les autres covariables. Notons  $\mathbf{Z} = (\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)})^t$  le vecteur des covariables, où  $\mathbf{Z}^{(2)}$  est le vecteur des  $p - 1$  covariables restantes.

On suppose qu'il n'existe pas d'interaction entre  $\mathbf{Z}^{(1)}$  et les autres covariables.

---

1. À condition que  $f$  n'ait pas de variation extrême et que la dépendance de  $\mathbb{E}(\mathbf{N} | \mathbf{X})$  par rapport à  $\mathbf{X}$  soit faible, comme par exemple en présence d'un taux modéré de censure.

**Première méthode** Supposons que  $\mathbf{Z}^{(1)}$  prenne  $K$  valeurs possibles. Si  $\mathbf{Z}^{(1)}$  est continue, nous stratifions les données suivant  $K$  strates, notées  $G_1, \dots, G_K$ . Si  $\mathbf{Z}^{(1)}$  est discrète, elle prend les valeurs  $1, 2, \dots, K$ .

Si nous utilisons un modèle de Cox stratifié, nous obtenons  $\hat{A}_{g0}(t)$ , ( $g = 1, \dots, K$ ), qui est le risque de base cumulé pour la  $g^e$  strate.

S'il y a proportionnalité des risques, le risque de base cumulé de chaque strate doit être un multiple des autres. Aussi, si nous traçons  $\ln [\hat{A}_{10}(t)], \dots, \ln [\hat{A}_{K0}(t)]$  en fonction de  $t$ , nous devons obtenir des courbes parallèles et l'écart entre ces courbes doit demeurer constant.

Nous pouvons aussi tracer  $\ln [\hat{A}_{g0}(t)] - \ln [\hat{A}_{10}(t)]$  ( $g = 2, \dots, K$ ) en fonction de  $t$ : nous devons alors obtenir des courbes globalement constantes.

**Deuxième méthode** Andersen (1982) propose, pour chaque  $t$ , de tracer  $\hat{A}_{g0}(t)$  ( $g = 2, \dots, K$ ) en fonction de  $\hat{A}_{10}(t)$ . En cas de proportionnalité, ces courbes doivent être des droites passant par l'origine.

De plus, si  $A_{g0}(t) = e^{\gamma g} A_{10}(t)$ , alors la pente de ces droites devrait approximativement être une estimation de  $e^{\gamma g}$ .

Gill et Schumacher (1987) ont montré que si le graphe de  $\hat{A}_{g0}(t)$  en fonction de  $\hat{A}_{10}(t)$  est convexe (respectivement concave), alors le rapport  $\alpha_{g0}(t)/\alpha_{10}(t)$  est une fonction croissante (resp. décroissante) de  $t$ .

Concernant ces deux premières méthodes, l'interprétation qui peut en être faite doit être considérée avec précaution, en raison des variances des courbes qui ne sont pas constantes au cours du temps.

**Troisième méthode** Arjas (1988) propose la méthode suivante.

Notons  $T_{(1)} < T_{(2)} < \dots < T_{(M)}$  ( $M \leq n$ , où  $n$  est le nombre de sujets) les temps de survenue de l'événement. Soit

$$\begin{aligned} \mathfrak{M}_i(k, \boldsymbol{\beta}) &= N_i(T_{(k)}) - \int_0^t \frac{Y_i(s) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{\sum_{j=1}^n Y_j \exp(\boldsymbol{\beta}^t \mathbf{Z}_j)} d \left[ \sum_{i=1}^n N_i(s) \right] \\ &= N_i(T_{(k)}) - \sum_{j \leq k} \frac{Y_i(T_{(j)}) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{\sum_{l=1}^n Y_l(T_{(j)}) \exp(\boldsymbol{\beta}^t \mathbf{Z}_l)}. \end{aligned}$$

Si nous notons

$$H(k, \boldsymbol{\beta}) = \sum_{i=1}^n \sum_{j \leq k} \frac{Y_i(T_{(j)}) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{\sum_{l=1}^n Y_l(T_{(j)}) \exp(\boldsymbol{\beta}^t \mathbf{Z}_l)},$$

alors

$$\begin{aligned} \bar{\mathfrak{M}}(k, \boldsymbol{\beta}) &= \sum_i \mathfrak{M}_i(k, \boldsymbol{\beta}) \\ &= k - H(k, \boldsymbol{\beta}) \end{aligned}$$

est une martingale.

Par suite, un graphe de  $H(k, \hat{\beta})$  en fonction de  $k$  peut être comparé à la première bissectrice.

**Quatrième méthode** Cette méthode est basée sur les **résidus du score**. Pour la  $k^e$  variable concernant le  $i^e$  sujet, la définition du résidu du score est, dans le contexte du modèle de Cox,

$$U_{ik}(\hat{\beta}, \infty) = \int_0^\infty [Z_{ik} - E_k(\hat{\beta}, s)] d\widehat{M}_i(s),$$

où  $E_j(\beta, s)$  est la  $j^e$  composante du vecteur  $\mathbf{E}(\cdot, \cdot)$  défini en page 32.

Le modèle de Cox avec les  $p$  covariables est ajusté. Quand toutes les covariables sont fixées au temps 0, le résidu du score à un temps de survenue de l'événement donné vaut

$$H_{ik}(t) = \delta_i Y_i(t) [Z_{ik} - E_k(\hat{\beta}, T_i)] - \sum_{t_{(i)} \leq t} [Z_{ik} - E_j(\hat{\beta}, t_{(i)})] Y_i(t_{(i)}) \exp(\hat{\beta}^t) [\hat{A}_0(t_{(i)}) - \hat{A}(t_{(i-1)})], \quad (2.13)$$

où  $0 = t_{(0)} < t_{(1)} < \dots < t_{(M)}$  sont, comme précédemment, les temps ordonnés de survenue de l'événement d'intérêt.

En utilisant les scores relatifs aux  $n$  individus, nous définissons un processus du score pour la  $k^e$  covariable égal à

$$U_k(t) = \sum_{i=1}^n H_{ik}(t).$$

Le processus des scores est la première dérivée partielle de la fonction de vraisemblance partielle du modèle de Cox ajusté, qui utilise uniquement l'information disponible au temps  $t$ . Il est clair que  $U_k(0) = 0$  et que  $\mathbf{U}(\infty) = 0$ , puisque la valeur de  $\beta$  utilisée lors de la construction des résidus du score est la solution du vecteur d'équation  $U_k(\infty) = 0$ ,  $k = 1, \dots, p$ .

Si l'adéquation est correcte, alors le processus

$$W_k(t) = U_k(t) \times \sqrt{\mathbb{V}(\hat{\beta}_k)}$$

converge vers un pont brownien fluctuant aux alentours de 0 (à condition que  $\text{Cov}(\hat{\beta}_k, \hat{\beta}_{k'}) = 0$  pour  $k \neq k'$ ). Par suite, un graphe de  $W_k(t)$  en fonction du temps devrait ressembler à une marche aléatoire fluctuant autour de 0. Si les risques pour différents niveaux de la covariable ne sont pas proportionnels, alors les graphes doivent avoir un maximum qui est trop grand en valeur absolue, à un instant donné.

L'utilisation des résidus du score pour s'assurer de la proportionnalité des risques présente deux avantages par rapport aux autres approches :

- 1° les covariables continues sont traitées naturellement et n'ont pas besoin d'être discrétisées ;
- 2° un seul modèle de Cox doit être ajusté pour vérifier la proportionnalité des risques concernant toutes les covariables du modèle.

Cependant, du fait que la puissance du processus du score à détecter une non-proportionnalité des risques n'a pas été comparée à celle des graphiques d'Andersen ou d'Arjas, il est recommandé de mettre en pratique toutes les méthodes possibles.

### Méthodes analytiques

La proportionnalité des risques concernant la covariable  $\mathbf{Z}_j$  ne sera pas vérifiée si la statistique

$$\sup_t \sum_i \int_0^t [Z_{ij} - E_j(\hat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s)$$

dépasse une certaine valeur (rappelons que  $E_j(\boldsymbol{\beta}, s)$  est la  $j^e$  composante du vecteur défini en page 32).

Cette statistique devrait être très sensible aux alternatives pour lesquelles les covariables ont un comportement (de croissance ou de décroissance) monotone au cours du temps, et tout spécialement pour les alternatives telles que

$$\frac{\lambda(t; \mathbf{Z} = x)}{\lambda(t; \mathbf{Z} = y)}$$

soit strictement croissant en fonction de  $t$  pour tout  $x < y$ , ou bien strictement décroissant pour tout  $x < y$ .

Pour obtenir la distribution de cette statistique, on démontre que, sous les deux hypothèses données ci-après,

$$\sup_t \sqrt{\mathcal{J}^{-1}(\hat{\boldsymbol{\beta}}, \infty)_{jj}} \sum_i \int_0^t [Z_{ij} - E_j(\hat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s)$$

suit asymptotiquement la distribution de

$$\sup_{0 \leq t \leq 1} W^0(t),$$

où  $W^0$  est un pont brownien.

Les deux hypothèses mentionnées ci-dessus sont les suivantes :

- 1° la  $j^e$  composante du vecteur des covariables satisfait l'hypothèse de proportionnalité des risques ;
- 2°  $\{V(t)\}_{jk} = 0$  pour tout  $t$ , où  $V(\cdot)$  est la covariance asymptotique de

$$\frac{1}{\sqrt{n}} \left( \sum_i \int_0^t [Z_{i1} - E_1(\hat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s), \dots, \sum_i \int_0^t [Z_{ip} - E_p(\hat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s) \right).$$

Cette dernière condition de nullité nécessite que la covariable  $\mathbf{Z}_j$  soit orthogonale aux autres covariables.

L'estimateur consistant  $n^{-1}\mathcal{J}^{-1}(\hat{\boldsymbol{\beta}}, t)$  de  $V(t)$  est la somme, au long des temps de décès survenus dans l'intervalle  $[0, t]$ , des covariances de  $\mathbf{Z}$  à chaque instant de décès.

Par exemple,  $\{V(t)\}_{jk} \approx 0$  dans les études d'intervention pour lesquelles la  $j^e$  covariable représente le traitement administré après « randomisation » – aussi longtemps que de fortes interactions entre traitement et facteur n'existent pas.

Les cas où cette hypothèse n'est pas vérifiée demeurent à l'étude.

#### 2.4.4 Justesse du modèle pour chaque sujet

L'usage graphique des résidus permet d'apprécier la pauvreté de la prédiction individuelle. La taille du résidu individuel  $\widehat{M}_i$  indique la justesse du modèle, avec une grande valeur positive pour un sujet qui « meurt trop tôt » – et inversement, une grande valeur négative pour un sujet qui « vit trop longtemps ».

Dans le modèle de Cox, les résidus martingales sont fortement étirés (*skewed*) et cet étirement – cette queue – déforme l'apparence du graphique standard des résidus. Notons que si la valeur maximale d'un tel résidu est finie et vaut 1, sa valeur minimale possible est  $-\infty$ .

Il est alors presque impossible de détecter les données aberrantes (*outliers*).

C'est pourquoi il est préférable d'utiliser les **résidus de déviance** (*deviance residuals*).

Ces résidus sont définis à partir des résidus martingales par

$$d_i(t) = \operatorname{sgn}(\widehat{M}_i(t)) \left\{ -2 \left[ \widehat{M}_i(t) + N_i(t) \log \left( \frac{N_i(t) - \widehat{M}_i(t)}{N_i(t)} \right) \right] \right\}^{\frac{1}{2}}.$$

Il s'agit d'une transformation empirique des résidus martingales. Dans cette expression, la racine carrée tend à diminuer les résidus martingales grandement négatifs, tandis que la transformation logarithmique accroît les résidus martingales qui sont proches de l'unité. Ainsi, la distribution des résidus de déviance est plus symétrique autour de zéro que celle des résidus martingales.

Dans le cas du modèle de Cox, les résidus de déviance s'écrivent

$$d_i(T_i) = \operatorname{sgn}(\widehat{M}_i(T_i)) \left\{ -2 \left[ \widehat{M}_i(T_i) + \delta_i \log \left( \delta_i - \widehat{M}_i(T_i) \right) \right] \right\}^{\frac{1}{2}}.$$

Le graphe des résidus de déviance en fonction des quantités

$$\sum_{k=1}^p \widehat{\beta}_k Z_{ik}$$

doit, lorsque la censure est modérée, ressembler à un échantillon de bruit distribué suivant une loi normale. En cas de censure sévère, une grande quantité de points proches de 0 va déformer l'approximation normale.

Dans tous les cas, les valeurs aberrantes potentielles auront des résidus de déviance dont les valeurs absolues seront trop importantes.

### 2.4.5 Concernant les « observations influentes »

#### Méthode graphique

Cette méthode consiste en la comparaison de l'estimation  $\hat{\beta}$  – obtenue en estimant  $\beta$  à partir de toutes les observations – et de l'estimation  $\hat{\beta}_{(i)}$  – obtenue en excluant des observations celle relative au  $i^e$  sujet.

Reprenons la définition des  $H_{ik}(t)$  (cf. p. 38), et notons  $H_{ik} = H_{ik}(\infty)$ . Nous avons maintenant

$$H_{ik} = \delta_i [Z_{ik} - E_k(\hat{\beta}, T_i)] - \sum_{t_b \leq T_i} [Z_{ik} - E_j(\hat{\beta}, t_b)] \exp(\hat{\beta}^t) [\hat{A}_0(t_b) - \hat{A}(t_{b-1})].$$

Le premier terme

$$\begin{aligned} R_{ik} &= \delta_i [Z_{ik} - E_k(\hat{\beta}, T_i)] \\ &= \delta_i \left[ Z_{ik} - \frac{\sum_{j=1}^n Y_j(t) Z_{jk} \exp(\hat{\beta} Z_j)}{\sum_{j=1}^n Y_j(t) \exp(\hat{\beta} Z_j)} \right] \end{aligned}$$

est le **résidu partiel de Schoenfeld** (Schoenfeld, 1982) : il s'agit de la différence entre la valeur observée de la covariable  $Z_{ik}$  à l'instant de survenue de l'événement, et la valeur obtenue par le modèle à ce même instant.

Il est possible de montrer que  $\hat{\beta} - \hat{\beta}_{(i)}$  est approximativement égal à

$$\mathcal{J}(\hat{\beta})(H_{i1}, \dots, H_{ip})^t,$$

où  $\mathcal{J}^{-1}(\hat{\beta})$  est la matrice d'information de Fisher observée.

Le graphe – pour chaque covariable – de ces résidus de Schoenfeld en fonction soit des temps de survenue de l'événement, soit de la covariable  $Z_{ik}$ , est utilisé pour établir l'influence de la  $i^e$  observation sur la  $k^e$  covariable (les  $R_{ik}$  en fonction des  $T_i$  doivent être centrés autour de 0).

#### Méthode analytique

L'influence d'une observation sur le modèle dépend à la fois du résidu obtenu après ajustement et de la valeur extrême de la covariable, soit grossièrement  $Z_i - E(\beta, t)$  que multiplie le résidu.  $E(\beta, t)$  est une fonction du temps : c'est la moyenne sur l'ensemble des individus à risque au temps  $t$ . Ceci suggère d'utiliser une valeur « moyenne au cours du temps » de  $Z_{ij} - E_j(\beta, t)$ , et finalement l'on parvient au résidu du score

$$\int_0^\infty (Z_{ij} - E_j(\hat{\beta}, s)) d\widehat{M}_i(s)$$

comme outil de mesure de l'influence d'une observation.

Une manière de formaliser ce résultat est d'ajouter des poids aux observations individuelles afin de donner une vraisemblance partielle et un vecteur de score pondérés. Ainsi<sup>1</sup>,

$$\begin{aligned} \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial w_i} &= \left( \frac{\partial \widehat{\boldsymbol{\beta}}}{\partial \mathbf{U}} \right) \left( \frac{\partial \mathbf{U}}{\partial w_i} \right) \\ &= \mathbf{J}(\widehat{\boldsymbol{\beta}})^{-1} \frac{\partial \mathbf{U}}{\partial w_i}, \end{aligned}$$

calculée au point  $w = 1$ , est l'estimateur *jackknife* infinitésimal de l'influence de la  $i^e$  observation sur  $\widehat{\boldsymbol{\beta}}$ . Une manipulation algébrique révèle que le second terme de l'équation ci-dessus est exactement le résidu du score, si bien que le vecteur d'influence du  $i^e$  sujet est

$$-\mathbf{J}(\widehat{\boldsymbol{\beta}})^{-1} \left( \int_0^\infty [Z_{i1} - E_1(\widehat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s), \dots, \int_0^\infty [Z_{ip} - E_p(\widehat{\boldsymbol{\beta}}, s)] d\widehat{M}_i(s) \right)^t.$$

Cette méthode sous-estime le *jackknife* réel, spécialement pour des valeurs extrêmes de  $\mathbf{Z}$ , puisque  $\mathbf{J}$  change également quand l'observation est modifiée.

Une autre méthode consiste à calculer le « premier pas de l'actualisation » (*1-step update*) de  $\widehat{\boldsymbol{\beta}}$  quand une seule covariable  $\mathbf{Z}_{p+1}$  est ajoutée, avec  $\mathbf{Z}_{p+1}$  valant 1 pour le sujet  $i$  et 0 pour tous les autres.

Ici, le changement au premier pas, au point  $(\widehat{\boldsymbol{\beta}}, 0)$  vaut

$$\Delta \widehat{\boldsymbol{\beta}}_{(i)} = \frac{-\mathbf{J}(\widehat{\boldsymbol{\beta}})^{-1} \gamma_i}{\eta_i - \gamma_i' \mathbf{J}(\widehat{\boldsymbol{\beta}})^{-1} \gamma_i} \widehat{M}_i,$$

où

$$\gamma_{ij} = \int_0^\infty Y_i(s) [Z_{ij} - E_j(\widehat{\boldsymbol{\beta}}, s)] \exp(\widehat{\boldsymbol{\beta}}^t \mathbf{Z}_i(s)) d\widehat{\Lambda}_0(s)$$

et

$$\eta_i = \int_0^\infty Y_i(s) [1 - E_{p+1}(\widehat{\boldsymbol{\beta}}, s)] \exp(\widehat{\boldsymbol{\beta}}^t \mathbf{Z}_i(s)) d\widehat{\Lambda}_0(s).$$

Cette expression est très similaire à celle des estimateurs *jackknife* du modèle linéaire, avec  $\widehat{M}_i$  considéré comme résidu.

**Remarque** — Les tests usuels d'adéquation – Kolmogorov-Smirnov, von Mises – peuvent être adaptés au cas où les covariables dépendent du temps (Marzec et Marzec, 1997).

## 2.5 Considération des ex-æquo

L'expression de la vraisemblance partielle de Cox (cf. p. 31) ne vaut que sous l'hypothèse de temps de survenue de l'événement distincts, c'est-à-dire lorsque  $\Delta N_i(t) = N_i(t) - N_i(t-)$  ne peut valoir que 0 ou 1, quel que soit  $i$ .

Cependant, des adaptations de cette vraisemblance partielle ont été conçues, afin de traiter des données de survie présentant des ex-æquo ; nous les donnons ci-dessous.

---

1. En reprenant la notation page 32.

### 2.5.1 Vraisemblance partielle de Breslow

Elle s'écrit

$$L(\boldsymbol{\beta}) = \prod_t \prod_i \frac{Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{[S^{(0)}(\boldsymbol{\beta}, t)]^{\Delta N_i(t)}}.$$

### 2.5.2 Vraisemblance partielle d'Efron

Elle est de la forme

$$L(\boldsymbol{\beta}) = \prod_t \prod_i \left\{ \frac{Y_i(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_i)}{\prod_{k=1}^{\Delta N_i(t)} \left[ \sum_j Y_j(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_j) - \frac{k-1}{\Delta N_i(t)} \sum_l Y_l(t) \delta_l \exp(\boldsymbol{\beta}^t \mathbf{Z}_l) \right]^{\Delta N_i(t)}} \right\}.$$

### 2.5.3 Vraisemblance partielle exacte

Son expression est

$$L(\boldsymbol{\beta}) = \prod_t \prod_i \left[ \int_0^{\tau} \prod_{j=1}^{\Delta N_i(t)} \left( 1 - \exp \left\{ - \left[ \frac{\exp(\boldsymbol{\beta}^t \mathbf{Z}_j)}{\sum_l Y_l(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_l)} \right] t \right\} \exp(-t) dt \right) \right].$$

## 2.6 Extensions du modèle

Le modèle de Cox peut être étendu à plusieurs cas, notamment ceux concernant une stratification des covariables, ou bien encore une dépendance de ces mêmes covariables vis-à-vis du temps.

Il est également possible, dans le cas de données emboîtées (où  $T_{ij}$  est le temps de survie de l'individu  $j$  appartenant au groupe  $i$ ), de diversifier la fonction de risque de base :  $\lambda_0(t)$  devient  $\lambda_{i0}(t)$ , fonction de risque de base propre au groupe  $i$ .

Il est enfin possible d'introduire une corrélation entre les données de survie : c'est précisément l'objet des chapitres à venir.

Notons que le modèle de Cox – tout comme un autre modèle, dit à *temps accélérés* – est un cas particulier d'un modèle ayant vu le jour en 1987 : le **modèle étendu de régression du risque instantané** (*extended hazard regression*) (Etezadi-Amoli et Ciampi). Ce modèle spécifie que le risque instantané s'écrit

$$\alpha(t) = g_1(\boldsymbol{\beta}^t \mathbf{Z}) \times \alpha_0[g_2(\boldsymbol{\gamma}^t \mathbf{Z})t],$$

où  $g_1(x)$  et  $g_2(x)$  sont des fonctions positives égales à 1 en 0,  $\alpha_0(t)$  est le risque de base, et  $\boldsymbol{\beta}$  et  $\boldsymbol{\gamma}$  sont les paramètres de la régression.



Deuxième partie

**MODÉLISATION DE LA SURVIE  
POUR DES DONNÉES  
CORRÉLÉES**



# Introduction

## Origine de la corrélation

La modélisation statistique s'est libérée, depuis une vingtaine d'années, de l'une des hypothèses les plus couramment posées : l'indépendance des données.

Dans certaines situations, il est en effet indispensable de remettre en cause cette hypothèse. Parmi celles-ci, citons le cas où une hétérogénéité inconnue opère de façon latente sur les variables à expliquer, ou bien encore le cas où les données sont susceptibles de présenter certaines associations.

Une illustration de ce dernier cas est l'étude d'animaux (niveau « individuel ») regroupés en troupeaux (niveau « groupe ») : de façon presque certaine, nous pouvons considérer que les données recueillies sur des animaux appartenant à un même troupeau ne seront pas indépendantes d'un point de vue statistique.

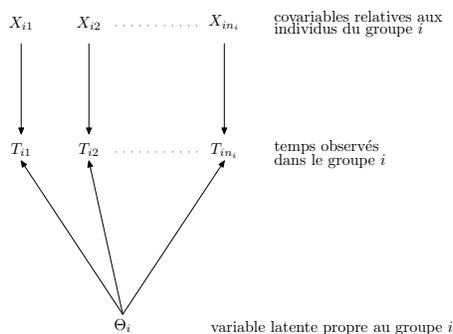


FIG. 2.2 – Schéma représentant la corrélation.

L'ignorance de cette corrélation entraîne un possible biais des paramètres du modèle et en particulier, dans le cas d'une corrélation positive, une sous-estimation de la variance des paramètres du niveau du groupe (*cluster*), et une surestimation de celle des variables intra-groupe.

Par la suite, nous nous placerons dans le cadre d'un modèle hiérarchique : nous supposons que les observations individuelles (d'un animal) sont corrélées au sein d'un même élevage (groupe).

Il est à noter que la littérature abonde en modèles pour données corrélées par paires (Clayton et Cuzick, 1985) ; les modèles pour données corrélées par groupes sont, quant à eux, plus rares.

En particulier, il s'avère souvent qu'une illustration numérique – au travers de simulations – d'un modèle pour données corrélées par groupes se fasse dans le cadre restreint de données corrélées par paires, la raison en étant la simplicité de l'écriture des distributions de probabilité qui, sous-tendant les simulations, doivent incorporer la corrélation.

Dans le domaine des modèles de survie, deux grandes approches de la corrélation sont fréquemment choisies : l'approche **marginale** et l'approche **mixte**. Elles seront étudiées tout au long de ce travail – d'abord d'un point de vue théorique, ensuite d'un point de vue pratique, notamment au travers de simulations ainsi que d'une application sur des données réelles.

## Approche marginale

### Modèle « classique »

L'approche marginale se caractérise par une modélisation de la réponse (d'un individu) conditionnellement aux covariables, mais pas aux autres réponses (individuelles). Le modèle marginal modélise donc une « réponse moyenne » ; il est qualifié de modèle de type **population-averaged**.

La corrélation, considérée comme une nuisance, est modélisée séparément. Ainsi, dans le cadre du modèle de Cox, le paramètre d'intérêt  $\beta$  peut, dans un premier temps, être estimé sous l'hypothèse d'indépendance. Puis, dans un second temps, un ajustement de sa variance permet d'inclure la dispersion supplémentaire (*overdispersion*).

Cependant, il est possible d'inclure la corrélation plus en amont : dans la fonction d'estimation de  $\beta$ . La nouvelle estimation du paramètre, ainsi que sa variance, sont donc corrigées grâce à l'introduction, dans la fonction d'estimation, d'un terme de pondération. Cette méthode est similaire aux **équations d'estimation généralisées** (*Generalized Estimating Equations* – GEE) pour l'analyse des données longitudinales (Liang *et al.*, 1992).

### Modèle de Liang, Self et Chang

Ce modèle (Liang *et al.*, 1993) base la fonction d'estimation sur une expression exploitant les comparaisons entre observations indépendantes.

Les auteurs démontrent la consistance, ainsi que la normalité asymptotique de l'estimateur du paramètre de régression.

Adapté au cas de données corrélées par paires, ce modèle se prête difficilement au cas de données corrélées par groupes ; aussi ne sera-t-il pas davantage développé ici.

### Modèle avec copule

Une troisième approche du modèle marginal est celle basée sur les **copules** (Liang *et al.*, 1995). Ce type de modèle s'attache à estimer aussi bien les paramètres de régression des covariables que la valeur de la dépendance entre les temps de survie. Ainsi, contrairement à l'approche « classique », la corrélation n'est plus traitée comme une simple nuisance, et est même estimée.

Cependant, le modèle marginal avec copule n'a pas été suffisamment développé; la prise en compte de données multivariées pose notamment problème. Aussi avons-nous laissé de côté ce type de modèle pour nous consacrer au modèle marginal « classique ».

## Approche mixte

Un modèle mixte modélise une réponse individuelle (emboîtée dans un groupe) conditionnellement aux covariables et aux réponses des autres individus de ce groupe. Il s'agit donc ici de la modélisation d'une réponse spécifique à un groupe: aussi qualifie-t-on ce type de modèle de *cluster-specific*.

La prise en compte de la corrélation se fait généralement *via* l'introduction, dans la fonction de risque, d'une variable aléatoire censée représenter l'« effet groupe ». Sa distribution de probabilité est généralement à valeurs dans  $\mathbb{R}^+$ .

Dans le cadre du modèle de Cox, cette variable aléatoire s'insère dans la partie régressive, devenant ainsi une covariable supplémentaire. Du fait de sa positivité, cette variable aléatoire a pour effet d'accroître le risque instantané de survenue de l'événement d'intérêt: c'est pourquoi elle est qualifiée de « fragilité », et ce type de modèle désigné sous les termes de **modèle de fragilité**<sup>1</sup>. L'étude de la significativité de cette nouvelle covariable permet d'établir s'il y a ou non, concernant une réponse (individuelle), influence de son appartenance à un groupe plutôt qu'à un autre.

Par la suite, les données d'un même groupe sont supposées être indépendantes conditionnellement à cette variable aléatoire.

Plusieurs études empiriques ont montré que l'ignorance de l'existence de la fragilité entraînait une sous-estimation des effets des covariables (Aalen, 1988; Klein, 1992; Henderson et Oman, 1999). Ce résultat a été conforté par une étude théorique de Lancaster (1990) concernant le modèle de fragilité de Weibull pour des données non censurées.

Le modèle mixte peut être envisagé de deux points de vue différents: soit d'un point de vue *fréquentiste*, soit d'un point de vue *bayésien*. Ces deux manières d'appréhender le modèle mixte seront présentées tour à tour.

## Approche pseudo-mixte

Le pseudo-modèle de fragilité (Mahé, 1998) est une combinaison des deux précédentes approches. Son avantage est de fournir à la fois une estimation de l'effet moyen des covariables (typique du modèle marginal) et une estimation de la corrélation intra-groupe (propre au modèle mixte).

Ce modèle, qui consiste à formuler le lien mathématique entre risque marginal et risque individuel, n'a pas été retenu dans les études à venir. Deux raisons peuvent être avancées, quant à cette décision d'écarter le pseudo-modèle de fragilité: d'une part le bénéfice attendu d'une telle approche, dans le cadre de nos travaux, n'était pas probant, et d'autre part nous souhaitions privilégier la comparaison de modèles plus « classiques » – ces seuls modèles nécessitant déjà de

1. Signalons que la dénomination de « modèle de fragilité partagée » est également usitée: elle permet d'insister sur le partage du caractère non observé (modélisé par la fragilité) entre tous les individus d'un même groupe.

nombreuses heures d'exécution informatique.

## Chapitre 3

# MODÈLE MARGINAL DE COX

### Contenu

---

<b>3.1</b>	<b>Estimation et tests dans le cadre général . . . . .</b>	<b>52</b>
3.1.1	Notations . . . . .	52
3.1.2	Estimation du risque de base cumulé . . . . .	52
3.1.3	Estimation du paramètre de la régression . . . . .	52
3.1.4	Tests . . . . .	54
<b>3.2</b>	<b>Estimations et tests dans un cadre restreint . . . . .</b>	<b>55</b>
<b>3.3</b>	<b>Critères d'adéquation . . . . .</b>	<b>55</b>
3.3.1	Test de la forme fonctionnelle d'une covariable . . . . .	57
3.3.2	Test de la fonction de lien . . . . .	57
3.3.3	Test de la proportionnalité des risques vis-à-vis d'une covariable . . . . .	58
3.3.4	Test de la proportionnalité des risques vis-à-vis de toutes les covariables . . . . .	58
3.3.5	Test concernant le choix même du modèle marginal . . . . .	58
<b>3.4</b>	<b>Coefficient de dispersion (<i>design effect</i>) . . . . .</b>	<b>59</b>

---

### 3.1 Estimation et tests dans le cadre général

#### 3.1.1 Notations

Le modèle que nous exposons ici est celui de Wei *et al.* (1989).

Considérons  $n$  groupes. Nous supposons que le groupe  $i$  ( $i = 1, \dots, n$ ) contient  $m_i$  individus.

Soient, pour  $(i, j) \in \{1, \dots, n\} \times \{1, \dots, m_i\}$ :

- $X_{ij}$  la date de survenue de l'événement chez l'individu  $j$  du groupe  $i$ ;
- $C_{ij}$  la date de censure correspondant;
- $T_{ij} = X_{ij} \wedge C_{ij}$ ;
- $\delta_{ij} = \mathbb{1}_{\{X_{ij} \leq C_{ij}\}}$ ;
- $\mathbf{Z}_{ij} = (Z_{1ij}, \dots, Z_{pij})$  le vecteur de dimension  $p$  des covariables;
- $Y_{ij}(t) = \mathbb{1}_{\{T_{ij} \geq t\}}$ ;
- $\mathbf{N}_i(t) = \{N_{ij}(t), 0 \leq t \leq \tau, j = 1, \dots, m_i\}$  le processus de comptage multivarié pour les  $m_i$  individus du groupe  $i$ ;
- $\mathbf{\Lambda}_i(t) = \{\Lambda_{ij}(t), 0 \leq t \leq \tau, j = 1, \dots, m_i\}$  le compensateur de  $N_i(t)$  par rapport à la filtration  $\mathcal{F}_t$ .

La fonction de **risque marginal** est

$$\alpha_{ij}(t) = \alpha_0(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ij}).$$

#### 3.1.2 Estimation du risque de base cumulé

L'estimateur du risque de base est déterminé sous hypothèse d'indépendance; il s'agit de l'estimateur – consistant et asymptotiquement normal – de Breslow (cf. p. 31):

$$\begin{aligned} \widehat{A}_0(t) &= \sum_{i=1}^n \sum_{j=1}^{m_i} \frac{\mathbb{1}_{\{T_{ij} \leq t\}} \delta_{ij}}{S^{(0)}(\widehat{\boldsymbol{\beta}}, T_{ij})} \\ &= \sum_{i=1}^n \sum_{j=1}^{m_i} \int_0^t \frac{dN_{ij}(u)}{S^{(0)}(\widehat{\boldsymbol{\beta}}, u)}, \end{aligned}$$

où

$$S^{(0)}(t) = \sum_i \sum_j Y_{ij}(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ij}).$$

#### 3.1.3 Estimation du paramètre de la régression

L'équation du score sous hypothèse d'indépendance, telle qu'en page 32, peut se réécrire matriciellement sous la forme

$$\sum_{i=1}^n \sum_{j=1}^{m_i} \int_0^\tau \mathbf{Z}_{ij} \widehat{M}_{ij}(dt) = 0$$

soit encore

$$\sum_{i=1}^n \int_0^\tau \mathbf{Z}_i^t \widehat{M}_i(dt) = 0,$$

avec

$$\begin{aligned}\widehat{M}_{ij}(t) &= N_{ij}(t) - \int_0^t Y_{ij}(s) e^{\beta^t \mathbf{Z}_{ij}(s)} \widehat{\Lambda}_0(ds), \\ \mathbf{Z}_i^t &= (\mathbf{Z}_{i1}, \dots, \mathbf{Z}_{in_i}), \\ \widehat{\mathbf{M}}_i^t(t) &= (\widehat{M}_{i1}(t), \dots, \widehat{M}_{in_i}(t)).\end{aligned}$$

De façon similaire aux équations d'estimations généralisées, et suivant les travaux de Binder (1992), cette équation du score peut être pondérée afin d'accroître l'efficacité de l'estimateur  $\widehat{\beta}$  :

$$\sum_{i=1}^n \int_0^\tau \mathbf{Z}_i^t \mathbf{W}_i(\beta, t) \widehat{\mathbf{M}}_i(dt) = 0, \quad (3.1)$$

où  $\mathbf{W}_i(\beta, t) = \{w_{ijk}(\beta, t), j, k = 1, \dots, n_i\}$  est la matrice de pondération, qui peut être définie de diverses façons – cette matrice peut, par exemple, être l'inverse de la matrice de corrélation de  $M_i^t(T_i) = (M_{i1}(T_{i1}), \dots, M_{in_i}(T_{in_i}))$ , ou une estimation de celle-ci.

Notons

$$\begin{aligned}\mathbf{S}^{(d)}(\beta, t) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} Y_{ij}(t) \mathbf{Z}_{ij}^d \exp(\beta^t \mathbf{Z}_{ij}) \quad d = 0, 1, \\ \mathbf{S}^{(d)}(\beta, t) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \sum_{k=1}^{n_i} \mathbf{Z}_{ij} w_{ijk}(\beta, t) Y_{ik}(t) \mathbf{Z}_{ij}^{d-2} \exp(\beta^t \mathbf{Z}_{ik}) \quad d = 2, 3, \\ \mathbf{S}^{(4)}(\beta, t) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \mathbf{Z}_{ij} Y_{ik}(t) \frac{\partial w_{ijk}(\beta, t)}{\partial \beta^t} \exp(\beta^t \mathbf{Z}_{ik})\end{aligned}$$

et

$$E(\beta, t) = \frac{\mathbf{S}^{(2)}(\beta, t)}{\mathbf{S}^{(0)}(\beta, t)}.$$

Notons enfin  $\widehat{\beta}_w$  et  $\widehat{\beta}_e$  les solutions de l'équation (3.1) respectivement dans le cas d'une matrice de pondération spécifiée dès l'origine et dans le cas d'une matrice de pondération estimée. Les résultats fondamentaux concernant ces deux estimateurs étant similaires, nous noterons  $\widehat{\beta}$  l'un ou l'autre de ces deux estimateurs.

Sous certaines conditions (cf. annexe § B.2 p. 164), l'énoncé suivant est vérifié :

$$\sqrt{n}(\widehat{\beta} - \beta_0) \rightsquigarrow \mathcal{N}(0, \Sigma),$$

où  $\Sigma$  peut être estimée par

$$\widehat{\Sigma} = \widehat{A}_w^{-1}(\widehat{\beta}) \widehat{\Sigma}_w \widehat{A}_w^{-1}(\widehat{\beta})$$

avec

$$\widehat{A}_w(\beta) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \left[ \frac{\widehat{\mathbf{S}}^{(3)}(\beta, t)}{\widehat{\mathbf{S}}^{(0)}(\beta, t)} - \frac{\widehat{\mathbf{S}}^{(2)}(\beta, t) \widehat{\mathbf{S}}^{(1)}(\beta, t)^t}{[\widehat{\mathbf{S}}^{(0)}(\beta, t)]^2} \right]$$

et

$$\widehat{\Sigma}_w(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \sum_{k=1}^{n_i} \widehat{G}_{ij}(\boldsymbol{\beta}) \widehat{G}_{ik}(\boldsymbol{\beta})^t ,$$

où

$$\begin{aligned} \widehat{G}_{ij}(\boldsymbol{\beta}) &= \int_0^\tau \left[ \sum_{k=1}^{n_i} \mathbf{Z}_{i\mathbf{k}} \widehat{w}_{ikj}(\boldsymbol{\beta}, t) - \widehat{E}(\boldsymbol{\beta}, t) \right] dN_{ij}(t) - \frac{1}{n} \sum_{l=1}^n \sum_{m=1}^{n_l} \int_0^\tau \frac{Y_{ij}(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ij})}{\widehat{\mathbf{S}}^{(0)}(\boldsymbol{\beta}, t)} \\ &\quad \times \left[ \sum_{k=1}^{n_i} \mathbf{Z}_{i\mathbf{k}} \widehat{w}_{ikj}(\boldsymbol{\beta}, t) - \widehat{E}(\boldsymbol{\beta}, t) \right] dN_{lm}(t) , \end{aligned}$$

avec  $\widehat{\mathbf{S}}^{(d)}$  ( $d = 0, \dots, 4$ ) valant exactement  $\mathbf{S}^{(d)}$  dans le cas d'une pondération non estimée, et valant sinon (c.-à-d. dans le cas d'une pondération estimée)  $\mathbf{S}^{(d)}$  avec  $w_{ikj}(\boldsymbol{\beta}, t)$  remplacé par  $\widehat{w}_{ikj}(\boldsymbol{\beta}, t)$ .

Différentes études (Cai et Prentice, 1995, 1997) ont montré que le bénéfice dû à la pondération de la fonction d'estimation est relatif, et notamment qu'il diminue à mesure que le taux de censure augmente.

### 3.1.4 Tests

Les adaptations des trois tests usuels (cf. § 2.3 p. 34) au cas d'un modèle marginal peuvent être entreprises (Cai, 1999).

**Test de Wald** Il se déduit du fait que

$$(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^t \widehat{\Sigma}(\widehat{\boldsymbol{\beta}}) (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \rightsquigarrow \chi^2(p) .$$

**Test du score** Il repose sur le résultat suivant :

$$\left[ U(\boldsymbol{\beta}_0, \tau)^t \widehat{\Sigma}^{-1}(\boldsymbol{\beta}_0) U(\boldsymbol{\beta}_0, \tau) \right] \rightsquigarrow \chi^2(p) ,$$

où  $U(\boldsymbol{\beta}, t)$  est l'équation du score (3.1).

**Test du rapport de vraisemblance** Il est basé sur le résultat suivant :

$$2 \left[ \log L(\widehat{\boldsymbol{\beta}}) - \log L(\boldsymbol{\beta}_0) \right] \rightsquigarrow \chi^2(p) ,$$

où  $L(\cdot)$  est la vraisemblance partielle de Cox sous l'hypothèse d'indépendance, et  $\widehat{\boldsymbol{\beta}}$  l'estimateur du maximum de cette même vraisemblance partielle.

### 3.2 Estimations et tests dans un cadre restreint

Reprenant le modèle précédent dans le cas particulier où la matrice de pondération est égale à la matrice identité, nous retrouvons le modèle exposé par Lin et Wei (1989).

L'estimateur du risque de base demeure celui de Breslow (cf. p. 31).

L'estimateur  $\hat{\beta}$ , calculé sous l'hypothèse d'indépendance des observations (cf. p. 32), demeure consistant et asymptotiquement distribué suivant une loi normale  $p$ -variée.

D'après les résultats de la section précédente, nous obtenons que, dans le cas d'une matrice de pondération égale à la matrice identité,

$$\sqrt{n}(\hat{\beta} - \beta) \rightsquigarrow \mathcal{N}(0, \mathcal{J}(\hat{\beta})^{-1} \Omega(\hat{\beta}) \mathcal{J}(\hat{\beta})^{-1}),$$

où

$$\mathcal{J}(\beta) = -\frac{\partial^2 L(\beta)}{\partial \beta^2}$$

et

$$\Omega(\beta) = \sum_{i=1}^n \sum_{j=1}^{m_i} \sum_{k=1}^{m_i} \Xi_{ij}(\hat{\beta}) \Xi_{ik}(\hat{\beta})^t \quad (3.2)$$

avec

$$\begin{aligned} \Xi_{ij}(\beta) &= \int_0^\tau \left[ \mathbf{Z}_{ij} - \frac{\bar{S}^{(1)}(\beta, t)}{\bar{S}^{(0)}(\beta, t)} \right] dN_{ij}(t) \\ &\quad - \sum_{k=1}^n \sum_{l=1}^{m_k} \int_0^\tau \left\{ \frac{Y_{ij}(t) \exp(\beta^t \mathbf{Z}_{ij})}{\bar{S}^{(0)}(\beta, t)} \times \left[ \mathbf{Z}_{ij} - \frac{\bar{S}^{(1)}(\beta, t)}{\bar{S}^{(0)}(\beta, t)} \right] \right\} dN_{kl}(t). \end{aligned}$$

Pour tester  $\beta = \beta_0$ , on utilise la statistique

$$(\hat{\beta} - \beta_0)^t D(\hat{\beta})^{-1} (\hat{\beta} - \beta_0) \rightsquigarrow \chi^2(p),$$

où

$$D(\hat{\beta}) = \mathcal{J}(\hat{\beta})^{-1} \Omega(\hat{\beta}) \mathcal{J}(\hat{\beta})^{-1}.$$

### 3.3 Critères d'adéquation

En reprenant les notations de la page 52, nous définissons les résidus martingales comme étant les quantités

$$\widehat{M}_{ik}(t) = N_{ik}(t) - \int_0^t Y_{ik}(s) e^{\hat{\beta}^t \mathbf{Z}_{ik}} d\widehat{A}_0(s, \hat{\beta}).$$

Ces résidus ont un comportement asymptotique similaire aux résidus ordinaires : leur espérance d'une part, leur somme totale d'autre part sont nulles.

L'étude des résidus martingales passe par celle des processus stochastiques suivants (dont les résidus martingales sont un cas particulier) :

$$W(t, z) = \sum_{i=1}^n \sum_{k=1}^K f(\mathbf{Z}_{ik}) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} \widehat{M}_{ik}(t), \quad (3.3)$$

où  $f$  est une fonction bornée et

$$\mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} = \mathbb{1}_{\{Z_{1ik} \leq z_1, \dots, Z_{pik} \leq z_p\}}.$$

D'après un développement de Taylor, il est possible de montrer (Speikerman et Lin, 1996) que le processus  $1/\sqrt{n} W(t, z)$  est asymptotiquement équivalent à un certain processus  $1/\sqrt{n} \widetilde{W}(t, z)$ , où

$$\begin{aligned} \widetilde{W}(t, z) &= \sum_{i=1}^n \sum_{k=1}^K \int_0^t \left[ f(\mathbf{Z}_{ik}) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} - \widetilde{G}(\boldsymbol{\beta}_0, s, z) \right] dM_{ik}(s) \\ &\quad + \widetilde{H}(\boldsymbol{\beta}_0, t, z) \widetilde{\mathcal{J}}^{-1}(\boldsymbol{\beta}_0) \sum_{i=1}^n \sum_{k=1}^K \int_0^\infty \left[ \mathbf{Z}_{ik} - \overline{Z}(\boldsymbol{\beta}_0, s) \right] dM_{ik}(s) \\ &= \sum_{i=1}^n \Psi_i(t, z), \end{aligned}$$

où  $\widetilde{G}(\boldsymbol{\beta}_0, s, z)$  et  $\widetilde{H}(\boldsymbol{\beta}_0, t, z)$  sont les limites respectives des processus

$$G(\boldsymbol{\beta}, t, z) = \sum_{i=1}^n \sum_{k=1}^K \frac{Y_{ik}(t) e^{\boldsymbol{\beta}^t \mathbf{Z}_{ik}} f(\mathbf{Z}_{ik}) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}}}{S^{(0)}(\boldsymbol{\beta}, t)}$$

et

$$H(\boldsymbol{\beta}, t, z) = \sum_{i=1}^n \sum_{k=1}^K \int_0^t Y_{ik}(s) e^{\boldsymbol{\beta}^t \mathbf{Z}_{ik}} f(\mathbf{Z}_{ik}) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} \left[ \mathbf{Z}_{ik} - \overline{Z}(\boldsymbol{\beta}, s) \right] d\widetilde{A}_0(s, \boldsymbol{\beta}).$$

Cette somme étant celle de  $n$  vecteurs i.i.d. et centrés, nous en déduisons que  $1/\sqrt{n} \widetilde{W}(t, z)$ , et par suite  $1/\sqrt{n} W(t, z)$ , sont asymptotiquement distribués suivant une loi normale centrée de matrice de covariance

$$\text{Cov} [(t, z), (t^b, z^b)] = \mathbb{E} [\Psi_1(t, z) \Psi_1(t^b, z^b)].$$

Cette fonction de covariance peut être estimée par  $1/n \sum_{i=1}^n \widehat{\Psi}_i(t, z) \widehat{\Psi}_i(t^b, z^b)^t$ , où

$$\begin{aligned} \widehat{\Psi}_i(t, z) &= \sum_{k=1}^K \left\{ \int_0^t \left[ f(\mathbf{Z}_{ik}) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} - G(\widehat{\boldsymbol{\beta}}, s, z) \right] d\widehat{M}_{ik}(s) \right. \\ &\quad \left. + H(\widehat{\boldsymbol{\beta}}, t, z) \widehat{\mathcal{J}}^{-1}(\widehat{\boldsymbol{\beta}}) \left[ \int_0^\infty \{ \mathbf{Z}_{ik} - \overline{Z}(\widehat{\boldsymbol{\beta}}, s) \} d\widehat{M}_{ik}(s) \right] \right\}, \end{aligned}$$

qui résulte du remplacement, dans l'expression initiale de  $\Psi_i$ , de toutes les quantités inconnues par leur estimations basées sur l'échantillon. Il est possible de démontrer que  $\widetilde{\Psi}_i$  est consistant.

Notons

$$\widehat{W}(t, z) = \sum_{i=1}^n \widehat{\Psi}_i(t, z) G_i,$$

où  $(G_1, \dots, G_n)$  sont des v.a. indépendantes entre elles, suivant des lois normales centrées et de variance 1, et indépendantes des données  $\{Y_{ik}, N_{ik}, \mathbf{Z}_{ik}\}$ .

Conditionnellement à  $\{Y_{ik}, N_{ik}, \mathbf{Z}_{ik}\}$ ,  $\widehat{W}(t, z)$  est la somme de  $n$  vecteurs indépendants et centrés, quel que soit le couple  $(t, z)$ . Le théorème de Lindeberg-Feller permet d'affirmer que les distributions conditionnelles de  $1/\sqrt{n} \widehat{W}(\cdot, \cdot)$  sont asymptotiquement gaussiennes et centrées.

Par ailleurs, la fonction de covariance de  $1/\sqrt{n} \widehat{W}(\cdot, \cdot)$ , conditionnellement à  $\{Y_{ik}, N_{ik}, \mathbf{Z}_{ik}\}$ , est simplement

$$\sum_{i=1}^n \widehat{\Psi}_i(t, z) \widehat{\Psi}_i(t^b, z^b)$$

et cette quantité converge vers la fonction de covariance déterministe de  $1/\sqrt{n} W(\cdot, \cdot)$ .

Ainsi,  $1/\sqrt{n} W(\cdot, \cdot)$  et  $1/\sqrt{n} \widehat{W}(\cdot, \cdot)$  convergent faiblement vers le même processus limite.

Pour approcher la distribution nulle de  $W(\cdot, \cdot)$ , on simule un certain nombre de réalisations à partir de  $\widehat{W}(\cdot, \cdot)$  en répétant la génération d'un échantillon aléatoire  $(G_1, \dots, G_n)$ , tandis que les données observées demeurent fixes.

### 3.3.1 Test de la forme fonctionnelle d'une covariable

On représente les sommes cumulées de  $\widehat{M}_{ik} = \widehat{M}_{ik}(\infty)$  en fonction des valeurs d'une covariable.

Le processus de somme partielle pour la  $j^{\text{e}}$  composante du vecteur des covariables est

$$W_j(x) = \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}_{\{Z_{jik} \leq x\}} \widehat{M}_{ik},$$

qui est un cas particulier de  $W(t, z)$ : ici,  $f(\cdot) = 1$ ,  $t = \infty$ ,  $\mathbf{Z}_j = x$  et  $\mathbf{Z}_k = 0$  pour tout  $k \neq j$ .

La distribution nulle des  $W_j$  peut être approchée au travers de la simulation du processus gaussien centré  $\widehat{W}_j$  correspondant. Afin d'étudier le profil résiduel sous le modèle adopté, on représente les résidus issus d'une vingtaine de réalisations de la distribution du processus  $\widehat{W}_j$ .

Pour accroître l'objectivité, on peut calculer un degré de signification (*p-value*) approximatif pour le test du maximum

$$\sup_x |W_i(x)|$$

en utilisant un grand nombre de réalisations de  $\widehat{W}_j$ .

### 3.3.2 Test de la fonction de lien

Pour  $p > 1$ , il peut être également souhaitable de tester la fonction de lien, c'est-à-dire la forme de régression exponentielle. Ce test peut être réalisé au travers de la quantité

$$\sup_x \left| \sum_{i=1}^n \sum_{k=1}^K W_i(x) \mathbb{1}_{\{\mathbf{Z}_{ik} \leq z\}} \widehat{M}_{ik} \right|.$$

Ce test, consistant, repose sur l'hypothèse alternative que

$$\lambda_k(t; \mathbf{Z}_k) = \lambda_0(t)g(\mathbf{Z}_k)$$

pour une certaine fonction  $g$ ; cependant, il n'existe pas de  $\beta$  pour lequel  $g(\mathbf{Z})/\exp(\beta^t \mathbf{Z})$  soit une constante quelles que soient les valeurs possibles de  $\mathbf{Z}$ .

### 3.3.3 Test de la proportionnalité des risques vis-à-vis d'une covariable

On considère le processus du score

$$U(\hat{\beta}, t) = \sum_{i=1}^n \sum_{k=1}^K \int_0^t [\mathbf{Z}_{ik} - \bar{\mathbf{Z}}(\hat{\beta}, s)] dN_{ik}(s),$$

qui est un autre cas spécial de  $W(t, z)$  avec  $z = \infty$  et  $f(x) = x$ .

Là encore, on simule les distributions des processus standardisés

$$\frac{1}{\sqrt{n}} \sqrt{\Omega^{-1}(\hat{\beta})_{jj}} U_j(\hat{\beta}, t) \quad (j = 1, \dots, p)$$

avec  $\Omega(\beta)$  définie en p. 55, et où  $U_j$  est la  $j^e$  composante de  $U$  et  $\Omega^{-1}(\hat{\beta})_{jj}$  le  $j^e$  élément diagonal de l'inverse de la matrice  $\Omega(\cdot)$ .

Une inspection, aussi bien graphique que numérique, de l'hypothèse de proportionnalité des risques en regard de la  $j^e$  composante du vecteur des covariables peut être ainsi menée.

### 3.3.4 Test de la proportionnalité des risques vis-à-vis de toutes les covariables

Le test

$$\sup_t \frac{1}{\sqrt{n}} \sum_{j=1}^p \sqrt{\Omega^{-1}(\hat{\beta})_{jj}} |U_j(\hat{\beta}, t)|$$

fournit un moyen de vérifier l'hypothèse globale de proportionnalité des risques, ce test étant consistant contre l'alternative

$$\lambda_k(t; \mathbf{Z}_k) = \lambda_0(t) e^{\theta(t)^t \mathbf{Z}_k} \quad (k = 1, \dots, K),$$

où  $\theta(t)$  est dépendante du temps.

### 3.3.5 Test concernant le choix même du modèle marginal

La statistique

$$\sup_{t, z} |W(t, z)|,$$

avec  $f(x) = 1$  (cf. p. 55), procure un test consistant contre l'hypothèse alternative de mauvaise spécification du modèle marginal de Cox.

### 3.4 Coefficient de dispersion (*design effect*)

L'expression **coefficient de dispersion** (*design effect*) – noté CD par la suite – désigne le rapport de la variance robuste d'un estimateur sur sa variance naïve. Dans le cadre du modèle linéaire, Scott et Holt (1982) ont développé une approche tendant à évaluer analytiquement le coefficient de dispersion relatif aux coefficients de régression.

Dans le cadre d'un modèle linéaire portant sur des groupes de même taille  $n$ , n'incluant qu'une covariable  $X$ , et dans lequel les variables réponses sont corrélées de la même manière à l'intérieur de chaque groupe (matrice d'équicorrélation (*exchangeable*) ou symétrique par composition (*compound symmetric*)), si nous notons  $\rho_Y$  le coefficient de corrélation de la variable réponse à l'intérieur d'un groupe (**coefficient de corrélation intraclasse**), soit

$$\rho_Y = \frac{\text{Cov}(Y_{ij}, Y_{ik})}{\mathbb{V}(Y_{ij})},$$

nous avons alors

$$\begin{aligned} \text{CD} &= \frac{\mathbb{V}(\hat{\beta})}{\mathbb{V}(\hat{\beta}_{\text{indep}})} \\ &= 1 + (n-1)\hat{\rho}_X\rho_Y, \end{aligned} \tag{3.4}$$

où  $\hat{\rho}_X$  est l'estimation du coefficient de corrélation intraclasse de la covariable  $X$ .

*Démonstration* — Notons  $\sigma^2\mathbf{V}$  la matrice diagonale par blocs de  $\mathbf{Y}$ , où chaque bloc correspond à un groupe. Nous savons que

$$\mathbb{V}(\hat{\beta}) = \sigma^2(\mathbf{X}^t\mathbf{X})^{-1}(\mathbf{X}^t\mathbf{V}\mathbf{X})(\mathbf{X}^t\mathbf{X})^{-1}$$

et que

$$\mathbb{V}(\hat{\beta}_{\text{indep}}) = \sigma^2(\mathbf{X}^t\mathbf{X})^{-1}.$$

Ainsi,

$$\mathbb{V}(\hat{\beta}) = \mathbb{V}(\hat{\beta}_{\text{indep}}) \mathbf{D}$$

avec

$$\mathbf{D} = (\mathbf{X}^t\mathbf{V}\mathbf{X})(\mathbf{X}^t\mathbf{X})^{-1}.$$

Dans le cas où les réponses sont corrélées de la même manière à l'intérieur de chaque groupe, l'expression de la matrice  $\mathbf{D}$  se simplifie, pour devenir

$$\mathbf{D} = \mathbf{I} + (\mathbf{M} - \mathbf{I})\rho_Y,$$

où  $\mathbf{M}$  est une matrice ne dépendant que de  $\mathbf{X}$ .

Par suite, dans le cas d'une unique covariable de paramètre  $\beta$ , et de groupes d'effectif commun  $n$ , nous obtenons

$$\mathbb{V}(\hat{\beta}) = \mathbb{V}(\hat{\beta}_{\text{indep}})[1 + (n-1)\hat{\rho}_X\rho_Y].$$

■

### 3.4. COEFFICIENT DE DISPERSION (DESIGN EFFECT)

---

Dans le cadre d'un modèle de régression linéaire portant sur  $n$  groupes de  $m$  individus chacun :

$$Y_{ij} = \alpha + \beta(x_{ij} - \bar{x}_{..}) + \epsilon_{ij} ,$$

si  $\text{Cov}(\epsilon) = \sigma^2 \mathbf{V}$  avec  $\mathbf{V}$  matrice diagonale par blocs, nous avons

$$\mathbf{M} = \begin{pmatrix} m & 0 \\ 0 & m \frac{K_x}{L_x} \end{pmatrix}$$

et

$$\hat{\rho}_x = \frac{1}{m-1} \left( m \frac{K_x}{L_x} - 1 \right)$$

où

$$K_x = m \sum_{i=1}^n (\bar{x}_{i.} - \bar{x}_{..})^2$$

et

$$L_x = \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - \bar{x}_{..})^2 .$$

Neuhaus et Segal (1993) ont démontré que l'équation (3.4) était encore valable dans le cadre d'un modèle linéaire généralisé pour données binaires.

Cependant aucun résultat analytique similaire n'a été démontré dans le cas d'un modèle de Cox.

## Chapitre 4

# MODÈLE DE FRAGILITÉ DE COX

### Contenu

---

<b>4.1</b>	<b>Écriture du modèle</b> . . . . .	<b>62</b>
<b>4.2</b>	<b>Concernant la fragilité</b> . . . . .	<b>63</b>
4.2.1	Différentes distributions pour la fragilité . . . . .	63
<b>4.3</b>	<b>Cas où la fragilité suit une loi gamma</b> . . . . .	<b>66</b>
4.3.1	Écriture du modèle . . . . .	66
4.3.2	Estimation par l'algorithme E.M. . . . .	68
<b>4.4</b>	<b>Tests portant sur la variance de la fragilité</b> . . . . .	<b>69</b>
4.4.1	Cas d'un effectif de groupes important . . . . .	69
4.4.2	Cas d'un effectif de groupes faible . . . . .	71
<b>4.5</b>	<b>Compléments</b> . . . . .	<b>71</b>
4.5.1	Concernant les modèles fréquentistes . . . . .	71
4.5.2	Concernant l'approche bayésienne . . . . .	72
<b>4.6</b>	<b>Prédiction de la fragilité</b> . . . . .	<b>73</b>
4.6.1	La théorie du meilleur prédicteur . . . . .	73
4.6.2	Tests d'hypothèses . . . . .	77
4.6.3	Dans le cadre du modèle de fragilité de Cox . . . . .	77
4.6.3.1	Cas d'une fragilité gamma . . . . .	78

---

## 4.1 Écriture du modèle

Les notations de la page 52 sont reprises ici.

Le **modèle de fragilité** (ou **modèle mixte**) de Cox spécifie que le risque instantané pour l'individu  $j$  du groupe  $i$  s'écrit

$$\begin{aligned}\tilde{\alpha}_{ij}(t) &= \Theta_i \times \alpha_{ij}(t) \\ &= \alpha_0(t) \Theta_i \exp(\beta^t \mathbf{Z}_{ij}),\end{aligned}$$

où  $\Theta_i$  est une v.a. représentant l'effet du groupe  $i$ .

L'introduction dans la fonction de risque du terme  $\Theta_i$  a généralement pour objectif d'accroître le risque de survenue de l'événement étudié chez un individu du groupe  $i$  (Vaupel *et al.*, 1979). La « fragilité » d'un tel individu étant alors plus grande, nous trouvons là, pour ce type de modèle, la qualification de « **modèle de fragilité** ».

La v.a.  $\Theta_i$  est désignée comme étant la composante de fragilité du modèle.

Une hypothèse fondamentale est posée : conditionnellement à  $\Theta_i$  ( $i \in \{1, \dots, n\}$ ), les observations du groupe  $i$  sont supposées indépendantes.

La vraisemblance vaut

$$\mathcal{L}(\beta) = \prod_{i,j} \left\{ \prod_t (\Theta_i Y_{ij}(t) \alpha_{ij}(t))^{\Delta N_{ij}(t)} \exp \left( - \Theta_i \int_0^\tau Y_{ij}(s) dA_{ij}(s) \right) \right\}, \quad (4.1)$$

et par suite

$$\mathcal{L}(\beta) = \prod_{i,j} \int_0^{+\infty} \left\{ \prod_t (\theta_i Y_{ij}(t) \alpha_{ij}(t))^{\Delta N_{ij}(t)} \exp \left( - \theta_i \int_0^\tau Y_{ij}(s) dA_{ij}(s) \right) \right\} p(\theta_i) d\theta_i, \quad (4.2)$$

où  $p$  est la densité de probabilité de  $\Theta_i$ .

*A priori*, n'importe quelle distribution peut être affectée à la fragilité. Cependant, cette distribution est généralement choisie parmi les distributions appartenant à la **famille exponentielle** – famille présentant de nombreuses propriétés mathématiques.

En particulier, les quatre distributions les plus fréquemment rencontrées dans la littérature sont la distribution **log-normale**, la distribution **inverse gaussienne**, la distribution **gamma** et la distribution **stable positive**.

Il est possible – comme dans le cadre d'une distribution log-normale de  $\Theta_i$  – d'introduire le terme de fragilité dans la partie exponentielle. Le modèle s'écrit alors

$$\tilde{\alpha}_{ij}(t) = \alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij} + \tilde{\Theta}_i),$$

avec  $\tilde{\Theta}_i = \log(\Theta_i)$  distribuée suivant une loi normale.

Le cas d'une distribution gamma sera développé plus loin.

Notons que l'identifiabilité du modèle nécessite l'adoption de la contrainte  $\mathbb{E}(\Theta_i) = 1$  (Elbers et Ridder, 1982).

**Remarque** — Dans ce même article, Elbers et Ridder (1982) ont montré qu'une fragilité de moyenne finie peut être identifiée à partir des données marginales ; pratiquement, dans le cas de données concernant des jumeaux, ceci signifie que nous pouvons estimer le paramètre  $\gamma$  de la distribution de la fragilité à partir de la connaissance d'un seul des deux jumeaux... Ainsi, la fragilité semble-t-elle décrire un peu plus que la seule corrélation existant entre deux (ou plusieurs) individus.

## 4.2 Concernant la fragilité

### 4.2.1 Différentes distributions pour la fragilité

Nous reprenons ici les notations et les termes employés par Hougaard (1984). En particulier, toute distribution considérant la fragilité  $\Theta$  comme fixe est appelée **distribution conditionnelle**, et toute distribution observée – c'est-à-dire après intégration suivant la loi de  $\Theta$  – est appelée **distribution marginale**.

La fonction de survie conditionnelle (sachant  $\Theta$ ) est

$$S(t|\Theta) = \exp\{-\Theta A(t)\}.$$

La fonction de survie de la population (*population survivor function*) – fonction de survie marginale d'un individu quelconque – est

$$\begin{aligned} S(t) &= \int \exp[-\theta A(t)] f(\theta) d\theta \\ &= \mathcal{L}_\theta[A(t)], \end{aligned}$$

où

$$\mathcal{L}_\theta(s) = \mathbb{E}[\exp(-s\theta)]$$

est la transformée de Laplace de la distribution de la fragilité.

Le risque cumulé de la population (*integrated population hazard*) est

$$H(t) = -g[A(t)], \tag{4.3}$$

où  $g(s) = \log[\mathcal{L}_\theta(s)]$ .

Le risque instantané de la population (*population hazard*) – le risque instantané marginal – est défini comme étant

$$h_{ij}(t) = \alpha_{ij}(t) \mathbb{E}(\Theta_i | T_{ij} > t),$$

que nous pouvons interpréter ainsi : le risque observé à l'instant  $t$  est égal au risque « moyenné » sur l'ensemble des personnes toujours en vie à cet instant  $t$ .

D'après (4.3), ce risque vaut

$$h(t) = -\alpha(t) g'[A(t)] .$$

Nous avons la relation<sup>1</sup>

$$S(t) = \exp [ - H(t) ] .$$

**Loi gamma** Si  $\Theta_i$  suit une loi  $\Gamma(\delta, \gamma)$ , sa distribution est

$$p(\theta_i) = \frac{\gamma^\delta}{\Gamma(\delta)} e^{-\gamma\theta_i} \theta_i^{\delta-1} \mathbf{1}_{\{\theta_i>0\}} .$$

Son espérance et sa variance valent respectivement  $\delta/\gamma$  et  $\delta/\gamma^2$ .

Le risque instantané marginal vaut dans ce cas

$$\frac{\delta\alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij})}{\gamma + A_0(t) \exp(\beta^t \mathbf{Z}_{ij})} .$$

Si, de plus, nous contraignons l'espérance de la distribution de la fragilité à valoir 1 (*i.e.*  $\delta = \gamma$ ), alors ce risque vaut

$$\frac{\alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij})}{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{ij})} .$$

Ces risques marginaux ne sont pas proportionnels; plus précisément, le rapport de deux d'entre eux vaudra bien, à l'origine, la valeur *correcte*<sup>2</sup>  $\exp[\beta^t(Z_{i1} - Z_{i2})]$ , mais ce rapport convergera par la suite de façon monotone vers 1 lorsque  $t$  tendra vers l'infini (ou plus justement lorsque  $A_0(t)$  tendra vers l'infini).

En effet, si  $\alpha_{i1}(t)/\alpha_{i2}(t) = c$ , où  $c$  est une constante positive, alors

$$\begin{aligned} \frac{h_{i1}(t)}{h_{i2}(t)} &= \left( \frac{\alpha_0(t) \exp(\beta^t \mathbf{Z}_{i1})}{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{i1})} \right) / \left( \frac{\alpha_0(t) \exp(\beta^t \mathbf{Z}_{i2})}{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{i2})} \right) \\ &= \frac{\alpha_{i1}(t)}{\alpha_{i2}(t)} \times \frac{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{i2})}{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{i1})} \\ &= c \times \frac{\gamma + A_0(t) \exp(\beta^t \mathbf{Z}_{i2})}{\gamma + A_0(t) \exp(\beta^t \mathbf{Z}_{i1})} , \end{aligned}$$

qui tend vers  $c/c = 1$  lorsque  $A_0(t)$  tend vers l'infini.

**Loi inverse gaussienne** Si  $\Theta_i$  suit une loi inverse gaussienne de paramètres  $\delta$  et  $\gamma$ , avec  $\gamma \geq 0$  et  $\delta > 0$ , sa densité vaut

$$g(z) = \sqrt{\frac{\delta}{\pi z^3}} \exp \left[ 2\sqrt{\delta\gamma} - \gamma z - \frac{\delta}{z} \right] \mathbf{1}_{\mathbb{R}^{++}}(z) .$$

---

1. Elle est le pendant de la relation valable en l'absence de corrélation, à savoir  $S(t) = \exp [ - A(t) ]$ .

2. Le terme « correcte » renvoie à la valeur du rapport des risques conditionnels.

L'espérance et la variance de cette loi valent respectivement  $\sqrt{\delta/\gamma}$  et  $\sqrt{\delta/(4\gamma^3)}$ .

Si nous contraignons l'espérance de cette loi à valoir 1 – ce qui entraîne  $\delta = \sqrt{\gamma}$  –, le risque instantané marginal vaut

$$\frac{\alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij})}{\sqrt{1 + \frac{1}{\gamma} A_0(t) \exp(\beta^t \mathbf{Z}_{ij})}}.$$

Là encore, ces risques marginaux ne sont pas proportionnels. Par un calcul similaire au précédent, nous obtenons une convergence monotone du rapport de deux risques marginaux vers  $\sqrt{c}$  lorsque  $A_0(t)$  tend vers l'infini.

**Loi positive stable** Supposons que  $\Theta_i$  suive une loi positive stable de paramètres  $\delta$  et  $\gamma$ ,  $0 < \delta < 1$ ,  $0 \leq \gamma \leq 1$ . Sa densité vaut

$$h(z) = -\frac{1}{\pi z} \sum_{k=1}^{\infty} \frac{\Gamma(k\gamma + 1)}{k!} \left( -\frac{\delta}{\gamma} z^{-\gamma} \right)^k \sin(\gamma k \pi) \mathbb{1}_{\mathbb{R}^{+*}}(z).$$

La moyenne de cette distribution est infinie ; par contre, pour  $\alpha < 1$ , les moments d'ordre  $q$  (pour  $q < \alpha$ ) existent et valent

$$\left( \frac{\delta}{\gamma} \right)^{\frac{q}{\gamma}} \frac{\Gamma(1 - \frac{q}{\gamma})}{\Gamma(1 - q)}.$$

Le cas standard – dans lequel nous nous placerons toujours par la suite – est celui où  $\delta = \gamma$ . Dans ce cas-là, le risque marginal vaut

$$\alpha_0(t) \exp(\delta \beta^t \mathbf{Z}_{ij}) \delta A_0(t)^{\delta-1},$$

fonction de risque proportionnel de la forme  $\omega_0(t) \exp(\tilde{\beta}^t \mathbf{Z}_{ij})$ , avec  $\omega_0(t) = \delta \alpha_0(t) A_0(t)^{\delta-1}$  et  $\tilde{\beta} = \delta \beta$ .

Ceci implique la consistance du modèle vis-à-vis des covariables inobservées, représentées par la fragilité. En revanche, la fonction de risque et les coefficients de régression ne sont pas consistants ; nous remarquons en particulier que les coefficients de régression sont numériquement plus faibles : ceci peut être considéré comme un biais (Hougaard (2000), p. 91).

Toutefois, en pratique, l'ajustement sera conduit dans l'ignorance de l'existence de ce biais, et le coefficient de régression estimé sera plus proche de 0 que celui du modèle de régression conditionnel.

Notons que la loi stable positive possède une propriété intéressante (Hougaard, 1986) pour la simulation de données corrélées, propriété qui sera énoncée lors de la comparaison des différents modèles pour données de survie corrélées (cf. p. 94).

### 4.3 Cas où la fragilité suit une loi gamma

#### 4.3.1 Écriture du modèle

$\Theta_i$  est supposée suivre une loi gamma  $\Gamma(\frac{1}{\gamma}, \frac{1}{\gamma})$ . En reprenant l'équation (4.2), nous pouvons écrire la vraisemblance sous la forme

$$\begin{aligned} \mathcal{L}(\beta, \gamma) &= \left\{ \prod_{i,j} \prod_t [Y_{ij}(t)\alpha_{ij}(t)]^{\Delta N_{ij}(t)} \right\} \\ &\quad \times \left\{ \prod_i \prod_t \int_0^{+\infty} \theta_i^{\sum_k \Delta N_{ik}(t)} \exp\left(-\theta_i \sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s)\right) p(\theta_i) d\theta_i \right\} \end{aligned} \quad (4.4)$$

$$\begin{aligned} &= \frac{1}{\gamma^{1/\gamma} \Gamma(\frac{1}{\gamma})} \times \left\{ \prod_{i,j} \prod_t [Y_{ij}(t)\alpha_{ij}(t)]^{\Delta N_{ij}(t)} \right\} \times \\ &\quad \left\{ \prod_i \prod_t \int_0^{+\infty} \theta_i^{\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma} - 1} \exp\left(-\theta_i \left[\sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right]\right) d\theta_i \right\} \end{aligned} \quad (4.5)$$

$$\begin{aligned} &= \left\{ \prod_{i,j} \prod_t [Y_{ij}(t)\alpha_{ij}(t)]^{\Delta N_{ij}(t)} \right\} \left\{ \prod_i \prod_t \frac{\Gamma(\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma})}{\gamma^{1/\gamma} \Gamma(\frac{1}{\gamma}) \left(\sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right)^{\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma}}} \right\} \\ &= \prod_i \prod_t \left\{ \frac{\Gamma(\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma})}{\gamma^{1/\gamma} \Gamma(\frac{1}{\gamma}) \left(\sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right)^{\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma}}} \prod_j [Y_{ij}(t)\alpha_{ij}(t)]^{\Delta N_{ij}(t)} \right\}. \end{aligned} \quad (4.6)$$

En tant que fonction de  $\theta_i$ , (4.5) est proportionnelle à

$$\prod_i \prod_t \left\{ \theta_i^{\sum_k \Delta N_{ik}(t) + \frac{1}{\gamma} - 1} \exp\left[-\theta_i \left(\sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right)\right] \right\},$$

soit à

$$\prod_i \left\{ \theta_i^{\sum_k N_{ik}(\tau) + \frac{1}{\gamma} - 1} \exp\left[-\theta_i \left(\sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right)\right] \right\}.$$

En d'autres termes, conditionnellement aux données, les v.a.  $\Theta_i$  sont indépendantes et distribuées suivant des lois

$$\Gamma\left(\sum_k N_{ik}(\tau) + \frac{1}{\gamma}, \sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}\right).$$

Posons

$$D_i = \sum_k N_{ik}(\tau),$$

$$F_i = D_i + \frac{1}{\gamma}$$

et

$$G_i = \sum_k \int_0^\tau Y_{ik}(s) dA_{ik}(s) + \frac{1}{\gamma}.$$

Supposons que les  $\Theta_i$  soient observées. D'après (4.1), la log-vraisemblance totale vaut (à un terme indépendant des paramètres d'intérêt près)

$$\mathcal{L}(\boldsymbol{\beta}, \alpha_0, \gamma; (\theta_i)_i) = \mathcal{L}_1(\gamma) + \mathcal{L}_2(\boldsymbol{\beta}, \alpha_0). \quad (4.7)$$

où

$$\mathcal{L}_1(\gamma) = -n \left[ \left( \frac{1}{\gamma} \right) \log(\gamma) + \log \left[ \Gamma \left( \frac{1}{\gamma} \right) \right] \right] + \sum_{i=1}^n \left\{ \left( \frac{1}{\gamma} + D_i - 1 \right) \log(\theta_i) - \frac{\theta_i}{\gamma} \right\}$$

et

$$\mathcal{L}_2(\boldsymbol{\beta}, \alpha_0) = \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \left\{ Y_{ij}(t) (\boldsymbol{\beta}' \mathbf{Z}_{ij} + \log[\alpha_0(t)]) - \theta_i A_0(t) \exp(\boldsymbol{\beta}' \mathbf{Z}_{ij}) \right\} dt.$$

Puisque les  $\Theta_i$  sont, conditionnellement aux données, distribuées suivant des lois  $\Gamma(F_i, G_i)$ , l'espérance de  $\mathcal{L}$ , étant données les valeurs des observations ainsi que les valeurs courantes de  $F_i$  et de  $G_i$ , vaut désormais  $\mathcal{L} = \mathcal{L}_1(\gamma) + \mathcal{L}_2(\boldsymbol{\beta}, \alpha_0)$  où

$$\mathcal{L}_1(\gamma) = -n \left[ \left( \frac{1}{\gamma} \right) \log(\gamma) + \log \left[ \Gamma \left( \frac{1}{\gamma} \right) \right] \right] + \sum_{i=1}^n \left\{ \left( \frac{1}{\gamma} + D_i - 1 \right) [\psi(F_i) - \log(G_i)] - \frac{F_i/G_i}{\gamma} \right\} \quad (4.8)$$

et

$$\mathcal{L}_2(\boldsymbol{\beta}, \alpha_0) = \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \left\{ Y_{ij}(t) [\boldsymbol{\beta}' \mathbf{Z}_{ij} + \log[\alpha_0(t)]] - \frac{F_i}{G_i} A_0(t) \exp(\boldsymbol{\beta}' \mathbf{Z}_{ij}) \right\} dt, \quad (4.9)$$

où  $\psi(\cdot)$  est la fonction digamma :  $\psi(x) = \partial \log [\Gamma(x)] / \partial x$ , où  $\Gamma(x)$  est la fonction gamma d'Euler.

### 4.3.2 Estimation par l'algorithme E.M.

L'estimation des paramètres nécessite le plus souvent l'emploi d'algorithmes, le calcul direct se révélant inexécutable. Dans le cadre du modèle de fragilité de Cox, l'algorithme le plus utilisé est l'algorithme E.M. (*Estimation - Maximization*) (Gill, 1985; Klein, 1992; Andersen *et al.*, 1991; Nielsen *et al.*, 1992). L'étape de maximisation (*M-step*) requière la maximisation de (4.8) et de (4.9) par rapport aux paramètres  $\gamma$  et  $\beta$ . Concernant l'étape d'estimation (*E-step*), la mise à jour de l'estimation de  $\gamma$  implique de maximiser numériquement (4.8). La mise à jour de l'estimation de  $\beta$  implique la maximisation de (4.9); cette log-vraisemblance contient  $\alpha_0$ , qui peut être considéré comme un paramètre de nuisance. Aussi, dans un premier temps, nous fixons la valeur de  $\beta$  et calculons l'estimateur de Breslow

$$\hat{A}_0(t) = \int_0^t \frac{\sum_{i,j} N_{ij}(u)}{\sum_{i,j} \hat{\theta}_i \exp(\beta' \mathbf{Z}_{ij})} . \quad (4.10)$$

Substituant (4.10) dans (4.9), nous obtenons comme log-vraisemblance

$$\mathcal{L}_3(\beta) = \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \left\{ Y_{ij}(t) [\beta' \mathbf{Z}_{ij} + \log[\tilde{\theta}_i]] - \log \left[ \sum_i \hat{\theta}_i \exp(\beta' \mathbf{Z}_{ij}) \right] \right\} dN_{ij}(t) . \quad (4.11)$$

Notons que cette log-vraisemblance partielle est celle d'un modèle de Cox ordinaire avec inclusion d'une covariable supplémentaire  $\log(\tilde{\theta}_i)$ . Ainsi, un programme ajustant le modèle de Cox permet d'obtenir une mise à jour de l'estimateur de  $\beta$ .

Résumons l'utilisation de l'algorithme E.M. :

- étape 1** utilisation d'un programme standard du modèle de Cox afin d'obtenir des estimations initiales de  $\beta$  et  $\alpha_0$  *via* respectivement les équations (4.11) et (4.10), avec  $\hat{\theta} = 1$  (*i.e.*  $\gamma = 0$ );
- étape 2** utilisation des valeurs courantes de  $\gamma$ ,  $\beta$  et  $\alpha_0$  afin d'obtenir  $F_i, G_i$  et  $\hat{\theta}_i$ ;
- étape 3** mise à jour de l'estimation de  $\gamma$  *via* (4.8) d'une part, de celles de  $\beta$  et de  $\alpha_0$  *via* (4.11) et (4.10) d'autre part;
- étape 4** itération des étapes 2 et 3 jusqu'à l'obtention de la convergence.

Une fois la convergence obtenue, la matrice d'information observée  $\mathcal{J}(\hat{\beta}, \hat{\gamma})$  peut être calculée. Nous donnons en annexe l'expression des composantes de cette matrice.

En négligeant la dépendance de  $\hat{\Theta}_i$  vis-à-vis de tous les paramètres, la procédure de l'algorithme E.M. telle qu'elle a été présentée plus haut tend à produire des estimations de la variance trop faibles. Andersen *et al.* (1997) ont, à cet effet, corrigé l'estimateur de la variance: le calcul de la matrice d'information porte alors sur  $\mathcal{J}(\hat{\beta}, \hat{\gamma}, \hat{\alpha}_0)$ , et non plus seulement sur  $\mathcal{J}(\hat{\beta}, \hat{\gamma})$ .

La base théorique sous-tendant ce modèle, qui concerne l'emploi de la matrice d'information pour l'estimation des paramètres de la variance, a été établie par Parner (1998) – qui reprenait les résultats de Murphy (1994; 1995) concernant la consistance et la normalité asymptotique des estimateurs des paramètres en l'absence de covariable. Dans ses travaux, Parner a démontré la consistance et la normalité asymptotique de  $\hat{A}_0(t)$ ,  $\hat{\gamma}$  et  $\hat{\beta}$ , ainsi que la consistance des estimateurs de la variance basés sur l'inverse de la matrice d'information, dans le cadre d'un modèle de Cox comportant des covariables.

Signalons enfin qu'il est également possible d'utiliser l'algorithme de Newton-Raphson lors de l'étape d'estimation (Shih et Louis, 1995).

**Remarque** — Pour une valeur de  $\gamma$  donnée, (4.4) est la vraisemblance totale ou partielle pour  $\beta$  basée sur  $(N_i, Y_i, \Theta_i)_i$  selon que conditionnellement à  $\Theta_i = \theta_i$ , la censure est non-informative ou informative pour  $\beta$ .

Par suite, si l'on intègre suivant  $\theta$  cette expression, on obtient respectivement soit la vraisemblance complète pour  $\beta$  basée sur  $(N_i, Y_i)_i$ , soit ce que l'on appellera la **vraisemblance partielle marginale**.

Sous l'hypothèse que « conditionnellement à  $\Theta_i = \theta_i$ , la censure est non-informative pour  $\theta_i$  », cette vraisemblance partielle marginale est égale à la **vraisemblance marginale partielle**, c'est-à-dire la vraisemblance partielle pour  $\beta$  basée sur  $(N_i, Y_i)_i$  dans le cadre « marginal » (ou de « données incomplètes ») quand  $\Theta_i$  n'est pas observée (*i.e.* le cas qui se présente en réalité).

Gill (1992) a démontré que de manière générale, il y a égalité entre vraisemblance partielle marginale et vraisemblance marginale partielle lorsque la censure n'est pas informative vis-à-vis de la variable suivant laquelle on intègre.

Ainsi pouvons-nous baser l'inférence concernant  $(\beta, \gamma)$  sur la vraisemblance marginale (non-conditionnelle à la fragilité) partielle (censure omise) des temps de survenue, qui est la même que la vraisemblance partielle marginale obtenue en écrivant la vraisemblance partielle usuelle pour  $(\beta, \gamma)$  valable dans le cas où  $\Theta_i$  est observable, puis intégrée suivant  $\theta_i$ .

La validité de cette égalité permet également de connaître la distribution de la fragilité conditionnellement aux données à partir de la vraisemblance partielle, et de faire usage de l'algorithme E.M. lors du calcul de l'estimateur de  $(\beta, \gamma)$  obtenu par la méthode du maximum de vraisemblance partielle.

## 4.4 Tests portant sur la variance de la fragilité

### 4.4.1 Cas d'un effectif de groupes important

Dans le cas où le nombre de groupes est important, il est possible (Gray, 1995 ; Commenges et Andersen, 1995) de tester l'effet groupe, en utilisant un test du score concernant l'hypothèse de nullité de la variance de la fragilité ; ce test fonctionne en l'absence de toute spécification de la distribution de la fragilité.

Notons

$$\alpha_0(s) Y_{ij}(s) \exp(\Theta_i + \beta^t \mathbf{Z}_{ij})$$

le processus d'intensité, où les  $\Theta_i$  sont i.i.d. de loi de distribution (inconnue)  $G$  centrée et de variance 1.

Notons également

$$\begin{aligned} N_{..} &= \sum_{i,j} N_{ij} , \\ S_i^{(0)}(\boldsymbol{\beta}, s) &= \sum_j Y_{ij}(s) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ij}) , \\ \bar{S}^{(0)}(\boldsymbol{\beta}, s) &= \sum_i S_i^{(0)}(\boldsymbol{\beta}, s) , \\ p_{ij}(\boldsymbol{\beta}, s) &= \frac{Y_{ij}(s) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ij})}{\bar{S}^{(0)}(\boldsymbol{\beta}, s)} \end{aligned}$$

et

$$p_i(\boldsymbol{\beta}, s) = \sum_j p_{ij}(\boldsymbol{\beta}, s) .$$

Désignons enfin par  $\hat{\boldsymbol{\beta}}$  l'estimateur du maximum de vraisemblance (partielle) sous l'hypothèse de nullité de la variance de la fragilité.

Pour tester cette hypothèse d'homogénéité, la statistique employée est celle du test du score ; elle est donnée par

$$T(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \left( \sum_{j=1}^{n_i} \widehat{M}_{ij}(\infty) \right)^2 - \bar{N}(\infty) + \int_0^\infty \sum_{i=1}^n p_i^2(\hat{\boldsymbol{\beta}}, s) d\bar{N}(s) , \quad (4.12)$$

avec les notations vues en page 52.

Remarquons que le premier terme de (4.12) est de la forme

$$\sum_i (\text{Observés}_i - \text{Attendus}_i)^2 ,$$

tandis que les deux autres termes ne servent qu'à assurer la nullité de l'espérance de la statistique.

Soit

$$H_i(\boldsymbol{\beta}, s) = 2 \left[ \widehat{M}_i(s) - \sum_{l=1}^n \widehat{M}_l(s-) p_l(\boldsymbol{\beta}, s) - p_i(\boldsymbol{\beta}, s) + \sum_{l=1}^n p_l^2(\boldsymbol{\beta}, s) \right] ,$$

où  $\widehat{M}_i(s) = \sum_j \widehat{M}_{ij}(s)$ .

La variance de  $T(\hat{\boldsymbol{\beta}})$  peut être estimée de façon consistante par

$$\hat{J}_c = \hat{J}(\hat{\boldsymbol{\beta}}) - \hat{J}(\hat{\boldsymbol{\beta}}) \mathcal{J}_{\hat{\boldsymbol{\beta}}}^{-1} \hat{J}(\hat{\boldsymbol{\beta}})^t ,$$

où  $\mathcal{J}_{\hat{\boldsymbol{\beta}}}^{-1}$  est la matrice d'information de Fisher relative à  $\hat{\boldsymbol{\beta}}$ ,

$$\hat{J}(\boldsymbol{\beta}) = \sum_{i=1}^n \int_0^\infty H_i^2(\boldsymbol{\beta}, s) p_i(\boldsymbol{\beta}, s) d\bar{N}(s) ,$$

et où

$$\widehat{J}(\widehat{\beta}) = \sum_{i=1}^n \int_0^{\infty} H_i(\beta, s) \left[ \sum_{j=1}^{n_i} z_{ij} p_{ij}(\beta, s) \right] d\overline{N}(s).$$

La statistique de l'hypothèse d'homogénéité est

$$H = \frac{T(\widehat{\beta})}{\sqrt{\widehat{J}_c}},$$

et elle suit asymptotiquement, sous l'hypothèse nulle, une loi normale standard.

#### 4.4.2 Cas d'un effectif de groupes faible

L'approche la plus naturelle est d'utiliser un test classique en considérant  $\Theta_i$  comme un paramètre ordinaire : le test peut être, par exemple, un test du score à  $n - 1$  degrés de liberté basé sur la vraisemblance partielle.

Une famille de tests du score, basés sur une vraisemblance pénalisée, est également exploitable dans le cas d'un nombre modéré ou faible de groupes (Gray, 1998).

### 4.5 Compléments

#### 4.5.1 Concernant les modèles fréquentistes

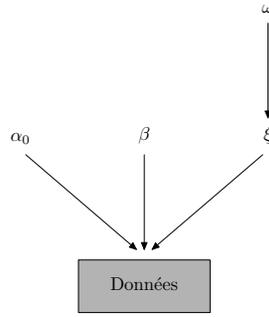
L'algorithme de Newton-Raphson peut être préféré à celui de l'algorithme E.M. lors de l'étape d'estimation : citons les travaux de MacGilchrist et Aisbett (1991) dans le cadre d'un modèle de fragilité log-normale, et ceux de Ha *et al.* (2001) dans le cas général.

Par ailleurs, des études comparatives portant sur les distributions de la fragilité ont été publiées (Paik *et al.*, 1994 ; Hougaard, 1995 ; Pickles et Crouchley, 1995). En particulier, Paik *et al.* ont proposé un modèle où la distribution de la fragilité est celle d'une loi *gamma par morceaux* : cette variante permet de modéliser une corrélation non constante. Notons que le modèle de Cox n'est pas l'unique modèle de fragilité : d'autres types de modèles supportent l'inclusion d'un terme de fragilité – parmi ceux-ci, citons le modèle de Weibull (Andersen *et al.*, 1997) et le modèle exponentiel (Sahu *et al.*, 1997).

Notons également que des modifications peuvent être portées au modèle de fragilité de Cox. Par exemple, le risque de base  $\alpha_0(t)$  peut être soit paramétrique – distribué suivant une loi de Weibull (Ducrocq et Casella, 1996), par exemple –, soit non-paramétrique mais, à la différence du modèle de Cox, spécifié tout de même en partie (par exemple constant par morceaux (Sinha, 1998)). Précisons que tous ces modèles n'ont, pour l'heure, porté que sur des données bivariées.

Enfin, nous n'avons considéré que le modèle incluant une **fragilité partagée**, c'est-à-dire que la fragilité est la même pour tous les individus d'un groupe – elle représente alors une source inobservable de variabilité *commune* à tous les sujets d'un même groupe. Or il est possible de modifier la définition de cette composante aléatoire, afin de prendre également en compte une variabilité inobservable *individuelle* : ainsi, la fragilité peut s'écrire sous la forme (Parner, 1998 ; Petersen, 1998)

$$\Theta_i^{(j)} = \Theta_{i0} + \Theta_{ij}$$

FIG. 4.1 – *Modèle bayésien.*

pour le  $j^{\text{e}}$  individu du  $i^{\text{e}}$  groupe. Les deux composantes de la fragilité peuvent, par exemple, être indépendantes, distribuées suivant des lois gamma de même paramètre d'échelle, mais de paramètres de forme différents. Dans ce cas, le problème d'identification du paramètre de surdispersion à partir d'un seul individu (cf. remarque p. 63) ne se pose plus.

#### 4.5.2 Concernant l'approche bayésienne

Cette approche est couramment rencontrée dans la bibliographie (Kalbfleisch, 1978; Clayton, 1991; Gamerman, 1991; Abrams *et al.*, 1996; Arjas et Liu, 1996; Ducrocq et Casella, 1996; Gustafson, 1997; Sargent, 1998).

##### Introduction

Reprenons (à la notation du terme de fragilité près) l'écriture de la fonction de risque du modèle de fragilité de Cox (cf. p. 62) :

$$\begin{aligned}\tilde{\alpha}_{ij}(t) &= \alpha_0 \Theta_i \exp(\beta^t \mathbf{Z}_{ij}) \\ &= \alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij} + \tilde{\Theta}_i),\end{aligned}$$

où  $\Theta_i$  est le terme de fragilité affecté au groupe  $i$  et  $\tilde{\Theta}_i = \log(\Theta_i)$ .

Par la suite, nous reprenons la notation  $\Theta_i$  pour la « log-fragilité » :

$$\tilde{\alpha}_{ij}(t) = \alpha_0(t) \exp(\beta^t \mathbf{Z}_{ij} + \Theta_i).$$

Nous supposons que  $\Theta_i$  suit une loi de densité  $g(\theta_i | \boldsymbol{\xi})$ , et que le couple des paramètres  $(\boldsymbol{\xi}, \boldsymbol{\beta})$  suit une loi de densité  $q(\boldsymbol{\xi}, \boldsymbol{\beta} | \omega)$ .

$\omega$ , désigné sous le nom d'**hyperparamètre**, est supposé connu.

Soit  $p(\boldsymbol{\xi}, \boldsymbol{\beta}, \theta | z, \omega)$  la **distribution a posteriori** des paramètres (c'est-à-dire une fois connues les valeurs des covariables).

Alors

$$p(\boldsymbol{\xi}, \boldsymbol{\beta}, \theta \mid z, \omega) \propto C(\boldsymbol{\beta} \mid z, \theta) \times g(\theta \mid \boldsymbol{\xi}) \times q(\boldsymbol{\xi}, \boldsymbol{\beta} \mid \omega), \quad (4.13)$$

où

$$C(\boldsymbol{\beta} \mid z, \theta) = \prod_t \prod_{i=1}^n \prod_{j=1}^{n_i} \left( \frac{y_{ij}(t) \exp(\boldsymbol{\beta}^t Z_{ij} + \theta_{ij})}{S^{(0)}(\boldsymbol{\beta}, t)} \right)^{\Delta N_{ij}(t)}$$

avec

$$S^{(0)}(\boldsymbol{\beta}, t) = \sum_{i=1}^n \sum_{j=1}^{n_i} y_{ij}(t) \exp(\boldsymbol{\beta}^t Z_{ij} + \theta_{ij}).$$

### Estimation par la méthode de Monte-Carlo par chaîne de Markov

La méthode de Monte-Carlo par chaîne de Markov (*Markov Chain Monte-Carlo (MCMC)*) permet d'obtenir des échantillons d'une densité  $\pi(\tau)$  qui peut être connue, mais difficile à échantillonner. Le principe consiste à créer une chaîne de Markov sur l'espace des états de  $\tau$  telle qu'elle admette comme distribution stationnaire  $\pi(\tau)$  (voir l'annexe 1).

L'échantillonnage des paramètres peut se faire au travers des **algorithmes de Metropolis, Hastings ou Gibbs**.

Sargent (1998) calcule le membre droit de la relation (4.13) de la façon suivante :

- il choisit l'approximation d'Efron pour  $C(\boldsymbol{\beta} \mid z, \theta)$  ;
- $g(\theta \mid \boldsymbol{\xi})$  est une loi normale  $\mathcal{N}(\mu, \nu)$  (*i.e.*  $\boldsymbol{\xi} = (\mu, \nu)$ ) – par suite, il suppose que  $\mu = 0$  (ce qui entraîne que  $\boldsymbol{\xi} = \nu$ ) ;
- il suppose que  $\boldsymbol{\xi} = \nu$  suit une inverse gamma de paramètres  $\alpha$  et  $\phi$  (*i.e.*  $1/\nu$  a pour espérance et variance respectivement  $\alpha/\phi$  et  $\alpha/\phi^2$ ), et que  $\boldsymbol{\beta}$  suit une loi de densité *a priori* uniforme (*flat*).

Des extensions de ce type de modèle sont possibles, telles que l'affectation de variances différentes pour chaque  $\Theta_i$ , ou encore l'attribution à  $\omega$  d'une loi de probabilité gaussienne.

## 4.6 Prédiction de la fragilité

### 4.6.1 La théorie du meilleur prédicteur

Considérant un modèle de régression linéaire mixte de la forme

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon} \quad (4.14)$$

où

- $\mathbf{Y}$  est le vecteur-colonne  $n \times 1$  des observations ;
- $\boldsymbol{\beta}$  est le vecteur  $p \times 1$  des effets fixes ;
- $\mathbf{u}$  est le vecteur  $q \times 1$  des effets aléatoires, de matrice de variance-covariance  $\mathbf{G}$  ;
- $\mathbf{X}$  est la matrice  $n \times p$  ;

- $\mathbf{Z}$  est la matrice  $n \times q$  ;
- $\boldsymbol{\epsilon}$  est le vecteur  $n \times 1$  des résiduelles, de matrice de variance-covariance  $\mathbf{R}$ .

Les variables aléatoires  $\mathbf{u}$  et  $\boldsymbol{\epsilon}$  étant par nature inobservables, il peut être utile d'en estimer les valeurs, connaissant les observations  $\mathbf{Y}$ . Ces estimations sont utilisées par certains tests portant sur les modèles mixtes, tests pour lesquels l'hypothèse de base est la nullité de combinaisons linéaires des effets fixes et aléatoires.

### Prédiction des seuls effets aléatoires

Nous reprenons le modèle (4.14), et nous supposons dans un premier temps que les effets fixes  $\boldsymbol{\beta}$  sont connus. Aucune hypothèse portant sur les distributions de  $\mathbf{u}$  et de  $\boldsymbol{\epsilon}$  n'est faite. Nous cherchons à prédire  $\mathbf{u}$ .

Nous notons

$$\begin{aligned} \mathbf{V} &= \mathbb{V}(\mathbf{Y}) \\ &= \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R}. \end{aligned}$$

Pour toute matrice symétrique définie-positive  $\mathbf{A}$ , si nous notons  $f(\mathbf{u}, \mathbf{Y})$  la fonction de densité conjointe de  $\mathbf{u}$  et de  $\mathbf{Y}$ , nous pouvons calculer la quantité suivante, qui est une généralisation des carrés moyens de l'erreur de prédiction sur  $\mathbf{u}$  :

$$\mathbb{E} [(\hat{\mathbf{u}} - \mathbf{u})' \mathbf{A} (\hat{\mathbf{u}} - \mathbf{u})] = \int \int (\hat{\mathbf{u}} - \mathbf{u})' \mathbf{A} (\hat{\mathbf{u}} - \mathbf{u}) f(\mathbf{u}, \mathbf{Y}) d\mathbf{Y} d\mathbf{u}. \quad (4.15)$$

Le vecteur  $\hat{\mathbf{u}}$  qui minimise cette quantité est appelé **meilleur prédicteur (BP – Best Predictor)** de  $\mathbf{u}$ . On montre que

$$\hat{\mathbf{u}} = \mathbb{E}(\mathbf{u} | \mathbf{Y}),$$

et que ce prédicteur ne dépend pas de la matrice  $\mathbf{A}$  choisie, ni de la forme de la densité  $f(\mathbf{u}, \mathbf{Y})$ .

Les propriétés de l'espérance conditionnelle permettent d'établir les résultats suivants :

$$\begin{aligned} \mathbb{E}_{\mathbf{Y}}(\hat{\mathbf{u}}) &= \mathbb{E}(\mathbf{u}) \\ \mathbb{V}(\hat{\mathbf{u}} - \mathbf{u}) &= \mathbb{E}_{\mathbf{Y}}[\mathbb{V}(\mathbf{u} | \mathbf{Y})] \\ \text{Cov}(\hat{\mathbf{u}}, \mathbf{u}^t) &= \mathbb{V}(\hat{\mathbf{u}}) \\ \text{Cov}(\hat{\mathbf{u}}, \mathbf{Y}^t) &= \text{Cov}(\mathbf{u}, \mathbf{Y}^t). \end{aligned}$$

Nous remarquons, en particulier, que le meilleur prédicteur de  $\mathbf{u}$  est sans biais.

Notons  $u_i$  le  $i^{\text{e}}$  élément de  $\mathbf{u}$ , et  $\hat{u}_i$  le  $i^{\text{e}}$  élément de  $\hat{\mathbf{u}}$ . Il est possible de montrer que le prédicteur  $\hat{u}_i = \text{BP}(u_i)$  est, parmi tous les prédicteurs de  $u_i$ , celui qui maximise la corrélation entre  $u_i$  et  $\hat{u}_i$ .

Posons

$$\begin{aligned} \mathbb{E} \left[ \begin{pmatrix} \mathbf{u} \\ \mathbf{Y} \end{pmatrix} \right] &= \begin{pmatrix} \mu_{\mathbf{u}} \\ \mu_{\mathbf{Y}} \end{pmatrix}, \\ \mathbb{V} \left[ \begin{pmatrix} \mathbf{u} \\ \mathbf{Y} \end{pmatrix} \right] &= \begin{pmatrix} \mathbf{G} & \mathbf{C} \\ \mathbf{C}^t & \mathbf{V} \end{pmatrix}. \end{aligned}$$

Si nous imposons maintenant au prédicteur  $\hat{\mathbf{u}}$  de  $\mathbf{u}$  d'être linéaire en  $\mathbf{Y}$ , c'est-à-dire d'être de la forme  $\hat{\mathbf{u}} = \mathbf{a} + \mathbf{B}\mathbf{Y}$ , la minimisation de (4.15) conduit au **meilleur prédicteur linéaire** (**BLP** – *Best Linear Predictor*)

$$\begin{aligned} \text{BLP}(\mathbf{u}) &= \hat{\mathbf{u}} \\ &= \mu_{\mathbf{u}} + \mathbf{C}\mathbf{V}^{-1}(\mathbf{Y} - \mu_{\mathbf{Y}}), \end{aligned}$$

dont les deux premiers moments sont

$$\mathbb{E}(\hat{\mathbf{u}}) = \mu_{\mathbf{u}}$$

et

$$\mathbb{V}(\hat{\mathbf{u}}) = \mathbf{C}\mathbf{V}^{-1}\mathbf{C}^t.$$

C'est un estimateur sans biais, qui ne dépend pas de la forme de  $f(\mathbf{u}, \mathbf{Y})$ .

Sous l'hypothèse de normalité, les deux estimateurs BP( $\mathbf{u}$ ) et BLP( $\mathbf{u}$ ) coïncident. Avec les notations de (4.14), on a  $\mu(\mathbf{u}) = 0$ ,  $\mu_{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta}$  et  $\mathbf{C} = \mathbf{G}\mathbf{Z}^t$ , d'où

$$\begin{aligned} \hat{\mathbf{u}} &= \text{BP}(\mathbf{u}) \\ &= \text{BLP}(\mathbf{u}) \\ &= \mathbf{G}\mathbf{Z}^t\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}). \end{aligned}$$

### Prédiction des effets aléatoires et estimation des effets fixes

Cette fois-ci, nous ne supposons plus  $\boldsymbol{\beta}$  connu. Nous nous intéressons donc conjointement à l'estimation des effets fixes et à la prédiction de l'effet aléatoire. Plus exactement, nous cherchons maintenant à estimer un vecteur  $w$  de la forme

$$w = \mathbf{L}^t\boldsymbol{\beta} + \mathbf{u},$$

où  $\mathbf{L}$  est une matrice connue telle que  $\mathbf{L}^t\boldsymbol{\beta}$  soit estimable. Par la suite, nous serons amenés à tester l'hypothèse de nullité de ce vecteur  $w$ .

Les deux premiers moments de  $w$  sont

$$\mathbb{E}(w) = \mathbf{L}^t\boldsymbol{\beta}$$

et

$$\begin{aligned} \mathbb{V}(w) &= \mathbb{V}(\mathbf{u}) \\ &= \mathbf{G} \end{aligned}$$

En outre,

$$\text{Cov}(w, \mathbf{Y}^t) = \mathbf{G}\mathbf{z}^t.$$

Nous ne pouvons plus prédire les effets aléatoires de la même manière que dans la section précédente, puisqu'il faut simultanément estimer  $\boldsymbol{\beta}$  et prédire  $u$ .

Nous cherchons un estimateur  $\hat{w}$  de  $w$  possédant les trois caractéristiques suivantes :

- « meilleur prédicteur » dans le sens de (4.15), *i.e.* quantité qui minimise

$$\mathbb{E} [(\hat{w} - w)^t A(\hat{w} - w)] ;$$

- linéaire en  $\mathbf{Y}$ , *i.e.* de la forme

$$\hat{w} = a + B\mathbf{Y} ,$$

où  $a$  et  $B$  ne dépendent pas de  $\beta$  ;

- sans biais :

$$\mathbb{E}(\hat{w}) = \mathbb{E}(w) .$$

La condition portant sur le biais de l'estimateur doit ici être explicitement posée et n'est plus une conséquence des conditions précédentes comme dans le cas du BP ou celui du BLP en raison de l'estimation de  $\beta$ . Un tel estimateur est appelé **meilleur prédicteur linéaire non biaisé (BLUP – Best Linear Unbiased Predictor)**.

La linéarité de l'estimateur recherché et le fait que celui-ci ne soit pas biaisé entraînent que l'égalité

$$a + B\mathbf{X}\beta = \mathbf{L}^t\beta$$

doit être vraie pour tout  $\beta$ , ce qui entraîne  $a = 0$  et  $B\mathbf{X} = \mathbf{L}^t$ .

Dans un article de 1975, Henderson reprend un résultat qu'il a obtenu en 1963 où il montre que  $\hat{w}$  est de la forme

$$\begin{aligned} \hat{w} &= \text{BLUP}(w) \\ &= \mathbf{L}^t\hat{\beta} + \mathbf{G}\mathbf{Z}^t\mathbf{V}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\beta}) , \end{aligned} \quad (4.16)$$

où  $\hat{\beta}$  est n'importe quelle solution de l'équation des moindres carrés généralisés

$$\mathbf{X}^t\mathbf{V}^{-1}\mathbf{X}\hat{\beta} = \mathbf{X}^t\mathbf{V}^{-1}\mathbf{Y} . \quad (4.17)$$

Dans ce même article, il reprend des résultats obtenus en 1959 et 1963, par lesquels il prouve que les  $\hat{\beta}$  et  $\hat{u}$  obtenus par les équations (4.16) et (4.17) sont aussi solutions du système

$$\begin{pmatrix} \mathbf{X}^t\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}^t\mathbf{R}^{-1}\mathbf{Z} \\ \mathbf{Z}^t\mathbf{R}^{-1}\mathbf{X} & \mathbf{Z}^t\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \hat{\beta} \\ \hat{u} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^t\mathbf{R}^{-1}\mathbf{Y} \\ \mathbf{Z}^t\mathbf{R}^{-1}\mathbf{Y} \end{pmatrix} , \quad (4.18)$$

qui est plus simple à calculer que celui constitué des équations (4.16) et (4.17), puisqu'il ne requiert pas l'inversion de  $\mathbf{V}$ , mais celles de  $\mathbf{R}$  et  $\mathbf{G}$  qui, en pratique, sont souvent faciles à calculer. Ce système d'équations est connu sous le nom d'**équations du modèle mixte**.

La démonstration de l'équivalence entre les solutions de (4.18) et celles de (4.16) et (4.17) utilise les égalités

$$\begin{aligned} \mathbf{R}^{-1} - \mathbf{R}^{-1}\mathbf{Z}(\mathbf{Z}^t\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1})^{-1}\mathbf{Z}^t\mathbf{R}^{-1} &= \mathbf{V}^{-1} \\ (\mathbf{Z}^t\mathbf{R}^{-1}\mathbf{Z} + \mathbf{G}^{-1})^{-1}\mathbf{Z}^t\mathbf{R}^{-1} &= \mathbf{G}\mathbf{Z}^t\mathbf{V}^{-1} \end{aligned}$$

et le fait que  $\mathbf{X}\hat{\beta}$  obtenu par (4.18) est aussi un estimateur des moindres carrés généralisés de  $\mathbf{X}\beta$ .

Le calcul de l'estimateur du BLUP suppose connues les matrices de variance-covariance  $\mathbf{R}$  et  $\mathbf{G}$ . En pratique, ces matrices doivent souvent être elles-mêmes estimées et l'on calcule alors un estimateur empirique du BLUP à partir des matrices estimées  $\hat{\mathbf{R}}$  et  $\hat{\mathbf{G}}$ . L'estimateur obtenu, appelé **meilleur prédicteur linéaire empirique (EBLUP – Empirical Best Linear Unbiased Predictor)**, perd alors les propriétés d'optimalité du BLUP.

### 4.6.2 Tests d'hypothèses

Nous considérons un modèle du type (4.14) sous les hypothèses habituelles de normalité et d'indépendance des vecteurs aléatoires  $\mathbf{u}$  et  $\boldsymbol{\epsilon}$ .

Nous supposons que les paramètres du modèle, soient  $\boldsymbol{\beta}, \mathbb{V}(\mathbf{u}) = \mathbf{G}$  et  $\mathbb{V}(\boldsymbol{\epsilon}) = \mathbf{R}$ , sont estimés par les méthodes du maximum de vraisemblance ou du maximum de vraisemblance restreinte.

En règle générale, il n'y a pas de solution explicite aux équations du maximum de vraisemblance, sauf si les données du dispositif sont équilibrées. De même, on ne dispose pas de test exact pour les paramètres du modèle mixte dans le cas général – seuls les cas équilibrés en bénéficient. On utilise alors des tests approchés.

### Test sur des combinaisons linéaires des effets fixes et aléatoires

Plusieurs types de tests sont possibles pour les effets fixes, parmi lesquels ceux que nous présentons ci-dessous.

Les hypothèses habituelles du modèle (4.14) peuvent s'écrire :

$$\begin{pmatrix} \mathbf{u} \\ \boldsymbol{\epsilon} \end{pmatrix} \rightsquigarrow \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{G} & 0 \\ 0 & \mathbf{R} \end{bmatrix} \right)$$

et elles entraînent donc :

$$\begin{aligned} \mathbb{E}(\mathbf{Y}) &= \mathbf{X}\boldsymbol{\beta}, \\ \mathbb{V}(\mathbf{Y}) &= \mathbf{Z}\mathbf{G}\mathbf{Z}^t + \mathbf{R} \end{aligned}$$

Nous décomposons  $\mathbf{V}$  sous la forme suivante :

$$\mathbf{V} = \sum_{t=1}^T \phi_t \mathbf{V}_t$$

où les matrices  $(\mathbf{V}_t)_t$  sont des matrices symétriques connues et les réels  $(\phi_t)_t$  les composantes d'un vecteur  $\boldsymbol{\phi} = (\phi_1, \dots, \phi_T)^t$ .

### 4.6.3 Dans le cadre du modèle de fragilité de Cox

En reprenant les notations de la page 62, il est possible d'exprimer la fonction de vraisemblance complète (inobservée) des triplets  $(N_{ij}, Y_{ij}, \Theta_i)$  en fonction de la transformée de Laplace de  $\Theta_i$  (Parner, 1997). Ainsi, cette fonction de vraisemblance complète vaut (cf. 62)

$$\begin{aligned} f(\mathbf{N}, \mathbf{Y}, \boldsymbol{\Theta}) &= \prod_{i=1}^n \prod_{j=1}^{m_i} \left\{ \prod_{t \in [0, \tau]} [\theta_i d\Lambda_{ij}(t)]^{\Delta N_{ij}(t)} \exp(-\theta_i \Lambda_{ij}(\tau)) \right\} p(\theta_i) \\ &= \left\{ \prod_{i=1}^n \theta_i^{N_{i+}(\tau)} \exp(-\theta_i \Lambda_{ij}(\tau)) p(\theta_i) \right\} \times \left\{ \prod_{i=1}^n \prod_{j=1}^{m_i} \prod_{t \in [0, \tau]} [d\Lambda_{ij}(t)]^{\Delta N_{ij}(t)} \right\}. \end{aligned} \tag{4.19}$$

La vraisemblance observée est obtenue en intégrant (4.19) suivant  $\Theta_i$  :

$$\begin{aligned} \int f(\mathbf{N}_i, \mathbf{Y}_i, \Theta_i) d\Theta_i &= \int \theta_i^{N_{i+}(\tau)} \exp(-\theta_i \Lambda_{i+}(\tau)) p(\theta_i) d\theta_i \times \prod_{j=1}^{m_i} \prod_{t \in [0, \tau]} [d\Lambda_{ij}(t)]^{\Delta N_{ij}(t)} \\ &= (-1)^{\varsigma_i} \mathcal{L}_\theta^{(\varsigma_i)}(\Lambda_{i+}(\tau)) \times \prod_{j=1}^{m_i} \prod_{t \in [0, \tau]} [d\Lambda_{ij}(t)]^{\Delta N_{ij}(t)}, \end{aligned} \quad (4.20)$$

où  $\mathcal{L}_\theta^{(\varsigma)}(\cdot)$  désigne la  $\varsigma^e$  dérivée de la transformée de Laplace de  $\Theta$ , et  $\varsigma_i = N_{i+}(\tau)$ .

Par suite, nous pouvons réécrire l'espérance conditionnelle de  $\Theta_i$  sachant  $(\mathbf{N}_i, \mathbf{Y}_i)_{i=1, \dots, n}$  sous la forme

$$\begin{aligned} \mathbb{E}[\Theta_i \mid (\mathbf{N}_i, \mathbf{Y}_i)] &= \int \theta_i f(\theta_i \mid \mathbf{N}_i, \mathbf{Y}_i) d\theta_i \\ &= \frac{\int \theta_i f(\mathbf{N}_i, \mathbf{Y}_i, \theta_i) d\theta_i}{\int f(\mathbf{N}_i, \mathbf{Y}_i, \theta_i) d\theta_i} \\ &= \frac{\mathcal{L}^{(\varsigma_i+1)}(\Lambda_{i+}(\tau))}{\mathcal{L}^{(\varsigma_i)}(\Lambda_{i+}(\tau))}. \end{aligned} \quad (4.21)$$

Lorsque la fragilité suit une loi log-normale, le calcul de la transformée de Laplace ne peut être explicitement fait. En revanche, lorsque la fragilité suit une loi gamma de paramètre  $\gamma$ , ce calcul peut être réalisé.

#### 4.6.3.1 Cas d'une fragilité gamma

Le calcul de la transformée de Laplace mène naturellement au résultat vu en page 66 : le meilleur prédicteur de  $\Theta_i$  ( $1 \leq i \leq n$ ) est

$$\begin{aligned} \widehat{\Theta}_i &= \mathbb{E}[\Theta_i \mid (N_{ij}(t), Y_{ij}(t), 1 \leq i \leq n, 0 \leq t \leq \tau)] \\ &= \frac{\sum_j N_{ij}(\tau) + \frac{1}{\gamma}}{\sum_k \int_0^\tau Y_{ik}(t) \alpha_0(t) \exp(\beta^t \mathbf{Z}_{ik}) dt + \frac{1}{\gamma}}. \end{aligned} \quad (4.22)$$

*Démonstration* — La transformée de Laplace d'une v.a.  $U$  suivant une loi gamma  $\Gamma(p, q)$  est

$$\mathcal{L}_U(s) = \left( \frac{p}{p+s} \right)^q$$

et sa dérivée d'ordre  $n$  vaut

$$\mathcal{L}_U^{(n)}(s) = (-1)^n p^q (p+s)^{-(q+n)} \frac{\Gamma(q+n)}{\Gamma(q)}.$$

D'après (4.21), et du fait que  $\Theta_i$  suit une loi  $\Gamma(\frac{1}{\gamma}, \frac{1}{\gamma})$ , nous avons

$$\begin{aligned} \mathbb{E} [\Theta_i \mid (\mathbf{N}_i, \mathbf{Y}_i)] &= \frac{(-1)^{\varsigma_i+1} \frac{1}{\gamma} [(\frac{1}{\gamma} + \Lambda_{i+}(\tau))^{-(\frac{1}{\gamma} + \varsigma_i + 1)} \frac{\Gamma(\frac{1}{\gamma} + \varsigma_i + 1)}{\Gamma(\frac{1}{\gamma})}]}{(-1)^{\varsigma_i} \frac{1}{\gamma} [\frac{1}{\gamma} + \Lambda_{i+}(\tau)]^{-(\frac{1}{\gamma} + \varsigma_i)} \frac{\Gamma(\frac{1}{\gamma} + \varsigma_i)}{\Gamma(\frac{1}{\gamma})}} \\ &= -\frac{\frac{1}{\gamma} + \varsigma_i}{\frac{1}{\gamma} + \Lambda_{i+}(\tau)} \end{aligned}$$

où  $\varsigma_i = N_{i+}(\tau)$ . ■

Concernant l'espérance de  $\widehat{\Theta}_i$ , nous la trouvons bien égale à l'unité comme l'hypothèse en avait été faite :

$$\begin{aligned} \mathbb{E}(\widehat{\Theta}_i) &= \mathbb{E}\left(\mathbb{E}[\Theta_i \mid (N_{ij}(t), Y_{ij}(t))]\right) \\ &= \mathbb{E}(\Theta_i) \\ &= 1. \end{aligned}$$

Quant à la variance et à la loi asymptotique de la fragilité, elles n'ont encore jamais été explicitées. Notons que (4.22) n'est pas le rapport d'une suite de martingales sur le processus de variation associé à cette suite ; par suite, le théorème de la limite centrale de Rebolledo ne peut être appliqué ici.

Tout au plus pouvons-nous écrire que

$$\begin{aligned} \widehat{\Theta}_i &= \frac{\sum_j N_{ij}(\tau) + \frac{1}{\gamma}}{\sum_k \int_0^\tau Y_{ik}(t) \alpha_0(t) \exp(\boldsymbol{\beta}^t \mathbf{Z}_{i\mathbf{k}}) dt + \frac{1}{\gamma}} \\ &= \frac{\sum_j M_{ij}(\tau) + \sum_j \Lambda_{ij}(\tau) + \frac{1}{\gamma}}{\sum_k \Lambda_{ik}(\tau) + \frac{1}{\gamma}} \\ &= \frac{\sum_j M_{ij}(\tau)}{\sum_k \Lambda_{ik}(\tau) + \frac{1}{\gamma}} + 1 \\ &= \frac{\sum_j M_{ij}(\tau)}{\sum_k \langle M_{ik} \rangle(\tau) + \frac{1}{\gamma}} + 1. \end{aligned}$$



Troisième partie

**DEUX ÉTUDES POUR UNE  
COMPARAISON EMPIRIQUE  
DES DIFFÉRENTS MODÈLES**



## Chapitre 5

# ÉTUDE PAR SIMULATIONS

### Contenu

---

<b>5.1</b>	<b>Présentation des procédures S-PLUS</b>	<b>84</b>
5.1.1	Concernant l'ajustement du modèle naïf	84
5.1.2	Concernant l'ajustement du modèle marginal	84
5.1.3	Concernant l'ajustement du modèle de fragilité	86
5.1.4	Types et fonctionnalités des résidus	92
<b>5.2</b>	<b>Présentation des simulations</b>	<b>92</b>
5.2.1	Objectifs	92
5.2.2	Modèle de fragilité gamma	93
5.2.3	Modèle de fragilité stable positive	94
<b>5.3</b>	<b>Paramétrage des simulations et choix des sorties</b>	<b>95</b>
<b>5.4</b>	<b>Résultats</b>	<b>97</b>
5.4.1	Estimations du paramètre de régression $\beta$	97
5.4.2	Estimations de la variance du paramètre de régression – coefficient de dispersion	100
5.4.3	Puissance du test du maximum de vraisemblance	102
5.4.4	Critère d'Akaike et taux de recouvrement de l'intervalle de confiance à 95 % de la variance de la fragilité	105
<b>5.5</b>	<b>Conclusions et perspectives</b>	<b>106</b>

---

## 5.1 Présentation des procédures S-PLUS

### 5.1.1 Concernant l'ajustement du modèle naïf

La procédure relative à l'ajustement du modèle de Cox est la suivante :

```
> coxph(Surv(temps,statut) ~ cov_1 + cov_2 + ... + factor(cov_i) + ...  
      + factor(cov_n) + cluster/frailty(groupe),  
      data = données, na.action=na.omit)
```

où :

- `coxph(Surv(...))` est la déclaration générale de la procédure ;
- `temps` est la variable correspondant aux temps de survie observés ;
- `statut` est la variable binaire indiquant les cas de censure ;
- `cov_1 + cov_2 + ...` est l'énumération des covariables quantitatives ;
- `factor(cov_i) + ... + factor(cov_n)` est l'énumération des covariables qualitatives ;
- `cluster(groupe)` est la déclaration de la variable induisant la corrélation dans le cadre du modèle marginal ;
- `frailty(groupe)` est la déclaration de la variable induisant la corrélation dans le cadre du modèle de fragilité ;
- `data = données` est la déclaration du nom du fichier contenant les données ;
- `na.action=na.omit` est une déclaration optionnelle permettant d'exclure les valeurs manquantes.

L'estimation se fait numériquement au travers de l'algorithme de Newton-Raphson.

### 5.1.2 Concernant l'ajustement du modèle marginal

Concernant le modèle marginal, le logiciel S-PLUS (MathSoft, Inc., 1999a; 1999b; 1999c) calcule une vraisemblance robuste similaire à celle du modèle de Lin et Wei (1989) (cf. p. 55), lorsque l'ajustement du modèle marginal de Cox est appelé au travers de l'option `cluster(...)` de la procédure `coxph`.

L'estimateur robuste de la variance – dit *estimateur sandwich* – est obtenu au travers d'une approche du type *jackknife*.

La procédure, qui a lieu lors des itérations suivant l'algorithme de Newton-Raphson, est la suivante :

1° le changement dans l'estimation du coefficient  $\beta$  est

$$\begin{aligned}\Delta\beta &= \mathbf{1}'(\mathbf{U}\mathbf{J}^{-1}) \\ &= \mathbf{1}'\mathbf{D},\end{aligned}$$

où  $\mathbf{U}$  est la matrice des résidus du score ; les changements concernant  $\hat{\beta}$  à chaque itération consistent en la somme des éléments d'une colonne de la matrice  $\mathbf{D}$ , définie comme les résidus du score, multipliée par la variance de  $\hat{\beta}$  ;

- 2° on retire le point  $i$  des données, et l'on recalcule  $\mathbf{U}$ , tandis que  $\mathbf{J}$  demeure inchangée ; cela correspond à la suppression de la colonne  $i$  de  $\mathbf{D}$  ;
- 3° on réitère une fois l'algorithme de Newton-Raphson, c'est-à-dire que l'on forme la nouvelle somme des éléments de la colonne qui doit être égale, par construction, à l'opposé de la colonne qui a été supprimée.

Lorsque les données sont constituées en groupes, la procédure consiste à ignorer à chaque étape un groupe entier, puis à calculer l'estimateur de la variance. Cette méthode est désignée sous le nom d'**estimation par la méthode du jackknife groupé** (*grouped jackknife*).

Plus précisément, supposons que l'estimateur sandwich soit de la forme

$$\mathbf{V} = \mathbf{A}\mathbf{B}^t\mathbf{A},$$

avec  $\mathbf{A} = \mathbf{J}^{-1}$  inverse de la matrice d'information et  $\mathbf{B}$  un terme correctif.

Sous certaines conditions,  $\hat{\beta}$  est asymptotiquement normal, de moyenne  $\beta$  et de matrice de variance-covariance  $\mathbf{V}$  telle que

$$\mathbf{A} = \left( \frac{\partial \mathbb{E} \left[ \sum_{i=1}^n \frac{\partial \log L_i(\beta)}{\partial \beta} \right]}{\partial \beta} \right)^{-1}$$

et

$$\mathbf{B} = \text{Cov} \left( \sum_{i=1}^n \frac{\partial \log L_i(\beta)}{\partial \beta} \right).$$

Si l'on note

$$u_i(\beta) = \frac{\partial \log L_i(\beta)}{\partial \beta},$$

on a

$$\mathbf{B} = \sum_{i=1}^n \mathbb{E} [u_i(\beta)^t u_i(\beta)] + \sum_{i \neq j} \mathbb{E} [u_i(\beta)^t u_j(\beta)]. \quad (5.1)$$

En cas d'indépendance, un estimateur naturel de  $\mathbf{B}$  est

$$\hat{\mathbf{B}} = \sum_{i=1}^n \mathbb{E} [u_i(\hat{\beta})^t u_i(\hat{\beta})].$$

En présence de corrélation (s'il y a  $k$  groupes), on élimine le terme croisé de (5.1) en écrivant

$$\mathbf{B} = \sum_{j=1}^k \tilde{u}(\beta)^t \tilde{u}(\beta),$$

où  $\tilde{u}$  est la somme des  $u_i$  sur l'ensemble des sujets appartenant au  $j^e$  groupe.

Ceci mène à l'**estimateur sandwich modifié**

$$\mathbf{A}\tilde{\mathbf{U}}^t\tilde{\mathbf{U}}\mathbf{A},$$

où la matrice  $\tilde{\mathbf{U}}$  est obtenue en remplaçant chaque groupe de colonnes de  $\mathbf{U} = (u_i)_i$  par la somme des colonnes de ce groupe.

### 5.1.3 Concernant l'ajustement du modèle de fragilité

Essentiellement pour des raisons numériques – les méthodes de calcul par l'algorithme E.M. étant instables et lentes –, la fonction `coxph` du logiciel S-PLUS utilise une procédure de vraisemblance pénalisée, proposée par Therneau et Grambsch (2000), lorsqu'un terme de fragilité est inclus dans le modèle (option `frailty(...)` de la procédure `coxph`).

#### Résolution numérique par pénalisation de la vraisemblance partielle de Cox

Soit le modèle

$$\begin{aligned}\alpha_{ij}(t) &= \alpha_0(t) \tilde{\Theta}_i \exp(\beta^t Z_{ij}) \\ &= \alpha_0(t) \exp(\beta^t Z_{ij} + \Theta_i) .\end{aligned}$$

Nous explicitons la procédure `coxph` dans les cas suivants où

- $\tilde{\Theta}_i$  suit une loi gamma d'espérance 1 et de variance  $\gamma$  ;
- $\Theta$  suit une loi gaussienne centrée de variance  $\gamma$ .

La log-vraisemblance partielle de Cox peut s'écrire

$$\begin{aligned}\log \mathcal{L}(\beta, \alpha_0, \gamma) &= \log \mathcal{L}(\beta, \alpha_0, \gamma; N_{ij}, Y_{ij}, \Theta_i) \\ &= \log \mathcal{L}(\beta, \alpha_0, \gamma; N_{ij}, Y_{ij} \mid \Theta_i) + \log \mathcal{L}(\Theta_i) \\ &= \sum_{i,j} \int_0^\tau \left\{ \alpha_0(t) Y_{ij}(t) (\beta^t Z_{ij} + \theta_i) - \log \left[ \sum_{k,l} \alpha_0(t) Y_{kl}(t) \exp(\beta^t Z_{kl} + \theta_k) \right] \right\} dN_{ij}(t) \\ &\quad + \log g(\theta_i; \gamma) \\ &= \mathfrak{L}(\beta, \alpha_0, \theta_i) + \log g(\theta_i; \gamma) ,\end{aligned}$$

où  $g(\theta_i; \gamma)$  est la vraisemblance associée à  $\Theta_i$ .

Quant à  $\mathfrak{L}(\beta, \alpha_0)$ , ce n'est en réalité que la vraisemblance partielle de Cox usuelle avec  $\theta$  (fixe – *offset*) incorporé dans la partie régressive du modèle.

Une écriture semblable peut être faite avec  $\tilde{\Theta}_i$ .

Nous définissons le logarithme de la fonction de **vraisemblance partielle pénalisée** comme étant

$$\text{PPL}(\beta, \alpha_0, \gamma) = \mathfrak{L}(\beta, \alpha_0, \theta_i) - h(\theta_i, \gamma) ,$$

où  $h(\cdot)$  est la fonction de pénalité. Cette définition demeure valable concernant  $\tilde{\theta}_i$ .

Dans le cas où  $\tilde{\Theta}_i$  suit une loi gamma,

$$h(\theta_i; \gamma) = -\frac{1}{\gamma} \sum_i [\tilde{\theta}_i - \exp(\tilde{\theta}_i)] \tag{5.2}$$

et dans le cas où  $\Theta_i$  suit une loi gaussienne,

$$h(\theta_i; \gamma) = -\frac{1}{2\gamma} \sum_i \theta_i^2 .$$

Ces deux expressions de  $h(\cdot; \gamma)$  sont justifiées ci-dessous.

**Estimation de  $\Theta_i$**  Pour estimer  $\Theta_i$ , nous résolvons

$$\begin{aligned} \frac{\partial \text{PPL}(\beta, \alpha_0, \theta_i, \gamma)}{\partial \theta_i} &= \frac{\partial \mathcal{L}_1(\beta, \alpha_0, \theta_i)}{\partial \theta_i} - \frac{\partial h(\theta_i, \gamma)}{\partial \theta_i} \\ &= 0. \end{aligned}$$

La dérivation se faisant suivant  $\theta_i$ , il s'ensuit que  $h(\theta_i; \gamma)$  n'est autre que la composante du logarithme de la vraisemblance relative à la fragilité – cette vraisemblance étant  $\prod_i g(\theta_i; \gamma)$  – qui dépend de  $\theta_i$ .

Ainsi, dans le cas d'une loi gamma,

$$g(\tilde{\theta}_i; \gamma) = \prod_i \frac{1}{\gamma^{\frac{1}{\gamma}} \Gamma(\frac{1}{\gamma})} \exp\left(-\frac{\tilde{\theta}_i}{\gamma}\right) \tilde{\theta}_i^{\frac{1}{\gamma}-1},$$

et avec le changement de variable  $\theta_i = \log(\tilde{\theta}_i)$ ,  $d\tilde{\theta}_i = \exp(\theta_i) d\theta_i$ ,

$$\begin{aligned} g(\theta_i; \gamma) &= \prod_i \frac{1}{\gamma^{\frac{1}{\gamma}} \Gamma(\frac{1}{\gamma})} \exp\left(-\frac{\exp(\theta_i)}{\gamma}\right) \exp\left[\left(\frac{1}{\gamma} - 1\right)\theta_i\right] \exp(\theta_i) \\ &= \prod_i K_1 \exp\left(-\frac{\exp(\theta_i)}{\gamma}\right) \exp\left(\frac{1}{\gamma}\theta_i\right), \end{aligned}$$

où  $K_1$  est indépendant de  $\theta_i$ . Il s'ensuit que

$$\log g(\theta_i; \gamma) = \sum_i \left\{ K_1 - \frac{1}{\gamma} [\exp(\theta_i) - \theta_i] \right\}$$

et par suite

$$h(\theta_i; \gamma) = \frac{1}{\gamma} \sum_i [\exp(\theta_i) - \theta_i]$$

Dans le cas d'une loi gaussienne,

$$g(\theta_i; \gamma) = \prod_i \frac{1}{\gamma \sqrt{2\pi}} \exp\left(-\frac{\theta_i^2}{2\gamma}\right),$$

d'où

$$\log g(\theta_i; \gamma) = \sum_i \left\{ K_2 - \frac{\theta_i^2}{2\gamma} \right\}$$

où  $K_2$  est indépendant de  $\theta_i$ , et finalement

$$h(\theta_i; \gamma) = \frac{1}{2\gamma} \sum_i \theta_i^2.$$

**Estimation de  $\beta$**  Elle est identique à celle d'un modèle de Cox ordinaire – avec  $\theta_i$  inclus en tant que terme fixe –, puisque  $\beta$  n'intervient pas dans la fonction de pénalisation.

La résolution numérique de l'ajustement d'un modèle de Cox avec pénalité se déroule alors de la façon itérative qui suit :

- 1° à l'appel initial de l'ajustement, une procédure `cfun` retourne une valeur initiale pour  $\gamma$ , paramètre de la variance de la fragilité ;
- 2° la vraisemblance partielle pénalisée, à  $\gamma$  fixé, est maximisée au travers de l'algorithme de Newton-Raphson (*inner loop* – boucle intérieure) ;
- 3° la dernière étape consiste à maximiser la log-vraisemblance par rapport à  $\gamma$  (*outer loop* – boucle extérieure).

**Remarque** — Cette procédure inclut des équations d'estimation simples, mais en contrepartie elle entraîne une sous-estimation des variances des paramètres des effets fixes, puisqu'elle ne tient pas compte de la variabilité  $\gamma$  des effets aléatoires lors des estimations de ces variances.

### Résolution numérique par l'algorithme E.M.

Une approche par l'algorithme EM est également possible. En effet, en reprenant le résultat (ainsi que les notations) de la page 78, il est possible d'exprimer l'espérance conditionnelle de  $\tilde{\Theta}_i$  sachant  $(\mathbf{N}_i, \mathbf{Y}_i)_{i=1, \dots, n}$  en fonction de la transformée de Laplace de  $\tilde{\Theta}_i$  :

$$\mathbb{E} [\tilde{\theta}_i \mid (\mathbf{N}_i, \mathbf{Y}_i)] = -\frac{\mathcal{L}^{(s_i+1)}(\Lambda_{i+}(\tau))}{\mathcal{L}^{(s_i)}(\Lambda_{i+}(\tau))}. \quad (5.3)$$

L'algorithme E.M. se déroule alors comme suit :

- étape de maximisation (*M-step*) : pour une valeur courante de  $\gamma$ , calcul de  $\hat{\beta}$  et  $\hat{\alpha}_0$  comme dans un modèle standard de Cox avec la fragilité (fixe) en *offset* ;
- étape d'estimation (*E-step*) : calcul de  $\hat{\theta}_i$  comme valeur attendue étant donnée  $\hat{\beta}$ ,  $\hat{\alpha}_0$  et les observations.

L'étape E, d'après (5.3), revient donc à calculer

$$\exp(\theta_i) = -\frac{\mathcal{L}^{(s_i+1)}(\hat{\Lambda}_{i+}(\tau))}{\mathcal{L}^{(s_i)}(\hat{\Lambda}_{i+}(\tau))}, \quad (5.4)$$

où

$$\hat{\Lambda}_{i+}(\tau) = \sum_j \int_0^\tau Y_{ij}(s) \hat{\alpha}_0(s) \exp(\beta^t Z_{ij}) ds.$$

Les estimateurs obtenus dépendent de  $\gamma$  et peuvent être notés  $\hat{\beta}(\gamma)$  et  $\hat{\alpha}_0(\gamma)$ . Aussi Parner suggère-t-il d'estimer  $\gamma$  à partir de la maximisation de la vraisemblance partielle

$$\mathcal{L}_m(\gamma) = \mathcal{L}_m(\hat{\beta}(\gamma), \hat{\alpha}_0(\gamma), \gamma). \quad (5.5)$$

où

$$\begin{aligned} \mathcal{L}_m(\beta, \alpha_0, \gamma) &= \sum_{i,j} \delta_{ij} \log \left[ \int_0^\tau \alpha_0(t) Y_{ij}(t) \exp(\beta^t Z_{ij}) dt \right] \\ &\quad + \sum_{i,j} \log \left[ (-1)^{\varsigma_i} \mathcal{L}_{\hat{\theta}}^{(\varsigma_i)} \left( \int_0^\tau \alpha_0(t) Y_{ij}(t) \exp(\beta^t Z_{ij}) dt \right) \right], \end{aligned} \quad (5.6)$$

d'après le résultat vu en page 78.

### Équivalence des estimations fournies par ces deux résolutions

**Lemme III.1** — *La solution du modèle de vraisemblance pénalisée – dans le cas d'une loi gamma comme dans le cas gaussien – coïncide avec la solution de l'algorithme E.M. de l'équation (4.20) avec  $\gamma$  fixé.*

*Démonstration* — Concernant  $\beta$ , la méthode par l'algorithme E.M. et celle par vraisemblance pénalisée ont la même équation de mise à jour de la variable, laquelle inclut  $\theta_i$  en temps qu'*offset* fixe. Par conséquent, si les solutions pour  $\theta_i$  sont les mêmes, alors celles concernant  $\beta$  le seront également.

Dans le cas où la fragilité suit une loi gamma de paramètre  $\gamma$ ,

$$\mathcal{L}^{(d)}(s) = (-\gamma)^d (1 + s\gamma)^{-(d+\frac{1}{\gamma})} \prod_{i=0}^{d-1} \left( \frac{1}{\gamma} + i \right)$$

et par suite,

$$\exp(\theta_i) = -\frac{\varsigma_i + \frac{1}{\gamma}}{\hat{A}_i + \frac{1}{\gamma}}. \quad (5.7)$$

Soit  $(\hat{\beta}, \hat{\theta}_i) = (\hat{\beta}(\gamma), \hat{\theta}_i(\gamma))$  une solution de l'algorithme E.M.  $\hat{\theta}_i$  satisfait exactement l'équation (5.7). Cette équation se réécrit

$$\hat{A}_i = \exp(-\hat{\theta}_i) \left( \varsigma_i + \frac{1}{\gamma} \right) - \frac{1}{\gamma}.$$

Si l'on introduit cette dernière équation dans l'équation du score pénalisé, et si l'on tient compte du fait que

$$\frac{\partial g(\theta_i; \gamma)}{\partial \theta_i} = 1 - \exp(\theta_i),$$

alors nous obtenons, au point  $(\hat{\beta}, \hat{\theta}_i)$ ,

$$\begin{aligned} \frac{\partial \text{PPL}(\beta, \alpha_0, \hat{\theta}_i, \gamma)}{\partial \hat{\theta}_i} &= \left\{ \varsigma_i - \exp(-\hat{\theta}_i) \left[ \varsigma_i + \frac{1}{\gamma} - \frac{1}{\gamma} \exp(\hat{\theta}_i) \right] \exp(\hat{\theta}_i) \right\} + \frac{1}{\gamma} (1 - \exp(\hat{\theta}_i)) \\ &= 0 \end{aligned}$$

pour tout  $i$ . Ceci montre que la solution de l'algorithme E.M. est aussi une solution des équations du score pénalisé.

Cependant, à  $\gamma$  fixé, la log-vraisemblance pénalisée et la log-vraisemblance des données observées (4.20), bien qu'elles aient la même solution, ne sont pas égales. ■

Nous pouvons écrire le profil de la log-vraisemblance suivant  $\gamma$  – c'est-à-dire l'équation (5.5) – comme un profil de la vraisemblance pénalisée auquel est ajouté un terme correctif qui ne dépend que de  $\gamma$  et de  $\varsigma_i$ . En utilisant le fait que la fragilité est la même pour tous les individus d'un groupe, nous obtenons que la vraisemblance partielle de Cox pour  $(\hat{\beta}, \hat{\theta}_i)$  doit être la même que celle relative au couple  $(\hat{\beta}, \hat{\theta}_i + c)$  pour toute constante  $c$ . Un calcul algébrique montre que la valeur  $c$  qui minimise la fonction de pénalité doit satisfaire

$$\sum_i \exp(\theta_i) = n. \quad (5.8)$$

Le lemme III.1 peut être obtenu directement à partir de la relation suivante (démontrée ci-dessous dans le cas d'une fragilité suivant une loi gamma) :

**Proposition III.1** — *Nous avons*

$$\mathcal{L}_m(\gamma) = PPL(\gamma) + \sum_i \left[ \frac{1}{\gamma} - \left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) + \frac{1}{\gamma} \log\left(\frac{1}{\gamma}\right) + \log\left(\frac{\Gamma\left(\frac{1}{\gamma} + \varsigma_i\right)}{\Gamma\left(\frac{1}{\gamma}\right)}\right) \right], \quad (5.9)$$

où  $PPL(\gamma) = PPL(\hat{\beta}(\gamma), \alpha_0(\gamma), \hat{\theta}_i(\gamma), \gamma)$  et où  $\mathcal{L}_m(\gamma)$  est défini par l'équation (5.5).

*Démonstration* — Le logarithme de la fonction de densité de  $\tilde{\Theta}_i$  est

$$\log [g(\tilde{\theta}_i)] = \left(\frac{1}{\gamma} - 1\right) \log(\tilde{\theta}_i) - \frac{1}{\gamma} \tilde{\theta}_i + \frac{1}{\gamma} \log\left(\frac{1}{\gamma}\right) - \log \left[\Gamma\left(\frac{1}{\gamma}\right)\right]$$

et la densité a une transformée de Laplace de la forme

$$\mathcal{L}(s) = (1 + s\gamma)^{-\frac{1}{\gamma}}$$

dont sa dérivée d'ordre  $d$  vaut

$$\mathcal{L}^{(d)}(s) = (-\gamma)^d (1 + s\gamma)^{-\left(\frac{1}{\gamma} + d\right)} \prod_{i=0}^{d-1} \left(\frac{1}{\gamma} + i\right).$$

L'équation (5.6) devient alors

$$\begin{aligned} \mathcal{L}_m(\beta, \alpha_0; \gamma) &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau Y_{ij}(t) \alpha_0(t) \exp(\beta^t Z_{ij}) dt \right) \\ &\quad + \sum_i \log \left[ \gamma^{\varsigma_i} (1 + \gamma A_i)^{-\left(\frac{1}{\gamma} + \varsigma_i\right)} \prod_{k=0}^{\varsigma_i-1} \left( \frac{1}{\gamma} + k \right) \right] \end{aligned} \quad (5.10)$$

$$\begin{aligned} &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau Y_{ij}(t) \alpha_0(t) \exp(\beta^t Z_{ij}) dt \right) \\ &\quad + \sum_i \left[ \varsigma_i \log(\gamma) - \left( \frac{1}{\gamma} + \varsigma_i \right) \log(1 + \gamma A_i) \right. \\ &\quad \left. + \log \left[ \Gamma\left(\frac{1}{\gamma} + \varsigma_i\right) - \log \left[ \Gamma\left(\frac{1}{\gamma}\right) \right] \right] \right]. \end{aligned} \quad (5.11)$$

Considérons cette équation restreinte à la courbe unidimensionnelle définie par les valeurs de maximisation  $\hat{\beta}(\gamma)$ ,  $\hat{\theta}_i(\gamma)$  et  $\hat{\alpha}_0(\gamma)$  pour chaque  $\gamma$ . D'après (5.7),

$$\hat{A}_i = \frac{d_i + \frac{1}{\gamma} - \frac{1}{\gamma} \exp(\hat{\theta}_i)}{\exp(\hat{\theta}_i)}.$$

En substituant cette expression dans la précédente, nous obtenons

$$\begin{aligned} \mathcal{L}_m(\beta, \alpha_0; \gamma) &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau \hat{\lambda}_{ij}(t) \exp(\hat{\beta}^t Z_{ij}) dt \right) + \sum_i \left[ -\varsigma_i \log\left(\frac{1}{\gamma}\right) - \left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) \right. \\ &\quad \left. + \left(\frac{1}{\gamma} + \varsigma_i\right) \log \left[ \frac{1}{\gamma} \exp(\hat{\theta}_i) \right] + \log \left[ \Gamma\left(\frac{1}{\gamma} + \varsigma_i\right) \right] \right] \\ &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau \hat{\lambda}_{ij}(t) \exp(\hat{\beta}^t Z_{ij}) dt \right) + \sum_i \left[ -\left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) \right. \\ &\quad \left. + \varsigma_i \log \left[ \exp(\hat{\theta}_i) \right] + \frac{1}{\gamma} \log \left[ \frac{1}{\gamma} \exp\left(\frac{1}{\gamma} \hat{\theta}_i\right) \right] - \log \left[ \Gamma\left(\frac{1}{\gamma}\right) \right] + \log \left[ \Gamma\left(\frac{1}{\gamma} + \varsigma_i\right) \right] \right] \\ &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau \hat{\lambda}_{ij}(t) \exp(\hat{\beta}^t Z_{ij} + \hat{\theta}_i) dt \right) + \sum_i \left[ -\left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) \right. \\ &\quad \left. + \frac{1}{\gamma} \log \left[ \frac{1}{\gamma} \exp(\hat{\theta}_i) \right] - \log \left[ \Gamma\left(\frac{1}{\gamma}\right) \right] + \log \left[ \Gamma\left(\frac{1}{\gamma} + \varsigma_i\right) \right] \right] \end{aligned}$$

En évaluant l'expression précédente en  $\hat{\beta}$ , après avoir ajouté la fonction de pénalité définie par (5.2), nous obtenons finalement

$$\begin{aligned} \mathcal{L}_m(\beta, \alpha_0; \gamma) &= \sum_{i,j} \delta_{ij} \log \left( \int_0^\tau \hat{\lambda}_{ij}(t) \exp(\hat{\beta}^t Z_{ij} + \hat{\theta}_i) dt \right) - h(\hat{\theta}_i; \gamma) + \sum_i \left[ -\frac{1}{\gamma} \hat{\theta}_i + \frac{1}{\gamma} \exp(\hat{\theta}_i) \right. \\ &\quad \left. - \left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) + \frac{1}{\gamma} \log \left[ \frac{1}{\gamma} \exp(\hat{\theta}_i) \right] - \log \left[ \Gamma\left(\frac{1}{\gamma} + \varsigma_i\right) \right] - \log \left[ \Gamma\left(\frac{1}{\gamma}\right) \right] \right] \\ &= \text{PPL}(\gamma) + \sum_i \left[ \frac{1}{\gamma} - \left(\frac{1}{\gamma} + \varsigma_i\right) \log\left(\frac{1}{\gamma} + \varsigma_i\right) + \frac{1}{\gamma} \log\left(\frac{1}{\gamma}\right) + \log \left( \frac{\Gamma\left(\frac{1}{\gamma} + \varsigma_i\right)}{\Gamma\left(\frac{1}{\gamma}\right)} \right) \right] \end{aligned}$$

où la dernière égalité fait appel à l'équation (5.8). ■

**Remarque** — Dans le cas d'une fragilité log-normale, il est possible de noter la similarité entre l'étape consistant à utiliser l'algorithme de Newton-Raphson lors de la résolution du modèle de Cox, et le calcul par les moindres carrés pondérés itérativement (*iteratively reweighted least-squares*) (MacGilchrist et Aisbett, 1991 ; MacGilchrist, 1993). Par suite, l'utilisation des méthodes d'estimation standard des problèmes gaussiens est envisageable. Ainsi,  $\gamma$  est choisi de telle manière qu'il satisfasse

$$\gamma = \frac{\sum_i \theta_i^2 + r}{n},$$

où  $r$  est un paramètre variant suivant la technique d'estimation : si l'on note  $H$  le minimum de la matrice hessienne de la vraisemblance pénalisée,

- pour l'estimateur du BLUP,  $r = 0$  ;
- pour l'estimateur du maximum de vraisemblance,  $r = \text{trace}[(H_{22})^{-1}]$  ;
- pour l'estimateur du maximum de vraisemblance restreinte,  $r = \text{trace}[(H^{-1})_{22}]$ .

### 5.1.4 Types et fonctionnalités des résidus

Les différents types de résidus fournis par S-PLUS sont :

- les résidus de Cox-Snell ;
- les résidus de déviance ;
- les résidus du score ;
- les résidus de Schoenfeld.

Quatre usages peuvent en être faits, qui ont été présentés au paragraphe 2.4 :

- la recherche de la forme fonctionnelle correcte concernant une covariable ;
- l'identification des sujets pour lesquels la justesse du modèle paraît douteuse ;
- l'identification des « observations influentes » ;
- la vérification de l'hypothèse de proportionnalité des risques.

## 5.2 Présentation des simulations

### 5.2.1 Objectifs

Nous souhaitons, au travers des simulations que nous entreprenons ici, répondre à deux interrogations : la première concerne la qualité des différents ajustements qui peuvent être menés sur un jeu de données de survie corrélées, ainsi que leur aptitude à prendre correctement en compte cette corrélation ; la seconde porte sur la dépendance de ces ajustements vis-à-vis de certains paramètres, parmi lesquels la taille des groupes, le taux de censure ou encore l'intensité de la corrélation.

Pour ce faire, nous avons procédé comme suit. À la première étape, nous choisissons la répartition en groupes des données de survie que nous allons générer ; deux cas de figures ont

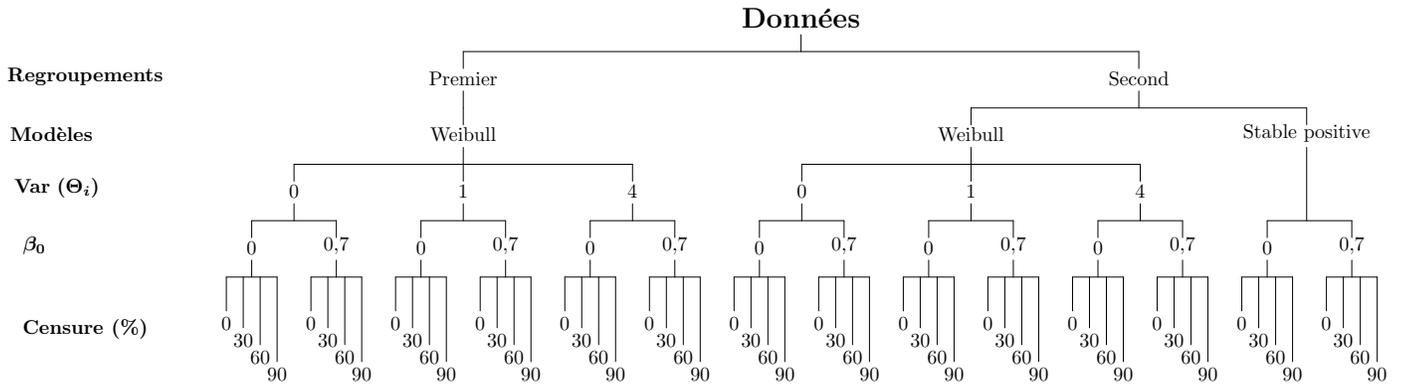


FIG. 5.1 – Mécanisme des simulations pas à pas.

été envisagés : 100 groupes de 10 données chacun (premier regroupement), ou 10 groupes de 100 temps chacun (second regroupement).

À la deuxième étape, deux modèles sont programmés pour la génération des temps de survie : un modèle de Weibull avec une fragilité suivant une loi gamma, et un modèle avec une fragilité stable positive<sup>1</sup>.

Lors de la troisième étape, nous affectons à la valeur initiale du coefficient de régression  $\beta$  les valeurs 0 et 0,7, et nous faisons varier la variance de la fragilité – qui peut prendre les valeurs 0, 1 et 4.

enfin, la quatrième étape concerne la censure des données, pour laquelle quatre cas de figures ont été sélectionnés : absence de censure, ou bien présence de censure à un taux de 30 %, 60 % ou 90 %. Une fois les données de survie générées, nous les avons triées par ordre croissant, afin de censurer les temps les plus longs.

Une fois les données générées, nous avons fait appel à cinq types d’ajustement par le modèle de Cox :

- le modèle naïf sans aucun « effet groupe » ;
- le modèle naïf comprenant un « effet groupe » fixe (l’identifiant du groupe étant inclus, lors de l’ajustement, sous la forme d’une covariable factorielle) ;
- le modèle robuste marginal ;
- le modèle robuste incluant une fragilité suivant une loi gamma ;
- le modèle robuste incluant une fragilité suivant une loi log-normale.

La figure 5.1 résume ces étapes successives.

### 5.2.2 Modèle de fragilité gamma

Nous nous servons d’un modèle de Weibull. Les densité et fonction de répartition d’une loi de Weibull  $W(a, b)$  sont respectivement

$$h(x) = ab^a t^{a-1} \exp[-(bt)^a] \mathbb{1}_{\{t>0\}}$$

1. Voir les deux sous-sections suivantes pour plus de détails concernant ces modèles.

et

$$H(x) = \left[ 1 - \exp(- (bt)^a) \right] \mathbf{1}_{\{t>0\}} .$$

Le risque instantané vaut

$$\alpha(t) = \theta_i \gamma \alpha_0^\gamma t^{\gamma-1} \exp(\boldsymbol{\beta}' \mathbf{Z})$$

où  $\theta_i$  est la fragilité.

Par suite, comme

$$f(t) = \alpha(t) \exp(- A(t)) ,$$

où

$$A(t) = \int_0^t \alpha(s) ds ,$$

nous avons

$$f(t) = \gamma \alpha_0^\gamma t^{\gamma-1} \theta_i \exp(\boldsymbol{\beta}' \mathbf{Z}) \exp \left[ - (\alpha_0 t)^\gamma \theta_i \exp(\boldsymbol{\beta}' \mathbf{Z}) \right] .$$

$f(t)$  est donc la densité d'une loi de Weibull  $W(\gamma, \alpha_0 \theta_i \exp[(\boldsymbol{\beta}' \mathbf{Z})/\gamma])$ .

### 5.2.3 Modèle de fragilité stable positive

Voici un résultat préliminaire (Chambers *et al.*, 1976 ; Hougaard, 1986).

Si  $\Theta_i, H_{i1}, \dots, H_{in_i}$  sont indépendantes, si  $\Theta_i$  suit une loi stable positive de paramètre  $\zeta$  (cf. p. 65), si les  $H_{ij}$  suivent des lois de Weibull de paramètre de forme commun  $\gamma$  et de paramètres d'échelle  $\exp(-\boldsymbol{\beta}^t \mathbf{Z}_{ij})$ , et si l'on pose  $T_{ij} = \Theta_i^{-1/\gamma} H_{ij}$ , alors  $T_{i1}, \dots, T_{in_i}$

- ont des distributions marginales de Weibull  $W(\zeta\gamma, \exp(-\zeta \boldsymbol{\beta}^t \mathbf{Z}_{ij}))$  ;
- suivent une loi de Weibull multivariée de fonction de survie jointe

$$\begin{aligned} \mathbb{P}(T_{i1} > t_1; \dots, T_{in_i} > t_{n_i}) &= \exp \left[ - (\Lambda_{i1}(t_{i1}) + \dots + \Lambda_{in_i}(t_{in_i}))^\zeta \right] \\ &= \exp \left[ - (\eta_{i1} t_1^\gamma + \dots + \eta_{in_i} t_{n_i}^\gamma)^\zeta \right] \end{aligned}$$

- sont telles que

$$\text{corr}(T_{ij}, T_{ik}) = 1 - \zeta \quad \forall j \neq k .$$

Nous adoptons un modèle de Weibull comportant un terme de fragilité distribué suivant une loi stable positive.

Hougaard (1986) a montré que si la distribution d'une v.a.  $T$  sachant  $\Theta_i$  (fragilité associée à  $T$ ) est une loi de Weibull  $W(\gamma, \omega \Theta_i)$  et si  $\Theta_i$  suit une loi stable positive  $P(\alpha)$ , alors la distribution de  $T$  est une loi de Weibull  $W(\alpha\gamma, \omega^\alpha)$ .

Considérons un modèle de fragilité de Weibull: le risque instantané pour l'individu  $j$  du groupe  $i$ , sachant  $\Theta_i$  (terme de fragilité), s'écrit

$$\alpha_{ij}(t) = \gamma t^{\gamma-1} \Theta_i e^{-\boldsymbol{\beta}^t \mathbf{Z}_{ij}} .$$

Par suite, la densité du temps de survie sachant  $\Theta_i$  vaut

$$\begin{aligned} f(t) &= \alpha_{ij}(t) \exp(-A_{ij}(t)) \\ &= \gamma t^{\gamma-1} \Theta_i e^{-\beta^t \mathbf{Z}_{ij}} e^{-t^\gamma \Theta_i e^{-\beta^t \mathbf{Z}_{ij}}} \mathbb{1}_{t \geq 0}. \end{aligned}$$

C'est la distribution d'une loi de Weibull  $W(\gamma, \Theta_i \exp(-\beta^t \mathbf{Z}))$ .

Supposons que la fragilité suive une loi stable positive de paramètre (index)  $0 < \zeta < 1$ .

D'après le résultat vu ci-dessus, nous obtenons que  $T_{ij}$  a une distribution marginale de Weibull  $W(\zeta\gamma, \exp(-\zeta\beta^t \mathbf{Z}_{ij}))$ .

De plus, la corrélation entre deux temps de survie d'un même groupe (mesuré par le tau de Kendall) vaut  $\tau = 1 - \zeta$ .

Ainsi, les simulations sont conduites de la façon suivante (Segal *et al.*, 1997) :

- tirage au sort de variables  $H_{ij}$  suivant des lois de Weibull  $W(\gamma, \exp(-\beta \mathbf{Z}_{ij}))$  ;
- tirage au sort (indépendamment des  $H_{ij}$ ) d'une variable  $\Theta_i$  suivant une loi stable positive ;
- calcul des temps de survie corrélés  $T_{ij} = Y_{ij} \Theta^{-1/\gamma}$ .

### 5.3 Paramétrage des simulations et choix des sorties

Nous avons adopté le schéma suivant :

- 1° les valeurs des paramètres ont été fixées comme suit :
  - le paramètre de la loi de Weibull  $\gamma$  prend la valeur 1,1, ce qui correspond à une fonction de risque croissante,
  - le risque de base  $\alpha_0$  vaut 0,01,
  - le paramètre de régression de la covariable,  $\beta$ , prend deux valeurs : 0 ou 0,7, ce qui correspond à des risques relatifs<sup>1</sup> valant respectivement 1 et (approximativement) 2,
  - $\mathbf{Z} = 0$  ou 1 avec la probabilité 1/2 ;
- 2° la fragilité  $\Theta_i$  suit une gamma d'espérance 1 et de variance  $\sigma^2$  valant 1 ou 4 ;
- 3° conformément à ce qui a été exposé à la section 5.2.2 p. 93, nous avons tiré les temps de survie au travers de l'inverse d'une distribution de Weibull  $W\left(\gamma, 1/(\alpha_0 \theta_i \exp[(\beta^t \mathbf{Z}/\gamma])\right)$  ;
- 4° deux situations ont été posées : une première comportant 100 groupes de 10 individus chacun, une seconde comprenant 10 groupes de 100 individus ; elles sont respectivement désignées sous les termes de **premier** et **second regroupement** ;
- 5° le taux de censure a pris successivement les quatre valeurs suivantes : 0 %, 30 %, 60 % et 90 % ;
- 6° 500 simulations ont été réalisées pour chacune des 96 configurations possibles, récapitulées dans le tableau 5.1.

---

1. Le risque relatif est le rapport  $R_1/R_0$  des risques de maladie chez les sujets exposés au facteur de risque étudié – indice 1 – et chez les sujets non-exposés – indice 0.

TAB. 5.1 – Récapitulatif des configurations choisies, pour chacune desquelles 500 simulations ont été réalisées.

Paramètre	Valeurs
Risque de base $\alpha_0(t)$	0,01
Temps de survie	inverse Weibull
Covariable $Z_{ij}$	dichotomique et individuelle
Coefficient de régression $\beta$	0 et 0,7
Distribution de la fragilité $\Theta_i$	log-normale et stable positive
Variance de la fragilité $\sigma^2$	0, 1 et 4
Taille des groupes	100 groupes de 10 individus et 10 groupes de 100 individus
Taux de censure	0 %, 30 %, 60 % et 90 %
<b>Nombre de configurations</b>	<b>56</b>

Nous nous sommes intéressés au biais<sup>1</sup> concernant l'estimation  $\hat{\beta}$  de  $\beta$  et à l'estimation de la variance de  $\hat{\beta}$ .

Plus précisément, si nous notons  $\hat{\beta}_i$  le  $i^e$  échantillon simulé, alors

$$\text{biais} = \frac{\sum_i \hat{\beta}_i}{500} - \beta_0 ,$$

où  $\beta_0$  est la valeur de  $\beta$  à la base des simulations.

Nous nous sommes également intéressés au coefficient de dispersion, qui vaut

$$\frac{[\sum_i (\hat{\beta}_i - \bar{\hat{\beta}})^2]_{\text{robuste}}}{[\sum_i (\hat{\beta}_i - \bar{\hat{\beta}})^2]_{\text{naïf}}} ,$$

où  $\bar{\hat{\beta}} = (\sum_i \hat{\beta}_i)/500$ .

Outre les estimations du paramètre de régression  $\beta$  et de sa variance, nous avons considéré les conclusions rendues par les trois tests usuels – tests du score, du rapport de vraisemblance et de Wald – ; étant donnée la similitude des résultats, nous nous sommes restreints à l'écriture des résultats d'un seul : le test du maximum de vraisemblance.

De même, nous n'avons retranscrit ici que les conclusions rendues par le test du maximum de vraisemblance lors d'un ajustement *via* un modèle de fragilité de loi gamma ; cependant, nous précisons que ces résultats demeurent quasi parfaitement les mêmes pour le modèle de fragilité de loi normale.

1. Nous employons abusivement le terme de « biais » concernant la différence entre d'une part la valeur  $\beta_0$  du paramètre de régression telle qu'elle est fixée lors des simulations, et d'autre part la valeur de l'estimation  $\hat{\beta}$  que fournit le modèle de Cox : en effet, l'estimation obtenue par ajustement du modèle de Cox est celle d'un paramètre qui n'est pas exactement  $\beta_0$ , étant donné que  $\beta_0$  est fixé avant l'introduction de la variabilité intra-groupe. Le biais, en statistique, mesure l'écart entre la valeur réelle et la valeur estimée d'un paramètre : ce qui fait ici défaut est la valeur *réelle* du paramètre, que nous ne connaissons pas, et qui se distingue de la valeur d'origine du paramètre, c'est-à-dire la valeur fixée lors de la mise en œuvre des simulations

Nous nous sommes intéressés au risque de première espèce et à la puissance du test du maximum de vraisemblance, ainsi qu'au taux de recouvrement de l'intervalle de confiance du paramètre de variance de la fragilité.

Concernant l'estimation du risque de première espèce, nous procédons comme suit (King *et al.*, 1996) :

- on se fixe  $\beta = \beta_0$  (valeur à la base des simulations) ;
- on simule  $k$  données et on réplique  $n$  fois ces  $k$  données ;
- on regarde le nombre  $\hat{n}$  de réplifications où, quand on compare  $\hat{\beta}$  à  $\beta_0$  via la statistique du test du maximum de vraisemblance, l'erreur de type I à 5 % (resp. à 1 %) est significative ;
- plus le taux  $\frac{\hat{n}}{n}$  est proche de 0,05 (resp. 0,01), meilleur est le test.

Concernant la puissance, qui mesure la précision de  $\hat{\beta}$ , la procédure est la suivante :

- on se fixe  $\beta = \beta_{(1)}$  différent de  $\beta_0$  ;
- on simule  $k$  données et on réplique  $n$  fois ces  $k$  données ;
- on regarde le nombre  $\hat{n}$  de réplifications où on rejette  $H_0 : \beta = \beta_0$  ;
- on réitère les trois étapes précédentes pour différentes valeur  $\beta_{(2)}, \beta_{(3)} \dots$  afin de tracer le graphe de la puissance.

Nous avons également calculé le taux de recouvrement de l'intervalle de confiance à 95 % du paramètre  $\gamma$ , soit la proportion de simulations telles que  $\gamma \in [\hat{\gamma} - 1,96\sqrt{\mathbb{V}(\hat{\gamma})}, \hat{\gamma} + 1,96\sqrt{\mathbb{V}(\hat{\gamma})}]$ .

Enfin, nous avons calculé le Critère d'Information d'Akaike (*Akaike's Information Criterion*) pour chacun des ajustements. Rappelons la définition de ce critère (Akaike, 1973) :

**Définition III.1** — Soit  $\{X_i\}_{i=1,\dots,n}$  un  $n$ -échantillon *i.i.d.* On dispose pour cet échantillon de  $M$  modèles, chacun caractérisé par un vecteur de  $k_m$  paramètres  $\beta_m$ . On note  $\hat{\beta}_m$  l'estimation de ce vecteur par le maximum de vraisemblance et  $f_{\hat{\beta}_m}$  la fonction de densité alors obtenue. La procédure de choix de modèle revient à maximiser, pour  $m = 1, \dots, M$ , le critère d'information d'Akaike défini par

$$\text{AIC}(m) = 2 \sum_{i=1}^n \log [f_{\hat{\beta}_m}(X_i)] - 2k_m .$$

Rappelons également que la procédure de maximisation du critère d'Akaike permet de comparer des modèles non nécessairement emboîtés – la seule condition étant que les paramètres de ces modèles soient estimés par le maximum de vraisemblance. Par ailleurs, il s'agit d'un critère de choix, et non d'un *test* statistique. Son utilisation conduit à opérer un classement des modèles candidats, mais on ne peut lui associer un risque de première ou de seconde espèce.

Enfin, pour tout complément d'information concernant les résultats des simulations, nous renvoyons le lecteur à l'annexe E, où se trouvent les tableaux exhaustifs des informations recueillies.

## 5.4 Résultats

### 5.4.1 Estimations du paramètre de régression $\beta$

Notons tout d'abord que les résultats concernant les deux ajustements par des modèles de fragilité – l'un avec une distribution gamma pour la fragilité, l'autre avec une distribution nor-

TAB. 5.2 – *Biais relatif (en %) de l'estimation du coefficient de régression pour  $\sigma^2 = 1$  et, entre parenthèses,  $\sigma^2 = 4$ .*

Reg.*	Censure	Naïf	Fixe	Mixte
1	0 %	-50 (-80)	8 (0)	0 (0)
	30 %	-38 (-73)	4 (3)	0 (0)
	60 %	-22 (-57)	9 (6)	1 (1)
	90 %	-4 (-16)	8 (8)	2 (2)
2	0 %	-42 (-67)	0 (0)	0 (0)
	30 %	-36 (-68)	0 (0)	0 (0)
	60 %	-30 (-52)	0 (0)	0 (0)
	90 %	-3 (-10)	2 (2)	2 (1)

\* : regroupement.

male – sont extrêmement similaires ; aussi nous contenterons-nous, par la suite, de retranscrire et de commenter les résultats du seul ajustement avec un modèle de fragilité gamma de Cox.

Le tableau 5.2 représente le biais relatif de l'estimation du coefficient de régression, soit le rapport du biais obtenu sur la valeur initiale de  $\beta$  – ici, la valeur initiale de  $\beta$  varie de 0 à 1, 2. Les évolutions des situations – vis-à-vis du taux de censure ou de la valeur initiale de  $\beta$  – dans le cas d'une variance de la fragilité égale à 1 ou dans celui d'une variance égale à 4, sont similaires.

**Comparaison des ajustements** Le modèle naïf sans « effet groupe » présente le biais le plus important, tandis que le biais le plus faible est obtenu avec l'ajustement mixte ; l'ajustement naïf avec un « effet groupe » fixe, tient une position intermédiaire, avec un biais légèrement plus important que celui du modèle mixte, mais qui demeure raisonnablement faible.

**Influence de la valeur initiale du paramètre de régression** Lorsque la valeur initiale de  $\beta$  augmente, le biais du modèle naïf simple augmente grandement : ainsi, dans le cas d'un taux de censure nul et d'une variance de la fragilité égale à 1, ce biais passe de la valeur 0,002 lorsque  $\beta = 0$  à la valeur (absolue) 0,260 lorsque  $\beta = 0,50$ .

Cette augmentation du biais à mesure que la valeur initiale de  $\beta$  augmente est également observée pour l'ajustement naïf avec « effet groupe » dans le cas du premier regroupement. En revanche, dans le cas du second regroupement, tout comme dans celui de l'ajustement avec un « effet groupe » aléatoire, cette augmentation du biais n'a pas lieu.

**Influence de la variance de la fragilité** Le passage de 1 à 4 de la variance (initiale) de la fragilité n'a que peu d'influence sur le biais de l'ajustement avec « effet groupe » aléatoire. Concernant le modèle naïf simple, l'augmentation de la variance initiale de la fragilité a pour effet d'accroître sensiblement le biais, qui passe par exemple de la valeur (absolue) 0,103 à la

valeur (absolue également) 0,289 – pour  $\beta = 0,50$ , un taux de censure de 60 % et un regroupement du premier type. *A contrario*, ce biais diminue lorsque nous incluons, dans le modèle naïf, un « effet groupe » fixe.

**Influence du type de regroupement** Le biais diminue lors du passage du premier type de regroupement (100 groupes de 10 individus) au second (10 groupes de 100 individus), dans le cas d'un ajustement naïf – qu'il intègre ou non un « effet groupe ». En revanche, l'ajustement avec un « effet groupe » aléatoire n'est pas sensible à cette distinction entre les différents regroupements.

Notons que la diminution du biais pour l'ajustement avec « effet groupe » fixe est très importante : ainsi, lorsque nous nous plaçons dans le cadre du second type de regroupement, nous pouvons affirmer que les biais des deux ajustements avec « effet groupe » sont extrêmement similaires, et très faibles.

**Influence du taux de censure** L'augmentation du taux de censure affecte de façon importante l'estimation du paramètre de régression dans le cadre du modèle naïf. En l'absence d'« effet groupe », le biais observé diminue considérablement à mesure que la censure se fait plus sévère : ainsi, ce biais passe de la valeur absolue 0,260 à la valeur absolue 0,103 (respectivement 0,025), lorsque le taux de censure passe de la nullité à la valeur de 60 % (resp. 90 %).

En présence d'un « effet groupe » (fixe), c'est le résultat inverse qui est observé, soit donc une augmentation du biais relatif à l'estimation du paramètre de régression, à mesure que le taux de censure croît.

Enfin, l'ajustement avec « effet groupe » aléatoire est indifférent à une variation du taux de censure, le biais demeurant dans tous les cas extrêmement faible.

En l'absence de corrélation lors des simulations des données de survie, nous obtenons logiquement le résultat inverse, à savoir un biais minimal (en valeur absolue) pour le modèle sans effet groupe : en effet, pour des taux de censure fixés à 0 %, 30 %, 60 % et 90 %, le biais maximal (en valeur absolue) est respectivement de

- 0,005, 0,005, 0,010 et 0,021 pour le modèle sans « effet groupe » ;
- 0,058, 0,041, 0,042 et 0,054 pour le modèle avec « effet groupe » fixe ;
- 0,006, 0,006, 0,010 et 0,022 pour le modèle avec « effet groupe » aléatoire.

En conclusion, nous pouvons souligner l'ordre d'importance des différents biais observés : le modèle sans « effet groupe » présente le biais le plus important, le modèle avec « effet groupe » aléatoire procurant, lui, l'estimation du paramètre de régression la moins « biaisée » ; enfin, en position intermédiaire, nous plaçons l'ajustement avec « effet groupe » fixe, dont le biais approche par valeurs supérieures celui du modèle avec « effet groupe » aléatoire (à l'exception du cas où le taux de censure est de 90 %).

TAB. 5.3 – Coefficient de dispersion et (entre parenthèses) critère d’Akaike pour  $\beta_0 = 0$ .

$\sigma^2$	Regroup.*	Censure	Effet fixe	Marginal	Fragilité
1	1	0 %	1,12 (1 036,7)	0,97 (6,4)	1,10 (835)
		30 %	1,10 (771,3)	0,99 (6,5)	1,08 (567)
		60 %	1,07 (492,6)	1,00 (6,4)	1,06 (286)
		90 %	1,05 (193,1)	0,98 (2,9)	1,03 (23)
	2	0 %	1,01 (816,3)	0,86 (8,5)	1,01 (795)
		30 %	1,01 (631,8)	0,90 (7,8)	1,01 (610)
		60 %	1,00 (372,8)	0,93 (6,6)	1,00 (404)
		90 %	1,00 (98,6)	0,87 (3,0)	1,00 (76)
4	1	0 %	1,12 (1 951)	0,94 (2)	1,11 (2 267)
		30 %	1,12 (1 444)	0,99 (2)	1,10 (1 580)
		60 %	1,10 (988)	0,99 (2)	1,08 (850)
		90 %	1,06 (309)	0,98 (2)	1,05 (112)
	2	0 %	1,02 (1 591)	0,86 (3)	1,01 (1 949)
		30 %	1,01 (1 432)	0,92 (3)	1,01 (1 556)
		60 %	1,01 (902)	0,93 (3)	1,01 (908)
		90 %	1,00 (201)	0,79 (2)	1,00 (198)

\* : regroupement.

### 5.4.2 Estimations de la variance du paramètre de régression – coefficient de dispersion

Lorsque  $\beta_0 = 0$ , la valeur moyenne de l’estimation de la variance du paramètre  $\beta$  est de l’ordre de (cf. annexe E)

- 0,06, 0,07, 0,10 et 0,20 pour des taux de censure fixés respectivement à 0 %, 30 %, 60 % et 90 % et pour le modèle sans effet groupe d’une part, les modèles avec effet groupe dans le cas du second regroupement d’autre part ;
- 0,07, 0,08, 0,10 et 0,21 pour les mêmes taux de censure respectifs et pour les modèles avec effet groupe (fixe comme aléatoire) dans le cas du premier regroupement.

Lorsque  $\beta_0 = 0,7$ , la valeur moyenne de l’estimation de la variance du paramètre  $\beta$  est de l’ordre de

- 0,06, 0,07, 0,10 et 0,22 pour des taux de censure fixés respectivement à 0 %, 30 %, 60 % et 90 % et pour le modèle sans effet groupe d’une part, les modèles avec effet groupe dans le cas du second regroupement d’autre part ;
- 0,07, 0,08, 0,11 et 0,22 pour les mêmes taux de censure respectifs et pour les modèles avec effet groupe (fixe comme aléatoire) dans le cas du premier regroupement.

Les tableaux 5.3 et 5.4 donnent, pour chacun des ajustements, le coefficient de dispersion – qui, rappelons-le, est le rapport de l’estimation de la variance robuste sur l’estimation de la variance naïve du paramètre de régression  $\beta$ .

Les coefficients de dispersion des modèles avec un effet groupe, qu’il soit fixe ou aléatoire, sont quasi égaux et supérieurs à 1. Par contre, le modèle marginal de Cox fournit un coefficient de dispersion inférieur à 1 dans le cas où  $\beta_0 = 0$ , et majoritairement supérieur à 1 dans le cas alternatif ( $\beta_0 = 0,7$ ).

TAB. 5.4 – Coefficient de dispersion et (entre parenthèses) critère d’Akaike pour  $\beta_0 = 0, 7$ .

$\sigma^{2*}$	Regroup.**	Censure	Effet fixe	Marginal	Fragilité
1	1	0 %	1,18 (1 141)	1,07 (53)	1,14 (892)
		30 %	1,14 (836)	1,06 (61)	1,11 (634)
		60 %	1,09 (543)	1,01 (54)	1,06 (337)
		90 %	1,05 (210)	0,98 (20)	1,02 (40)
	2	0 %	1,06 (882)	1,14 (73)	1,05 (860)
		30 %	1,04 (692)	1,33 (67)	1,04 (670)
		60 %	1,01 (424)	1,08 (56)	1,02 (403)
		90 %	1,00 (115)	0,90 (20)	1,00 (94)
4	1	0 %	1,19 (1 935)	0,96 (10)	1,13 (2 300)
		30 %	1,18 (1 458)	1,01 (12)	1,13 (1 618)
		60 %	1,14 (1 034)	1,03 (18)	1,11 (868)
		90 %	1,07 (324)	0,98 (15)	1,05 (127)
	2	0 %	1,07 (1 916)	0,96 (24)	1,05 (1 982)
		30 %	1,07 (1 446)	1,29 (19)	1,05 (1 607)
		60 %	1,05 (944)	1,30 (22)	1,04 (940)
		90 %	1,00 (168)	0,80 (16)	1,00 (214)

\* : variance de la fragilité; \*\* : regroupement.

Dans tous les cas, ce coefficient de dispersion demeure faible – de valeurs maximales 1,12 lorsque  $\beta_0 = 0$ , et 1,19 lorsque  $\beta_0 = 0, 7$ .

- Quelle que soit la valeur de  $\beta_0$ , le passage du premier regroupement au second a pour effet
- de réduire le coefficient de dispersion pour les ajustements avec effet groupe, ainsi que celui de l’ajustement marginal lorsque  $\beta_0 = 0$ ;
  - d’accroître le coefficient de dispersion pour l’ajustement marginal dans le cas où  $\beta_0 = 0, 7$ , excepté lorsque le taux de censure est de 90 %.

Quant au passage de la variance de la fragilité de 1 à 4, il semble n’avoir aucun effet sur la valeur du coefficient de dispersion.

Enfin, l’augmentation du taux de censure entraîne

- une diminution du coefficient de dispersion pour les deux ajustements avec effet groupe ;
- une augmentation du coefficient de dispersion pour l’ajustement marginal et pour un taux de censure passant de 0 % à 30 %, ou de 30 % à 60 % – au-delà de 60 %, le coefficient de dispersion diminue.

Ces résultats surprennent par le fait que la covariable simulée est de type « intra-groupe » ; sa variance, dans le cas d’un ajustement naïf, devrait donc être surestimée. Or seul le modèle marginal semble corriger dans le sens attendu l’estimation de la variance de  $\beta$  – c’est-à-dire que nous obtenons, dans ce cas-là, un coefficient de dispersion inférieur à 1.

## 5.4.3 Puissance du test du maximum de vraisemblance

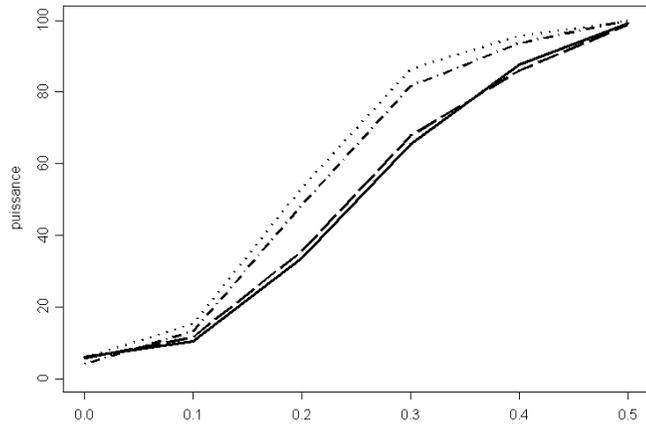
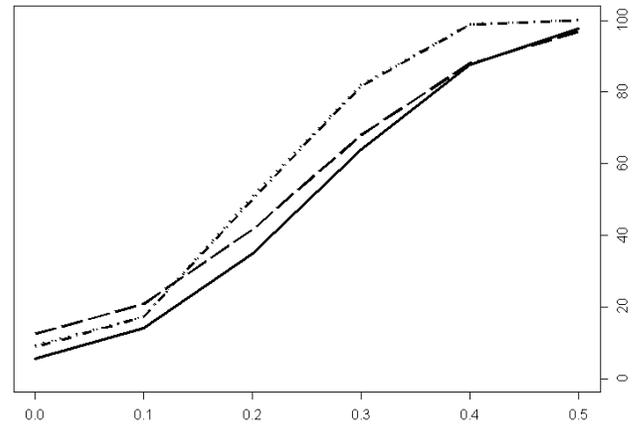
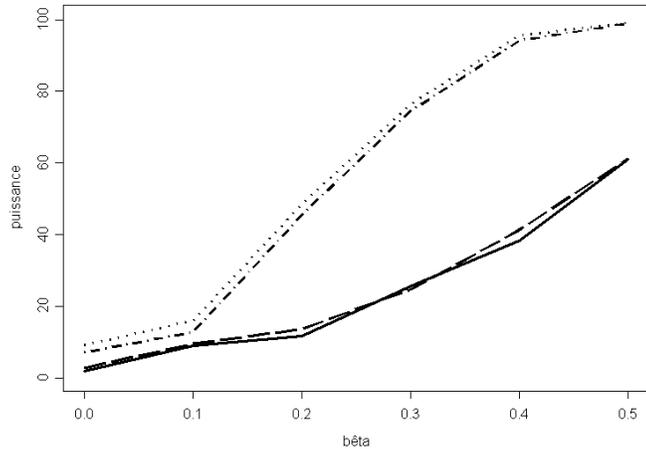
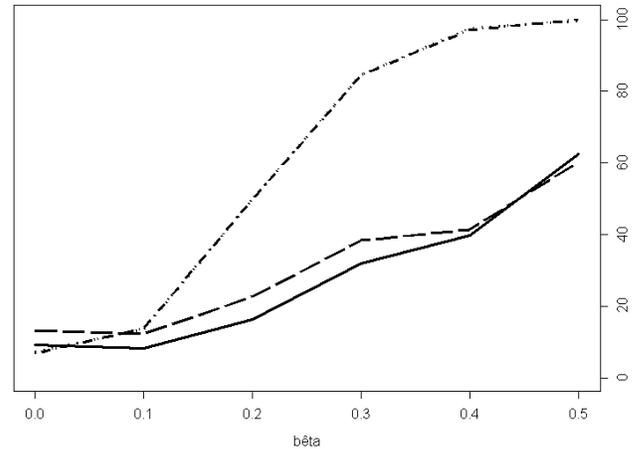
(a) Premier regroupement et  $\sigma^2 = 1$ (b) Second regroupement et  $\sigma^2 = 1$ (c) Premier regroupement et  $\sigma^2 = 4$ (d) Second regroupement et  $\sigma^2 = 4$ 

FIG. 5.2 – Puissance du test du maximum de vraisemblance au seuil de 5 %, pour  $\beta_0 = 0$  et pour un taux de censure de 60 % : ajustement naïf sans « effet groupe » (ligne pleine), ajustement naïf avec « effet groupe » fixe (ligne pointillée), ajustement avec fragilité (ligne hachurée-pointillée) et ajustement marginal (ligne hachurée).

TAB. 5.5 – Puissance du test du maximum de vraisemblance au seuil de 5 %, pour un taux de censure de 60 %.

Regroup.*	$\beta_0^{**}$	$\sigma^{2b}$	$\beta_{(k)}^{bb}$	Sans effet groupe	Avec effet fixe	Marginal	Avec effet aléatoire
1 <sup>er</sup>	0,00	1	0,00	6,0	5,6	5,6	4,0
			0,50	99,2	99,6	98,8	100
	4	0,00	1,9	9,0	2,8	7,1	
			0,50	61,1	99,0	61,1	99,0
	0,70	1	0,70	40,4	10,4	40,4	6,0
			1,20	56,8	100	54,4	99,6
	4	0,70	97,1	8,2	96,6	3,8	
			1,20	45,5	100	41,1	99,5
2 <sup>nd</sup>	0,00	1	0,00	5,6	9,2	12,4	8,8
			0,50	97,6	100	96,8	100
	4	0,00	9,2	7,2	13,2	6,8	
			0,50	62,6	99,6	60,2	100
	0,70	1	0,70	30,8	6,0	34,0	6,0
			1,20	63,6	99,6	54,8	99,6
	4	0,70	88,6	6,0	82,1	5,2	
			1,20	34,4	99,2	20,0	100

\* : regroupement ; \*\* : valeur initiale du paramètre de régression ; <sup>b</sup> : variance de la fragilité ; <sup>bb</sup> : valeur du paramètre de régression pour l'hypothèse alternative.

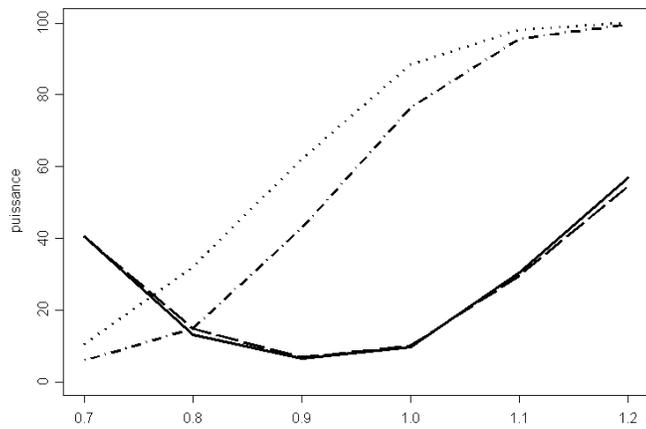
La puissance du test du maximum de vraisemblance est maximale pour les deux ajustements avec effet groupe : les représentations graphiques de la puissance pour ces deux ajustements se chevauchent lorsque  $\beta_0 = 0$  (fig. 5.2). Lorsque  $\beta_0 = 0,7$ , ce chevauchement est préservé lors du second regroupement (fig. 5.3(b) et 5.3(d)), mais non lors du premier (fig. 5.3(a) et 5.3(c)) : c'est alors le modèle de Cox avec effet fixe qui fournit au test du maximum de vraisemblance une puissance maximale.

La puissance observée avec les deux autres ajustements – ajustement sans effet groupe et ajustement marginal – est faible lorsque  $\beta_0 = 0$ , et même aberrante lorsque nous avons conjointement  $\beta_0 = 0,7$  et un taux de censure de 60 % (fig. 5.3(c) et 5.3(d)).

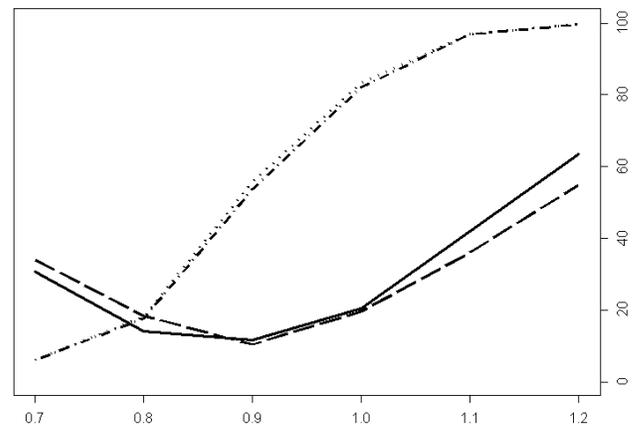
L'augmentation du taux de censure entraîne une diminution de la puissance du test du maximum de vraisemblance, quel que soit l'ajustement adopté.

Le passage de la variance (initiale) de la fragilité de 1 à 4 a lui aussi pour conséquence une diminution de la puissance du test. Le tableau 5.5 fournit une illustration de cette diminution.

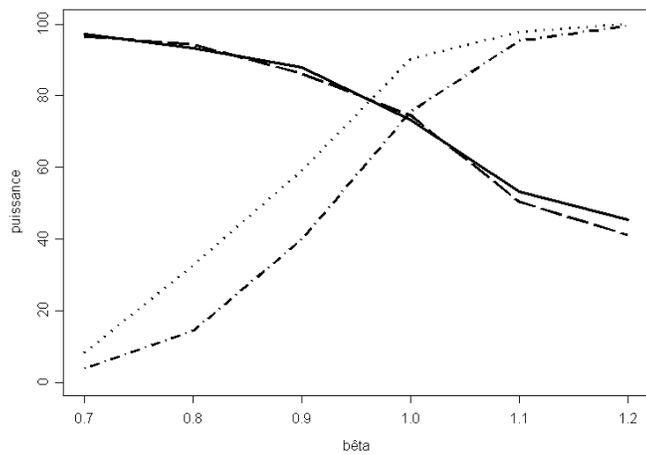
À l'inverse, le passage du premier au second regroupement a pour effet d'accroître la puissance de ce test – ici encore, nous observons ce résultat pour l'ensemble des ajustements.



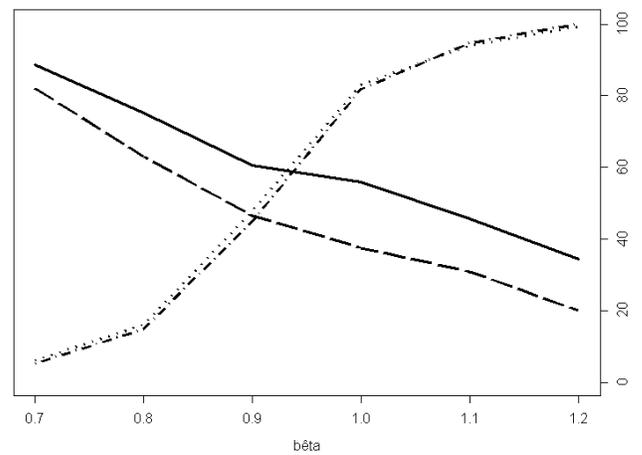
(a) Premier regroupement et  $\sigma^2 = 1$



(b) Second regroupement et  $\sigma^2 = 1$



(c) Premier regroupement et  $\sigma^2 = 4$



(d) Second regroupement et  $\sigma^2 = 4$

FIG. 5.3 – Puissance du test du maximum de vraisemblance au seuil de 5 %, pour  $\beta_0 = 0,7$  et pour un taux de censure de 60 % : ajustement naïf sans « effet groupe » (ligne pleine), ajustement naïf avec « effet groupe » fixe (ligne pointillée), ajustement avec fragilité (ligne hachurée-pointillée) et ajustement marginal (ligne hachurée).

TAB. 5.6 – Étude du taux de recouvrement à 95 % de  $\sigma^2$ .

Regroup.*	$\sigma^{2**}$	Censure	$\beta_0 = 0$		$\beta_0 = 0,7$	
			Gamma <sup>b</sup>	Gaussienne <sup>bb</sup>	Gamma <sup>b</sup>	Gaussienne <sup>bb</sup>
1	1	0 %	99,6	98,1	99,4	99,5
		30 %	97,2	99,8	97,3	99,8
		60 %	87,1	100	87,2	100
		90 %	67,1	100	67,9	100
	4	0 %	51,7	0,0	57,6	1,0
		30 %	62,8	7,5	70,5	59,7
		60 %	73,6	77,5	79,3	85,3
		90 %	46,6	61,5	44,9	99,6
2	1	0 %	84,3	93,7	99,1	100
		30 %	99,3	95,9	99,3	100
		60 %	99,2	95,6	99,1	100
		90 %	95,0	98,6	95,4	100
	4	0 %	96,6	17,3	95,6	64,9
		30 %	98,7	28,4	98,0	80,1
		60 %	97,2	51,6	98,2	96,5
		90 %	81,7	83,4	82,7	94,9

\* : regroupement ; \*\* : variance de la fragilité ; <sup>b</sup> : distribution gamma affectée à la fragilité lors de l'ajustement ; <sup>bb</sup> : distribution gaussienne affectée à la fragilité lors de l'ajustement.

#### 5.4.4 Critère d'Akaike et taux de recouvrement de l'intervalle de confiance à 95 % de la variance de la fragilité

La valeur du critère d'Akaike, pour chacun des ajustements, est présentée dans les tableaux 5.3 et 5.4. D'après la définition de ce critère (cf. page 97), nous obtenons une seule et même valeur pour le modèle sans effet groupe et pour le modèle marginal.

Les résultats d'un tableau à l'autre – c'est-à-dire pour les deux valeurs de  $\beta_0$  choisies lors de la réalisation des simulations – sont extrêmement proches. Nous observons que

- le modèle sans effet groupe (ou le modèle marginal) présente toujours le critère le plus faible ;
- lorsque le taux de censure augmente, la valeur du critère d'Akaike diminue pour tous les ajustements réalisés ;
- le passage du premier au second regroupement a pour effet de
  - réduire la valeur du critère pour l'ajustement avec effet groupe fixe, ainsi que pour l'ajustement avec effet groupe aléatoire en l'absence de censure,
  - augmenter la valeur du critère pour l'ajustement sans effet groupe ou pour l'ajustement marginal, ainsi que pour l'ajustement avec effet groupe aléatoire lorsque le taux de censure est de 30 %, 60 % ou 90 % ;
- pour une valeur (initiale) de la fragilité de 1, c'est l'ajustement avec effet groupe fixe qui fournit le critère le plus important, et donc qui s'impose comme le meilleur (notamment vis-à-vis de l'ajustement avec effet aléatoire),
- pour une valeur (initiale) de la fragilité de 4, et pour un taux de censure de 60 % ou de 90 %, c'est là encore le modèle avec effet fixe qui semble le plus adapté aux données ;
- pour une valeur (initiale) de la fragilité de 4, et pour un taux de censure de 0 % ou de 30 %, c'est l'ajustement avec effet groupe aléatoire qui présente alors le critère maximal.

Concernant le taux de recouvrement de l'intervalle de confiance à 95 % de la variance de la fragilité (tab. 5.6), lorsque la variance (initiale) de la fragilité est fixée à 1 :

- ce taux est excellent pour l'ajustement incluant un terme de fragilité suivant une loi normale, quels que soient la valeur initiale de  $\beta$ , le regroupement et le taux de censure ;
- ce taux est également excellent pour l'ajustement incluant un terme de fragilité suivant une loi gaussienne – là encore quels que soient la valeur initiale de  $\beta$  et le regroupement –, mais il décroît sensiblement à mesure que le taux de censure augmente dans le cas du premier regroupement (passant, par exemple lorsque  $\beta_0 = 0,7$ , de 100 % à 70 %, le taux de censure valant respectivement 0 % ou 90 %).

Lorsque la variance (initiale) de la fragilité est fixée à 4, nous observons les résultats suivants :

- le taux de recouvrement est moyennement satisfaisant – se situant globalement autour de 50 % – pour l'ajustement avec une fragilité suivant une loi gamma, lors du premier regroupement, quels que soient la valeur de  $\beta_0$  et le taux de censure,
- pour ce même ajustement, le taux de recouvrement est excellent lors du second regroupement, quels que soient la valeur de  $\beta_0$  et le taux de censure – cependant qu'il (le taux de recouvrement) accuse une légère diminution lorsque la censure est fixée à 90 %,
- le taux de recouvrement pour l'ajustement avec une fragilité suivant une loi gaussienne est beaucoup moins bon que celui obtenu pour le modèle de fragilité alternatif ; de plus, ce taux est plus élevé lorsque  $\sigma^2$  vaut 4 (plutôt que 1) et lorsque les données sont constituées suivant le second regroupement (plutôt que le premier) – notons l'extrême faiblesse que peut atteindre le taux de recouvrement : lors du premier regroupement par exemple, et en l'absence de censure, ce taux est nul.

Rappelons que la variabilité censée représenter l'« effet élevage » a été introduite, lors des simulations, sous la forme d'une variable aléatoire suivant une loi gamma. Avec ce choix de distribution de probabilité, et du strict point de vue des résultats concernant le taux de recouvrement, nous devons conclure à un meilleur ajustement du modèle de fragilité de Cox où le terme aléatoire suit une loi gamma, plutôt qu'à celui où la fragilité est distribuée suivant une loi gaussienne.

## 5.5 Conclusions et perspectives

« Quelles doivent être les caractéristiques d'un jeu de données pour qu'un ajustement du modèle de Cox puisse en tirer le maximum d'informations ? »

Si la question se pose en ces termes, nous sommes en mesure d'y apporter deux éléments de réponse, concernant deux de ces caractéristiques du jeu de données, visibles *a priori*, à savoir :

- la constitution des données sous forme de groupes : dans notre étude, il apparaît qu'un jeu de données constitué d'une dizaine de groupes comportant chacun une centaine d'individus (second regroupement) « se prêtera mieux » à un ajustement du modèle de Cox qu'un jeu de données constitué d'une centaine de groupes comportant chacun une dizaine d'individus (premier regroupement), c'est-à-dire que le biais du paramètre de régression sera moindre, tandis que la puissance des tests statistiques réalisés sera meilleure ;

TAB. 5.7 – Tableau récapitulatif des résultats pour chaque paramètre mutatis mutandis, lors des ajustements avec le modèle sans effet groupe (SEG), le modèle avec effet groupe fixe (EGF) et le modèle avec effet groupe aléatoire (EGA).

	Regroupement*			Variance**			Censure <sup>b</sup>			$\beta^{\text{bb}}$		
	SEG	EGF	EGA	SEG	EGF	EGA	SEG	EGF	EGA	SEG	EGF	EGA
« Biais »	~	-	-	+	~	~	-	-	-	+	~	~
CD <sup>•</sup>	-	-	-	~	~	~	~	-	-	+	+	+
Puissance	+	+	+	-	-	-	+	+	+	-	+	+
AIC	~	-	~	-	+	+	~	-	-	+	+	~

\* : passage du premier au second regroupement ; \*\* : passage de la variance de la fragilité de 1 à 4 ; <sup>b</sup> : passage du taux de censure de la nullité à 90 % ; <sup>bb</sup> : passage de  $\beta$  de la nullité à la valeur 0,7 ; <sup>•</sup> : coefficient de dispersion.

- le taux de censure : nous concluons de notre étude qu’un taux de censure de 30 %, voire de 60 %, peut améliorer l’estimation de  $\beta$ , ainsi que la puissance des tests statistiques réalisés – cette amélioration est relative au cas d’absence de censure –, tandis qu’un taux de censure de 90 % a, lui, tendance à rendre très imprécis les résultats statistiques d’un ajustement du modèle de Cox (nous pouvions nous attendre à ce résultat, puisqu’un taux de censure aussi élevé – nécessairement consécutif d’une anomalie dans la constitution de l’enquête épidémiologique – traduit un manque d’information certain).

Les autres caractéristiques, invisibles lors de l’élaboration du jeu de données – le degré de corrélation ou la valeur de l’estimation du paramètre de régression –, sont donc des caractéristiques *a posteriori* qui influent sur l’ajustement du modèle de Cox de la manière suivante :

- une variabilité « inter-groupe » importante (se situant, lors de nos simulations, dans la variance initiale de la fragilité, que nous avons fixée à 1 ou 4) a un effet néfaste d’une part sur la puissance des tests statistiques (qui diminue quel que soit l’ajustement réalisé), d’autre part sur l’estimation du paramètre de régression (dont le biais croît) lors de l’ajustement naïf sans effet groupe ;
- l’augmentation de la valeur « initiale<sup>1</sup> » du paramètre de régression (c’est-à-dire la valeur  $\beta_0$  lors de nos simulations) entraîne un accroissement de la puissance des tests, à condition qu’un effet groupe soit pris en compte (sinon c’est le résultat contraire qui est obtenu), tandis qu’à cette même condition, elle n’a aucun effet sur le biais de l’estimation finale de  $\beta$  (en dehors de cette condition, le biais se trouve au contraire accru).

Ces conclusions concernent les quatre types d’ajustement du modèle de Cox étudiés – ajustement naïf sans et avec effet groupe fixe, ajustement marginal et ajustement avec effet groupe aléatoire.

Le tableau 5.7 récapitule l’influence de chacun des paramètres, *toutes choses étant égales par ailleurs*, sur les différents ajustements du modèle de Cox étudiés au travers de notre travail.

1. Ce terme étant en réalité impropre, puisque l’étude d’un jeu de données de survie ne laisse présager d’aucune valeur particulière concernant le paramètre de régression et, par suite, concernant son estimation...



## Chapitre 6

# ÉTUDE ÉPIDÉMIOLOGIQUE DES GASTRO-ENTÉRITES NÉONATALES CHEZ LE VEAU

### Contenu

---

<b>6.1</b>	<b>Introduction</b>	<b>110</b>
<b>6.2</b>	<b>Protocole</b>	<b>111</b>
6.2.1	Population cible	111
6.2.2	Constitution des données	111
6.2.3	Informations relatives à l'entretien du bâtiment	112
6.2.4	Informations relatives à l'alimentation	116
6.2.5	Informations relatives aux mesures de prophylaxie	116
6.2.6	Informations relatives au vêlage et aux soins apportés au veau	116
6.2.7	Informations relatives à la létalité et à la mortalité	117
<b>6.3</b>	<b>Statistiques descriptives</b>	<b>117</b>
<b>6.4</b>	<b>Statistiques analytiques</b>	<b>119</b>
6.4.1	Sélection des variables	119
6.4.2	Comparaisons des courbes de survie	120
6.4.3	Estimations des paramètres	121
6.4.4	Concernant la fragilité	123
6.4.5	Justesse des ajustements	126
6.4.6	Conclusions	128
6.4.7	Discussion	131

---

## 6.1 Introduction

Les gastro-entérites néonatales (GENN) sont l'une des pathologies néonatales prépondérantes du veau. Les taux d'incidence<sup>1</sup> de diarrhée varient, selon les études, entre 15 % et 20 %. Le taux de létalité varie, quant à lui, entre 1,5 % et 8 %. Le risque de diarrhée est maximal durant la première semaine de vie du veau, puis il décroît à mesure que l'âge de l'animal augmente.

Ces morbidité et mortalité précoces constituent une source importante de perte financière pour l'éleveur – due d'une part à la perte effective des animaux, et d'autre part à la diminution du potentiel génétique d'amélioration de l'élevage. Le coût moyen (par veau et par an associé à la prévention et à la mortalité des gastro-entérites est estimé à 30,78 €. Une étude américaine (Franck et Kaneene, 1993) avance le chiffre de 230 millions d'euros par an, concernant le coût total de cette pathologie supporté par l'industrie laitière.

En France, dans la région Midi-Pyrénées, seule une proportion de 20 % (respectivement 80 %) d'élevages (resp. de veaux) est estimée indemne de diarrhée néonatale. Dans 50 % des cas, l'issue de cette diarrhée est fatale. Le coût de cette pathologie – dans cette même région – est évalué à 64,4 € par veau et par an. Plus généralement, 21 % des interventions vétérinaires chirurgicales sont consacrées au traitement des entérites.

L'identification des causes de diarrhée est difficile : d'une part, la liste des étiologies possibles est longue, et d'autre part les méthodes de diagnostic sont limitées. Des agents infectieux peuvent être à l'origine de ces affections : citons rotavirus, coronavirus, *Escherichia coli* et *Cryptosporidium*. Ces agents, dits de type « zoonotique » (c'est-à-dire pouvant être à l'origine de zoonoses), ont été associés à l'apparition de pathologies humaines consécutives à une ingestion d'aliments d'origine animale.

La lutte contre la diarrhée du veau se résume essentiellement à l'emploi d'antibiotiques, constituant ainsi un risque sanitaire potentiel pour l'homme, en termes de résidus médicamenteux – au travers de la possible nocivité de ces résidus d'une part, et du développement de résistances aux bactéries d'autre part.

Les diarrhées du veau illustrent parfaitement l'exemple d'une maladie complexe et multifactorielle. Elles résultent d'interactions entre des facteurs propres à l'animal (technique d'élevage, alimentation, logement, hygiène, pathologies concomitantes telles que malformation congénitale ou dyspnée<sup>2</sup>, etc.) et des agents pathogènes. Toutes ces composantes doivent être étudiées simultanément, afin de pouvoir déterminer l'importance de chacun de ces facteurs et, par la suite, pouvoir proposer une gestion globale des risques.

L'objectif de l'étude prospective menée en région Midi-Pyrénées (Bendali *et al.*, 1999a,b) était de décrire cette pathologie de la diarrhée, ainsi que d'évaluer l'impact que peuvent avoir sur son incidence des facteurs aussi divers que l'âge de l'animal, la saison de sa naissance ou les caractéristiques d'élevage.

---

1. Pour les termes relatifs au domaine de l'épidémiologie, consulter le glossaire présent en page 197.  
2. Difficulté de la respiration.

## 6.2 Protocole

### 6.2.1 Population cible

Une étude de cohorte longitudinale prospective a été menée dans 92 élevages bovins allaitants de 8 départements du Sud-Ouest de la France. L'inclusion dans l'étude des élevages reposait sur une technique d'échantillonnage aléatoire stratifié en grappes – la stratification, basée sur le département, était ajustée sur le nombre de fermes dans chaque département (suivant une méthode d'échantillonnage proportionnel). Un algorithme d'échantillonnage aléatoire a permis d'identifier les élevages que l'étude devait retenir. Cet algorithme intégrait le critère en vertu duquel, si une ferme était sélectionnée, alors tous les animaux la composant l'étaient également.

La population d'origine consistait en l'intégralité des fermes de la région Midi-Pyrénées. Le critère d'inclusion était le suivant : un élevage était retenu si et seulement s'il contenait plus de 20 vaches adultes, au moins trois quarts de veaux au moment du sondage, et s'il présentait les cinq races les plus importantes de la région (Charolaise, Limousine, Gascogne, Aubrac et Blonde d'Aquitaine). Au total, 7 000 élevages – représentant 30 % de l'effectif (en élevages) de la région – étaient candidats à l'inclusion dans l'étude.

Sur les 95 élevages bovins de Midi-Pyrénées sélectionnés de la façon aléatoire décrite précédemment, 92 furent finalement retenus pour l'étude ; ils furent suivis sur une période allant de décembre 1995 à avril 1996.

Dans ces 92 élevages étudiés, 3 047 veaux – nés pendant la période d'observation – ont été retenus, soit 33 veaux en moyenne par élevage.

### 6.2.2 Constitution des données

Des questionnaires concernant la santé, la démographie et les conditions d'élevage pouvant influencer sur l'apparition de cas de diarrhées ont été constitués. Chaque élevage a reçu 9 visites d'au moins un vétérinaire et un technicien vétérinaire, chargés de la collecte des données. Les visites étaient espacées de 8 à 12 jours – suivant la période de vêlage. Ces visites régulières ont permis de suivre minutieusement l'évolution des « données continues », ainsi que d'enregistrer les variations dans le temps de certains paramètres (température, humidité, densité de résidence des veaux).

Lors de la première visite, datée de décembre 1995, les questionnaires furent enrichis d'un entretien avec l'éleveur, afin de déterminer les pratiques d'élevage (notamment le type de bâtiment), les conditions d'alimentation et le degré de propreté.

Par ailleurs, différents examens ont été effectués, afin de déterminer les agents pathogènes susceptibles d'être à l'origine de cas de diarrhée. Des prélèvements fécaux ont ainsi été réalisés, tant au sein d'élevages touchés par la pathologie étudiée, qu'au sein d'élevages indemnes. Les laboratoires vétérinaires des 8 départements possèdent une méthode d'analyse standard (homogène d'un laboratoire à l'autre) permettant d'identifier rotavirus, coronavirus, *Escherichia coli* et *Cryptosporidium*. Cependant, les résultats de ces examens n'ont pas été retenus dans notre analyse, principalement en raison d'un nombre important de données manquantes.

L'ensemble des informations obtenues à partir des questionnaires est présenté dans les tableaux 6.1a, 6.1b et 6.1c. Pour chaque variable, nous indiquons :

- les modalités prises par cette variable ;
- le niveau auquel elle intervient – niveau individuel (celui de l’animal) ou niveau collectif (celui de l’élevage) – ;
- la répartition (effectifs et pourcentages) de ses différentes modalités, relativement au niveau auquel cette variable s’applique ;
- son type, c’est-à-dire la catégorie à laquelle elle se rattache (prophylaxie, entretien du bâtiment, etc.) ;
- le nombre de cas de diarrhée pour chacune de ses modalités.

### 6.2.3 Informations relatives à l’entretien du bâtiment

Nous avons observé les mêmes proportions entre stabulations libres et stabulations entravées (respectivement 51 % et 49 %). Dans 53 % des cas, l’aération et le renouvellement d’air ont été jugés satisfaisants.

Peu d’élevages ont été déclarés possesseur et/ou utilisateur d’une infirmerie pour isoler les animaux malades – soient respectivement 27 % et 34 % pour les adultes et pour les jeunes animaux.

Le nettoyage systématique du bâtiment, avant ou après la saison de vêlage, ne semble pas être une pratique très courante. La majorité des éleveurs – entre 74 % et 78 % d’entre eux – n’ont pas adopté de nettoyage systématique et régulier. En revanche, le paillage et la désinfection ont été davantage employés : de ce fait, les animaux ont généralement présenté une bonne note de propreté, notamment les veaux dont 85 % de l’effectif étaient relativement propres (les adultes, quant à eux, n’étaient propres qu’à 52 %).

Le confort thermique des animaux, notamment des veaux, a été estimé par un indice synthétique ; cet indice a montré que près de la moitié des veaux (44 % des élevages) ont été logés dans des conditions de confort thermique satisfaisantes.

La concentration des veaux a été relativement bien respectée dans la plupart des cas (64 % des élevages). Les vêlages n’étaient pas regroupés au sein d’une période déterminée : en effet, environ 55,3 % d’entre eux ont enregistré des naissances étalées sur toute l’année, probablement en raison d’un grand nombre d’animaux dans l’exploitation, ou du type de production (tel que le veau d’Aveyron, dont la naissance peut avoir lieu à tout moment de l’année).

Le box de vêlage n’a pas été une pratique générale : il n’a été disponible que dans 33 % des cheptels. Son entretien (nettoyage et désinfection) a été généralement bien suivi (83 % des cas).

TAB. 6.1a – Variables explicatives étudiées dans le cadre des GENN: description, modalités, effectifs et niveau hiérarchique de ces variables.

Description de la variable	Modalités	Effectifs*	%**	Évé.♦	Type <sup>b</sup>	Niveau <sup>bb</sup>
Nombre de veaux dans la portée	un seul	2 944	96,6	429	N	I
	plus d'un	103	3,4	11		
État du nouveau-né (prématurité)	à terme	3 027	99,3	437	N	I
	prématuré	20	0,7	3		
État du nouveau-né (malformation(s) congénitale(s))	normal	3 030	99,4	437	N	I
	malformation	17	0,6	3		
Respiration du veau à la naissance	normale	2 832	92,9	401	N	I
	dyspnée	220	7,1	39		
Utilisation d'orexygènes et stimulants à titre curatif	non	3 012	98,9	437	N	I
	oui	35	1,1	3		
Utilisation d'orexygènes et stimulants à titre préventif	non	2 291	98,2	433	N	I
	oui	56	1,8	7		
Désinfection de l'ombilic	oui	1 860	61,0	271	N	I
	non	1 187	39,0	169		
Produit de désinfection de l'ombilic	aucun	1 207	40,0	169	N	I
	iodine	880	29,0	146		
	autre	960	31,0	125		
Première prise de colostrum assistée par l'éleveur	oui	2 250	74,0	335	N	I
	non	797	26,0	105		
Délai de la première prise de colostrum	< 6h.	2 684	88,0	392	N	I
	> 6h.	363	12,0	48		
Fréquence d'alimentation de veau pendant le premier jour (nombre de tétées)	suffisante ( $\geq 3$ )	2 549	84,0	353	N	I
	insuffisante ( $< 3$ )	498	16,0	87		
Mois de naissance des veaux	déc. 95	434	14,3	76	N	I
	janv. 96	754	24,7	67		
	fév. 96	963	31,6	131		
	mars 96	896	29,4	166		
Condition de vêlage	non assisté	1 785	58,6	226	N	I
	aide facile	795	26,1	138		
	aide difficile	467	15,3	76		
Fréquence d'alimentation du veau pendant les 3 premiers jours (nombre de tétées)	suffisante ( $\geq 3$ )	2 553	84,0	361	N	I
	insuffisante ( $< 3$ )	494	16,0	79		
Vaccination de la mère contre coronavirus	oui	809	27,0	110	P	I
	non	2 238	73,0	330		
Vaccination de la mère contre rotavirus	oui	807	26,0	110	P	I
	non	2 240	74,0	330		
Vaccination de la mère contre <i>E. coli</i>	oui	777	26,0	114	P	I
	non	2 270	74,0	326		
Vaccination de la mère contre d'autres agents que GENN	oui	751	25,0	81	P	I
	non	2 296	75,0	359		
Parité de la mère	primipare	620	20,0	93	C	I
	multipare	2 427	80,0	347		
Sexe du veau	mâle	1 490	50,0	235	D	I
	femelle	1 463	50,0	202		

\* : répartition des modalités de la variable suivant que les pratiques sont recueillies à l'échelle de l'élevage ( $n = 92$ ) ou à celle du veau ( $n = 3\ 047$ ); \*\* : pourcentage de distribution de la variable par rapport à la population générale; ♦ : nombre de survenues de diarrhée; <sup>b</sup> : C = conduite du troupeau, D = divers, N = vêlage, P = prophylaxie; <sup>bb</sup> : I = variable au niveau individuel (animal).

## 6.2. PROTOCOLE

TAB. 6.1b – Variables explicatives étudiées dans le cadre des GENN: description, modalités, effectifs et niveau hiérarchique de ces variables.

Description de la variable	Modalités	Effectifs*	%**	Évé. •	Type <sup>b</sup>	Niveau <sup>bb</sup>
Distribution de concentré	oui	31	33,7	123	A	G
	non	61	66,3	323		
Aliment à base d'ensilage de maïs	oui	41	44,5	189	A	G
	non	51	55,5	257		
Aliment à base d'ensilage d'herbe	oui	57	61,9	270	A	G
	non	35	38,1	176		
Alimentation des animaux rationnée	non	22	23,9	107	A	G
	oui	70	76,1	339		
Alimentation adaptée entre gestante et allaitante	oui	42	45,6	178	A	G
	non	50	54,4	268		
Alimentation adaptée entre primipares et multipares	oui	29	31,5	153	A	G
	non	63	68,5	293		
Pratique du flushing	oui	16	17,4	76	A	G
	non	76	82,6	370		
Type de stabulation	libre	47	51,1	193	B	G
	entravée	45	48,9	253		
Aération correcte du bâtiment	oui	49	53,2	206	B	G
	non	43	46,8	240		
Odeur d'ammoniac perceptible dans le bâtiment	non	86	93,4	394	B	G
	oui	6	6,6	52		
Humidité perceptible dans le bâtiment	non	86	93,4	396	B	G
	oui	6	6,6	50		
Existence d'une infirmerie pour adultes	oui	25	27,1	80	B	G
	non	67	72,9	366		
Existence d'une infirmerie pour veaux	oui	31	33,7	115	B	G
	non	61	66,3	331		
Densité des veaux	correcte	61	66,3	258	C	G
	forte	31	33,7	188		
Animaux groupés en lots	oui	74	80,4	375	C	G
	non	18	19,6	71		
Parcage des veaux	box collectif	10	10,8	54	C	G
	box individuel	67	72,8	319		
	pas de box	15	16,4	73		
Vêlages groupés en hiver	oui	42	45,6	255	C	G
	non	50	54,4	191		
Incidence de la diarrhée la saison précédente	< 5 %	28	30,4	78	D	G
	≥ 5 %	64	69,6	368		

\* : répartition des modalités de la variable suivant que les pratiques sont recueillies à l'échelle de l'élevage ( $n = 92$ ) ou à celle du veau ( $n = 3\ 047$ );

\*\* : pourcentage de distribution de la variable par rapport à la population générale; • : nombre de survenues de diarrhée; <sup>b</sup> : A = alimentaire, B = bâtiment, C = conduite du troupeau, D = divers; <sup>bb</sup> : G = variable au niveau groupe (élevage).

TAB. 6.1c – Variables explicatives étudiées dans le cadre des GENN: description, modalités, effectifs et niveau hiérarchique de ces variables.

Description de la variable	Modalités	Effectifs*	%**	Évé. •	Type <sup>b</sup>	Niveau <sup>bb</sup>
Paillage suffisant et régulier pour les adultes	oui	39	42,4	148	E	G
	non	53	57,6	298		
Paillage suffisant et régulier pour les veaux	oui	31	33,7	191	E	G
	non	61	66,3	255		
Nettoyage des bâtiments des veaux avant la saison des vêlages	oui	24	26,1	162	E	G
	non	68	73,9	284		
Nettoyage des bâtiments des veaux après la saison des vêlages	oui	20	21,7	71	E	G
	non	72	78,3	375		
Nettoyage suite à chaque épisode diarrhéique	oui	4	4,3	30	E	G
	non	88	95,7	416		
Désinfection régulière et correcte chez les veaux	oui	39	42,4	189	E	G
	non	53	57,6	257		
Propreté des veaux	satisfaisante	78	84,7	374	E	G
	faible	14	15,3	72		
Propreté des vaches	satisfaisante	52	56,5	212	E	G
	faible	40	43,5	234		
Désinfection du lieu de vêlage	oui	15	16,3	95	E	G
	non	77	83,7	351		
Nettoyage du local de vêlage	fréquent	76	82,6	390	E	G
	non fréquent	16	17,4	56		
Utilisation d'un local de vêlage	oui	31	33,7	141	N	G
	non	61	66,3	305		
Vaccination de la majorité des vaches	oui	22	23,9	113	N	G
	non	70	76,1	333		
Vaccination de la majorité des veaux contre au moins un des agents pathogènes	oui	21	22,8	113	P	G
	non	71	77,2	333		
Supplémentation en vitamines et minéraux des vaches	oui	55	59,8	264	P	G
	non	37	40,2	182		
Supplémentation en vitamines et minéraux des veaux	oui	40	43,4	255	P	G
	non	52	56,6	191		

\* : répartition des modalités de la variable suivant que les pratiques sont recueillies à l'échelle de l'élevage ( $n = 92$ ) ou à celle du veau ( $n = 3\ 047$ ); \*\* : pourcentage de distribution de la variable par rapport à la population générale; • : nombre de survenues de diarrhée; <sup>b</sup> : E = entretien, N = vêlage, P = prophylaxie; <sup>bb</sup> : G = variable au niveau groupe (élevage).

### 6.2.4 Informations relatives à l'alimentation

Les critères alimentaires ont été recueillis selon la nature du rationnement.

Deux types d'ensilage ont généralement été distribués aux animaux : l'ensilage de maïs dans 44 % des exploitations, et l'ensilage d'herbe dans 63 %. Un tiers des éleveurs ont fourni aux animaux une alimentation à base de concentré.

Bien que deux aliments de même nature ne soient pas obligatoirement de même qualité (du fait, par exemple, d'un mode de conservation différent), les investigations relatives à l'alimentation n'ont concerné que la nature du produit et non sa qualité. C'est pour cette raison qu'aucune analyse supplémentaire d'échantillons d'aliment n'a été effectuée.

En outre, une alimentation à volonté n'a eu lieu que dans 23 élevages – la majorité des éleveurs rationnant l'aliment. La ration a été rarement ajustée suivant le stade de la gestation, ou suivant la parité (51 % et 64 % des élevages n'ont présenté aucun ajustement, dans les deux cas de figure respectifs). Enfin, seuls 16 élevages sur les 92 retenus ont pratiqué le flushing<sup>1</sup>.

### 6.2.5 Informations relatives aux mesures de prophylaxie

En plus de l'entretien et des soins réguliers (nettoyage, désinfection...), les pratiques vaccinales et les supplémentations en vitamines et minéraux ont été étudiées. De façon générale, la vaccination n'a été réalisée que dans peu d'élevages (moins du quart des cheptels). Par ailleurs, au-delà de la pratique générale de la vaccination, nous avons disposé, pour chaque veau, des renseignements concernant le statut vaccinal de sa mère (nature et date des vaccinations). À partir de ces données, nous avons constaté qu'environ 25 % des animaux avaient été réellement vaccinés contre les principaux agents pathogènes.

L'introduction d'animaux étrangers (achat, remplacement) a très rarement été observé : un unique cas de remplacement et cinq cas d'achat ont ainsi été enregistrés.

Contrairement à la vaccination, la supplémentation en minéraux et vitamines a été davantage observée. Les adultes ont reçu un complément en minéraux ou vitamines dans 60 % des cas, et les veaux dans 43 % des cas.

### 6.2.6 Informations relatives au vêlage et aux soins apportés au veau

Les résultats ont montré que plus de la moitié des veaux sont nés sans aucune aide (57 %) ; une simple intervention a suffi pour le quart des naissances (26 %), alors que 15 % des vêlages ont été dystociques, c'est-à-dire laborieux. 45 veaux ont nécessité à leur naissance une césarienne, ce qui représente 1,5 % des naissances. La majorité des mères étaient multipares (80 %, contre 20 % de primipares).

La quasi-totalité des naissances concernaient des veaux uniques (96,5 %). Parmi les autres, nous avons enregistré 53 naissances gémellaires (3,4 %) et un triplet. La proportion des malformations congénitales ou des prématurités n'a pas dépassé 0,6 %.

Dès sa naissance, le veau requiert la plus grande attention des éleveurs. La désinfection du nombril a été réalisée dans plus de la moitié des cas (61 %). Par ailleurs, dans 88 % des cas, la première prise de colostrum a été observée dans les 6 premières heures de vie. Enfin, 83 % des veaux ont été suffisamment alimentés pendant les premiers jours de leur vie.

---

1. Supplémentation alimentaire des vaches pendant la période de reproduction

TAB. 6.2 – Tableau récapitulatif des types d’observations rencontrés lors de l’étude.

Type de l’observation	Effectif	Proportion
Diarrhée	440	14,4 %
Censure		
perdition de vue	127	4,2 %
exclusion-vivant	2 480	81,4 %

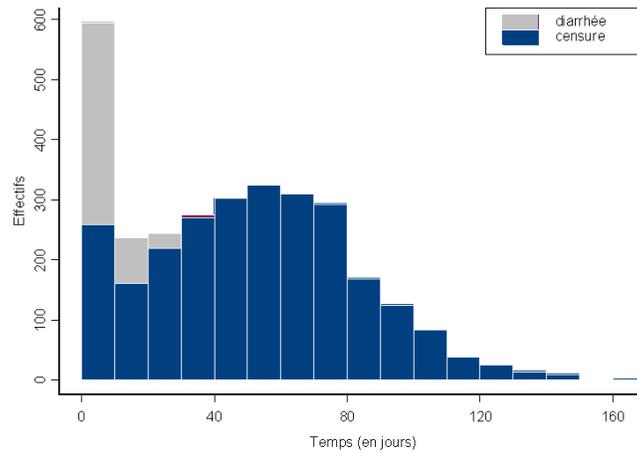


FIG. 6.1 – Répartition conjointe des cas de censure et des cas de survenue de diarrhée.

### 6.2.7 Informations relatives à la létalité et à la mortalité

Nous nous sommes intéressés à la survenue de la diarrhée dans les 31 premiers jours de vie. Nous avons enregistré au total 440 cas de diarrhées pendant ce premier mois de vie, ce qui représente une proportion de 14,4 % des veaux. Par ailleurs, 113 veaux sont morts en période néonatale (soit 3,7 % de l’effectif total).

Enfin, 2 607 veaux ont été *censurés*. Parmi ceux-ci, 113 ont été perdus de vue pour cause de mort, et 14 ont été perdus de vue pour une autre raison (par exemple pour cause de vente). Les 2 480 animaux restants ont été exclus-vivants, soit parce qu’ils n’ont pas présenté de diarrhée lors de leur premier mois de vie, soit parce qu’ils sont nés durant les 30 derniers jours de la période de suivi. Le tableau 6.2 et la figure 6.1 récapitulent les informations concernant ces données.

## 6.3 Statistiques descriptives

Le tableau 6.3 recense les temps de survenue des cas de diarrhée : ces temps, par définition même du phénomène étudié, sont tous inférieurs ou égaux à 31. Pour chacun de ces temps, nous donnons la valeur de la fonction de survie, ainsi qu’un intervalle de confiance à 95 % pour cette valeur. Les ex-æquo, qui apparaissent en grand nombre ici, sont traités au travers de

TAB. 6.3 – Temps de survenue de diarrhées et valeur de la fonction de survie pour chacun de ces temps.

Temps (en jours)	Nombre d'animaux à risque	Nombre de cas de diarrhée	Estimation de la survie $\hat{S}(t)$	Écart-type $\sigma(\hat{S}(t))$	Intervalle de confiance à 95 % de $\hat{S}(t)$
1	3 038	7	0,998	0,0008	0,996 - 0,999
2	2 974	36	0,986	0,0021	0,981 - 0,990
3	2 908	41	0,972	0,0030	0,966 - 0,978
4	2 838	29	0,962	0,0035	0,955 - 0,969
5	2 784	33	0,950	0,0040	0,943 - 0,958
6	2 735	45	0,935	0,0045	0,926 - 0,944
7	2 671	42	0,920	0,0050	0,910 - 0,930
8	2 609	51	0,902	0,0055	0,891 - 0,913
9	2 532	33	0,890	0,0058	0,879 - 0,902
10	2 486	19	0,884	0,0060	0,872 - 0,895
11	2 452	10	0,880	0,0060	0,868 - 0,892
12	2 432	10	0,876	0,0061	0,864 - 0,888
13	2 406	12	0,872	0,0062	0,860 - 0,884
14	2 387	8	0,869	0,0063	0,857 - 0,881
15	2 365	16	0,863	0,0064	0,851 - 0,876
16	2 331	8	0,860	0,0065	0,848 - 0,873
17	2 309	6	0,858	0,0065	0,845 - 0,871
18	2 285	5	0,856	0,0066	0,843 - 0,869
19	2 260	1	0,856	0,0066	0,843 - 0,869
21	2 215	2	0,855	0,0066	0,842 - 0,868
22	2 190	5	0,853	0,0066	0,840 - 0,866
23	2 167	1	0,853	0,0066	0,840 - 0,866
24	2 138	1	0,852	0,0066	0,839 - 0,865
25	2 111	4	0,851	0,0067	0,837 - 0,864
26	2 088	3	0,849	0,0067	0,836 - 0,863
27	2 061	2	0,848	0,0067	0,835 - 0,862
28	2 038	3	0,847	0,0068	0,834 - 0,861
29	2 013	1	0,847	0,0068	0,834 - 0,860
30	1 993	2	0,846	0,0068	0,833 - 0,859
31	1 971	4	0,844	0,0068	0,831 - 0,858

l'approximation d'Efron (cf. p. 43).

L'histogramme 6.1 représente la distribution des censures et des cas de diarrhées observés.

Nous remarquons, à partir de ce tableau et de cette figure, que les cas de diarrhée surviennent le plus fréquemment durant les 10 premiers jours de vie du veau.

La figure 6.2 représente quant à elle l'estimation de Nelson-Aalen de la courbe de survie générale, c'est-à-dire lorsque tous les animaux sont considérés comme appartenant à un même troupeau, et qu'aucune distinction n'est faite quant aux valeurs prises par les différentes covariables. Plus exactement, cette figure ne représente la fonction de survie que pour les 31 premiers jours – la courbe présentant, à partir du 31<sup>e</sup> jour, un plateau correspondant à une stationnarité de la valeur (égale à 0,844) de la fonction de survie.

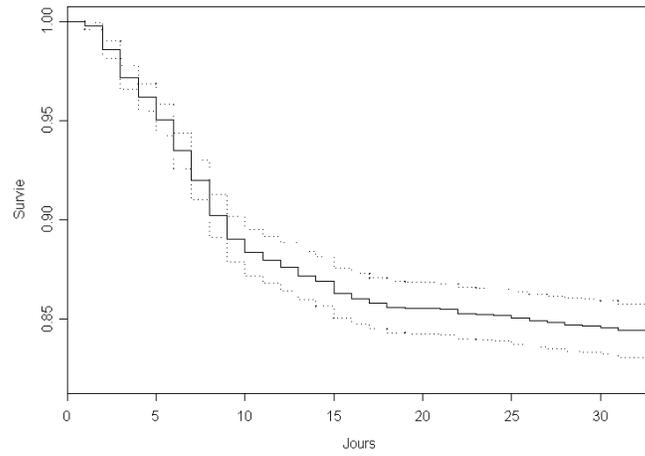


FIG. 6.2 – Estimation de Nelson-Aalen de la courbe de survie, avec un intervalle de confiance à 95 %.

## 6.4 Statistiques analytiques

### 6.4.1 Sélection des variables

Après vérification, transformation et combinaison des variables, nous avons adopté trois méthodes descendantes de sélection de ces variables, la première ignorant tout « effet élevage », les deux autres prenant en compte cet effet (l'une au travers d'un modèle marginal de Cox, l'autre au travers d'un modèle de fragilité gamma de Cox). Ces trois méthodes, similaires quant à leur mode de fonctionnement, s'entreprennent de la façon suivante :

- 1° toutes les variables explicatives sont incluses dans le modèle de Cox ;
- 2° au vu des résultats concernant la significativité des effets de ces variables, nous excluons celles dont l'effet est jugé non significatif ;
- 3° nous réitérons les deux premières étapes, jusqu'à ne retenir que les variables dont l'effet est jugé significatif.

À l'issue de la deuxième étape, et dans le cas de la suppression d'une (ou plusieurs) variable(s), nous pratiquons le test du maximum de vraisemblance afin de tester la meilleure qualité de l'ajustement retenu.

Au terme de cette méthode de sélection, nous avons retenu 15 variables : elles sont présentées dans le tableau 6.4a.

Afin de comparer les estimations des paramètres et les conclusions tirées des différents ajustements, nous avons appliqué à ces 15 variables les modèles de Cox usuels – soient le modèle naïf avec et sans « effet élevage » (fixe), le modèle marginal et le modèle de fragilité (gamma).

Lors de l'introduction des interactions entre les 15 variables finalement retenues et la variable « mois de naissance », nous avons adopté la méthode itérative de sélection des variables suivante :

- 1° introduction des 45<sup>1</sup> termes d'interaction ;
- 2° analyse par le modèle naïf de Cox ;
- 3° au vu des résultats de l'analyse :
  - si une variable et ses interactions avec le mois de naissance ont des effets non significatifs, nous retirons dans un premier temps les interactions,
  - si une variable a un effet non significatif, mais qu'en revanche l'effet (d'au moins une) de ses interactions est jugé significatif – cas de figure non rencontré –, alors nous conservons la variable et ses interactions,
  - si une variable a un effet significatif, mais qu'en revanche les effets de toutes ses interactions sont jugés non significatifs, alors nous retirons les termes d'interaction,
  - si une variable et les termes d'interaction qui lui sont adjoints ont des effets jugés significatifs, nous les conservons ;
- 4° répétition de l'étape 2 ;

---

1. À chaque variable – hormis la variable « mois de naissance » – sont rattachés 3 termes d'interactions, à raison d'un terme par mois de naissance ; la variable « condition de vélage » étant à trois classes, elle produit un terme d'interaction supplémentaire, et le total est ainsi de  $13 \times 3 + 2 \times 3 = 45$  termes d'interactions.

- 5° répétition de l'étape 3, enrichie de la considération suivante : une variable dont nous avons supprimé les interactions à l'issue de l'étape 3, et dont l'effet est jugé non significatif au sortir de l'étape 4, est retirée du modèle ;
- 6° répétition des étapes 2, 3, 4 et 5 jusqu'à la conservation des seules variables et termes d'interactions ayant un effet jugé significatif.

Notons que, là encore, sitôt que nous retirons une (ou plusieurs) variable(s) du modèle, nous pratiquons le test du maximum de vraisemblance afin de tester la meilleure qualité de l'ajustement obtenu.

Enfin, parallèlement à l'introduction des termes interactifs, nous avons créé une variable relative au rang de naissance du veau au sein de son élevage : cet indice (à trois classes) a été testé – aussi bien seul qu'en compagnie des autres variables –, et son effet n'a pas été jugé statistiquement significatif.

De même, une étude a porté sur l'importance du rang des cas de diarrhée à l'intérieur de chaque troupeau : nous souhaitons déterminer si l'apparition précoce d'un cas de diarrhée, à l'intérieur d'un troupeau, exposait davantage ce troupeau à l'apparition de nouveaux cas de diarrhée. La réponse apportée à cette interrogation est négative.

#### 6.4.2 Comparaisons des courbes de survie

Le test du log-rank de comparaison des courbes de survie a été effectué séparément pour chacune des 15 variables finalement retenues. Rappelons que ce test consiste à comparer le nombre de cas de diarrhée observés à celui que nous devrions théoriquement obtenir (cf. p. 27).

Une différence significative est ressortie de ces comparaisons pour 11 des 15 variables : les 4 variables pour lesquelles ce test n'a pas conclu à une différence significative des courbes de survie, sont la condition de vêlage, la vaccination de la mère contre *E. coli*, le nettoyage du bâtiment après chaque épisode diarrhéique, et enfin la supplémentation en vitamines et minéraux des veaux à leur naissance.

#### 6.4.3 Estimations des paramètres

Les risques relatifs obtenus avec le modèle naïf sont similaires à ceux obtenus lors de l'étude originelle<sup>1</sup> (Bendali *et al.*, 1999b) : nous retrouvons, parmi les variables présentant les risques relatifs les plus élevés :

- la condition du vêlage, dont le risque relatif (RR) vaut 1,50 pour la modalité « aide facile » (contre la valeur 1,46 dans l'étude originelle) ;
- la vaccination de la mère contre d'autres agents que GENN (RR = 1,81 contre 2,01 à l'origine) ;
- la présence d'une odeur d'ammoniac (RR = 1,62 contre 1,46 à l'origine) ;
- le nettoyage des bâtiments des veaux après la saison des vêlages (RR = 1,59 contre 1,92 à l'origine).

De même, nous retrouvons, parmi les variables présentant les risques relatifs les plus faibles :

1. Toutes proportions gardées... puisque les fichiers de données des deux études diffèrent quelque peu.

TAB. 6.4a – Comparaison des différents ajustements sans la prise en compte de possibles interactions : modèle naïf et modèle marginal.

Variable	Modalité	Naïf sans « effet élevage »				Marginal		
		RR <sup>*</sup>	IC <sup>**</sup>	EC <sup>♭</sup>	p <sup>♭♭</sup>	IC <sup>**</sup>	EC <sup>♭</sup>	p <sup>♭♭</sup>
Mois de naissance	janv.	0,42	0,30 - 0,60	0,17	< 0,001	0,22 - 0,82	0,33	0,011
	fév.	0,68	0,50 - 0,92	0,15	0,015	0,36 - 1,31	0,33	NS
	mars	1,42	1,04 - 1,92	0,15	0,024	0,78 - 2,63	0,31	NS
Condition de vêlage	aide facile	1,50	1,20 - 1,87	0,11	< 0,001	1,13 - 1,97	0,14	0,005
	aide difficile	1,39	1,05 - 1,83	0,14	0,019	1,06 - 1,93	0,15	0,018
Vaccination de la mère contre <i>E. coli</i>	non	0,82	0,72 - 0,94	0,06	0,004	0,65 - 1,02	0,11	NS
Vaccination de la mère contre d'autres agents que GENN	non	1,81	1,53 - 2,14	0,08	< 0,001	1,33 - 2,47	0,15	< 0,001
Distribution de concentré	non	1,39	1,24 - 1,56	0,06	< 0,001	1,14 - 1,70	0,10	0,001
Aliment à base d'ensilage de maïs	non	0,69	0,61 - 0,78	0,06	< 0,001	0,53 - 0,90	0,13	0,006
Odeur d'ammoniac perceptible dans le bâtiment	oui	1,62	1,35 - 1,94	0,09	< 0,001	1,18 - 2,23	0,16	0,003
Existence d'une infirmerie pour veaux	non	1,13	1,01 - 1,26	0,05	0,033	0,90 - 1,41	0,11	NS
Densité des veaux	forte	1,32	1,19 - 1,47	0,05	< 0,001	1,08 - 1,62	0,10	0,006
Incidence de la diarrhée la saison précédente	forte	1,37	1,20 - 1,56	0,06	< 0,001	1,02 - 1,83	0,15	0,036
Nettoyage des bâtiments des veaux avant la saison des vêlages	non	0,79	0,71 - 0,88	0,05	< 0,001	0,64 - 0,98	0,11	0,037
Nettoyage des bâtiments des veaux après la saison des vêlages	non	1,59	1,37 - 1,84	0,07	< 0,001	1,24 - 2,04	0,12	< 0,001
Nettoyage suite à chaque épisode diarrhéique	non	0,53	0,42 - 0,67	0,12	< 0,001	0,32 - 0,87	0,25	0,013
Propreté des vaches	faible	1,32	1,18 - 1,47	0,05	< 0,001	1,04 - 1,67	0,12	0,023
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	0,87	0,78 - 0,97	0,05	0,013	0,69 - 1,09	0,11	NS

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>♭</sup> : écart-type de  $\hat{\beta}$ ; <sup>♭♭</sup> : degré de significativité, avec NS = non significatif.

- la vaccination de la mère contre *E. coli* (RR = 0,82 contre 0,47 à l'origine) ;
- le nettoyage des bâtiments des veaux avant la saison des vêlages (RR = 0,79 contre 0,54 à l'origine) ;
- la supplémentation en vitamines et minéraux des veaux à leur naissance (RR = 0,87 et valant anciennement 0,68).

**Passage du modèle naïf au modèle marginal** Trois variables, dont les effets étaient jugés significatifs lors de l'analyse naïve, voient cette significativité de leur effet contestée par l'analyse marginale. Il s'agit de la vaccination de la mère contre *E. coli*, de l'existence d'une infirmerie pour veaux, et de la supplémentation en vitamines et minéraux des veaux à leur naissance.

Par ailleurs, le passage de l'ajustement naïf à l'ajustement marginal se traduit – à l'exception de la variable « condition de vêlage » – par un doublement de l'estimation de la variance des paramètres de régression : ainsi, le coefficient de dispersion varie entre les valeurs extrêmes 1,14 (« condition de vêlage ») et 6,25 (« incidence de la diarrhée la saison précédente »).

**Passage du modèle naïf au modèle mixte** Parmi les 15 variables retenues, 5 voient leur effet perdre toute significativité statistique lors du passage de l'ajustement naïf à l'ajustement mixte : il s'agit de la vaccination de la mère contre *E. coli*, de l'odeur d'ammoniac perceptible dans le bâtiment, de l'existence d'une infirmerie pour veaux, du nettoyage des bâtiments des veaux avant la saison des vêlages et de la supplémentation en vitamines et minéraux des veaux à leur naissance.

Concernant les estimations des paramètres de régression, nous obtenons des valeurs dif-

TAB. 6.4b – Comparaison des différents ajustements sans la prise en compte de possibles interactions : modèle naïf avec « effet élevage » fixe et modèle de fragilité.

Variable	Modalité	Naïf avec « effet élevage »				Mixte			
		RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>	RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>
Mois de naissance	janv.	0,41	0,28 - 0,61	0,19	< 0,001	0,43	0,29 - 0,62	0,19	< 0,001
	fév.	0,78	0,54 - 1,12	0,18	NS	0,77	0,54 - 1,09	0,17	NS
	mars	2,00	1,36 - 2,92	0,19	< 0,001	1,83	1,28 - 2,61	0,18	< 0,001
Condition de vêlage	aide facile	1,17	0,91 - 1,49	0,12	NS	1,25	0,98 - 1,59	0,12	NS
	aide difficile	1,39	1,03 - 1,88	0,15	0,032	1,41	1,05 - 1,89	0,15	0,023
Vaccination de la mère contre <i>E. coli</i>	non	1,00	0,76 - 1,31	0,13	NS	0,95	0,76 - 1,18	0,11	NS
Vaccination de la mère contre d'autres agents que GENN	non	1,63	1,23 - 2,17	0,14	< 0,001	1,73	1,37 - 2,19	0,12	< 0,001
Distribution de concentré	non	39,9	10 <sup>-6</sup> - 10 <sup>8</sup>	8,26	NS	1,50	1,15 - 1,95	0,12	0,001
Aliment à base d'ensilage de maïs	non	1,14	0,07 - 17,4	1,39	NS	0,67	0,52 - 0,86	0,12	0,001
Odeur d'ammoniac perceptible dans le bâtiment	oui	10 <sup>5</sup>	10 <sup>-71</sup> - 10 <sup>80</sup>	89,0	NS	1,51	0,93 - 2,46	0,24	NS
Existence d'une infirmerie pour veaux	non	0,85	0,25 - 2,85	0,61	NS	1,04	0,81 - 1,34	0,13	NS
Densité des veaux	forte	1,47	0,08 - 24,6	1,43	NS	1,36	1,06 - 1,75	0,13	0,017
Incidence de la diarrhée la saison précédente	forte	0,89	0,21 - 3,72	0,72	NS	1,42	1,10 - 1,85	0,13	0,007
Nettoyage des bâtiments des veaux avant la saison des vêlages	non	0,77	0,18 - 3,20	0,72	NS	0,83	0,63 - 1,09	0,14	NS
Nettoyage des bâtiments des veaux après la saison des vêlages	non	1,67	0,17 - 1,58	1,14	NS	1,41	1,04 - 1,90	0,15	0,023
Nettoyage suite à chaque épisode diarrhéique	non	0,58	0,02 - 15,8	1,68	NS	0,47	0,26 - 0,83	0,29	0,009
Propreté des vaches	faible	1,86	0,22 - 15,2	1,07	NS	1,33	1,06 - 1,68	0,12	0,015
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	0,68	10 <sup>-14</sup> - 10 <sup>13</sup>	15,8	NS	0,81	0,63 - 1,03	0,12	NS

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>b</sup> : écart-type de  $\hat{\beta}$ ; <sup>bb</sup> : degré de significativité, avec NS = non significatif.

féralant très faiblement d'un ajustement à l'autre, ces différences pouvant jouer dans un sens (augmentation du risque relatif) comme dans l'autre (diminution du risque relatif). Le sens de ces variations n'est pas lié à l'appartenance des variables au « niveau élevage » ou au niveau individuel. Notons que l'écart le plus important survient pour la première modalité de la variable « condition de vêlage », dont le risque relatif passe de la valeur 1,50 à la valeur 1,25, lors de la prise en compte d'un « effet élevage » aléatoire.

Concernant les estimations de la variance de ces paramètres de régression, là encore, le passage du modèle naïf au modèle robuste se traduit par une augmentation sensible de ces estimations – à l'exception des variables « mois de naissance » et « condition de vêlage ». Ainsi, le coefficient de dispersion présente comme valeurs extrêmes la valeur 1,19 (« condition de vêlage ») et la valeur 7,11 (« nettoyage suite à chaque épisode diarrhéique »).

**Comparaison des modèles marginal et mixte** Les conclusions rendues par les deux ajustements robustes diffèrent sensiblement. Rappelons que ce sont respectivement 3 et 5 variables qui voient la significativité de leur effet rejetée lors de la prise en compte d'un « effet élevage », respectivement dans le cas d'un ajustement marginal et dans celui d'un ajustement mixte. Comme cela était souhaitable, les 3 variables ayant perdu la significativité de leur effet dans le cadre de l'ajustement marginal font partie des 5 ayant perdu la leur dans le cadre de l'ajustement mixte.

TAB. 6.5 – Ajustements ne prenant en compte que les variables individuelles.

Variable	Modalité	Naïf				p <sup>bb</sup>				
		RR*	IC**	EC <sup>b</sup>						
Mois de naissance	janv.	0,47	0,34 - 0,66	0,16	< 0,001					
	fév.	0,76	0,58 - 1,02	0,14	NS					
	mars	1,25	0,95 - 1,64	0,14	NS					
Condition de vêlage	aide facile	1,47	1,18 - 1,81	0,11	< 0,001					
	aide difficile	1,40	1,08 - 1,82	0,13	0,01					
Vac. <i>E. coli</i> <sup>o</sup>	non	0,93	0,83 - 1,04	0,06	NS					
Vac. autres <sup>oo</sup>	non	1,24	1,09 - 1,42	0,06	< 0,001					
		Marginal				« Effet élevage » fixe				
Mois de naissance	janv.	0,47	0,24 - 0,93	0,34	0,029	0,41	0,28 - 0,61	0,19	< 0,001	
	fév.	0,76	0,42 - 1,38	0,30	NS	0,78	0,54 - 1,12	0,18	NS	
	mars	1,25	0,70 - 2,23	0,29	NS	1,99	1,36 - 2,92	0,19	< 0,001	
Condition de vêlage	aide facile	1,47	1,07 - 2,01	0,16	0,017	1,16	0,91 - 1,49	0,12	NS	
	aide difficile	1,40	1,01 - 1,94	0,16	0,040	1,39	1,02 - 1,88	0,15	0,032	
Vac. <i>E. coli</i> <sup>o</sup>	non	0,93	0,71 - 1,22	0,13	NS	1,01	0,76 - 1,31	0,13	NS	
Vac. autres <sup>oo</sup>	non	1,24	0,90 - 1,72	0,16	NS	1,63	1,23 - 2,17	0,14	< 0,001	
		Naïf avec strates*				Robuste avec fragilité*				
Mois de naissance	janv.	0,41	0,28 - 0,61	0,19	< 0,001	0,43	0,29 - 0,62	0,19	< 0,001	
	fév.	0,75	0,52 - 1,08	0,18	NS	0,77	0,54 - 1,10	0,18	NS	
	mars	2,03	1,38 - 2,98	0,19	< 0,001	1,78	1,24 - 2,56	0,18	0,001	
Condition de vêlage	aide facile	1,14	0,89 - 1,45	0,12	NS	1,21	0,95 - 1,54	0,12	NS	
	aide difficile	1,34	0,99 - 1,81	0,15	NS	1,39	1,03 - 1,88	0,15	0,028	
Vac. <i>E. coli</i> <sup>o</sup>	non	1,01	0,77 - 1,33	0,14	NS	0,99	0,79 - 1,24	0,11	NS	
Vac. autres <sup>oo</sup>	non	1,59	1,20 - 2,10	0,14	0,001	1,47	1,16 - 1,85	0,12	0,001	

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>b</sup> : écart-type de  $\hat{\beta}$ ; <sup>bb</sup> : degré de significativité, avec NS = non significatif; \* : stratification par élevages; <sup>o</sup> : vaccination de la mère contre *E. coli*; <sup>oo</sup> : vaccination de la mère contre d'autres agents que GENN.

**Concernant d'éventuelles interactions** Le tableau F.2, présenté dans l'annexe F p. 193, fournit les résultats des ajustements naïf et marginal à l'issue de la sixième étape de la méthode de sélection descendante (cf. p. 120). Pour 7 variables seulement, un au moins des 3 termes d'interaction qui sont associés à chacune d'entre elles ont un effet jugé significatif. Ajoutons que 5 variables parmi ces 7 ont un effet jugé non significatif. Enfin, notons que le modèle de fragilité de Cox n'a pu être ajusté, à l'aide du logiciel S-Plus, et en présence des éventuelles interactions, qu'en excluant les termes d'interaction relatifs à la variable « vaccination de la mère contre *E. coli* » (tab. F.3 p. 196).

Ainsi la prise en compte des différentes interactions n'apparaît-elle pas, au travers de ces résultats – similaires à ceux de l'étude antérieure – indispensable.

#### 6.4.4 Concernant la fragilité

##### Étude du rôle de la fragilité

Nous avons réalisé les deux ajustements robustes usuels – modèles de Cox marginal et mixte – sur les seules variables individuelles, afin d'étudier la prise en compte d'un excédent de

variabilité « exhaustif ». Les résultats sont présentés dans le tableau 6.5.

Par « exhaustif », nous entendons le fait que, dans le cadre du modèle de fragilité de Cox, la fragilité intègre maintenant toute la variabilité qui n'est pas expliquée par les seules variables individuelles – ce qui est le rôle exact théoriquement prêté à la fragilité. À la section précédente, le modèle de fragilité portait à la fois sur des variables propres aux individus, et sur des variables propres aux troupeaux. Ici, nous excluons de l'analyse les variables relatives aux troupeaux, afin de rendre à la fragilité sa fonction naturelle : celle de représenter toute la variabilité qui n'est pas due aux variables individuelles, et qui traduit l'effet de l'appartenance à un même troupeau.

Concernant l'ajustement marginal, si nous comparons les tableaux 6.4a et 6.5, à l'exception de la variable « vaccination de la mère contre d'autres agents que GENN », nous retrouvons des résultats similaires, tant en regard des estimations des risques relatifs que des estimations de leurs écart-types (et par conséquent des résultats concernant la significativité des effets individuels).

La variable « vaccination de la mère contre d'autres agents que GENN » est la seule variable pour laquelle le jugement statistique concernant son effet diffère : antérieurement hautement significatif ( $p < 0,001$ ), cet effet ne l'est maintenant plus.

Concernant l'ajustement mixte (tableaux 6.4b et 6.5), l'exclusion des variables relatives au niveau collectif (« niveau troupeau ») ne modifie en rien les résultats concernant les paramètres de régression (risques relatifs et écart-types). Par contre, cette même exclusion entraîne une modification importante de la valeur de l'estimation de la variance de la fragilité, qui passe de 0,76 à 1,35 ; ce résultat ne nous surprend pas, dans la mesure où la suppression de variables explicatives entraîne nécessairement une augmentation de la variabilité résiduelle.

Ainsi le doublement de la variabilité non observée (contenue toute entière dans le terme de fragilité), lors du passage de l'ajustement portant sur toutes les variables explicatives à celui ne portant plus que sur les variables explicatives du niveau individuel, n'affecte en rien les estimations relatives aux paramètres de régression.

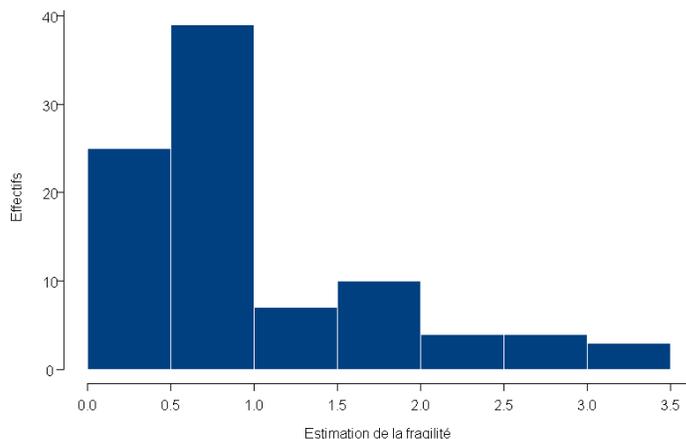
Ajoutons que concernant la fragilité elle-même, son effet est jugé hautement significatif ( $p < 0,001$ ).

Enfin, notons l'augmentation attendue de l'estimation de la variance des différents paramètres, lors de la prise en compte de la possible corrélation des données. Nous voyons que cette augmentation est bien plus marquée dans le cas de l'ajustement marginal que dans celui de l'ajustement mixte. Ainsi, concernant le mois de naissance par exemple, le coefficient de dispersion atteint la valeur 4,51 dans le cadre de l'ajustement marginal, alors qu'il n'est que de 1,41 dans celui de l'ajustement mixte. L'explication de cette très nette différence est fournie par la remarque faite en page 88, concernant la sous-estimation des variances des paramètres des effets fixes dans le cas d'un ajustement par le modèle de fragilité de Cox.

### Étude des estimations de la fragilité

Les valeurs des estimations de la fragilité sont représentées graphiquement au travers de l'histogramme 6.3, et analytiquement dans le tableau 6.6, qui contient les valeurs extrêmes de ces estimations.

La moyenne des estimations de la fragilité vaut 1. La moyenne des variances de l'estimation

FIG. 6.3 – Valeurs des estimations de la fragilité  $\exp(\hat{\Theta})$ .

TAB. 6.6 – Étude des élevages présentant les valeurs extrêmes des estimations de la fragilité.

Élevage	Effectif	% C   % D <sup>b</sup>	« Effet élevage » fixe				« Effet élevage » aléatoire			
			cov. ind.*		ttes cov.**		cov. ind.*		ttes cov.**	
			$\exp(\hat{\beta})$	$\sigma(\hat{\beta})$	$\hat{\beta}$	$\sigma(\hat{\beta})$	$\exp(\hat{\Theta})$	$\sigma(\hat{\Theta})$	$\exp(\hat{\Theta})$	$\sigma(\hat{\Theta})$
63	29	100   0	1,02	0,01	1,01	0,51	0,08	1,43	<b>0,10</b>	0,85
16	56	96,4   3,6	1,02	0,04	1,14	2,27	0,23	0,38	<b>0,18</b>	0,43
25	66	92,4   7,6	0,95	0,04	1,28	1,46	0,88	0,21	<b>0,20</b>	0,49
59	23	60,8   39,2	1,02	0,02	1,19	0,31	2,27	0,11	<b>3,03</b>	0,14
11	24	37,5   62,5	1,06	0,08	0,14	13,06	3,73	0,07	<b>3,09</b>	0,12
70	26	57,7   42,3	1,02	0,02	1,01	0,02	4,49	0,09	<b>3,24</b>	0,11

\* : ajustement pour les seules covariables individuelles ; \*\* : ajustement pour toutes les covariables ; <sup>b</sup> : pourcentage du nombre de cas de censures | pourcentage du nombre de cas de diarrhées.

de la fragilité vaut 0,40.

Rappelons que dans le cas du modèle de fragilité sur l'ensemble des covariables (tab. 6.4b), l'estimation de la variance de la fragilité vaut 0,76, tandis que dans celui du modèle de fragilité sur les covariables individuelles (tab. 6.5), l'estimation de la variance de la fragilité vaut 1,35. Ainsi, la prise en compte de variables explicatives au niveau du troupeau réduit – comme attendu – la surdispersion, dans une proportion de moitié.

Au sein des élevages présentant les plus faibles valeurs pour  $\exp(\hat{\Theta})$ , la répartition des censures et des cas de diarrhée est fortement déséquilibrée : le pourcentage de censures y est toujours supérieur à 90 %.

Par contre, lorsque nous nous intéressons aux troupeaux présentant les plus fortes estimations pour la fragilité, nous constatons que cette répartition entre cas de censures et cas de diarrhée est moins marquée : tantôt un plus grand nombre de censure que de diarrhée prévaut à l'intérieur du troupeau, tantôt ce sont les cas de diarrhée qui sont les plus nombreux.

Cette répartition, au sein de chaque troupeau, entre cas de censures et cas de diarrhée ne saurait expliquer la valeur (ni même l'amplitude de cette valeur) prise par l'estimation de la fragilité. Ainsi, à titre de comparaison, citons :

- l'élevage 1, dont la répartition en pourcentages des cas de censures / diarrhées vaut 100 / 0, dont l'effectif total est de 62 animaux, et dont l'estimation de la fragilité est de 0,29 ;
- l'élevage 18, dont les mêmes quantités valent respectivement 100 / 0, 27 et 0,32 ;
- l'élevage 7, présentant les valeurs respectives de 32 / 68, 22 et 2,00 ;
- l'élevage 27, dont les caractéristiques sont de 29,1 / 70, 9, 24 et 1,29.

#### 6.4.5 Justesse des ajustements

**Test de la proportionnalité des risques** Le tableau 6.7 nous informe sur la qualité des différents ajustements. Le test de la proportionnalité des risques est effectué pour chacune des covariables ; la justesse de l'ajustement adopté est testée au travers d'un test global, qui porte sur l'ensemble des covariables.

Ainsi, quel que soit l'ajustement choisi, la proportionnalité des risques est rejetée lorsqu'elle est testée sur l'ensemble des covariables (test *global*). C'est ce test global qui permet d'apprécier la qualité de l'ajustement – en particulier, nous ne pouvons en aucun cas conclure à une meilleure qualité de l'ajustement naïf par la seule considération du fait que pour un plus grand nombre de covariables, l'hypothèse de base n'est pas rejetée.

**Sujets médiocrement prédits** Concernant la figure 6.4, les animaux présentant un résidu important sont médiocrement prédits par le modèle. Nous distinguons clairement, sur la figure, deux ensembles de points : un premier nuage, situé vers le bas, comprenant les animaux dont la valeur du résidu de la déviance est faible, et un second nuage comprenant, au contraire, les animaux pour lesquels le résidu a une valeur importante.

Nous avons cherché à savoir quel était le mécanisme de répartition des individus entre ces deux ensembles de points. Il est apparu que les animaux pour lesquels nous observons un résidu de la déviance important sont ceux qui ont présenté un épisode diarrhéique.

#### 6.4.6 Conclusions

##### En période *peripartum*

Parmi les facteurs étudiés, nous avons constaté que l'apparition des GENN dépend essentiellement des conditions dans lesquelles se déroule le vêlage, ainsi que de l'état du veau à sa naissance. En effet, les conditions de mise bas sont étroitement liées à l'apparition de la diarrhée : les veaux nés d'un vêlage difficile ont 41 fois plus de (mal)chances d'avoir la diarrhée que ceux nés sans aide (RR = 1,41, p = 0,023). De plus, le vêlage dystocique affaiblit le nouveau-né, rendu de la sorte sensible aux infections ; par suite, le veau ne pourra se pourvoir suffisamment en colostrum, et il demeurera plus longtemps exposé aux agents infectieux contenus dans la litière. Par ailleurs, le risque d'inoculation d'agents pathogènes, lors d'interventions humaines, ne doit pas non plus être négligé.

TAB. 6.7 – Test de la proportionnalité des risques pour les différents ajustements : degré de significativité\* pour le test de l’hypothèse nulle  $H_0$  : « les risques sont proportionnels ».

Variable	Modalités	Naïf	Marginal	Mixte
Mois de naissance	janv.	NS	NS	0,005
	fév.	0,033	0,009	0,015
	mars	< 0,001	< 0,001	< 0,001
Condition de vêlage	aide facile	NS	NS	NS
	aide difficile	0,028	< 0,001	0,036
Vaccination de la mère contre <i>E. coli</i>	non	< 0,001	< 0,001	< 0,001
Vaccination de la mère contre d’autres agents que GENN	non	NS	< 0,001	< 0,001
Distribution de concentré	non	NS	< 0,001	< 0,001
Aliment à base d’ensilage de maïs	non	NS	0,042	< 0,001
Odeur d’ammoniac perceptible dans le bâtiment	oui	0,003	< 0,001	< 0,001
Existence d’une infirmerie pour veaux	non	NS	0,001	0,014
Densité des veaux	forte	NS	NS	NS
Incidence de la diarrhée la saison précédente	forte	NS	0,008	< 0,001
Nettoyage des bâtiments des veaux avant la saison des vêlages	non	0,007	< 0,001	< 0,001
Nettoyage des bâtiments des veaux après la saison des vêlages	non	NS	NS	0,002
Nettoyage suite à chaque épisode diarrhéique	non	NS	< 0,001	< 0,001
Propreté des vaches	faible	NS	< 0,001	< 0,001
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	NS	NS	NS
<b>Test global</b>		< 0,001	< 0,001	< 0,001

\* : degré de significativité, avec NS = non significatif.

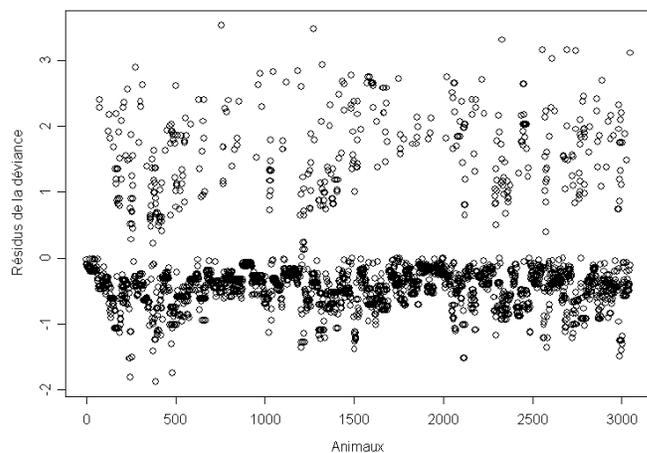


FIG. 6.4 – Résidus de la déviance pour le modèle de fragilité.

Ces résultats suggèrent de limiter les interventions aux cas de dystocie où l’aide de l’éleveur s’avère réellement indispensable – cependant, dans la pratique, il paraît difficile de refuser à l’éleveur qui assiste au vêlage d’y participer...

La dyspnée du veau à la naissance n’a pas présenté de lien statistiquement significatif avec

la diarrhée – dans le cadre de notre étude, comme dans celle d’origine. Ce résultat pourrait découler d’une association entre certaines variables : les difficultés respiratoires ne seraient alors que la conséquence d’une parturition difficile.

Concernant les naissances gémellaires, 50 couples de jumeaux et un triplet sont nés pendant la période du suivi. Or la faiblesse de cette proportion (3,4 %) pourrait être à l’origine de l’absence de lien significatif entre le nombre de veaux dans la portée d’une part, et la survenue de la diarrhée d’autre part. Il pourrait en être de même concernant la prématurité (proportion de 0,7 %) et les malformations congénitales (0,6 %) : dans tous ces cas, le faible nombre de cas observés ne permettrait pas de distinguer les veaux exposés au facteur au risque des veaux qui ne le sont pas.

Par ailleurs, aucun facteur lié à la prise de colostrum n’a pu être significativement associé à la survenue de la diarrhée – ce qui est contraire aux conclusions de nombreuses autres études : il est en effet reconnu que ce nutriment est un élément capital de l’élaboration de la protection du veau vis-à-vis des infections.

Toutefois, 84 % des veaux étaient supposés avoir reçu suffisamment de colostrum pendant leur premier jour de vie ; de façon similaire, concernant la précocité des veaux, 88 % d’entre eux auraient tété lors de leurs 6 premières heures de vie. Or, dans le même temps, ce sont 74 % des éleveurs qui déclarent avoir assisté à la première prise de colostrum. Ces différences de proportions, qui mettent à mal la concordance entre les dires et les faits relatifs à la présence de l’éleveur au sein de son troupeau, pourraient entraîner une surestimation des proportions de veaux recevant précocement et/ou en quantité suffisante le colostrum. Par suite, cette surestimation serait susceptible de masquer toute association entre la prise suffisante et/ou précoce de colostrum d’une part, et la survenue de la diarrhée d’autre part. (Signalons également qu’il est difficile d’évaluer exactement la quantité absorbée par le veau s’il tète sa mère...

Enfin, notons qu’une interaction entre les conditions du vêlage et l’absorption de colostrum ne peut être exclue, un veau fatigué du fait d’un vêlage dystocique étant susceptible de téter peu (ou pas) sa mère.

D’autres facteurs relatifs aux soins du veau à sa naissance ne semblent pas liés aux GENN, d’après les résultats de notre étude : il s’agit de la désinfection de l’ombilic et de l’utilisation d’un local de vêlage.

### **Bâtiment et ambiance**

Le type de stabulation ne semble pas lié à la survenue de la diarrhée.

Par contre, la présence d’odeur d’ammoniac apparaît significativement liée à la survenue de la diarrhée dans le cadre des modèles naïf ( $RR = 1,62$ ,  $p < 0,001$ ) et marginal, mais pas dans celui du modèle mixte ( $RR = 1,51$ , non significatif). Une odeur d’ammoniac a été ressentie dans 6,6 % des élevages. Ce facteur nous renseigne indirectement sur l’aération et l’entretien de la litière. Concernant l’aération elle-même, nous avons observé que 53,2 % des bâtiments présentaient une aération correcte ; ce facteur n’a pas, lors de l’étude statistique, présenté de lien significatif avec la morbidité.

L’humidité du bâtiment, ainsi que l’existence d’infirmier pour les adultes, n’ont pas non plus présenté de lien significatif avec la survenue de la diarrhée. Seule, l’existence d’une infirme-

rie pour veaux a paru liée aux GENN, dans le cas d'une approche naïve : sitôt l'« effet élevage » pris en compte lors de l'ajustement statistique, la significativité de ce lien disparaît.

Il ressort de l'ensemble de ces résultats que le bâtiment n'aurait pas d'effet direct sur l'apparition des diarrhées. En revanche, comme nous allons le voir dans la section suivante, d'autres paramètres semblent conditionner la survenue des gastro-entérites ; or ces paramètres sont indirectement liés au bâtiment (il s'agit, par exemple, de son entretien et de son hygiène).

### Entretien de la stabulation

L'entretien du bâtiment regroupe plusieurs pratiques qui semblent fortement liées à l'apparition des GENN.

Nous avons constaté que le nettoyage des locaux *après* la saison des vélages réduisait presque de moitié les cas de survenue de diarrhées (RR = 1,41,  $p = 0,023$  avec l'ajustement mixte). Ce résultat peut s'interpréter ainsi : cette pratique de nettoyage permettrait d'assainir les locaux (c'est-à-dire de réduire la concentration des agents infectieux), et d'accueillir ainsi dans des conditions d'hygiène satisfaisantes les nouveaux-nés.

À l'inverse, et contre toute attente, nous avons constaté que le nettoyage du bâtiment *avant* la saison des vélages constituait un facteur de risque de survenue de la diarrhée pour les veaux, dans le cadre des ajustements naïf et marginal (RR = 0,79 avec  $p < 0,001$  et  $p = 0,037$  lors d'ajustements respectivement naïf et marginal). En revanche, cette significativité n'apparaît plus lors de l'ajustement du modèle de Cox avec fragilité.

Concernant ce phénomène inattendu, nous pouvons avancer l'explication suivante : la majorité des élevages (73,9 %) ne pratiquant pas ce nettoyage, il se pourrait que les 24 élevages qui y ont recours présentent, à l'origine, des problèmes particuliers de GENN, ce qui entraînerait de leur part une pratique accrue de ce type de nettoyage – dans le but de diminuer l'incidence de la maladie. Pour appuyer cette explication, soulignons que les 24 élevages concernés représentent 36,3 % des cas de diarrhée.

La propreté des vaches est un autre facteur qui, d'après notre étude, est significativement lié à la survenue des GENN (RR = 1,33 avec  $p = 0,015$  dans le cadre du modèle mixte) : une propreté insatisfaisante des animaux augmente ainsi de 33 % le risque de voir apparaître des cas de diarrhée.

Les autres facteurs relatifs à l'entretien de la stabulation (tels la quantité ou le rythme de paillage) n'ont pas été conservés lors de l'étape de sélection des variables.

### Facteurs alimentaires

Les facteurs caractérisant l'alimentation des animaux adultes ne paraissent avoir aucune répercussion sur la survenue de diarrhée chez les veaux.

La seule variable étant apparue significativement liée à la morbidité est la distribution d'aliments concentrés aux vaches : une privation de ce type d'aliments entraîne une augmentation du risque de GENN d'environ 50 % (RR = 1,50 avec  $p = 0,001$  dans le cadre de l'ajustement mixte). Ces résultats sont conformes à ceux rencontrés dans la littérature vétérinaire.

Nous devons cependant admettre que la quantification précise de l'impact de l'absorption d'aliment sur l'incidence de la maladie n'est pas aisée. En effet, la nature de l'aliment – qui est une information de première importance – n'est pas le seul élément qui doit être connu.

La qualité de l'ensilage, par exemple, peut ainsi considérablement varier d'une étable à l'autre, en fonction de sa préparation ou de sa conservation. En toute rigueur, des analyses physico-chimiques auraient dû être menées afin d'évaluer précisément l'effet des aliments sur la survenue de diarrhée.

### Conduite et gestion du troupeau

Parmi les paramètres étudiés, la parité des vaches n'a pas paru significativement liée à la morbidité – ce qui est contraire aux conclusions d'autres études épidémiologiques (Bendali, 1999). D'après ces études, les veaux nés de primipares seraient davantage exposés à la maladie, en raison de la faible production de colostrum et de sa moindre qualité d'une part, et des complications lors du vêlage (naturellement plus fréquentes chez les primipares) d'autre part. Cependant, nous pouvons nous demander si l'effet de la parité ne serait pas, en réalité, masqué par celui – significatif – des conditions du vêlage.

D'autres facteurs, tels que le groupement des adultes en lots, le groupement des vêlages en hiver ou le parcage des veaux, ne semblent pas liés à la survenue des GENN.

En revanche, la densité des veaux est un facteur de risque significatif (RR = 1,36 avec  $p = 0,017$  lors de l'ajustement mixte) : cet accroissement du risque est sans doute lié à la pression microbienne, à la contamination transversale entre animaux et au stress supplémentaire, consécutifs de cette promiscuité accrue.

Enfin, l'effet du mois de naissance est apparu significativement associé à l'incidence de la diarrhée : le risque de GENN pour un veau né au mois de mars (RR = 1,83 avec  $p < 0,001$ ) serait ainsi presque deux fois supérieur à celui d'un veau né en décembre.

### Prophylaxie

Les variables relevant du domaine prophylactique ont été les plus nombreuses, au sein de notre étude.

La supplémentation vitaminique et minérale des les veaux, pratiquée dans 43,4 % des élevages, semble protéger les veaux des diarrhées, d'après les résultats de l'analyse naïve (RR = 0,87 avec  $p = 0,013$ ) ; mais la significativité de cet effet ne résiste pas à la prise en compte d'un « effet élevage », qu'elle soit entreprise dans le cadre d'un modèle marginal ou d'un modèle mixte.

La vaccination, quant à elle, constitue une part importante des mesures préventives contre les maladies chez les jeunes bovins. Les résultats de notre étude ont montré que la vaccination de la mère contre d'autres agents que les GENN avait un effet bénéfique sur le veau : il est naturel de penser que cette pratique vaccinale témoigne d'une bonne politique de gestion de l'élevage en terme de prophylaxie et de prévention.

Paradoxalement, la vaccination de la parturiente contre *E. coli*, quand elle se révèle être un facteur de risque significatif (c'est le cas lors de l'approche naïve), semble favoriser la survenue de diarrhée. Diverses raisons pourraient expliquer ce phénomène inattendu ; nous en avançons trois :

- 1° les élevages à forte incidence auraient tendance à présenter une pratique vaccinale plus importante que les autres ;

- 2° la vaccination serait réalisée à l'aide de souches qui diffèreraient de celles à l'origine des GENN ;
- 3° la vaccination s'effectuerait dans de mauvaises conditions (elle entraînerait un stress supplémentaire chez l'animal, ou bien elle serait elle-même mal réalisée).

Cependant, la vaccination constituant l'un des agents de lutte contre les maladies les plus efficaces, il ne serait d'aucun intérêt d'en décider la suspension. Aussi, au vu de nos résultats, devons-nous préférentiellement suggérer d'examiner et d'améliorer les conditions et le respect des indications d'emploi qui entourent l'administration de vaccins, plutôt que de suggérer l'abandon de cette pratique.

### 6.4.7 Discussion

À l'éleveur qui interroge le statisticien sur les bénéfices à attendre d'une enquête statistique multifactorielle, portant sur les différentes conditions dans lesquelles cet éleveur tient son troupeau, nous pouvons apporter quelques éléments de réponse.

Tout d'abord, la récolte des données constituée d'informations recueillies sur des animaux regroupés en troupeaux, doit inciter le statisticien à tenir compte d'un « effet élevage », afin de corriger les erreurs de jugement statistique qui découleraient de mauvaises estimations des différents paramètres entrant dans le modèle mathématique. Cette certitude, le statisticien doit faire en sorte que l'éleveur la partage avec lui<sup>1</sup>.

À partir de cette situation, deux types de modèle se présentent au statisticien : le modèle marginal, qui modélise le risque relatif obtenu « en moyenne » sur l'ensemble de l'échantillon, et le modèle de fragilité, qui modélise pour sa part un risque relatif « intra-élevage ». Le statisticien désireux de mettre en pratique ces deux types de modèle robuste – auquel, à titre de comparaison, il peut adjoindre le modèle naïf – pourra conclure de la sorte :

- les estimations des risques relatifs diffèrent très peu d'une approche à l'autre ;
- la prise en compte d'un « effet élevage » ampute 3 variables de leur significativité : il s'agit des variables « vaccination de la mère contre *E. coli* », « existence d'une infirmerie pour veaux » et « supplémentation en vitamines et minéraux des veaux à leur naissance » ;
- le choix entre approche marginale et approche mixte, parallèlement à l'interprétation que l'éleveur désirera faire du risque relatif obtenu, entraîne les remarques suivantes :

- 1° 2 variables supplémentaires perdent leur significativité lors du passage du modèle marginal au modèle mixte, qui sont le « nettoyage des bâtiments des veaux avant la saison des vélages » et l'« odeur d'ammoniac perceptible dans le bâtiment »,
- 2° ces deux variables, relatives au niveau « élevage », et qui conservent des risques relatifs similaires à ceux de l'approche marginale, doivent le changement de statut de leur significativité à l'accroissement de variabilité obtenue lors de l'introduction d'un effet aléatoire dans le modèle (ainsi, pour la seconde de ces deux variables, l'écart-type vaut respectivement 0,16 et 0,24 suivant que l'approche est de type marginal ou mixte).

---

1. À condition qu'il en soit lui-même pénétré... ce qui n'est pas encore un fait établi, si l'on s'en réfère au nombre d'études portant sur des données potentiellement corrélées, et ne tenant aucun compte de cette corrélation éventuelle.

TAB. 6.8 – Valeur de la log-vraisemblance pour les différents ajustements.

Modèle	Log $\mathcal{L}^*$
Naïf sans covariable	-3 461,65
Fragilité (gamma) sans covariable	-3 215,76
Fragilité (gaussienne) sans covariable	-3 221,27
Naïf avec covariables	-3 316,49
Naïf avec covariables et « effet élevage » fixe	-3 141,28
Marginal avec covariables	-3 316,49
Fragilité (gamma) avec covariables	-3 163,01
Fragilité (gaussienne) avec covariables	-3 160,19

\* : log-vraisemblance du modèle.

Dans le cas présent, concernant le choix du modèle robuste finalement retenu, deux remarques supplémentaires peuvent être faites.

La première consiste à rappeler (cf. p. 129) l'existence d'une certaine redondance entre les informations apportées par le facteur « odeur d'ammoniac » et celles que procurent d'autres facteurs présents dans le modèle robuste final (par exemple la densité des veaux, la propreté des vaches ou encore le nettoyage des bâtiments). Aussi, face au dilemme de l'éleveur qui ne sait lequel des deux modèles robustes il doit privilégier, ni non plus quantifier la perte liée à l'exclusion de ce facteur explicatif du modèle, le statisticien doit-il pouvoir raisonnablement assurer que la décision d'opter pour le modèle mixte ne sacrifie en rien une partie de l'information.

La seconde concerne l'autre variable dont la significativité disparaît lors du passage du modèle marginal au modèle mixte : la variable « nettoyage des bâtiments des veaux avant la saison des vêlages ». Le risque relatif associé à ce facteur ne laissait pas de nous surprendre, puisqu'il se trouvait aller à l'encontre du sens commun : du fait de son infériorité à 1 (RR = 0,79 dans le cadre du modèle marginal), ce risque amenait à déconseiller le nettoyage des bâtiments avant la saison des vêlages. Or le passage au modèle mixte a pour conséquence non seulement de réviser la conclusion statistique relative à la significativité de l'effet de ce facteur, mais encore d'augmenter légèrement la valeur du risque relatif (qui vaut 0,83 dans le cadre du modèle mixte), qui se rapproche ainsi de l'unité. Il se pourrait ainsi que le modèle mixte *corrige* les conclusions – sinon fausses, du moins surprenantes – tirées de l'ajustement marginal, et qui semblaient liées à une sorte de biais dans la représentation des élevages vis-à-vis de ce facteur (cf. p. 129).

En outre, l'étude de la fragilité apporte une information d'importance : elle permet d'estimer, pour chaque élevage, la quantité de variabilité due au seul « effet élevage ». Ainsi, quand une valeur de 1 signifie un effet nul du facteur « élevage », nous voyons (cf. fig. 6.3) qu'un nombre non négligeable d'élevages présentent des valeurs supérieures à 2 concernant l'estimation de la fragilité. À partir de cette observation, il est possible de déterminer quels sont les élevages présentant les variabilités les plus fortes, et par suite d'étudier plus précisément les facteurs non encore observés qui peuvent être à l'origine de cette variabilité.

Nous donnons, à titre indicatif, les valeurs de la log-vraisemblance obtenues pour les différents ajustements (tab. 6.8). Bien que statistiquement non (toutes) comparables entre elles, ces valeurs nous éclairent sur l'apport d'information que constitue l'inclusion, dans le modèle

TAB. 6.9 – Les différentes procédures algorithmiques pour l’ajustement du modèle de fragilité.

	(RE)ML*		AIC		CAIC		DL**	
	Gamma	Normale	Gamma	Normale	Gamma	Normale	Gamma	Normale
Variance <sup>b</sup>	0,76	0,93	1,76	1,34	1,80	1,13	1,07	1,13
Degrés de liberté <sup>b</sup>	55,7	60,48	64,1	56,1	63,9	54,8	60	55
Log-vraisemblance	-3 163,0	-3 160,2	-3 162,4	-3 155,5	-3 150,9	-3 157,4	-3 157,1	-3 157,2
Outer*	8	5	10	6	10	6	6	8
Newton-Raphson**	41	24	155	43	106	40	18	26

\* : ML (*Maximum Likelihood*) pour la fragilité gamma, REML (*REstricted ML*) pour la fragilité gaussienne; \*\* : degrés de libertés – associés à la fragilité – fixés à 60 pour la fragilité gamma, et à 55 pour la fragilité gaussienne; <sup>b</sup> : concernant la fragilité; \* : nombre d’itérations de la boucle extérieure (maximisation de la log-vraisemblance par rapport à  $\gamma$ ); \*\* : nombre d’itérations de la boucle intérieure (maximisation de la vraisemblance partielle pénalisée – à  $\gamma$  fixé).

statistique, de covariables ou d’un « effet élevage ».

Ainsi, outre le bénéfice certain tiré de l’inclusion des variables explicatives dans le modèle, nous retrouvons un résultat similaire à l’un de ceux que nous avons énoncés lors de l’analyse par simulations : il s’agit de la forte similarité entre les valeurs de la log-vraisemblance dans les cas d’un modèle avec « effet troupeau » fixe et aléatoire. Cependant, dans le cas présent, la préférence concernant le meilleur ajustement va nécessairement à celui incluant un « effet groupe » aléatoire, du fait des estimations aberrantes obtenues dans le cadre du modèle de Cox avec l’« effet troupeau » introduit comme paramètre fixe.

Quant au tableau 6.9, il nous renseigne sur les résultats obtenus avec les différentes options du logiciel S-Plus, disponible lors d’un ajustement mixte et concernant la méthode d’estimation de la variance de la fragilité. Rappelons ces procédures :

- dans le cas d’une fragilité gamma : méthodes du maximum de vraisemblance (ML), du Critère d’Information d’Akaike<sup>1</sup> (AIC), du Critère d’Information d’Akaike corrigé<sup>1</sup> (CAIC) et du nombre de degrés de liberté fixés à l’avance<sup>2</sup> ;
- dans le cas d’une fragilité gaussienne : méthodes du maximum de vraisemblance résiduelle (REML), du Critère d’Information d’Akaike (AIC), du Critère d’Information d’Akaike corrigé (CAIC) et du nombre de degrés de liberté fixés à l’avance.

Le statisticien, quant à lui, s’interroge naturellement sur la qualité des outils qui sont à sa disposition, ainsi que sur l’utilisation la plus judicieuse qu’il peut faire de ceux-ci. Dans le cas présent, et malgré le type de données étudié (la survenue de la diarrhée dans les 31 premiers jours de vie), type approprié aux modèles pour données de survie, force est de constater que le choix du modèle statistique – celui de Cox – n’a peut-être pas été le meilleur : les deux études portant sur la qualité de l’ajustement de ce modèle ont révélé sa mauvaise adéquation aux données. En particulier, la proportionnalité des risques ne s’est pas vérifiée. Dans ce cas-là, c’est le choix même du modèle de Cox qui doit être remis en cause.

1. Ce critère consiste ici à minimiser le Critère d’Information d’Akaike, ou ce même critère corrigé.

2. L’utilisateur spécifie le nombre de degré de libertés, qui varie de 0 à  $n$  (où  $n$  est le nombre de groupes) : ceci est rendu possible du fait de la considération de l’effet aléatoire comme version pénalisée d’une variable factorielle. Voir Therneau et Grambsch (2000) p. 233 pour plus de détails.

Notons également l'importance du taux de censure présent dans notre jeu de données. Par sa valeur élevée (85,6 %), ce taux peut entraîner des estimations faussées et des erreurs concernant les conclusions apportées par certains tests (voir, en particulier, les résultats obtenus par simulations avec un taux de censure fixé à 90 %).

Cependant, l'analyse originelle (Bendali, 1999) – qui utilisait déjà le modèle de Cox – tenait compte de ce taux de censure élevé. Aussi, dans le souci de poursuivre au plus près la première étude qui avait été menée sur les facteurs de risques relatifs aux GENN, avons-nous tenu à reprendre – et dans les mêmes conditions – l'outil analytique de Cox.



Quatrième partie

# CONCLUSIONS



Les méthodes classiques d'analyse de survie supposent l'indépendance des temps de survenue de l'événement d'intérêt ; or cette hypothèse ne peut plus être raisonnablement posée lors de l'étude de données de survie groupées. Les modèles robustes, qu'ils soient de type marginal ou de type mixte, permettent alors de traiter ces données hétérogènes.

Nous avons souhaité, dans un premier temps, réunir toutes les informations relatives au modèle de Cox et à ses extensions concernant la prise en compte de la possible corrélation des données de survie. Nous avons exposé les différents modèles, leur signification, ainsi que leurs critères d'adéquation.

Nous avons ensuite présenté une longue étude par simulations, ayant pour objectif de comparer les différents ajustements qui pouvaient être entrepris à partir de données de survie corrélées. Nous avons montré par ces simulations que le modèle de fragilité de Cox semblait le plus à même de prendre en compte cette surdispersion des données de survie, due à leur regroupement. Cependant, un modèle de Cox avec « effet groupe » fixe a également prouvé sa robustesse, notamment lors de l'étude de la puissance des tests statistiques pratiqués dans cette étude par simulations.

Concernant le modèle de fragilité, nous nous sommes toutefois heurtés au problème de la sous-estimation de la variance du paramètre de régression, lors de l'ajustement de ce type de modèle au travers du logiciel S-Plus.

Nous avons également analysé les effets, sur les estimations des différents paramètres, du taux de censure et des effectifs des données de survie au sein des différents groupes. Le rôle extrême d'un taux de censure de 90 % – diminuant de façon importante, tant la précision des estimations des paramètres de régression, que la puissance des tests statistiques pratiqués –, à l'opposé du faible rôle que joue ce taux de censure ramené à 60 % – et *a fortiori* à 30 % –, a été illustré au travers des simulations. Quant à la répartition des données entre les différents groupes, nous avons vu que la situation d'un faible nombre de groupes, chaque groupe contenant de nombreux individus, était préférable à celle d'un nombre important de groupes d'effectifs réduits.

Par ailleurs, la particularité essentielle que présentait le modèle de fragilité de Cox, à savoir la quantification de l'hétérogénéité, a été mise en évidence. Ainsi, outre le test appelé à décider de la significativité ou non de l'« effet groupe », il est possible de donner une valeur explicite de cette surdispersion – au travers des estimations de la fragilité et de sa variance.

Ce résultat important prend tout son sens lors d'une étude épidémiologique au sein de laquelle le statisticien désire connaître les groupes « les plus fragiles ». Ainsi, lors de notre étude sur les GENN chez le veau, nous avons pu déterminer quels étaient les troupeaux qui présentaient les risques instantanés de diarrhée les plus importants. Un prolongement naturel de cette étude<sup>1</sup> eût été l'analyse plus particulière de ces groupes « fragilisés », afin de déterminer les causes de ce phénomène de fragilisation.

En effet, l'épidémiologiste conçoit l'idéal que serait la détermination de *tous* les facteurs de risques : il ne lui serait plus nécessaire alors de recourir à l'introduction d'un effet aléatoire dans le modèle statistique, puisque cet effet – censé modéliser la variabilité non observable – serait nul...

Les conclusions que nous avons tirées de cette étude concernant les facteurs de risques liés à la survenue de GENN sont apparues en accord avec celles de l'étude antérieure (Bendali,

---

1. Naturel mais irréalisable ici, car il dépasse le cadre de cette thèse.

1999) : il est possible d'associer à cette pathologie une dizaine de facteurs de risque, influant soit au niveau individuel soit au niveau collectif, certains favorisant fortement la survenue de la diarrhée (l'absence de vaccination de la mère contre d'autres agents que GENN, par exemple, dont le risque relatif s'élève à 1,73), d'autres au contraire prouvant leur effet protecteur (la non-alimentation à base d'ensilage de maïs, dont le risque relatif est de 0,67).

Cependant, des limites à l'étendue de ces études ont clairement été posées.

Concernant l'étude par simulations tout d'abord, nous n'avons introduit qu'une seule variable explicative (binaire), agissant au niveau individuel. Il eût été intéressant de complexifier ce modèle et d'y inclure, notamment, des variables explicatives quantitatives, agissant tant au niveau individuel qu'au niveau collectif.

Nous aurions pu, également, renoncer à l'équilibre de la répartition des données entre les différents groupes, et observer l'influence de ce déséquilibre sur les résultats statistiques.

Au cours de l'étude des GENN, nous n'avons pu occulter le problème que posait l'ampleur de la censure sur l'ensemble de l'échantillon (85 % des données). Cependant, nous avons tenu à reprendre le travail qui avait été entrepris et publié auparavant, afin d'y ajouter l'analyse au travers du modèle mixte qui avait fait défaut jusque-là.

Notons aussi qu'il aurait été possible d'introduire dans le modèle statistique une dépendance vis-à-vis du temps concernant certaines covariables ; nous aurions pu également modéliser la censure sous une forme « informative », c'est-à-dire que l'indépendance entre le phénomène de survenue de la diarrhée et celui de la censure aurait été rejetée. L'étude portant sur la survenue de diarrhée dans le premier mois de vie du veau, nous aurions pu enfin adopter une autre stratégie d'analyse, et opter pour un mécanisme de troncature des données – touchant tous les temps de survenue de diarrhée supérieurs à 31 jours.

Nous avons également souhaité développer les connaissances en prédiction, dans le cadre du modèle de fragilité de Cox. Cependant nous nous sommes heurtés à certaines difficultés concernant la loi asymptotique de la fragilité : les outils habituels propres aux processus de comptage (tel le théorème de Rebolledo) n'ont pas permis de déterminer les propriétés asymptotiques du meilleur prédicteur de la fragilité.

Nous nous sommes intéressés dans ce travail à une seule source d'hétérogénéité dans la population, celle provenant de facteurs de risque non observés communs à un groupe d'individus (tels que des facteurs de risques génétiques ou environnementaux). Or il peut exister une autre source d'hétérogénéité dans la population qui peut provenir de variables individuelles négligées. Le modèle de fragilité corrélée permet alors de traiter ces deux sources d'hétérogénéité ; il serait souhaitable d'entreprendre son application dans une étude future.

Ainsi les modèles robustes de Cox sont-ils prometteurs, et leurs applications épidémiologiques nombreuses. Nous nous sommes appuyés sur un problème épidémiologique de données groupées pour illustrer cette méthodologie statistique ; cependant il est possible d'en faire l'application à une structure d'échantillon présentant une potentielle – et quelconque – corrélation des données.

Cinquième partie

**BIBLIOGRAPHIE**



- O. O. AALEN  
*Non-Parametric Inference for a Family of Counting Processes*  
Annals of Statistics, vol. 6 p. 701–726, 1978.
- O. O. AALEN  
*Heterogeneity in Survival Data Analysis*  
Statistics in Medicine, vol. 7 p. 1121–1137, 1988.
- K. ABRAMS, D. ASHBY ET D. ERRINGTON  
*A Bayesian Approach to Weibull Survival Models — Application to a Cancer Clinical Trial*  
Lifetime Data Analysis, vol. 2 p. 159–174, 1996.
- H. AKAIKE  
*Information Theory and an Extension of the Maximum Likelihood Principle*  
In *Proceedings of International Symposium on Information Theory*. B.N. Petrov et F. Czaki (éd.), Budapest, 1973.
- P. K. ANDERSEN, O. BORGAN, R. D. GILL ET N. KEIDING  
*Statistical Models Based on Counting Processes*  
Springer-Verlag, 1991.
- P. K. ANDERSEN, J. P. KLEIN, K. M. KNUDSEN ET R. T. Y PALACIOS  
*Estimation of Variance in Cox's Regression Model with Shared Gamma Frailties*  
Biometrics, vol. 53 p. 1475–1484, 1997.
- P. K. ANDERSEN, J. P. KLEIN ET M.-J. ZHANG  
*Testing for Centre Effects in Multi-Center Survival Studies: a Monte Carlo Comparison of Fixed and Random Effects Tests*  
Statistics in Medicine, vol. 18 p. 1489–1500, 1999.
- P. K. ANDERSEN  
*Testing Goodness of Fit of Cox's Regression and Life Model*  
Biometrics, vol. 38 p. 67–77, 1982.
- E. ARJAS ET L. LIU  
*Non-Parametric Bayesian Approach to Hazard Regression: a Case Study with a Large Number of Missing Covariate Values*  
Statistics in Medicine, vol. 15 p. 1757–1770, 1996.
- E. ARJAS  
*A Graphical Method for Assessing Goodness of Fit in Cox's Proportional Hazards Model*  
Journal of the American Statistical Association, vol. 83 p. 204–212, 1988.
- F. BENDALI, H. BICHET, F. SCHELCHER ET M. SANAA  
*Pattern of Diarrhoea in Newborn Beef Calves in South-West France*  
Veterinary Research, vol. 30 p. 61–74, 1999a.
- F. BENDALI, M. SANAA, H. BICHET ET F. SCHELCHER  
*Risk Factors with Diarrhoea in Newborn Calves*  
Veterinary Research, vol. 30 p. 509–522, 1999b.

- F. BENDALI  
*Étude Épidémiologique des Gastro-Entérites Néonatales chez le Veau*  
Thèse de Doctorat, Université de Franche-Comté, Faculté de Médecine et de Pharmacie de Besançon, 1999.
- D. A. BINDER  
*Fitting Cox's Proportional Hazards Models from Survey Data*  
Biometrika, vol. 79 p. 139–147, 1992.
- P. E. BÖHMER  
*Theorie der unabhängigen Wahrscheinlichkeiten Rapports*  
Mémoires et Procès-verbaux du 7<sup>e</sup> Congrès International d'Actuaires, vol. 2 p. 327–343, 1912.
- N. BRESLOW  
*A Generalized Kruskal-Wallis Test for Comparing K Samples Subject to Unequal Patterns of Censorship*  
Biometrika, vol. 57 p. 579–594, 1970.
- N. BRESLOW  
*Covariance Analysis of Censored Survival Data*  
Biometrics, vol. 30 p. 89–99, 1974.
- J. CAI ET R. L. PRENTICE  
*Estimating Equations for Hazard Ratio Parameters Based on Correlated Failure Time Data*  
Biometrika, vol. 82 p. 151–164, 1995.
- J. CAI ET R. L. PRENTICE  
*Regression Estimation Using Multivariate Failure Time Data and a Common Baseline Hazard Function Model*  
Lifetime Data Analysis, vol. 3 p. 197–213, 1997.
- J. CAI  
*Hypothesis Testing of Hazard Ratio Parameters in Marginal Models for Multivariate Failure Time Data*  
Lifetime Data Analysis, vol. 5 p. 39–53, 1999.
- J. M. CHAMBERS, C. L. MALLOWS ET B. W. STUCK  
*A Method for Simulating Stable Random Variables*  
Journal of the American Statistical Association, vol. 71 p. 340–344, 1976.
- D. G. CLAYTON  
*A Monte-Carlo Method for Bayesian Inference in Frailty Models*  
Biometrics, vol. 47 p. 467–485, 1991.
- D. CLAYTON ET J. CUZICK  
*Multivariate Generalizations of the Proportional Hazards Model*  
Journal of the Royal Statistical Society, Series A, vol. 148 p. 82–117, 1985.
- D. COMMENGES ET P. K. ANDERSEN  
*Score Test of Homogeneity for Survival Data*  
Lifetime Data Analysis, vol. 1 p. 145–159, 1995.

- D. R. COX ET E. J. SNELL  
*A General Definition of Residuals (with Discussion)*  
Journal of the Royal Statistical Society, Series B, vol. 30 p. 248–275, 1968.
- D. R. COX  
*Regression Models and Life Tables (with Discussion)*  
Journal of the Royal Statistical Society, Series B, vol. 74 p. 187–220, 1972.
- D. DACUNHA-CASTELLE ET M. DUFLO  
*Probabilités et Statistiques – tomes 1 et 2 (2<sup>e</sup> édition)*  
Masson, 1993.
- J.-J. DROESBEKE, B. FICHET ET P. TASSI  
*Analyse Statistique des Durées de Vie – Modélisation des données censurées*  
Economica, 1989.
- V. DUCROCQ ET G. CASELLA  
*A Bayesian Analysis of Mixed Survival Models*  
Genetics, Selection, Evolution, vol. 28 p. 505–529, 1996.
- C. ELBERS ET G. RIDDER  
*True and Spurious Duration Dependence: the Identifiability of the Proportional Hazard Model*  
Review of Economics Studies, vol. 49 p. 403–409, 1982.
- J. ETEZADI-AMOLI ET A. CIAMPI  
*Extended Hazard Regression for Censored Survival Data with Covariates: a Spline Approximation for the Baseline Hazard Function*  
Biometrics, vol. 43 p. 181–192, 1987.
- R. V. FOUTZ  
*On the Unique Consistent Solution to the Likelihood Equations*  
Journal of the American Statistical Association, vol. 72 p. 147–148, 1977.
- N. A. FRANCK ET J. B. KANEENE  
*Management Risk Factors Associated with Calf Diarrhoea in Michigan Herds*  
Journal of Dairy Science, vol. 76 p. 1313–1323, 1993.
- D. GAMERMAN  
*Dynamic Bayesian Models for Survival Data*  
Applied Statistics, vol. 40 p. 63–79, 1991.
- E. GEHAN  
*A Generalized Wilcoxon Test for Comparing Arbitrarily Singly-Censored Data*  
Biometrika, vol. 52 p. 202–223, 1965.
- R. D. GILL  
*Discussion of Multivariate Generalizations of the Proportional Hazards Model, by Clayton and Cuzick*  
Journal of the Royal Statistical Society, Series A, vol. 148 p. 108–109, 1985.

- R. D. GILL  
*On Estimating Transition Intensities of a Markov Process with Aggregated Data of a Certain Type: 'Occurrences but no Exposures'*  
Scandinavian Journal of Statistics, vol. 13 p. 113–134, 1986.
- R. D. GILL  
*Marginal Partial Likelihood*  
Scandinavian Journal of Statistics, vol. 19 p. 133–137, 1992.
- R. GILL ET M. SCHUMACHER  
*A Simple Test of the Proportional Hazards Model*  
Biometrika, vol. 74 p. 289–300, 1987.
- R. J. GRAY  
*Tests for Variation Over Groups in Survival Data*  
Journal of the American Statistical Society, vol. 90 p. 198–203, 1995.
- R. J. GRAY  
*On Tests for Group Variation With a Small to Moderate Number of Groups*  
Lifetime Data Analysis, vol. 4 p. 139–148, 1998.
- M. GREENWOOD  
*The Natural Duration of Cancer*  
Reports on Public Health and Medical Subjects, vol. 33 p. 1–26, 1926.
- P. GUSTAFSON  
*Large Hierarchical Bayesian Analysis of Multivariate Survival Data*  
Biometrics, vol. 53 p. 230–242, 1997.
- I. D. HA, Y. LEE ET J.-K. SONG  
*Hierarchical Likelihood Approach for Frailty Models*  
Biometrika, vol. 88 p. 233–243, 2001.
- C. R. HENDERSON  
*Best Linear Unbiased Estimation and Prediction Under a Selection Model*  
Biometrics, vol. 31 p. 423–447, 1975.
- R. HENDERSON ET P. OMAN  
*Effect of Frailty on Marginal Regression Estimates in Survival Analysis*  
Journal of the Royal Statistical Society, Series B, vol. 61 p. 367–379, 1999.
- P. HOUGAARD  
*Life Table Methods for Heterogeneous Populations: Distributions Describing the Heterogeneity*  
Biometrika, vol. 71 p. 75–83, 1984.
- P. HOUGAARD  
*Survival Models for Heterogeneous Populations Derived from Stable Distributions*  
Biometrika, vol. 73 p. 387–393, 1986.
- P. HOUGAARD  
*Frailty Models for Survival Data*  
Lifetime Data Analysis, vol. 1 p. 255–273, 1995.

- P. HOUGAARD  
*Analysis of Multivariate Survival Data*  
Statistics for Biology and Health. Springer, 2000.
- M. P. JONES ET J. CROWLEY  
*A General Class of Nonparametric Tests for Survival Analysis*  
Biometrics, vol. 45 p. 157–170, 1989.
- J. D. KALBFLEISCH ET R. L. PRENTICE  
*The Statistical Analysis of Failure Time Data*  
Wiley, 1980.
- J. D. KALBFLEISCH  
*Non-parametric Bayesian Analysis of Survival Time Data*  
Journal of the Royal Statistical Society, Series A, vol. 40 p. 214–221, 1978.
- E. L. KAPLAN ET P. MEIER  
*Nonparametric Estimation from Incomplete Observations*  
Journal of the American Statistical Association, vol. 53 p. 457–481, 1958.
- T. M. KING, T. H. BEATY ET K.-Y. LIANG  
*Comparison of Methods for Survival Analysis of Dependent Data*  
Genetic Epidemiology, vol. 13 p. 139–158, 1996.
- J. P. KLEIN  
*Semiparametric Estimation of Random Effect Using the Cox Model Based on the EM Algorithm*  
Biometrics, vol. 48 p. 795–806, 1992.
- T. LANCASTER  
*The Econometric Analysis of Transition Data*  
Cambridge University Press, 1990.
- K.-Y. LIANG, L. G. SELF ET Y.-C. CHANG  
*Modelling Marginal Hazards in Multivariate Failure-Time Data*  
Journal of the Royal Statistical Society, Series B, vol. 55 p. 441–453, 1993.
- K.-Y. LIANG, S. G. SELF, K. J. BANDEEN-ROCHE ET S. L. ZEGER  
*Some Recent Developments for Regression Analysis of Multivariate Failure Time Data*  
Lifetime Data Analysis, vol. 1 p. 403–415, 1995.
- K.-Y. LIANG, S. L. ZEGER ET B. QAQISH  
*Multivariate Regression Analyses for Categorical Data*  
Journal of the Royal Statistical Society, Series B, vol. 54 p. 3–40, 1992.
- D. Y. LIN ET L. J. WEI  
*The Robust Inference for the Cox Proportional Hazards Model*  
Journal of the American Statistical Association, vol. 84 p. 1074–1078, 1989.
- C. A. MACGILCHRIST ET C. W. AISBETT  
*Regression with Frailty in Survival Analysis*  
Biometrics, vol. 47 p. 461–466, 1991.

- C. A. MACGILCHRIST  
*REML Estimation for Survival Models with Frailty*  
 Biometrics, vol. 49 p. 221–225, 1993.
- C. MAHÉ  
*Analyse Statistique de Délais d'Événements Corrélés*  
 Thèse de Doctorat, Université Paris 7 – Denis Diderot, 1998.
- N. MANTEL  
*Evaluation of Survival Data and Two New Rank Order Statistics Arising in its Consideration*  
 Cancer Chemotherapy Reporting, vol. 50 p. 163–170, 1966.
- L. MARZEC ET P. MARZEC  
*On Fitting Cox's Regression Model with Time-Dependent Coefficients*  
 Biometrika, vol. 84 p. 901–908, 1997.
- MATHSOFT, INC.  
*S-Plus 2000, Guide to Statistics, vol.1*  
 Data Analysis Products Division, 101 Main Street Cambridge, MA 02142 USA edition, 1999a.
- MATHSOFT, INC.  
*S-Plus 2000, Guide to Statistics, vol.2*  
 Data Analysis Products Division, 101 Main Street Cambridge, MA 02142 USA edition, 1999b.
- MATHSOFT, INC.  
*S-Plus 2000, User's Guide*  
 Data Analysis Products Division, 101 Main Street Cambridge, MA 02142 USA edition, 1999c.
- S. A. MURPHY  
*Consistency in a Proportional Hazards Model Incorporating a Random Effect*  
 The Annals of Statistic, vol. 22 p. 712–731, 1994.
- S. A. MURPHY  
*Asymptotic Theory for the Frailty Model*  
 The Annals of Statistic, vol. 23 p. 182–198, 1995.
- W. B. NELSON  
*Theory and Applications of Hazard Plotting for Censored Data*  
 Technometrics, vol. 14 p. 945–965, 1972.
- J. M. NEUHAUS ET M. R. SEGAL  
*Design Effects for Binary Regression Models Fitted to Dependent Data*  
 Statistics in Medicine, vol. 12 p. 1259–1268, 1993.
- G. G. NIELSEN, R. D. GILL, P. K. ANDERSEN ET T. H. I. A. SØRENSEN  
*A Counting Process Approach to Maximum Likelihood Estimation in Frailty Models*  
 Scandinavian Journal of Statistics, vol. 19 p. 25–43, 1992.
- M. C. PAIK, W.-Y. TSAI ET R. OTTMAN  
*Multivariate Survival Analysis Using Piecewise Gamma Frailty*  
 Biometrics, vol. 50 p. 975–988, 1994.

- E. PARNER  
*Inference in Semiparametric Frailty Models*  
Thèse de Doctorat, Université d'Aarhus (Danemark), 1997.
- E. PARNER  
*Asymptotic Theory for the Correlated Gamma-Frailty Model*  
The Annals of Statistics, vol. 26 p. 183–214, 1998.
- J. H. PETERSEN  
*An Additive Frailty Model for Correlated Life Times*  
Biometrics, vol. 54 p. 646–661, 1998.
- R. PETO ET J. PETO  
*Asymptotically Efficient Rank-Invariant Test Procedures*  
Journal of the Royal Statistical Society, Series A, vol. 135 p. 185–206, 1972.
- A. PICKLES ET R. CROUCHLEY  
*A Comparison of Frailty Models for Multivariate Survival Data*  
Statistics in Medicine, vol. 14 p. 1447–1461, 1995.
- R. L. PRENTICE  
*Linear Rank Tests with Right-Censored Data*  
Biometrika, vol. 65 p. 167–179, 1978.
- S. K. SAHU, D. K. DEY, H. ASLANIDOU ET D. SINHA  
*A Weibull Regression Model with Gamma Frailties for Multivariate Survival Data*  
Lifetime Data Analysis, vol. 3 p. 123–137, 1997.
- D. J. SARGENT  
*A General Framework for Random Effects Survival Analysis in the Cox Proportional Hazards Setting*  
Biometrics, vol. 54 p. 1486–1497, 1998.
- M. D. SCHLUCHTER  
*Methods for the Analysis of Informatively Censored Longitudinal Data*  
Statistics in Medicine, vol. 11 p. 1861–1870, 1992.
- D. SCHOENFELD  
*Partial Residuals for the Proportional Hazards Regression Model*  
Biometrika, vol. 69 p. 239–241, 1982.
- A. J. SCOTT ET D. HOLT  
*The Effect of Two-Stage Sampling on Ordinary Least Squares Methods*  
Journal of the American Statistical Association, vol. 77 p. 485–497, 1982.
- M. R. SEGAL, J. M. NEUHAUS ET I. R. JAMES  
*Dependance Estimation for Marginal Models of Multivariate Survival Models*  
Lifetime Data Analysis, vol. 3 p. 251–268, 1997.
- J. SHIH ET T. A. LOUIS  
*Assessing Gamma Frailty Models for Clustered Failure Time Data*  
Lifetime Data Analysis, vol. 1 p. 205–220, 1995.

- D. SINHA  
*Posterior Likelihood Methods for Multivariate Survival Data*  
Biometrics, vol. 54 p. 1463–1474, 1998.
- C. F. SPIEKERMAN ET D. Y. LIN  
*Checking the Marginal Cox Model for Correlated Failure Time Data*  
Biometrika, vol. 83 p. 143–156, 1996.
- Y. SUN ET M. SHERMAN  
*Some Permutation Tests for Survival Data*  
Biometrics, vol. 52 p. 87–97, 1996.
- R. TARONE ET J. WARE  
*On Distribution-Free Test for Equality of Survival Distributions*  
Biometrika, vol. 64 p. 156–160, 1977.
- T. M. THERNEAU, P. M. GRAMBSCH ET P. R. FLEMING  
*Martingale-Based Residuals for Survival Models*  
Biometrika, vol. 77 p. 147–160, 1990.
- T. M. THERNEAU ET P. M. GRAMBSCH  
*Modeling Survival Data – Extending the Cox Model*  
Springer, 2000.
- J. W. VAUPEL, K. G. MANTON ET E. STALLARD  
*The Impact of Heterogeneity in Individual Frailty on the Dynamics of Mortality*  
Demography, vol. 16 p. 439–454, 1979.
- L. J. WEI, D. Y. LIN ET L. WEISSFELD  
*Regression Analysis of Multivariate Incomplete Failure Time Data by Modelling Marginal Distributions*  
Journal of the American Statistical Association, vol. 84 p. 1065–1073, 1989.

Sixième partie

**ANNEXES**



# Annexe A

## THÉORIE STATISTIQUE

### A.1 Processus aléatoires et intégrales stochastiques

**Définition VI.1** — Soit un espace probabilisé  $(\Omega, \mathcal{A}, \mathbb{P})$  où  $\mathbb{P}$  est la mesure de probabilité sur  $(\Omega, \mathcal{A})$ . Un **processus aléatoire**, ou encore une **fonction aléatoire réelle** (f.a.r.) est une fonction à deux variables :  $t$  – le temps – et  $\omega$  – le hasard –, et elle est notée  $X(t, \omega)$ , avec  $t \in [0, \infty[$  et  $\omega \in \Omega$ .

À  $t$  fixé, la fonction  $X_t : \omega \mapsto X(t, \omega)$  est appelée **coordonnée** à l'instant  $t$  (c'est donc une v.a.). La **trajectoire** est  $\omega \mapsto (X(t, \omega), t \geq 0)$ .

Une f.a.r. à trajectoire continue (f.a.r.c.) est une application

$$\begin{aligned} X : [0, \infty[ \times \Omega &\rightarrow \mathbb{R} \\ (t, \omega) &\mapsto X(t, \omega) \end{aligned}$$

telle que

- a) pour presque tout  $\omega$ ,  $t \mapsto X(t, \omega)$  est continue,
- b) pour tout  $t \geq 0$ ,  $X_t : \omega \mapsto X(t, \omega)$  est une v.a. réelle.

**Définition VI.2** — Soit un espace probabilisé  $(\Omega, \mathcal{A}, \mathbb{P})$ ,  $t \in \mathbb{N}$  ou  $\mathbb{R}^+$ . Une **filtration** est une famille  $\mathcal{F}_t$  de tribus,  $t \in \mathbb{N}$  ou  $\mathbb{R}^+$ , telle que

$$\mathcal{F}_s \subset \mathcal{F}_t \subset \mathcal{A},$$

$\forall s \leq t$ .

**Définition VI.3** — Soient un espace probabilisé  $(\Omega, \mathcal{A}, \mathbb{P})$  et  $\mathcal{F}_t$  une filtration. Un processus  $X = X(t, \omega)$  est dit  **$\mathcal{F}_t$ -adapté** si  $\forall t$ ,  $X_t$  est  $\mathcal{F}_t$ -mesurable.

**Définition VI.4** — Soit un espace probabilisé  $(\Omega, \mathcal{A}, \mathbb{P})$ . Soit  $(M_t)_t$ ,  $t \in \mathbb{N}$  ou  $\mathbb{R}^+$  un processus réel défini sur  $\Omega$ . Soit  $(\mathcal{F}_t)_t$  une filtration sur  $\Omega$ .

$(M_t)_t$  est une  **$\mathcal{F}_t$ -martingale** si

- (i)  $\forall t$ ,  $M_t$  est  $\mathcal{F}_t$ -adaptée et  $M_t \in \mathcal{L}^1$  (i.e. est intégrable),
- (ii)  $\forall s, t$ ,  $0 \leq s \leq t$ ,  $\mathbb{E}(M_t | \mathcal{F}_s) = M_s$  p.s.

**Définition VI.5** — Soit un espace probabilisé  $(\Omega, \mathcal{A}, \mathbb{P})$  et une filtration  $\mathcal{F} = (\mathcal{F}_t)_{t \in \mathbb{R}^+}$ .

Un  **$\mathcal{F}$ -processus croissant**  $A$  est un processus adapté à  $\mathcal{F}$  à valeurs réelles satisfaisant la propriété suivante : pour tout  $\omega \in \Omega$ , les trajectoires  $t \mapsto A_t(\omega)$  sont croissantes, continues à droite et nulles en 0.

**Définition VI.6** — La tribu  $P$  de  $]0, \infty[ \times \Omega$  engendrée par les ensembles  $]s, t] \times \Gamma$ , où  $0 < s < t$  et  $\Gamma \in \mathcal{F}_s$ , est la **tribu des ensembles prévisibles**. Une v.a.  $C$  définie sur  $]0, \infty[ \times \Omega, P$  est un **processus prévisible**. On note  $C(t, \omega) = C_t(\omega)$ .

**Proposition VI.1** — Soient  $A$  un processus croissant et  $C_s$  un processus prévisible. Alors, pour tout  $t$ ,  $\int_0^t C_s dA_s$  est une v.a.

## A.2 Processus de comptage

**Définition VI.7** — Un **processus de comptage**  $N$  est un processus cadlag (continu à droite avec une limite à gauche), adapté, nul en zéro, croissant et ayant des sauts d'amplitude 1.

**Proposition VI.2** — Soit  $N(t)$  un processus de comptage. Il existe un processus  $\Lambda(t)$  prévisible, croissant, continu à droite et nul en zéro tel que

$$M(t) = N(t) - \Lambda(t)$$

soit une martingale.

$\Lambda(t)$  s'appelle le **compensateur** de  $N(t)$ , ou encore son **processus d'intensité cumulée**.

**Proposition VI.3** — Soient  $N$  un processus ponctuel de dimension 1, et  $\Lambda$  son compensateur. Si  $N$  est absolument continu, alors  $N$  possède une **intensité**  $\lambda$ , i.e. il existe un processus prévisible  $\lambda$  tel que

$$\Lambda(t) = \int_0^t \lambda(s) ds$$

pour tout  $t$ .

L'intensité est définie par

$$\lambda(s) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \mathbb{P}(N(s + \epsilon) - N(s) \geq 1 \mid \mathcal{F}_s),$$

où  $\mathcal{F}_s$ ,  $s \in \mathcal{S}$ , est la filtration naturelle, c'est-à-dire l'ensemble des événements observables à l'instant  $s$ .

**Proposition VI.4** — Soit  $M(t)$  une martingale. Alors  $M^2(t)$  est une sous-martingale et

$$M^2(t) = M_t + \langle M \rangle_t,$$

avec  $\langle M \rangle_t$  défini par

$$\langle M \rangle_t = \lim_{|\delta| \rightarrow 0} \sum_i \mathbb{E} [(M_{t_{i+1}} - M_{t_i})^2 | \mathcal{F}_{t_i}].$$

$\langle M \rangle_t$  est appelé le **processus prévisible croissant** associé à  $M(t)$ .

**Proposition VI.5** — Soit  $N$  un processus de comptage et  $\Lambda$  son compensateur. Le processus prévisible associé à la martingale locale de carré intégrable  $M = N - \Lambda$  (ou encore le compensateur de  $M^2$ ) vaut

$$\begin{aligned} \langle M \rangle &= \Lambda - \int \Delta \Lambda \, d\Lambda \\ &= \int (1 - \Delta \Lambda) \, d\Lambda, \end{aligned}$$

et en particulier, si  $\Lambda$  est continu,

$$\langle M \rangle = \Lambda.$$

**Théorème VI.1 (Théorème de l'innovation)** — Soit  $N$  un processus de comptage adapté par rapport à deux filtrations  $(\mathcal{F}_t)_t$  et  $(\mathcal{G}_t)_t$  telles que  $\mathcal{F}_t \subseteq \mathcal{G}_t$ . On suppose que  $N$  a pour intensité  $\lambda$  par rapport à  $(\mathcal{G}_t)_t$ .

Alors il existe un processus  $\tilde{\lambda}$  prévisible par rapport à  $(\mathcal{F}_t)_t$  et tel que

$$\tilde{\lambda}(t) = \mathbb{E} [\lambda(t) | \mathcal{F}_{t-}].$$

$\tilde{\lambda}$  est le processus d'intensité de  $N$  par rapport à  $(\mathcal{F}_t)_t$ .

**Définition VI.8** — Un processus de comptage  $r$ -dimensionnel  $N = \{N_i : i = 1, \dots, r\}$  est appelé **processus de comptage multivarié** si chacune de ses composantes est un processus de comptage univarié et s'il ne peut y avoir simultanément des sauts de deux (ou plus) de ses composantes.

**Proposition VI.6** — Soit  $M$  une martingale localement de carré intégrable, et  $T$  un temps d'arrêt. On définit le processus d'arrêt  $M^T$  par

$$M^T(t) = M(t \wedge T).$$

$M^T$  est localement de carré intégrable si et seulement si

$$\mathbb{E} \langle M \rangle(T) < \infty$$

et, également, si et seulement si

$$\mathbb{E} [M](T) < \infty.$$

De plus, si  $M$  est une martingale locale à variation finie, alors  $M^T$  est une martingale uniformément intégrable, et de variation intégrable si et seulement si

$$\mathbb{E} \int_0^T |dM(s)| < \infty.$$

**Proposition VI.7** — Soient  $N$  un processus de comptage d'intensité  $\lambda$ ,  $M = N - \int \lambda$  et  $H$  un processus prévisible localement borné. Alors  $M$  et  $\int HdM$  sont des martingales localement de carré intégrable et

$$\begin{aligned} \langle M \rangle &= \text{diag} \int \lambda, \\ [M] &= \text{diag} N, \\ \left\langle \int HdM \right\rangle &= \int H \text{diag} \lambda H^t, \\ \left[ \int HdM \right] &= \int H d(\text{diag} N) H^t, \end{aligned} \tag{A.1}$$

où  $\text{diag} N$ , pour un vecteur  $N$ , est la matrice diagonale formée des composantes du vecteur. Nous pouvons réécrire ces équations sous la forme

$$\begin{aligned} \langle M_h, M_l \rangle &= \delta_{hl} \int \lambda_h, \\ [M_h, M_l] &= \delta_{hl} N_h, \\ \left\langle \sum_h \int H_{jh} dM_h, \sum_l \int H_{kl} dM_l \right\rangle &= \sum_h \int H_{jh} H_{kh} \lambda_h, \\ \left[ \sum_h \int H_{jh} dM_h, \sum_l \int H_{kl} dM_l \right] &= \sum_h \int H_{jh} H_{kh} dN_h. \end{aligned} \tag{A.2}$$

**Proposition VI.8** — Soit une suite de processus  $X^{(n)}$  tels que

$$X^{(n)}(s) \xrightarrow{\mathbb{P}} f(s) \quad (n \rightarrow \infty)$$

pour presque tout  $s \in [0, \tau]$ , où la fonction déterministe  $f$  satisfait

$$\int_0^\tau |f(s)| ds < \infty.$$

Supposons de plus qu'il existe  $k_\delta$  tel que  $\int_0^\tau k_\delta < \infty$  et tel que

$$\liminf_{n \rightarrow \infty} \mathbb{P} \left( \left| X^{(n)}(s) \right| \leq k_\delta, \forall s \right) \geq 1 - \delta.$$

Alors

$$\sup_t \left| \int_0^t X^{(n)}(s) ds - \int_0^t f(s) ds \right| \xrightarrow{\mathbb{P}} 0. \tag{A.3}$$

### A.3 Théorème de la limite centrale

**Théorème VI.2 (Rebolledo)** — Si  $M_n$  est une suite de martingales, et si

(i)  $\langle M_n \rangle_t$  converge en probabilité vers  $v_t$  déterministe,

(ii)  $\forall \epsilon, \exists M_{n,\epsilon}$  suite de martingales telles qu'aucune différence  $M_n - M_{n,\epsilon}$  n'ait une amplitude supérieure à  $\epsilon$ ,

alors  $M_n(t)$  a une limite  $M(t)$  de processus croissant  $v_t$ , et  $M(t)$  est un processus gaussien :

$$\frac{M_n(t)}{v_t} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1).$$

## A.4 Produit infini (ou intégral)

**Définition VI.9** — Soit  $X(s)$  un processus cadlag, nul en 0, et à variation bornée. On obtient une mesure additive en posant

$$X(]s, t]) = X(t) - X(s).$$

Soit une partition  $t_0 = s < t_1 < \dots < t_n = t$ . Son **pas** est

$$|\delta| = \sup_i |t_i - t_{i-1}|.$$

On appelle **produit-intégrale** (ou **produit infini**)

$$\begin{aligned} \mathcal{P}_s^t(1 + dX) &= \mathcal{P}_{]s,t]}(1 + dX) \\ &= \lim_{|\delta| \rightarrow 0} \prod_{i=1}^n \left[ 1 + X(]t_{i-1}, t_i]) \right], \end{aligned}$$

qui est indépendante de la suite des  $(\delta)$ .

**Propriété VI.1** — Si  $X(t)$  est continu, alors

$$\mathcal{P}_{]0,t]}(1 + dX) = e^{X(t)}.$$

**Théorème VI.3 (Duhamel)** — Soient  $Y = \mathcal{P}(1 + dX)$  et  $Y' = \mathcal{P}(1 + dX')$ . Alors

$$Y(t) - Y'(t) = \int_{s \in [0,t]} \mathcal{P}_{]0,s]}(1 + dX) (X(ds) - X'(ds)) \mathcal{P}_{(s,t]}(1 + dX').$$

Si  $Y'(t)$  est régulière, alors

$$\begin{aligned} \frac{Y(t)}{Y'(t)} - 1 &= \int_{s \in [0,t]} \mathcal{P}_{]0,s]}(1 + dX) (X(ds) - X'(ds)) \left[ \mathcal{P}_{]0,s]}(1 + dX') \right]^{-1} \\ &= \int_0^t \frac{Y(s-)}{Y'(s)} [X(ds) - X'(ds)]. \end{aligned}$$

**Théorème VI.4 (Jacod)** — Soit un processus de comptage multivarié  $N = (N_1, \dots, N_k)$  sur  $[0, \tau]$ .

Soient  $\mathbb{P}$  et  $\tilde{\mathbb{P}}$  deux mesures de probabilité sur deux espaces tels que  $N$  ait pour compensateur dans chacun d'eux respectivement  $\Lambda$  et  $\tilde{\Lambda}$ . On suppose que  $\tilde{\mathbb{P}}$  est absolument continue par rapport à  $\mathbb{P}$  :  $\tilde{\mathbb{P}} \ll \mathbb{P}$ . Alors

$$\tilde{\Lambda}_h \ll \Lambda_h \quad \text{pour tout } h, \mathbb{P} - p.s.$$

et

$$\begin{aligned} \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} &= \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} \bigg|_{\mathcal{F}_0} \frac{\mathcal{P}_{t \in [0, \tau]} \left\{ \prod_h d\tilde{\Lambda}_h(t)^{\Delta N_h(t)} (1 - d\tilde{\Lambda}_+(t))^{1 - \Delta N_+(t)} \right\}}{\mathcal{P}_{t \in [0, \tau]} \left\{ \prod_h d\Lambda_h(t)^{\Delta N_h(t)} (1 - d\Lambda_+(t))^{1 - \Delta N_+(t)} \right\}} \\ &= \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} \bigg|_{\mathcal{F}_0} \prod_t \prod_h \left( \frac{d\tilde{\Lambda}_h(t)}{d\Lambda_h(t)} \right)^{\Delta N_h(t)} \frac{\mathcal{P}_{t \in [0, \tau] : \Delta N_+(t) \neq 1} (1 - d\tilde{\Lambda}_+(t))}{\mathcal{P}_{t \in [0, \tau] : \Delta N_+(t) \neq 1} (1 - d\Lambda_+(t))}. \end{aligned}$$

**Théorème VI.5** — Si  $\Lambda$  et  $\tilde{\Lambda}$  sont presque sûrement continues, alors

$$\frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} = \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} \bigg|_{\mathcal{F}_0} \frac{\prod_t \prod_h d\tilde{\Lambda}_h(t)^{\Delta N_h(t)} \exp[-\tilde{\Lambda}_+(\tau)]}{\prod_t \prod_h d\Lambda_h(t)^{\Delta N_h(t)} \exp[-\Lambda_+(\tau)]}$$

et si  $\Lambda$  et  $\tilde{\Lambda}$  sont presque sûrement absolument continues, alors

$$\frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} = \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}} \bigg|_{\mathcal{F}_0} \frac{\prod_t \prod_h \tilde{\lambda}_h(t)^{\Delta N_h(t)} \exp[-\tilde{\Lambda}_+(\tau)]}{\prod_t \prod_h \lambda_h(t)^{\Delta N_h(t)} \exp[-\Lambda_+(\tau)]}.$$

**Théorème VI.6** — La vraisemblance vaut

$$d\mathbb{P} = d\mathbb{P} \bigg|_{\mathcal{F}_0} \mathcal{P}_{t \in [0, \tau]} \left( \prod_h \lambda_h(t)^{\Delta N_h(t)} [1 - \lambda_+(t)]^{1 - \Delta N_+(t)} \right)$$

et dans le cas continu

$$d\mathbb{P} = d\mathbb{P} \bigg|_{\mathcal{F}_0} \prod_t \prod_h \lambda_h(t)^{\Delta N_h(t)} \exp[-\Lambda_+(\tau)]. \quad (\text{A.4})$$

## A.5 Vraisemblances complète et partielle

La vraisemblance pour un processus de comptage  $N(t)$  *self-exciting* – i.e. dont la filtration associée est exactement (ou augmentée d'ensembles nuls)  $\sigma(N(s), s \leq t)$  pour tout  $t$  – vaut

$$d\mathbb{P} = \mathcal{P}_{t \in [0, \tau]} \left( d\Lambda(t)^{\Delta N(t)} [1 - d\Lambda_+(t)]^{1 - \Delta N_+(t)} \right). \quad (\text{A.5})$$

Notons que, même si le processus n'est pas *self-exciting* – i.e. si la filtration associée à ce processus est plus grande que celle générée par  $N$  –, les rapports de vraisemblance et fonctions de score obtenus avec la partie droite de (A.5) conservent des propriétés de martingales, ce qui peut être utile pour toute inférence basée sur la méthode du maximum de vraisemblance.

Revenons au cas d'un processus *self-exciting*. La vraisemblance s'écrit sous la forme

$$\mathcal{P}_{t \in [0, \tau]} \left( \prod_{x \in E} \Lambda(dt, dx)^{N(dt, dx)} [1 - \Lambda(dt, E)]^{1 - N(dt, E)} \right), \quad (\text{A.6})$$

où

- $N$  est le processus de marques avec ses marques dans  $(E, \mathcal{E})$ ;
- $\Lambda$  est le compensateur de  $N$ ;
- $N$  et  $\Lambda$  sont des mesures aléatoires sur  $(\mathcal{T} \times E, \mathcal{B}(\mathcal{T}) \otimes \mathcal{E})$ .

L'interprétation de (A.6) demeure la même :

- conditionnellement à ce qui s'est passé avant l'instant  $t$ , un événement survient
  - dans l'intervalle de temps  $[t, t + dt)$  – que nous noterons  $dt$  pour simplifier –,
  - avec une marque dans  $dx$ ,
 avec la probabilité  $\Lambda(dt, dx)$  : nous noterons  $N(dt, dx) = 1$  ;
- conditionnellement à ce qui s'est passé avant l'instant  $t$ , aucun événement ne survient
  - dans l'intervalle de temps  $dt$ ,
  - sans marque,
 avec la probabilité  $1 - \Lambda(dt, E)$  : nous noterons  $N(dt, E) = 0$ .

Nous allons formuler la vraisemblance avec des conditionnalités, puis supprimer les termes conditionnels sous réserve qu'ils ne dépendent pas du paramètre d'intérêt – ce qui, finalement, nous donne une vraisemblance non pas partielle mais complète vis-à-vis de ce paramètre d'intérêt.

Soit  $\emptyset$  un point en dehors de  $E$  que nous appellerons « marque vide », et qui représente l'absence d'événement. Soit  $\bar{E} = E \cup \{\emptyset\}$ . Soit  $G$  un espace de marques tel que  $\emptyset \notin G$ , et soit  $\bar{G} = G \cup \{\emptyset\}$ .

Soit  $g : \bar{E} \rightarrow \bar{G}$  une application mesurable telle que  $g(\emptyset) = \emptyset$ . Soit  $N^g$  le processus de marques d'espace  $G$  défini par

$$N^g((0, t] \times A) = N((0, t] \times g^{-1}(A)).$$

Le compensateur  $\Lambda^g$  de  $N^g$  est défini par

$$\Lambda^g((0, t] \times A) = \Lambda((0, t] \times g^{-1}(A)).$$

La distribution conditionnelle, étant donnée  $\mathcal{F}(t-)$ , peut être construite de la manière suivante :

- 1° avec la probabilité  $\Lambda^g(dt, dy)$ , le processus réduit a un événement dans  $dt \times dy$  (*i.e.*  $N^g(dt, dy) = 1$ ) ; avec la probabilité  $1 - \Lambda^g(dt, G)$ , il n'a aucun événement, *i.e.*  $N^g(dt, G) = 0$  ;
- 2° étant donné que le processus réduit a un événement dans  $dt \times dy$ , *i.e.*  $N^g(dt, dy) = 1$ , le processus originel en a également un dans  $dt \times dx$ , *i.e.*  $N(dt, dx) = 1$  [pour  $x$  t.q.  $g(x)=y$ ], avec la probabilité  $\Lambda(dt, dx)/\Lambda^g(dt, dy)$  ;
- 3° mais étant donné que le processus n'a pas d'événement dans  $dt \times dx$ , c.-à-d. que  $N^g(dt, G) = 0$ ,
  - le processus originel a un événement dans  $dt \times dx$  (*i.e.*  $N(dt, dx) = 1$  [pour  $x$  t.q.  $g(x)=y$ ]) avec la probabilité conditionnelle  $\Lambda(dt, dx)/[1 - \Lambda^g(dt, G)]$  ;

- le processus originel n'a pas d'événement (*i.e.*  $N(dt, E) = 1$  avec la probabilité conditionnelle complémentaire  $1 - \Lambda(dt, g^{-1}(\emptyset))/(1 - \Lambda^g(dt, G))$ ).

Combinant ces possibilités dans le même ordre que celui de leur description, nous parvenons à une réécriture de (A.6) sous la forme

$$\begin{aligned}
d\mathbb{P} &= \mathcal{P}_{t \in [0, \tau]} \left( \prod_{y \in G} \Lambda^g(dt, dy)^{N^g(dt, dy)} [1 - \Lambda^g(dt, G)]^{1 - N^g(dt, G)} \right) \\
&\times \prod_{y \in G} \left( \prod_{x: g(x)=y} \left( \frac{\Lambda(dt, dx)}{\Lambda^g(dt, dy)} \right)^{N(dt, dx)} \right)^{N^g(dt, dy)} \\
&\times \left[ \prod_{x: g(x)=\emptyset} \left( \frac{\Lambda(dt, dx)}{1 - \Lambda^g(dt, G)} \right)^{N(dt, dx)} \right. \\
&\times \left. \left( 1 - \frac{\Lambda(dt, g^{-1}(\emptyset))}{1 - \Lambda^g(dt, G)} \right)^{1 - N(dt, E)} \right]^{1 - N^g(dt, G)}. \tag{A.7}
\end{aligned}$$

L'interprétation probabiliste se fait comme suit :

- 1° la première ligne de (A.7) peut être mathématiquement interprétée telle quelle : elle donne la vraisemblance partielle basée sur  $N^g$  (avec ignorance du reste de l'information basée sur  $N$ ) ; cette vraisemblance partielle a exactement la même forme que la vraisemblance basée sur  $N^g$  seul dans le cas d'un processus *self-exciting* ;
- 2° la deuxième ligne donne la contribution  $\Lambda(dt, dx)/\Lambda^g(dt, dy)$  pour  $t, x$  et  $y = g(x)$ , qui correspond aux événements (en nombre fini) communs à  $N$  et  $N^g$ . Puisque  $\Lambda^g$  est une « marginalisation » de  $\Lambda$ , on peut « désintégrer »  $\Lambda$  [restreint à  $\mathcal{F} \times E \setminus g^{-1}(\emptyset)$ ] en un produit de l'image  $\Lambda^g$  de  $\Lambda$  et d'une mesure de probabilité de transition  $\Lambda(dx | t, y)$  sur  $\{x : g(x) = y\}$  ; ainsi, on peut écrire

$$\Lambda(dt, dx) = \Lambda^g(dt, dy)\Lambda(dx | t, y),$$

dans ce sens que l'intégrale suivant  $t$  et  $x$  du membre de gauche est égale à la triple intégrale suivant  $t, y$  et  $x = g(y)$  du membre de droite. Intuitivement, pour  $y = g(x)$ , nous écrivons la probabilité d'avoir une marque dans  $dx$  dans la période  $dt$ , étant donné le passé, que multiplie la probabilité d'avoir une marque dans  $dx$  étant donnée une marque réduite  $y$  au temps  $t$ . Ainsi,  $\Lambda(dt, dx)/\Lambda^g(dt, dy)$  a une interprétation mathématique tout comme  $\Lambda(dx | t, y)$ , et le rapport de deux tels termes (de compensateurs  $\Lambda$  et  $\tilde{\Lambda}$  suivant deux mesures de probabilité différentes  $\mathbb{P}$  et  $\tilde{\mathbb{P}}$ ) peut être mathématiquement interprété comme la dérivée de Radon-Nikodym, au point  $(t, y)$ , de la mesure de probabilité de transition  $\Lambda(\cdot | t, y)$  par rapport à  $\tilde{\Lambda}(\cdot | t, y)$  sur  $\{x : g(x) = y\}$  ;

- 3° la troisième ligne n'intervient qu'un nombre fini de fois ; cela suggère que  $\Lambda(dt, dx)/(1 - \Lambda^g(dt, G))$  peut être interprété mathématiquement comme  $\left[1 - \Lambda^g(\{t\} \times G)\right]^{-1} \times \Lambda(dt, dx)$ , qui est le rapport qui peut être interprété comme la dérivée de Radon-Nikodym  $(d\Lambda/d\tilde{\Lambda})(t, x)$  que multiplie la fonction  $\left[1 - \tilde{\Lambda}^g(\{t\} \times G)\right]/\left[1 - \Lambda^g(\{t\} \times G)\right]$  ;

4° la quatrième ligne participe pour chaque  $t$  pour lequel  $N$  n'a pas d'événement dans  $dt$ . Parce que

$$\begin{aligned} 1 - \frac{\Lambda(dt, g^{-1}(\emptyset))}{1 - \Lambda^g(dt, G)} &= \frac{1 - \Lambda(dt, g^{-1}(\emptyset)) - \Lambda^g(dt, G)}{1 - \Lambda^g(dt, G)} \\ &= \frac{1 - \Lambda(dt, E)}{1 - \Lambda^g(dt, G)}, \end{aligned}$$

nous pouvons interpréter un produit, sur l'intervalle des temps, de termes tels que celui-ci comme un rapport de produits-intégrales.

Pour les additionner, nous réécrivons (A.7) sous la forme

$$\begin{aligned} d\mathbb{P} &= \mathcal{P}_{t \in [0, \tau]} \left[ \left( \prod_{y \in G} \Lambda^g(dt, dy)^{N^g(dt, dy)} [1 - \Lambda^g(dt, G)]^{1 - N^g(dt, G)} \right) \right. \\ &\quad \times \prod_{x : g(x) \neq \emptyset} \Lambda(dx | t, g(x))^{N(dt, dx)} \prod_{x : g(x) \neq \emptyset} \left( \frac{\Lambda(dt, dx)}{1 - \Lambda^g(\{t\} \times G)} \right)^{N(dt, dx)} \\ &\quad \left. \times \left( \frac{1 - \Lambda(dt, E)}{1 - \Lambda^g(dt, G)} \right)^{1 - N(dt, E)} \right], \end{aligned} \quad (\text{A.8})$$

qui peut être interprétée mathématiquement après calcul des rapports et formation des dérivées de Radon-Nikodym d'une part, des produits-intégrales d'autre part, comme suit : la première ligne est la vraisemblance partielle basée sur  $N^g$ , le produit des autres lignes forme la vraisemblance partielle basée sur le reste de  $N$ .

Un cas spécial – celui où  $E$  et  $G$  sont dénombrables et où  $\Lambda$  est absolument continue sur  $\mathcal{F} \times E$  par rapport à la mesure de Lebesgue que multiplie la mesure de comptage – est intéressant. Dans ce cas,  $N$  est un processus de comptage dans le sens usuel (c.-à-d. avec un nombre dénombrable de composantes) et  $N^g$  est une agrégation de  $N$ . Ces deux processus ont chacun une intensité – respectivement  $\lambda_x(x \in E)$  et  $\lambda_y^g(y \in G)$ , nous avons

$$\lambda_y^g(t) = \sum_{x : g(x)=y} \lambda_x(t).$$

La mesure de transition  $\Lambda(dx | t, g(x))$  est une mesure de probabilité sur l'ensemble fini  $\{x : g(x) = y\}$  et ses atomes sont simplement  $\lambda_x(t)/\lambda_y^g(t)$ .

La partie atomique  $1 - \Lambda^g(\{t\} \times g)$  disparaît et la factorisation devient

$$\begin{aligned} d\mathbb{P} &\propto \mathcal{P}_{t \in [0, \tau]} \left[ \left( \prod_{y \in G} \lambda_y^g(t)^{N_y^g(dt)} [1 - \lambda_+^g(t)]^{1 - N_+^g(dt)} \right) \right. \\ &\quad \times \prod_{x : g(x) \neq \emptyset} \left( \frac{\lambda_x(t)}{\lambda_{g(x)}^g(t)} \right)^{N_x(dt)} \prod_{x : g(x) \neq \emptyset} \lambda_x(t)^{N_x(dt)} \left( \frac{1 - \lambda_+(t)dt}{1 - \lambda_+^g(t)dt} \right)^{1 - N_+(dt)} \\ &\propto \prod_{t, y} \lambda_y^g(t)^{N_y^g(dt)} \exp \left( - \int_0^\tau \lambda_+^g(t) dt \right) \prod_{t, x : g(x) \neq \emptyset} \left( \frac{\lambda_x(t)}{\lambda_{g(x)}^g(t)} \right)^{N_x(dt)} \\ &\quad \times \prod_{t, x : g(x) \neq \emptyset} \lambda_x(t)^{N_x(dt)} \exp \left( - \int_0^\tau [\lambda_+(t) - \lambda_+^g(t)] \right). \end{aligned} \quad (\text{A.9})$$

On démontre facilement que cette quantité est réellement une factorisation de la vraisemblance totale :

$$d\mathbb{P} \propto \mathcal{P}_t \left[ \prod_x \lambda_x(t)^{N_x(dt)} [1 - \lambda_+(t)dt]^{1-N_+(dt)} \right].$$

## A.6 Résultats complémentaires

**Théorème VI.7 (Glivenko-Cantelli)** — a) Pour une suite  $(X_n)$  de vecteurs aléatoires à valeurs dans  $\mathbb{R}^k$ , indépendants, de loi  $F$ , les répartitions empiriques  $(\bar{F}_n(\omega, \cdot))$  où

$$\bar{F}_n(\omega, \cdot) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i(\omega)}$$

convergent étroitement vers  $F$  pour presque tout  $\omega$ .

b) Pour  $k = 1$ , la convergence des fonctions de répartition est p.s. uniforme : pour presque tout  $\omega$ ,

$$\sup_x |\bar{F}_n(\omega, x) - F(x)| \rightarrow 0.$$

## Annexe B

# DÉMONSTRATIONS COMPLÉMENTAIRES

### B.1 Concernant le test non-paramétrique d'égalité des fonctions de risque

La démonstration de la proposition I.6 p. 27 est issue de Andersen *et al.* (1991) (p. 360-362). Nous reprenons ici les notations de la page 26.

Notons  $\mathbf{N}^{(n)} = (N_1^{(n)}, \dots, N_k^{(n)})$  une suite de  $n$  ( $n = 1, 2, \dots$ ) processus de comptage  $k$ -variés, d'intensités  $\boldsymbol{\lambda}^{(n)} = (\lambda_1^{(n)}, \dots, \lambda_k^{(n)})$  de la forme  $\lambda_h^{(n)} = \alpha_h(t) Y_n^{(n)}(t)$ , où  $\alpha_h$  est le même pour tout  $n$ .

**Théorème VI.8** — *Supposons qu'il existe une suite de constantes positives  $(c_n)$  et de fonctions positives  $y_1, \dots, y_n, \kappa$  telles que  $\kappa^2 \alpha y_+$  (où  $y_+ = \sum_h y_h$ ) est intégrable sur  $[0, t]$ . Soit*

$$\sigma_{hj}(t) = \int_0^t \kappa^2(s) \frac{y_h(s)}{y_+(s)} \left( \delta_{hj} - \frac{y_j(s)}{y_+(s)} \right) \alpha(s) y_+(s) ds \quad (\text{B.1})$$

pour tout  $t \in [0, t]$  et tous  $h, j$ , et supposons de plus que sous l'hypothèse nulle (1.7), nous avons

(A) pour tous  $h, j$  et pour tout  $t \in [0, t]$ ,

$$c_n^2 \int_0^t [K^{(n)}(s)]^2 \frac{Y_h^{(n)}(s) Y_j^{(n)}(s)}{Y_+^{(n)}(s)} \alpha(s) ds \xrightarrow{\mathbb{P}} \int_0^t \kappa^2(s) \frac{y_h(s) y_j(s)}{y_+(s)} \alpha(s) ds \quad (n \rightarrow \infty);$$

(B) pour tout  $t \in [0, t]$  et tout  $\epsilon > 0$ ,

$$c_n^2 \int_0^t [K^{(n)}(s)]^2 \mathbf{1}_{\{|c_n K^{(n)}(s)| > \epsilon\}} Y_+^{(n)}(s) \alpha(s) ds \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty).$$

Alors, sous l'hypothèse nulle (1.7),

$$c_n(\Upsilon_1^{(n)}, \dots, \Upsilon_k^{(n)}) \xrightarrow{\mathcal{L}} (U_1, \dots, U_k) \quad (n \rightarrow \infty)$$

dans  $D([0, \tau])^k$ , où les  $U_h$  sont des martingales gaussiennes centrées avec

$$U_h(0) = 0$$

et

$$\text{Cov}[U_h(s), U_j(t)] = \sigma_{hj}(s \wedge t).$$

Par suite, sous l'hypothèse nulle, pour tous  $h, j$ ,

$$\sup_{t \in [0, \tau]} |c_n^2 \hat{\sigma}_{hj}(t) - \sigma_{hj}(t)| \xrightarrow{\mathbb{P}} 0 \quad (n \rightarrow \infty),$$

où  $\hat{\sigma}_{hj}(t)$  est défini par (1.11).

*Démonstration* — Notons en premier lieu que  $\sigma_{hj}(t)$ , donné par (B.1), est bien défini d'une part par le fait que  $\kappa^2 \alpha y$  est intégrable, et d'autre part par l'inégalité de Cauchy-Schwartz. D'après (1.9), nous pouvons réécrire

$$c_n \Upsilon_h^{(n)}(t) = \sum_{l=1}^k \int_0^t H_{hl}^{(n)}(s) dM_l^{(n)}(s),$$

où

$$H_{hl}^{(n)}(s) = c_n K^{(n)}(s) \left\{ \delta_{hl} - \frac{Y_h^{(n)}(s)}{Y_+^{(n)}(s)} \right\}$$

et le résultat de convergence monotone résulte du théorème central limite pour les martingales (th. VI.2), en utilisant (1.10) et les conditions A et B.

En appliquant le même théorème aux martingales

$$c_n W_h^{(n)} = c_n \sum_{l=1}^k \int_0^t K^{(n)}(s) \frac{Y_h^{(n)}(s)}{Y_+^{(n)}(s)} dM_l^{(n)}(s)$$

pour  $h = 1, 2, \dots, k$ , nous obtenons

$$\begin{aligned} c_n^2 [W_h^{(n)}, W_j^{(n)}] &= c_n^2 \int_0^t [K^{(n)}(s)]^2 \frac{Y_h^{(n)}(s) Y_j^{(n)}(s)}{Y_+^{(n)2}(s)} dN_+^{(n)}(s) \\ &\xrightarrow{\mathbb{P}} \int_0^t \kappa^2(s) \frac{y_h(s) y_j(s)}{y_+(s)} \alpha(s) d(s) \end{aligned} \quad (\text{B.2})$$

uniformément sur  $[0, \tau]$  quand  $n \rightarrow \infty$  pour tous  $h, j$ .

L'uniforme consistance de  $c_n^2 \hat{\sigma}_{hj}$  dérive alors de (1.11) et (B.1). ■

Concernant la vérification des conditions (A) et (B), nous pouvons nous aider de la proposition suivante.

**Proposition VI.9** — *Supposons que :*

(i) *il existe une suite croissante de constantes  $(a_n)$  tendant vers  $+\infty$  lorsque  $n \rightarrow \infty$ , ainsi qu'une suite de constantes  $(d_n)$ , telles que pour tout  $h$  et pour presque tout  $t \in [0, \tau]$ , nous ayons, sous l'hypothèse nulle,*

$$\frac{Y_h^{(n)}(s)}{a_n^2} \xrightarrow{\mathbb{P}} y_h(s) < \infty \quad (\text{B.3})$$

et

$$\frac{K^{(n)}(s)}{d_n} \xrightarrow{\mathbb{P}} \kappa(s) < \infty \quad (\text{B.4})$$

quand  $n \rightarrow \infty$  ;

(ii) *pour tout  $\delta > 0$ , il existe une fonction  $k_\delta$  telle que*

$$\liminf_{n \rightarrow \infty} \mathbb{P} \left( \left[ \frac{K^{(n)}(s)}{(a_n d_n)} \right]^2 Y_+^{(n)}(s) \leq k_\delta(s), \forall s \in [0, \tau] \right) \geq 1 - \delta \quad (\text{B.5})$$

*sous l'hypothèse nulle.*

Alors, à condition que  $\alpha k_\delta$  soit intégrable sur  $[0, \tau]$  pour tout  $\delta > 0$ , les conditions (A) et (B) du théorème VI.8 sont vérifiées, avec  $c_n = (a_n d_n)^{-1}$ .

*Démonstration* — Les quantités à intégrer des conditions (A) et (B) sont bornées : elles sont majorées par

$$c_n^2 [K^{(n)}(s)]^2 \alpha(s) Y_+(s).$$

Par suite, la proposition VI.8 est applicable, et le résultat annoncé est obtenu. ■

En notant  $\Upsilon^{(n)} = (\Upsilon_1^{(n)}(t), \dots, \Upsilon_k^{(n)}(t))^t$ , nous avons pour tout  $t \in [0, \tau]$  que

$$c_n \Upsilon^{(n)}(t) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma(t)) \quad (\text{B.6})$$

quand  $n \rightarrow \infty$ , sous les conditions du théorème VI.8, où  $\Sigma(t)$  est la matrice singulière  $k \times k$  dont les éléments  $\sigma_{hj}(t)$  sont définis par (B.1).

$\Sigma(t)$  est de rang  $k - 1$  (c.-à-d. de rang maximum) sous certaines conditions, telles que l'existence pour chaque couple  $(h, j)$  de valeurs dans  $\{1, 2, \dots, k\}$  d'un ensemble  $B_{hj} \subset [0, t]$ , de mesure de Lebesgue positive et tel que

$$\kappa^2(s) \frac{y_h(s) y_j(s)}{y_+(s)} \alpha(s) > 0$$

pour  $s \in B_{hj}$  (cf. Gill (1986) p. 114).

Par suite, d'après (B.2),  $c_n^2 \widehat{\Sigma}(t)$  est un estimateur consistant de  $\Sigma(t)$ , et de plus, la probabilité que  $c_n^2 \widehat{\Sigma}(t)$  soit de rang  $k - 1$  tend vers 1 lorsque  $n$  tend vers l'infini.

Sous l'hypothèse nulle (1.7), la statistique donnée par (1.12) suit asymptotiquement un chi-deux à  $k - 1$  degrés de liberté.

## B.2 Concernant le modèle marginal de Cox

Nous exposons ci-dessous la démonstration du résultat de la page 53, telle qu'elle est énoncée par Cai et Prentice (1997).

Soient les notations suivantes :

$$\begin{aligned} \mathbf{S}_k^{(d)}(\boldsymbol{\beta}, t) &= \frac{1}{n} \sum_{i=1}^n Y_{ik}(t) \mathbf{Z}_{ik}^d \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ik}) \quad d = 0, 1, \\ \mathbf{S}_k^{(d)}(\boldsymbol{\beta}, t) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \mathbf{Z}_{ik} w_{ijk}(\boldsymbol{\beta}, t) Y_{ik}(t) \mathbf{Z}_{ik}^{d-2} \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ik}) \quad d = 2, 3, \\ \mathbf{S}_k^{(4)}(\boldsymbol{\beta}, t) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \mathbf{Z}_{ij} Y_{ik}(t) \frac{\partial w_{ijk}(\boldsymbol{\beta}, t)}{\partial \boldsymbol{\beta}^t} \exp(\boldsymbol{\beta}^t \mathbf{Z}_{ik}), \\ \mathbf{S}^{(d)}(\boldsymbol{\beta}, t) &= \sum_{k=1}^{n_i} \mathbf{S}_k^{(d)}(\boldsymbol{\beta}, t) \quad d = 0, 1, 2, 3, \\ \mathbf{S}^{(4)}(\boldsymbol{\beta}, t) &= \sum_{k=1}^{n_i} \mathbf{S}_k^{(4)}(\boldsymbol{\beta}, t), \\ E(\boldsymbol{\beta}, t) &= \frac{\mathbf{S}^{(2)}(\boldsymbol{\beta}, t)}{\mathbf{S}^{(0)}(\boldsymbol{\beta}, t)} \end{aligned}$$

et

$$V(\boldsymbol{\beta}, t) = \frac{\mathbf{S}^{(3)}(\boldsymbol{\beta}, t)}{\mathbf{S}^{(0)}(\boldsymbol{\beta}, t)} - \frac{\mathbf{S}^{(2)}(\boldsymbol{\beta}, t) \mathbf{S}^{(1)}(\boldsymbol{\beta}, t)^t}{[\mathbf{S}^{(0)}(\boldsymbol{\beta}, t)]^2}.$$

Notons  $\boldsymbol{\beta}_0$  la vraie valeur du paramètre  $\boldsymbol{\beta}$ .

Pour tout  $k = 1, \dots, n_i$ , on suppose que :

A.  $\int_0^\tau \lambda_0(t) dt < \infty$  ;

B. il existe un voisinage  $\mathfrak{B}$  de  $\boldsymbol{\beta}_0$  qui satisfait les trois conditions suivantes :

- (i) il existe une fonction (scalaire, vectorielle ou matricielle, selon l'exposant entre parenthèses)  $s_k^{(d)}$ ,  $d = 0, \dots, 4$ , définie sur  $\mathfrak{B} \times [0, \tau]$  telle que

$$\sup_{t, \boldsymbol{\beta}} \left\| \mathbf{S}_k^{(d)}(\boldsymbol{\beta}, t) - s_k^{(d)}(\boldsymbol{\beta}, t) \right\| \xrightarrow{\mathbb{P}} 0$$

$$\forall k = 1, \dots, n_i,$$

- (ii)

$$\frac{1}{\sqrt{n}} \sup_{t, \boldsymbol{\beta}, i, j, k} \left\| \mathbf{Z}_{ij} \frac{\partial w_{ijk}(\boldsymbol{\beta}, t)}{\partial \boldsymbol{\beta}^t} \right\| \xrightarrow{\mathbb{P}} 0,$$

(iii) il existe une matrice  $\Sigma_w = \Sigma_w(\beta_0)$  telle que

$$\frac{1}{n} \sum_{i=1}^n \mathbb{V}(D_i) \longrightarrow \Sigma_w ,$$

où

$$D_i = \sum_{j=1}^{n_i} \int_0^\tau \left\{ \sum_{k=1}^{n_i} \mathbf{Z}_{i\mathbf{k}} w_{ikj}(\beta_0, t) - e(\beta_0, t) \right\} M_{ij}(du)$$

et

$$e(\beta, u) = \frac{\sum_j s_j^{(2)}(\beta, u)}{\sum_j s_j^{(0)}(\beta, u)} ;$$

C. soit

$$v(\beta, u) = \frac{\sum_j s_j^{(3)}(\beta, u)}{\sum_j s_j^{(0)}(\beta, u)} - \frac{[\sum_j s_j^{(2)}(\beta, u)] [\sum_j s_j^{(1)}(\beta, u)]}{[\sum_j s_j^{(0)}]^2} ,$$

alors, pour tout  $\beta \in \mathfrak{B}$ ,  $t \in [0, \tau]$  et  $j = 1, \dots, n_i$ ,

$$s_j^{(1)}(\beta, t) = \frac{\partial s_j^{(0)}(\beta, t)}{\partial \beta}$$

et

$$s_j^{(3)}(\beta, t) = \frac{\partial s_j^{(2)}(\beta, t)}{\partial \beta^t} - s_j^{(4)}(\beta, t) ;$$

de plus,  $s_k^{(0)}$  est bornée autour de 0 (*bounded away from 0*) sur  $\mathfrak{B} \times [0, 1]$ , et la matrice

$$\int_0^\tau v(\beta_0, t) \left[ \sum_j s_j^{(0)}(\beta, t) \right] \lambda_0(t) dt$$

est définie positive ;

enfin,  $\partial w_{ijk}(\beta, t)/\partial \beta$  existe pour tout  $t, \beta, i, j, k$ , et est une fonction continue de  $\beta \in \mathfrak{B}$  ;

D. la condition de Lindeberg est vérifiée :

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ \|D_i\|^2 \mathbf{1}_{\|D_i\| > \epsilon \sqrt{n}} \right] \xrightarrow{\mathbb{P}} 0 ;$$

E. concernant la matrice de pondération :

(i)  $\widehat{W}_i(\beta, t) - W_i(\beta, t) \xrightarrow{\mathbb{P}} 0$  uniformément en  $t \in [0, 1]$  et  $\beta \in \mathfrak{B}$  ,

(ii)  $\partial \widehat{w}_{ijk}(\beta, t)/\partial \beta$  est une fonction continue de  $\beta$  pour tout  $(i, j, k)$ ,  $t \in [0, 1]$ ,  $\beta \in \mathfrak{B}$ ,

- (iii)  $\partial \widehat{w}_{ijk}(\boldsymbol{\beta}, t) / \partial \boldsymbol{\beta} - \partial w_{ijk}(\boldsymbol{\beta}, t) / \partial \boldsymbol{\beta} \xrightarrow{\mathbb{P}} 0$  uniformément en  $t \in [0, 1]$  et  $\boldsymbol{\beta} \in \mathfrak{B}$ , et pour tout  $(i, j, k)$ .

Notons que les énoncés sont relatifs à la filtration étendue qui, concernant le  $j^{\text{e}}$  élément du  $i^{\text{e}}$  groupe, inclut toute l'information portant sur les autres éléments du vecteur réponse. Les composantes de la condition E constituent des conditions de régularité et de stabilité concernant les poids estimés.

**Théorème VI.9** — *Sous les conditions A-D,  $\widehat{\boldsymbol{\beta}}_w$  vérifiant (3.1) est un estimateur consistant de  $\boldsymbol{\beta}_0$ .*

*$\sqrt{n}(\widehat{\boldsymbol{\beta}}_w - \boldsymbol{\beta}_0)$  est asymptotiquement normal, de moyenne nulle et de matrice de variance-covariance*

$$\Sigma = A_w^{-1}(\boldsymbol{\beta}_0) \Sigma_w A_w^{-1}(\boldsymbol{\beta}_0),$$

où

$$A_w(\boldsymbol{\beta}) = \int_0^\tau v(\boldsymbol{\beta}, u) \left[ \sum_j s_j^{(0)}(\boldsymbol{\beta}, u) \right] \lambda_0(u) du.$$

**Théorème VI.10** — *Sous les conditions A-E,  $\widehat{\boldsymbol{\beta}}_e$  vérifiant*

$$\sum_{i=1}^n \int_0^\tau \mathbf{Z}_i^t \widehat{\mathbf{W}}_i(\boldsymbol{\beta}, t) \widehat{\mathbf{M}}_i(dt) = 0 \tag{B.7}$$

*est un estimateur consistant de  $\boldsymbol{\beta}_0$  et  $\sqrt{n}(\widehat{\boldsymbol{\beta}}_e - \boldsymbol{\beta}_0)$  a la même distribution asymptotique que  $\sqrt{n}(\widehat{\boldsymbol{\beta}}_w - \boldsymbol{\beta}_0)$ .*

*Démonstration* — On note  $R_n(\boldsymbol{\beta})$  le membre de gauche de (3.1) que multiplie  $1/n$ .

$\widehat{\boldsymbol{\beta}}_w$  vérifie  $R_n(\widehat{\boldsymbol{\beta}}_w) = 0$ .

Basée sur une extension de Foutz (Foutz, 1977), on peut montrer que  $\widehat{\boldsymbol{\beta}}_w$  est consistant vis-à-vis de  $\boldsymbol{\beta}$  à condition que :

- (i)  $\partial R_n(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}$  existe et soit continue dans un voisinage ouvert  $\mathcal{B}$  de  $\boldsymbol{\beta}_0$  ;
- (ii)  $\partial R_n(\boldsymbol{\beta}_0) / \partial \boldsymbol{\beta}_0^t$  soit définie négative avec une probabilité 1 à mesure que  $n \rightarrow \infty$  ;
- (iii)  $\partial R_n(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}^t$  converge en probabilité vers une fonction fixe,  $\sigma(\boldsymbol{\beta})$ , uniformément dans un voisinage ouvert de  $\boldsymbol{\beta}_0$  ;
- (iv)  $R_n(\boldsymbol{\beta})$  tende vers 0 en probabilité.

On peut écrire

$$\frac{\partial R_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^t} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau H_{ij}(\boldsymbol{\beta}, u) M_{ij}(du) \quad (\text{B.8a})$$

$$- \frac{1}{n} \sum_{l=1}^n \sum_{j=1}^{n_l} \int_0^\tau \left[ \frac{\mathbf{S}^{(4)}(\boldsymbol{\beta}, u)}{\mathbf{S}^{(0)}(\boldsymbol{\beta}, u)} \right] M_{lj}(du) \quad (\text{B.8b})$$

$$+ \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \left[ \frac{\mathbf{S}^{(2)}(\boldsymbol{\beta}, u) \mathbf{S}^{(1)}(\boldsymbol{\beta}, u)^t}{[\mathbf{S}^{(0)}(\boldsymbol{\beta}, u)]^2} \right. \quad (\text{B.8c})$$

$$\left. - \frac{\mathbf{S}^{(3)}(\boldsymbol{\beta}, u)}{\mathbf{S}^{(0)}(\boldsymbol{\beta}, u)} \right] N_{ij}(du), \quad (\text{B.8d})$$

où

$$H_{ij}(\boldsymbol{\beta}, t) = \sum_{k=1}^{n_i} \mathbf{Z}_{ik} \frac{\partial w_{ikj}(\boldsymbol{\beta}, t)}{\partial \boldsymbol{\beta}^t}.$$

(B.8b) est asymptotiquement équivalent à

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau \frac{\sum_l s_l^{(4)}(\boldsymbol{\beta}, t)}{\sum_l s_l^{(4)}(\boldsymbol{\beta}, t)} M_{ij}(dt),$$

qui est une martingale locale de carré intégrable. Aussi, d'après l'inégalité de Lengart, on a que, pour tout  $\Delta$ ,  $\mu > 0$ , il existe  $n_0$  tel que pour  $n \geq n_0$ ,

$$\mathbb{P} \left\{ \left| \int_0^\tau \left[ \frac{\sum_k \mathbf{s}_k^{(4)}(\boldsymbol{\beta}, u)}{\sum_k \mathbf{s}_k^{(0)}(\boldsymbol{\beta}, u)} \right]_{ll'} \frac{1}{n} \sum_{i=1}^n M_{ij}(du) \right| > \Delta \right\} \leq \frac{\mu}{\Delta^2} \\ + \mathbb{P} \left\{ \frac{1}{n} \int_0^\tau \left[ \frac{\sum_k \mathbf{s}_k^{(4)}(\boldsymbol{\beta}, u)}{\sum_k \mathbf{s}_k^{(0)}(\boldsymbol{\beta}, u)} \right]_{ll'}^2 \mathbf{S}_j^{(0)}(\boldsymbol{\beta}, u) \lambda_0(u) du > \mu \right\},$$

où l'indice  $ll'$  indique l'élément  $(l, l')$  de la matrice en question.

Sur la base des conditions A, B(i) et C, le second terme du membre de droite de cette dernière expression tend vers 0 en probabilité, uniformément en  $\boldsymbol{\beta} \in \mathcal{B}$ , et par conséquent (B.8b) converge de la même manière.

Un argument similaire permet de montrer que (B.8a) tend également vers 0 en probabilité, uniformément en  $\boldsymbol{\beta} \in \mathcal{B}$ .

Par ailleurs, en utilisant les conditions A, B et C, on peut montrer que

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{n_i} \int_0^\tau V(\boldsymbol{\beta}, u) N_{ij}(du) \xrightarrow{\mathbb{P}} \int_0^\tau v(\boldsymbol{\beta}, u) \left[ \sum_k s_k^{(0)}(\boldsymbol{\beta}, u) \right] \lambda_0(u) du$$

uniformément en  $\boldsymbol{\beta} \in \mathcal{B}$ .

En conséquence, (i) est satisfait sur la base de (B.8) et de la continuité des dérivées partielles des poids (condition C). De même, (ii) et (iii) sont satisfaites, puisque

$$\frac{\partial R_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^t} \xrightarrow{\mathbb{P}} - \int_0^\tau v(\boldsymbol{\beta}, u) \left[ \sum_k s_k^{(0)}(\boldsymbol{\beta}, u) \right] \lambda_0(u) du \\ = -A_w(\boldsymbol{\beta}) \quad (\text{B.9})$$

uniformément en  $\beta \in \mathcal{B}$ .

Pour (iv), on peut écrire

$$R_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{n_i} \int_0^\tau \left[ \sum_{k=1}^{n_i} \mathbf{Z}_{ij} w_{ijk}(\beta, u) - E(\beta, u) \right] M_{ik}(du). \quad (\text{B.10})$$

En utilisant l'inégalité de Lengart, le membre de droite de (B.10) est asymptotiquement équivalent à

$$\frac{1}{n} \sum_{i=1}^n D_i, \quad (\text{B.11})$$

où  $D_i$  est défini dans la condition B. Puisque chaque  $\mathbf{Z}_{ij}$  est prévisible, l'espérance de (B.11) est nulle, et sa variance converge vers 0 quand  $n$  tend vers l'infini, comme l'indique la condition B(iii).

Par suite, l'inégalité de Chebichev implique que  $R_n(\beta)$  converge vers 0 en probabilité. Il s'ensuit qu'il existe  $\beta_w$  tel que  $R_n(\beta_w) = 0$  avec la probabilité de converger quand  $n$  tend vers l'infini égale à 1. De plus,  $\beta_w$  converge en probabilité vers  $\beta_0$ .

Considérons maintenant la convergence asymptotique de  $\beta_w$ . L'inégalité de Lengart permet d'affirmer que  $\sqrt{n}R_n(\beta_0)$  est asymptotiquement équivalent à  $\sqrt{n} \sum_{i=1}^n D_i$ , qui est une somme de  $n$  v.a. indépendantes (mais non identiquement distribuées) de moyennes nulles et de matrices de variance  $\mathbb{V}(D_i)$ . D'après les conditions B(iii) et D, et selon le théorème central limite multivarié, il vient que la distribution de  $\sqrt{n}R_n(\beta_0)$  converge vers une distribution normale centrée de variance  $\Sigma_w$ .

De plus, compte tenu de (B.8) et de la condition C, nous pouvons dire que

$$\frac{1}{n} \frac{\partial R_n(\beta^*)}{\partial \beta^*} \longrightarrow A_w(\beta^*)$$

uniformément en  $\beta \in \mathcal{B}$ , si bien que

$$\frac{1}{n} \frac{\partial R_n(\beta^*)}{\partial \beta^*} \longrightarrow A_w(\beta_0),$$

à condition que  $\beta^*$  soit consistant pour  $\beta_0$ .

Par un développement de Taylor de  $R_n(\beta)$  autour de  $\hat{\beta}_w$ , et d'après la consistance de  $\hat{\beta}_w$ , nous avons que  $\sqrt{n}(\hat{\beta} - \beta_0)$  est asymptotiquement normal, centré et de variance

$$A_w^{-1}(\beta_0) \Sigma_w(\beta_0) A_w^{-1}(\beta_0)$$

comme annoncé. ■

La preuve du second théorème est similaire à la précédente : elle repose sur la réécriture du

membre de gauche de (B.7), des dérivées ainsi que des éléments de  $\widehat{\Sigma}$  sous la forme

$$\begin{aligned} & \frac{1}{\sqrt{n}} \sum_{i=1}^n \int_0^t H_i(\boldsymbol{\beta}, u) \widehat{w}_{ijk}(\boldsymbol{\beta}, u) M_{ij}(du) \\ & + \frac{1}{\sqrt{n}} \sum_{i=1}^n \int_0^t H_i(\boldsymbol{\beta}, u) [\widehat{w}_{ijk}(\boldsymbol{\beta}, u) - w_{ijk}(\boldsymbol{\beta}, u)] M_{ij}(du) , \end{aligned} \quad (\text{B.12})$$

où  $H_i$  est un processus prévisible.

Le premier terme de (B.12) est le même que celui qui intervient dans le cas d'une pondération connue (comme cela a été vu plus haut), tandis que le second peut s'écrire

$$\begin{aligned} & \frac{1}{\sqrt{n}} \sum_{i=1}^n \int_0^t H_i(\boldsymbol{\beta}, u) [\widehat{w}_{ijk}(\boldsymbol{\beta}, u) - w_{ijk}(\boldsymbol{\beta}, u)] \mathbb{E}_{\mathcal{F}_i(u)|\widetilde{\mathcal{F}}_i(u)}[\widetilde{M}_{ij}(du)] \\ = & \mathbb{E}_{\mathcal{F}_i(u)|\widetilde{\mathcal{F}}_i(u)} \left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n \int_0^t H_i(\boldsymbol{\beta}, u) [\widehat{w}_{ijk}(\boldsymbol{\beta}, u) - w_{ijk}(\boldsymbol{\beta}, u)] \widetilde{M}_{ij}(du) \right] , \end{aligned} \quad (\text{B.13})$$

où  $\mathcal{F}_i(u)$  est la filtration naturelle et  $\widetilde{\mathcal{F}}_i(u)$  est la filtration plus large qui inclut toute l'information concernant le processus de comptage avant l'instant  $u$ , et où  $\widetilde{M}_{ij}$  est la martingale associée au processus  $N_{ij}$ , relative à la plus large des deux filtrations – son existence est assurée par le théorème de l'innovation.

(B.13) est l'espérance de la somme d'intégrales stochastiques de processus prédictibles suivant des martingales. Par conséquent, (B.13) est l'espérance d'une martingale, et il est possible de montrer que le processus prévisible de variation converge vers 0 quand  $n \rightarrow \infty$  sous la condition E.

Par suite, (B.13) converge vers 0 en probabilité, et seuls les termes semblables au premier de l'équation (B.12) contribueront à la distribution asymptotique de  $\widehat{\beta}_e$  et à la consistance de  $\widehat{\Sigma}$  en tant qu'estimateur de  $\Sigma$ .



## Annexe C

# MATRICE D'INFORMATION DE FISHER DANS LE MODÈLE DE FRAGILITÉ GAMMA DE COX

Concernant les notations, nous rappelons celles de la page 67, auxquelles nous en ajoutons quelques-unes :

$$\begin{aligned}\hat{\alpha}_0(t) &= \frac{\sum_{i,j} dN_{ij}(t)}{\sum_{i,j} \hat{\theta}_i \exp(\beta' \mathbf{Z}_{ij})}, \\ \hat{\alpha}_{0k}(t) &= \frac{\partial \hat{\alpha}_0(t)}{\partial \beta_k}, \\ \hat{\alpha}_{0kl}(t) &= \frac{\partial^2 \hat{\alpha}_0(t)}{\partial \beta_k \partial \beta_l}, \\ \hat{A}_0(t) &= \int_0^t \hat{\alpha}_0(u) du, \\ \hat{A}_{0k}(t) &= \int_0^t \hat{\alpha}_{0k}(u) du\end{aligned}$$

et

$$\hat{A}_{0kl}(t) = \int_0^t \hat{\alpha}_{0kl}(u) du,$$

pour  $k = 1, \dots, p$ .

Reprenant l'expression de la vraisemblance observable (4.6), nous obtenons que la log-

vraisemblance vaut

$$\begin{aligned}
 & \log \mathcal{L}(\boldsymbol{\beta}, \gamma, \alpha_0) \\
 &= \sum_{i=1}^n \left( \log \left[ \Gamma \left( D_i + \frac{1}{\gamma} \right) \right] + \frac{1}{\gamma} \log \left( \frac{1}{\gamma} \right) - \log \left[ \Gamma \left( \frac{1}{\gamma} \right) \right] \right. \\
 & \quad \left. - \left( D_i + \frac{1}{\gamma} \right) \log \left( \sum_{j=1}^{n_i} \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t) + \frac{1}{\gamma} \right) + \int_0^\tau \sum_{j=1}^{n_i} \left[ Y_{ij}(t) \left\{ \log(\alpha_0(t)) + \boldsymbol{\beta}' Z_{ij} \right\} \right] dt \right) \\
 &= \sum_i \left( \log \left[ \Gamma \left( D_i + \frac{1}{\gamma} \right) \right] + D_i \log \left( \frac{1}{\gamma} \right) - \log \left[ \Gamma \left( \frac{1}{\gamma} \right) \right] - \left( D_i + \frac{1}{\gamma} \right) \log \left( \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t) + 1 \right) \right. \\
 & \quad \left. + \int_0^\tau \sum_j \left[ \delta_{ij} \left\{ \log(\alpha_0(t)) + \boldsymbol{\beta}' Z_{ij} \right\} \right] dt \right). \tag{C.1}
 \end{aligned}$$

En substituant (4.10) (cf. p. 68) dans (C.1), nous obtenons les dérivées suivantes :

$$\begin{aligned}
 & \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\beta}, \gamma)}{\partial \gamma^2} \\
 &= \sum_i \left\{ \frac{D_i}{\gamma^2} + \frac{2\psi(\frac{1}{\gamma})}{\gamma^3} + \frac{\psi'(\frac{1}{\gamma})}{\gamma^4} - \frac{2\psi(\frac{1}{\gamma} + D_i)}{\gamma^3} - \frac{\psi'(\frac{1}{\gamma} + D_i)}{\gamma^4} - \frac{2 \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)}{\gamma^2 [1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)]} \right. \\
 & \quad \left. - \frac{(\frac{1}{\gamma} + D_i) [\sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)]^2}{[1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)]^2} + \left( \frac{2}{\gamma^3} \right) \log \left[ \gamma \left( 1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t) \right) \right] \right\},
 \end{aligned}$$

$$\begin{aligned}
 & \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\beta}, \gamma)}{\partial \gamma \partial \beta_k} \\
 &= \sum_i \left\{ \frac{D_i \sum_j \int_0^\tau Y_{ij}(t) \exp(\boldsymbol{\beta} Z_{ij}) [Z_{ijk} \hat{A}_0(t) + \hat{A}_{0k}(t)] dt}{1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)} \right. \\
 & \quad \left. - \frac{[1 + \gamma D_i] [\sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)] \sum_j \int_0^\tau Y_{ij}(t) \exp(\boldsymbol{\beta} Z_{ij}) [Z_{ijk} \hat{A}_0(t) + \hat{A}_{0k}(t)] dt}{[1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)]^2} \right\},
 \end{aligned}$$

$$\begin{aligned}
 & \frac{\partial^2 \log \mathcal{L}(\boldsymbol{\beta}, \gamma)}{\partial \beta_k \partial \beta_l} \\
 &= \sum_i \left[ \frac{(1 + \gamma D_i) \sum_j \int_0^\tau Y_{ij}(t) \exp(\boldsymbol{\beta}' Z_{ij}) [Z_{ijk} Z_{ijl} \hat{A}_0(t) + Z_{ijk} \hat{A}_{0k}(t) + Z_{ijl} \hat{A}_{0l}(t) + \hat{A}_{0kl}(t)] dt}{1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)} \right. \\
 & \quad \left. - (\gamma + \gamma^2 D_i) \left\{ \frac{\sum_j \int_0^\tau Y_{ij}(t) \exp(\boldsymbol{\beta} Z_{ij}) [Z_{ijk} \hat{A}_0(t) + \hat{A}_{0k}(t)] dt}{1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)} \right\} \right. \\
 & \quad \left. \times \left\{ \frac{\sum_j \int_0^\tau Y_{ij}(t) \exp(\boldsymbol{\beta} Z_{ij}) [Z_{ijl} \hat{A}_0(t) + \hat{A}_{0l}(t)] dt}{1 + \gamma \sum_j \int_0^\tau Y_{ij}(t) d\hat{A}_{ij}(t)} \right\} + \sum_j \delta_{ij} \frac{\hat{\alpha}_{0k}(t) \hat{\alpha}_{0l}(t) - \hat{\alpha}_{0kl}(t) \hat{\alpha}_0(t)}{\hat{\alpha}_0^2(t)} \right].
 \end{aligned}$$

# Annexe D

## PROCÉDURE DES SIMULATIONS SOUS S-PLUS

Les lignes débutant par un dièse sont des commentaires.

```
1 "simulations"<-
2 function(x)
3 { A <- matrix(NULL,ncol=31)
4   for (i in 1:x)
5     #boucle de 1 à x, où "x" est le nombre de simulations (x=500)
6     {
7       M <- cbind(matrix(rep(1:100,rep(10,100)),ncol=1),matrix(nrow=1000,ncol=4))
8       #matrice M dont la première colonne est la variable "groupe"
9       #ici, 100 groupes de 10 individus chacun
10      M[,5] <- rep(rgamma(20,0.25,0.25),rep(50,20))
11      #dans la cinquième colonne de M, la variable "fragilité" (esp.=1, var.=4)
12      M[,2] <- rbinom(1000,1,0.5)
13      #dans la deuxième colonne de M, génération de la covariable binaire
14      M[,3] <- rweibull(1000,1.1,1/(0.01 *M[,5]* exp((0.3 * M[,2] )/1.1)))
15      #dans la troisième colonne de M, génération du temps de survie suivant une Weibull
16      M <- M[rev(order(M[,3])),1:5]
17      #tri
18      M[1:300,4] <- 0
19      M[301:1000,4] <- 1
20      #dans la quatrième colonne de M, variable de censure (au taux de 30 % ici)
21      dimnames(M) <- list(NULL, c("groupe","cov", "temps","cens","fragil"))
22      #attribution des noms
23      M <- data.frame(M)
24      #création d'un "data.frame"
25      nule <- 0
26      #hypothèse de base pour le test de Wald
27      try({
28        #fonction permettant de passer à la simulation suivante si erreur
29        fit1 <- coxph(Surv(temps,cens) ~ cov,data= M)
30        #ajustement du modèle de Cox naïf
31        fit2 <- coxph(Surv(temps,cens) ~ cov + groupe,data= M)
32        #ajustement du modèle de Cox avec "effet groupe" fixe
33        fit3 <- coxph(Surv(temps,cens) ~ cov + cluster(groupe),data= M)
34        #ajustement du modèle marginal de Cox
35        fit4 <- coxph(Surv(temps,cens) ~ cov + frailty(groupe),data= M)
36        #ajustement du modèle de fragilité (gamma) de Cox
37        fit5 <- coxph(Surv(temps,cens) ~ cov + frailty(groupe,distrib="gauss"),data= M)
38        #ajustement du modèle de fragilité (gaussienne) de Cox
39      })
}
```

## ANNEXE D. PROCÉDURE DES SIMULATIONS SOUS S-PLUS

```

40   A <- rbind(A,c(fit1$coef,sqrt(fit1$var),1-pchisq(((fit1$coef - nule)/ ((fit1$var[1,1])^(0.5)))^2,1),
   fit1$loglik,afit2$coef[1],sqrt(afit2$var[1,1]),afit2$coef[2],sqrt(afit2$var[2,2]),
   afit2$loglik,1-pchisq((afit2$coef[1]/((afit2$var[1,1])^(0.5)))^2,1),
   1-pchisq((afit2$coef[2]/((afit2$var[2,2])^(0.5)))^2,1),sqrt(fit3$var),
   1-pchisq(((fit3$coef - nule)/ ((fit3$var[1,1])^(0.5)))^2,1),
   fit3$loglik,fit4$coef,sqrt(fit4$var),1-pchisq(((fit4$coef - nule)/ ((fit4$var[1,1])^(0.5)))^2,1),
   mean(fit4$frail),mean(fit4$history[[1]]$theta),fit4$loglik,fit5$coef,sqrt(fit5$var),
   1-pchisq(((fit5$coef - nule)/ ((fit5$var[1,1])^(0.5)))^2,1),mean(fit5$frail),
   mean(fit5$history[[1]]$theta),fit5$loglik),1-pchisq(-2*(afit1$loglik[2]-afit2$loglik[2]),19),
   2*(afit1$loglik[2]-afit1$loglik[1]-sum(afit1$df)),2*(afit2$loglik[2]-afit2$loglik[1]-sum(afit2$df)),
   2*(afit3$loglik[2]-afit3$loglik[1]-sum(afit3$df)),2*(afit4$loglik[2]-afit4$loglik[1]-sum(afit4$df)),
   2*(afit5$loglik[2]-afit5$loglik[1]-sum(afit5$df)))
41   #récupération des valeurs des différentes estimations
42   }
43   A
44   AA <- A[,c(1,2,6,7,11,15,16,22,23,18,19,25,26,4,5,8,9,13,14,20,21,27,28,29,30,31,32,33,34)]
45   MOYENNE <- apply(AA,2,mean)
46   #moyenne des différentes estimations (contenues dans AA) sur l'ensemble des simulations
47   names(MOYENNE) <- c("coeff", "var naiv","coeff 1","var coeff 1","coeff 2","var coeff 2","var rob (marginal)"
   ,"coeff rob (fragil gamma)","var rob (fragil gamma)","coeff rob (fragil gauss)", "var rob (fragil gauss)",
   "mean fragil (gamma)","var fragil (gamma)","mean fragil (gauss)","var fragil (gauss)","likeli (naive) 1",
   "likeli (naive) 2","likeli (fixe) 1","likeli (fixe) 2","likeli (marginal) 1","likeli (marginal) 2",
   "likeli (fragil gamma) 1","likeli (fragil gamma) 2","likeli (fragil gauss) 1","likeli (fragil gauss) 2",
   "test likeli fixe","aic naif","aic fixe","aic marg","aic frail gam","aic frail gaus")
48   MOY <- data.frame(MOYENNE)
49   VARIANCE <- apply(AA,2,var)
50   #variance des différentes estimations (contenues dans AA) sur l'ensemble des simulations
51   E <- cbind(MOY,VARIANCE)
52   F <- data.frame(A[,c(19,26)])
53   INTERVAL <- matrix(cbind(F[,1] - (1.96)*sqrt(E[13,2]),F[,1] + (1.96)*sqrt(E[13,2]),
   F[,2] - (1.96)*sqrt(E[15,2]),F[,2] + (1.96)*sqrt(E[15,2])),ncol=4)
54   INTERVAL2 <- apply(INTERVAL,2,mean)
55   #intervalles nécessaires au calcul du taux de recouvrement
56   TRGAMMA <- (apply(rbind(INTERVAL[INTERVAL[,1]<= teta & INTERVAL[,2]>= teta,],c(0,0,0,0)),2,length) -1)/x
57   #taux de recouvrement pour la variance de la fragilité gamma
58   TRGAUSS <- (apply(rbind(INTERVAL[INTERVAL[,3]<= teta & INTERVAL[,4]>= teta,],c(0,0,0,0)),2,length) -1)/x
59   #taux de recouvrement pour la variance de la fragilité gaussienne
60   CI <- apply(INTERVAL,2,function(y) length(y[y<0.05]))/x
61   AAA <- A[,c(3,10,12,17,24)]
62   INTCINQ <- apply(AAA,2,function(y) length(y[y<0.05])/length(y))
63   #puissance des tests statistiques au seuil de 5 %
64   names(INTCINQ) <- c('<0.05 : wald naif','wald coeff 1','wald coeff 2','wald marg','wald frail ga','wald frail no')
65   TESTCINQ <- data.frame(INTCINQ)
66   TESTUN <- apply(AAA,2,function(y) length(y[y<0.01])/length(y))
67   #puissance des tests statistiques au seuil de 1 %
68   EE <- cbind(TESTCINQ,TESTUN)
69   print(E)
70   print(EE)
71   names(INTERVAL2) <- c("moy gamma rec1","moy gamma rec2","moy gauss rec1","moy gauss rec2")
72   names(TRGAMMA) <- c("dedans (gamma)","dedans (gamma)","dedans (gamma)","dedans (gamma)")
73   names(TRGAUSS) <- c("dedans (gaus)","dedans (gaus)","dedans (gaus)","dedans (gaus)")
74   print(INTERVAL2)
75   print(TRGAMMA)
76   print(TRGAUSS)
77   }

```

**Annexe E**

**RÉSULTATS COMPLETS DES  
SIMULATIONS**

TAB. E.1 – Estimations des paramètres, puissance du test du maximum de vraisemblance et AIC – premier regroupement,  $\beta_0 = 0$  et absence de corrélation.

Censure	$\beta$	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne						
		$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*			
0 %	0,00	-0,004	0,063	6,0	0,8	1,0	-0,004	0,070	7,6	1,6	110,7	0,063	7,6	1,6	-0,004	0,063	6,0	0,8	-0,3	-0,004	0,064	6,4	0,8	-0,4
	0,10	0,102	0,063	34,0	17,6	3,6	0,112	0,070	37,2	17,2	116,0	0,063	36,0	16,8	0,102	0,063	34,0	17,6	2,4	0,103	0,064	34,4	17,6	2,5
	0,20	0,201	0,063	89,6	72,4	10,9	0,224	0,071	85,6	67,6	122,3	0,063	89,2	72,0	0,202	0,064	89,6	72,0	9,7	0,203	0,064	89,2	71,6	9,7
	0,30	0,298	0,064	99,6	98,8	22,4	0,326	0,071	99,6	98,0	132,0	0,063	99,6	98,8	0,298	0,064	99,6	98,8	21,0	0,299	0,064	99,6	98,8	20,9
	0,40	0,394	0,064	100	100	38,1	0,434	0,072	100	100	149,0	0,064	100	100	0,395	0,064	100	100	36,8	0,397	0,065	100	100	36,8
0,50	0,504	0,065	100	100	60,1	0,558	0,072	100	100	173,1	0,065	100	100	0,506	0,065	100	100	59,1	0,508	0,065	100	100	59,1	
30 %	0,00	-0,001	0,075	4,8	0,4	0,9	-0,002	0,081	6,0	3,2	104,3	0,075	4,4	0,8	-0,001	0,075	4,8	0,4	-0,6	-0,001	0,076	4,8	0,8	-0,5
	0,10	0,105	0,076	30,0	15,2	3,1	0,112	0,081	30,8	14,8	108,8	0,075	30,4	14,8	0,106	0,076	29,2	14,8	1,9	0,106	0,076	29,2	14,8	1,9
	0,20	0,202	0,076	79,2	53,6	8,1	0,215	0,082	74,8	51,6	112,6	0,075	80,0	57,6	0,203	0,076	79,6	54,0	6,7	0,203	0,076	78,8	53,6	6,7
	0,30	0,301	0,076	97,2	90,0	16,7	0,324	0,082	96,8	89,6	122,8	0,075	96,8	90,4	0,302	0,076	97,2	90,0	15,4	0,303	0,076	97,2	90,0	15,5
	0,40	0,403	0,076	100	100	28,8	0,436	0,082	100	100	135,7	0,075	100	100	0,404	0,076	100	100	27,7	0,405	0,076	100	100	27,8
0,50	0,497	0,076	100	100	43,4	0,541	0,082	100	100	148,5	0,075	100	100	0,498	0,076	100	100	42,2	0,499	0,077	100	100	42,3	
60 %	0,00	-0,010	0,100	6,8	2,0	1,1	-0,008	0,106	8,4	2,4	107,1	0,100	6,8	2,4	-0,010	0,100	6,8	2,0	-0,2	-0,010	0,100	6,8	2,0	0,02
	0,10	0,102	0,100	18,4	6,0	2,1	0,106	0,107	19,6	6,4	109,5	0,099	20,0	7,2	0,102	0,100	18,4	6,0	0,8	0,103	0,100	18,4	6,0	1,0
	0,20	0,194	0,100	47,6	23,6	4,8	0,207	0,107	49,6	29,6	111,2	0,099	48,8	25,6	0,195	0,100	47,6	24,0	3,6	0,195	0,100	48,0	24,0	3,7
	0,30	0,296	0,100	82,8	62,8	9,6	0,315	0,107	84,0	62,0	114,7	0,100	84,4	62,0	0,297	0,101	83,2	62,8	8,3	0,297	0,101	83,2	62,4	8,4
	0,40	0,403	0,101	96,8	91,2	17,0	0,431	0,108	96,4	91,2	123,1	0,100	96,8	91,6	0,403	0,101	96,8	91,2	15,4	0,404	0,101	96,8	90,8	15,6
0,50	0,501	0,102	99,6	99,2	25,6	0,542	0,109	99,6	99,6	132,7	0,101	99,6	99,6	0,502	0,102	99,6	99,2	24,3	0,503	0,102	99,6	99,6	24,5	
90 %	0,00	0,006	0,201	5,6	1,2	1,0	0,015	0,212	8,4	1,6	115,8	0,198	6,0	0,8	0,006	0,201	5,6	1,2	-0,7	0,006	0,201	5,6	1,2	-0,2
	0,10	0,112	0,201	9,6	2,4	1,3	0,115	0,212	11,6	3,6	118,7	0,199	9,6	3,2	0,112	0,201	10,0	2,4	-0,3	0,112	0,202	10,0	2,4	0,3
	0,20	0,195	0,202	19,2	7,2	2,0	0,209	0,213	19,2	9,2	117,6	0,201	18,8	7,6	0,196	0,202	19,2	7,2	0,2	0,196	0,202	18,8	7,6	0,6
	0,30	0,282	0,202	27,6	8,4	2,8	0,299	0,214	27,2	12,0	120,1	0,201	28,0	10,4	0,283	0,203	27,2	9,2	1,1	0,283	0,203	27,2	9,2	1,8
	0,40	0,421	0,205	54,8	27,2	5,2	0,454	0,217	53,6	32,4	122,8	0,204	54,8	29,2	0,422	0,205	54,4	27,6	3,5	0,423	0,206	54,0	27,2	4,1
0,50	0,510	0,207	72,0	42,4	7,2	0,547	0,218	69,2	48,4	123,1	0,205	71,2	45,6	0,511	0,207	71,6	42,8	5,4	0,512	0,207	71,6	43,2	5,9	

\* : estimation de l'écart-type de  $\beta$ ; <sup>a</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; • : critère d'information d'Akaike.

TAB. E.2 – Estimations des paramètres, puissance du test du maximum de vraisemblance et AIC – second regroupement,  $\beta_0 = 0$  et absence de corrélation.

Censure	$\beta$	Modèle naïf sans effet groupe					Modèle naïf avec effet groupe					Modèle marginal					Modèle de fragilité gamma					Modèle de fragilité gaussienne				
		$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	AIC*	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	AIC*	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	AIC*	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	AIC*		
0 %	0,00	0,003	0,063	4,0	0,4	0,9	0,003	0,064	4,0	0,4	10,0	0,059	10,4	4,4	0,003	0,063	4,0	0,4	-0,5	0,003	0,063	4,0	0,4	0,4	-0,9	
	0,10	0,101	0,063	35,6	15,6	3,5	0,102	0,064	34,8	16,0	12,5	0,059	41,2	24,0	0,102	0,063	35,2	15,6	2,0	0,102	0,063	34,4	15,6	1,6	1,6	
	0,20	0,202	0,063	90,0	71,2	10,9	0,204	0,064	90,4	72,8	19,9	0,058	91,6	74,0	0,202	0,063	90,0	72,0	9,3	0,202	0,064	90,8	72,0	9,0	9,0	
	0,30	0,293	0,064	99,2	96,4	21,8	0,297	0,064	99,6	97,2	31,1	0,059	98,8	97,2	0,294	0,064	99,2	96,8	20,3	0,294	0,064	99,2	96,8	20,0	20,0	
	0,40	0,401	0,064	100	100	39,2	0,405	0,065	100	100	48,5	0,059	100	100	0,401	0,064	100	100	37,7	0,402	0,064	100	100	37,4	37,4	
0,50	0,500	0,065	100	100	59,3	0,505	0,066	100	100	68,4	0,060	100	100	0,501	0,065	100	100	57,7	0,501	0,065	100	100	57,3	57,3		
30 %	0,00	-0,002	0,075	6,0	1,6	1,0	-0,002	0,076	6,0	2,0	10,2	0,068	12,8	4,4	-0,002	0,075	6,0	1,6	-0,3	-0,002	0,075	6,0	1,6	1,6	-0,5	
	0,10	0,107	0,075	30,4	12,8	3,0	0,108	0,076	28,4	12,4	12,3	0,069	36,8	21,6	0,107	0,075	30,4	13,2	1,7	0,107	0,075	30,0	12,8	1,4	1,4	
	0,20	0,204	0,076	80,0	53,6	8,2	0,206	0,076	80,0	54,8	16,8	0,069	81,2	64,4	0,205	0,076	80,0	54,0	6,6	0,205	0,076	80,0	54,4	6,3	6,3	
	0,30	0,300	0,076	98,8	92,8	16,6	0,301	0,076	98,8	92,8	25,5	0,070	97,6	91,6	0,300	0,076	98,8	92,8	15,1	0,300	0,076	98,8	92,8	14,8	14,8	
	0,40	0,402	0,076	100	100	28,9	0,404	0,076	100	100	38,6	0,070	100	99,2	0,402	0,076	100	100	27,7	0,402	0,076	100	100	27,4	27,4	
0,50	0,501	0,076	100	100	44,3	0,505	0,077	100	100	53,3	0,070	100	100	0,502	0,076	100	100	42,9	0,502	0,076	100	100	42,6	42,6		
60 %	0,00	-0,003	0,100	5,2	0,4	0,9	-0,004	0,100	5,2	0,4	9,6	0,090	8,8	4,0	-0,003	0,100	5,2	0,4	-0,7	-0,003	0,100	5,2	0,4	0,4	-0,8	
	0,10	0,093	0,100	16,8	6,0	1,9	0,095	0,100	15,2	6,0	10,7	0,094	21,2	9,6	0,094	0,100	16,8	6,0	0,5	0,099	0,100	16,8	6,0	0,3	0,3	
	0,20	0,204	0,100	53,6	30,4	5,2	0,206	0,101	54,4	31,6	14,4	0,094	56,0	38,8	0,204	0,100	54,0	30,4	3,7	0,204	0,100	54,8	30,4	3,6	3,6	
	0,30	0,299	0,101	84,0	67,6	9,8	0,300	0,101	84,4	67,6	18,6	0,094	86,0	71,2	0,299	0,101	84,0	67,6	8,3	0,299	0,101	84,0	68,0	8,1	8,1	
	0,40	0,399	0,101	96,8	91,6	16,6	0,400	0,102	96,4	92,0	25,7	0,091	98,0	92,8	0,399	0,101	96,8	91,8	15,2	0,399	0,101	96,8	92,0	15,0	15,0	
0,50	0,494	0,102	100	99,6	24,8	0,497	0,102	100	99,2	34,2	0,091	100	98,4	0,494	0,102	100	99,6	23,5	0,494	0,102	100	99,6	23,4	23,4		
90 %	0,00	0,003	0,201	3,6	0,0	0,9	0,001	0,202	4,0	0,0	9,8	0,189	7,2	2,4	0,003	0,201	3,6	0,0	-0,5	0,003	0,201	3,6	0,0	0,0	-0,6	
	0,10	0,093	0,201	6,8	1,6	1,2	0,095	0,202	8,0	1,6	10,4	0,187	12,4	6,0	0,093	0,201	6,8	1,6	-0,1	0,093	0,201	7,2	1,6	1,6	-0,1	
	0,20	0,214	0,202	18,0	4,8	2,1	0,214	0,203	18,4	5,6	11,0	0,183	27,2	10,4	0,214	0,202	17,6	4,8	0,6	0,214	0,202	17,6	4,8	0,6	0,6	
	0,30	0,322	0,203	34,4	14,0	3,4	0,324	0,204	35,2	14,8	12,2	0,191	41,2	21,2	0,323	0,203	34,4	14,0	1,8	0,323	0,203	34,2	14,0	1,8	1,8	
	0,40	0,385	0,204	50,4	24,0	4,5	0,387	0,205	49,6	24,4	14,0	0,186	54,0	33,6	0,386	0,204	50,4	24,0	3,0	0,386	0,204	50,4	23,6	3,0	3,0	
0,50	0,489	0,206	66,8	45,6	6,7	0,491	0,207	66,4	45,2	16,2	0,191	70,0	53,2	0,489	0,206	67,2	45,6	5,3	0,489	0,206	67,2	45,6	5,3	5,3		

\* : estimation de l'écart-type de  $\hat{\beta}$ ; <sup>a</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; • : critère d'information d'Akaike.

TAB. E.3 – Estimations des paramètres, puissance du test du maximum de vraisemblance et AIC – premier regroupement,  $\beta_0 = 0,7$  et absence de corrélation.

Censure	Modèle naïf sans effet groupe					Modèle naïf avec effet groupe					Modèle marginal					Modèle de fragilité gamma					Modèle de fragilité gaussienne				
	$\beta$	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*
0 %	0,70	0,700	0,067	4,8	0,8	109,5	0,767	0,075	16,0	6,4	220,7	0,066	5,6	0,4	0,702	0,067	5,6	0,8	108,4	0,704	0,067	5,6	0,4	108,4	
	0,80	0,800	0,068	28,8	8,0	139,2	0,878	0,074	67,6	41,6	250,0	0,067	32,0	10,4	0,802	0,068	31,2	8,8	137,9	0,805	0,068	32,4	9,2	137,9	
	0,90	0,897	0,069	83,2	62,0	170,0	0,980	0,076	93,6	87,2	278,5	0,068	82,0	62,4	0,899	0,069	84,0	62,0	168,6	0,901	0,069	84,4	63,6	168,5	
	1,00	0,997	0,070	98,4	94,8	203,5	1,088	0,077	100	99,2	311,8	0,069	98,4	95,2	0,999	0,070	98,4	94,8	202,1	1,002	0,070	99,2	95,6	202,0	
	1,10	1,104	0,071	100	100	241,5	1,204	0,079	100	100	352,2	0,070	100	100	1,107	0,071	100	100	240,3	1,111	0,072	100	100	240,3	
	1,20	1,202	0,072	100	100	276,5	1,310	0,080	100	100	386,2	0,071	100	100	1,205	0,073	100	100	275,2	1,208	0,073	100	100	275,2	
30 %	0,70	0,697	0,077	4,4	1,6	83,7	0,754	0,084	8,0	3,6	188,4	0,077	4,8	2,0	0,699	0,077	4,4	1,2	82,4	0,700	0,077	4,8	1,6	82,5	
	0,80	0,792	0,078	23,2	6,8	107,3	0,862	0,084	53,2	27,6	212,9	0,077	22,8	6,4	0,794	0,078	23,6	6,8	106,0	0,796	0,078	24,0	6,8	106,0	
	0,90	0,901	0,078	72,4	51,2	137,2	0,984	0,085	89,6	74,8	244,8	0,078	72,0	53,6	0,904	0,078	73,2	54,0	136,2	0,907	0,079	73,6	54,8	136,3	
	1,00	1,000	0,079	98,8	88,0	167,0	1,087	0,086	99,6	95,6	272,8	0,078	98,8	89,6	1,003	0,079	98,8	89,6	165,6	1,005	0,079	98,8	89,6	165,7	
	1,10	1,106	0,080	100	98,8	201,7	1,205	0,087	100	100	308,7	0,079	100	100	1,110	0,080	100	98,8	200,5	1,112	0,080	100	98,8	200,6	
	1,20	1,201	0,080	100	100	235,1	1,307	0,088	100	100	340,2	0,080	100	100	1,204	0,081	100	100	233,6	1,206	0,081	100	100	233,6	
60 %	0,70	0,694	0,103	6,4	1,2	47,5	0,742	0,110	7,6	2,8	153,5	0,102	7,6	0,8	0,696	0,104	6,0	1,2	46,3	0,697	0,104	5,6	1,2	46,5	
	0,80	0,810	0,105	14,8	4,0	63,4	0,868	0,112	28,8	12,8	170,3	0,105	14,0	4,0	0,812	0,105	15,2	4,4	62,2	0,814	0,105	16,0	4,4	62,5	
	0,90	0,901	0,106	45,2	20,0	77,6	0,966	0,113	64,8	38,8	182,4	0,105	47,6	22,0	0,903	0,106	46,0	20,8	76,1	0,904	0,106	46,4	21,2	76,3	
	1,00	0,985	0,107	75,2	52,8	91,8	1,053	0,114	86,0	69,2	198,9	0,107	76,4	55,6	0,987	0,108	76,0	53,6	90,5	0,989	0,108	76,0	53,6	90,8	
	1,10	1,102	0,109	94,8	88,8	112,7	1,182	0,116	98,0	96,0	219,9	0,108	94,8	88,8	1,104	0,110	94,8	89,2	111,4	1,106	0,110	94,8	89,2	111,6	
	1,20	1,219	0,112	100	99,2	135,5	1,305	0,119	100	99,2	242,0	0,111	99,6	99,2	1,222	0,112	100	99,2	134,2	1,224	0,112	100	99,2	134,3	
90 %	0,70	0,720	0,213	3,6	0,4	13,0	0,761	0,224	4,8	0,4	129,9	0,210	3,6	0,4	0,722	0,213	3,6	0,4	11,3	0,722	0,213	3,6	0,4	11,9	
	0,80	0,804	0,216	7,2	1,6	15,8	0,851	0,227	10,4	3,2	131,3	0,214	7,2	2,4	0,806	0,216	7,2	1,6	14,1	0,806	0,216	7,2	1,6	14,6	
	0,90	0,907	0,220	13,2	3,6	19,5	0,955	0,231	18,4	4,4	134,8	0,219	12,8	4,8	0,909	0,220	13,6	3,6	17,7	0,909	0,221	13,2	3,6	18,1	
	1,00	1,011	0,225	25,6	6,8	23,5	1,073	0,236	36,8	14,0	141,0	0,223	28,4	8,4	1,013	0,225	26,0	6,8	21,8	1,014	0,225	26,0	6,8	22,4	
	1,10	1,110	0,230	35,2	18,0	27,6	1,177	0,241	47,2	24,8	144,9	0,229	37,2	17,6	1,114	0,230	36,0	18,4	25,9	1,115	0,231	36,0	18,4	26,4	
	1,20	1,216	0,236	58,4	31,6	32,0	1,280	0,247	66,0	40,0	147,9	0,234	58,8	32,2	1,218	0,236	58,8	32,0	30,3	1,218	0,236	58,8	32,0	30,8	

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : critère d'information d'Akaike.

TAB. E.4 – Estimations des paramètres, puissance du test du maximum de vraisemblance et AIC – second regroupement,  $\beta_0 = 0,7$  et absence de corrélation.

Censure	Modèle naïf sans effet groupe					Modèle naïf avec effet groupe					Modèle marginal					Modèle de fragilité gamma					Modèle de fragilité gaussienne				
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	AIC*
0 %	0,70	0,701	0,067	6,4	1,2	110,0	0,707	0,067	5,2	1,2	119,0	0,061	8,0	3,6	108,5	0,702	0,067	6,4	1,2	108,5	0,702	0,067	6,4	1,2	108,1
	0,80	0,804	0,068	33,2	13,6	140,7	0,810	0,068	34,4	16,8	149,7	0,062	39,2	20,0	139,1	0,805	0,068	33,6	13,6	139,1	0,805	0,068	33,6	13,2	138,7
	0,90	0,909	0,069	87,2	68,4	174,0	0,917	0,069	89,6	71,6	183,0	0,064	88,4	72,8	172,6	0,911	0,069	87,6	69,2	172,6	0,911	0,069	87,6	70,0	172,2
	1,00	0,998	0,070	99,2	96,4	203,8	1,006	0,070	98,8	97,2	212,9	0,065	99,2	96,4	202,3	0,999	0,070	99,2	96,4	202,3	1,000	0,070	99,2	87,4	202,0
	1,10	1,097	0,071	100	100	238,7	1,105	0,072	100	100	247,5	0,065	100	100	237,1	1,098	0,071	100	100	237,1	1,099	0,071	100	100	236,7
1,20	1,204	0,072	100	100	277,8	1,212	0,073	100	100	287,0	0,066	100	100	276,4	1,205	0,073	100	100	276,4	1,206	0,073	100	100	276,0	
30 %	0,70	0,698	0,077	3,2	0,8	83,9	0,703	0,078	4,4	0,8	93,0	0,070	8,4	3,2	82,4	0,699	0,077	3,2	0,8	82,4	0,699	0,077	3,2	0,8	82,1
	0,80	0,803	0,078	23,2	9,6	109,9	0,808	0,078	25,2	9,6	119,0	0,069	35,6	18,0	108,5	0,804	0,078	23,2	9,6	108,5	0,804	0,078	23,2	9,6	108,2
	0,90	0,903	0,078	75,6	53,2	137,7	0,909	0,079	77,2	55,2	147,3	0,071	78,4	62,8	136,3	0,904	0,078	76,4	53,6	136,3	0,904	0,078	76,4	53,6	136,1
	1,00	1,004	0,079	97,6	91,2	168,3	1,011	0,079	97,6	92,0	177,5	0,073	98,0	91,6	166,8	1,005	0,079	97,6	91,2	166,8	1,006	0,079	97,6	91,6	166,6
	1,10	1,098	0,080	100	99,2	199,1	1,105	0,080	100	99,2	208,1	0,074	99,6	99,2	197,7	1,099	0,080	100	99,2	197,7	1,100	0,080	100	99,2	197,4
1,20	1,200	0,081	100	100	234,3	1,209	0,081	100	100	243,5	0,073	100	100	232,8	1,201	0,081	100	100	232,8	1,202	0,081	100	100	232,6	
60 %	0,70	0,702	0,104	4,4	0,8	48,4	0,706	0,104	4,0	0,8	57,4	0,097	8,0	3,6	46,9	0,703	0,104	4,4	0,8	46,9	0,703	0,104	4,4	0,8	46,8
	0,80	0,806	0,105	16,4	6,0	63,0	0,810	0,105	17,6	5,6	72,4	0,098	22,4	10,4	61,7	0,807	0,105	17,2	6,0	61,7	0,807	0,105	17,2	5,6	61,6
	0,90	0,895	0,106	42,0	24,4	76,8	0,901	0,107	45,6	24,0	86,2	0,098	51,6	28,4	75,5	0,896	0,106	42,8	24,0	75,5	0,896	0,106	42,8	24,0	75,4
	1,00	1,008	0,108	80,8	65,2	95,9	1,014	0,108	83,2	65,6	104,8	0,098	84,8	64,0	94,5	1,009	0,108	80,8	65,2	94,5	1,009	0,108	81,2	65,6	94,3
	1,10	1,089	0,109	95,2	80,0	110,4	1,095	0,110	96,1	80,8	119,5	0,102	95,2	80,4	108,9	1,090	0,109	95,2	80,4	108,9	1,091	0,109	95,2	80,4	108,8
1,20	1,196	0,111	99,6	98,4	130,8	1,204	0,112	99,6	98,4	140,2	0,102	99,6	99,2	129,5	1,198	0,111	99,6	99,6	129,5	1,198	0,111	99,6	98,4	129,4	
90 %	0,70	0,687	0,212	4,8	0,8	12,0	0,690	0,213	5,2	0,4	21,3	0,195	10,8	3,6	10,5	0,688	0,212	4,8	0,8	10,5	0,688	0,212	4,8	0,8	10,5
	0,80	0,799	0,212	6,8	1,6	15,7	0,802	0,217	7,6	1,6	24,7	0,196	12,4	5,6	14,1	0,799	0,216	7,2	1,6	14,1	0,800	0,216	7,2	1,6	14,1
	0,90	0,918	0,220	14,0	4,4	19,9	0,924	0,221	15,2	4,8	28,9	0,203	22,4	9,2	18,3	0,919	0,220	15,2	4,8	18,3	0,919	0,220	15,2	4,8	18,3
	1,00	1,005	0,224	26,0	7,2	23,3	1,009	0,225	25,2	7,6	32,3	0,213	27,6	11,6	21,7	1,006	0,224	26,0	7,2	21,7	1,006	0,224	26,0	7,2	21,7
	1,10	1,133	0,231	46,4	19,6	28,4	1,137	0,232	49,2	20,4	37,1	0,215	53,6	30,8	20,0	26,9	1,134	0,231	46,4	20,0	26,9	1,134	0,231	46,8	20,0
1,20	1,215	0,236	57,2	31,2	32,0	1,220	0,237	57,6	30,0	41,1	0,216	63,6	41,6	31,2	30,7	1,216	0,236	58,0	31,2	30,7	1,216	0,236	58,0	31,2	30,6

\* : estimation de l'écart-type de  $\beta$  ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 % ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 % ; \* : critère d'information d'Alkaïe.



TAB. E.6 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0$  et fragilité gamma de variance  $\sigma^2 = 1$ .

Censure	Modèle naïf sans effet groupe					Modèle naïf avec effet groupe					Modèle marginal					Modèle de fragilité gamma					Modèle de fragilité gaussienne				
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	ESp.*	Var.**	$\hat{\beta}$	EC*	5 % <sup>a</sup>	1 % <sup>b</sup>	ESp.*	Var.**	
0 %	0,00	-0,001	0,063	4,4	2,0	0,000	0,064	4,8	4,0	0,053	11,2	4,8	0,000	0,064	4,4	0,0	-0,648	1,230	0,000	0,064	4,4	0,0	10 <sup>-18</sup>	1,929	
	0,10	0,054	0,063	12,4	4,8	0,094	0,064	35,2	12,8	0,052	25,2	12,8	0,093	0,064	34,8	12,4	-0,616	1,196	0,093	0,064	35,2	12,4	10 <sup>-18</sup>	1,758	
	0,20	0,113	0,063	42,4	22,0	0,205	0,064	89,6	72,8	0,057	53,2	41,2	0,204	0,064	89,6	71,6	-0,623	1,205	0,204	0,064	89,6	71,6	10 <sup>-18</sup>	1,798	
	0,30	0,174	0,063	77,2	54,0	0,305	0,064	100	99,6	0,057	83,6	69,6	0,303	0,064	100	99,6	-0,619	1,197	0,303	0,064	100	99,6	10 <sup>-18</sup>	1,728	
	0,40	0,227	0,063	90,0	78,4	0,406	0,065	100	100	0,057	94,0	81,6	0,404	0,065	100	100	-0,642	1,220	0,403	0,065	100	100	10 <sup>-18</sup>	1,836	
0,50	0,289	0,064	98,4	95,2	0,509	0,065	100	100	0,058	97,6	94,8	0,506	0,065	100	100	-0,606	1,190	0,505	0,065	100	100	10 <sup>-18</sup>	1,724		
30 %	0,00	0,002	0,075	7,2	2,0	0,002	0,076	7,2	0,8	0,067	14,4	6,0	0,003	0,076	6,8	0,8	-0,606	1,185	0,003	0,076	7,2	0,8	10 <sup>-18</sup>	1,677	
	0,10	0,061	0,075	13,2	4,8	0,097	0,076	28,4	12,4	0,070	21,2	8,8	0,097	0,076	27,6	12,4	-0,597	1,176	0,097	0,076	27,6	12,4	10 <sup>-18</sup>	1,608	
	0,20	0,130	0,075	43,2	22,0	0,205	0,076	75,6	51,6	0,070	48,0	32,4	0,204	0,076	74,4	51,2	-0,622	1,212	0,204	0,076	75,6	51,6	10 <sup>-18</sup>	1,733	
	0,30	0,195	0,075	66,8	48,4	0,306	0,076	94,4	90,4	0,073	69,6	52,8	0,304	0,076	94,0	90,4	-0,605	1,199	0,304	0,076	94,0	90,4	10 <sup>-18</sup>	1,711	
	0,40	0,268	0,076	88,4	74,4	0,407	0,077	99,6	99,2	0,075	87,2	76,4	0,405	0,077	99,6	99,2	-0,582	1,149	0,404	0,077	99,6	99,2	10 <sup>-19</sup>	1,574	
0,50	0,318	0,076	96,4	91,6	0,509	0,077	100	100	0,079	96,0	86,8	0,506	0,077	100	100	-0,615	1,202	0,505	0,077	100	100	10 <sup>-18</sup>	1,697		
60 %	0,00	-0,001	0,100	5,6	1,2	0,003	0,100	9,2	1,6	0,088	12,4	3,6	0,003	0,101	8,8	1,6	-0,586	1,160	0,003	0,101	8,8	1,6	10 <sup>-18</sup>	1,499	
	0,10	0,071	0,100	14,0	5,6	0,104	0,100	17,2	8,0	0,094	20,8	9,6	0,103	0,101	17,2	8,0	-0,636	1,243	0,103	0,101	16,0	8,0	10 <sup>-19</sup>	1,645	
	0,20	0,153	0,100	34,8	14,8	0,193	0,101	50,8	28,8	0,092	41,6	24,8	0,192	0,101	50,0	28,4	-0,560	1,142	0,192	0,101	50,4	28,4	10 <sup>-18</sup>	1,414	
	0,30	0,233	0,100	64,0	39,6	0,301	0,101	82,0	65,6	0,093	68,0	50,0	0,299	0,101	81,6	64,4	-0,617	1,220	0,299	0,101	81,2	64,4	10 <sup>-18</sup>	1,561	
	0,40	0,329	0,101	87,6	75,6	0,412	0,101	98,8	94,8	0,095	88,0	75,6	0,410	0,102	98,8	94,8	-0,592	1,164	0,410	0,102	98,8	94,8	10 <sup>-18</sup>	1,490	
0,50	0,389	0,101	97,6	89,2	0,500	0,102	100	99,6	0,098	96,8	88,0	0,497	0,102	100	99,6	-0,597	1,172	0,497	0,102	100	99,6	10 <sup>-19</sup>	1,542		
90 %	0,00	0,007	0,201	4,4	0,4	0,001	0,202	7,6	0,8	0,178	13,2	6,4	0,002	0,202	6,8	0,8	-0,501	1,135	0,001	0,202	6,8	0,8	10 <sup>-18</sup>	1,149	
	0,10	0,093	0,201	6,8	0,4	0,098	0,202	7,6	0,4	0,179	13,6	7,2	0,098	0,202	7,6	0,4	-0,481	1,088	0,098	0,202	7,6	0,4	10 <sup>-19</sup>	1,098	
	0,20	0,193	0,202	16,8	5,2	0,211	0,202	17,2	6,8	0,170	27,2	14,4	0,209	0,202	16,8	6,8	-0,503	1,121	0,209	0,202	16,8	6,8	10 <sup>-18</sup>	1,156	
	0,30	0,300	0,203	33,6	12,0	0,315	0,204	36,0	13,6	0,180	42,4	27,6	0,314	0,204	35,2	12,8	-0,507	1,139	0,314	0,204	35,2	12,8	10 <sup>-18</sup>	1,163	
	0,40	0,381	0,204	39,2	21,2	0,405	0,205	41,6	24,8	0,176	49,6	35,2	0,403	0,205	42,0	24,0	-0,490	1,105	0,403	0,205	42,0	24,0	10 <sup>-19</sup>	1,136	
0,50	0,491	0,206	64,0	37,2	0,513	0,207	66,4	44,0	0,182	69,2	51,6	0,511	0,207	66,0	44,0	-0,504	1,136	0,511	0,207	66,0	44,0	10 <sup>-18</sup>	1,151		

\* : estimation de l'écart-type de  $\hat{\beta}$ ; <sup>a</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \*\* : variance de la fragilité.

ANNEXE E. RÉSULTATS COMPLETS DES SIMULATIONS

TAB. E.7 – AIC et (entre parenthèses) taux de recouvrement de l'intervalle de confiance à 95 % de  $\hat{\sigma}^2$  pour  $\beta_0 = 0$ .

(a)  $\sigma^2 = 1$  et premier regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,00	1	1 024	822 (99,1)	817 (90,4)
	0,10	1	1 037	834 (99,5)	829 (99,5)
	0,20	3	1 038	841 (100)	835 (100)
	0,30	6	1 042	840 (100)	834 (100)
	0,40	11	1 025	824 (99,1)	818 (99,1)
	0,50	15	1 052	850 (100)	844 (100)
30 %	0,00	1	762	556 (96,8)	553 (99,5)
	0,10	1	769	566 (97,7)	563 (100)
	0,20	1	764	558 (96,8)	556 (100)
	0,30	6	777	572 (97,2)	569 (100)
	0,40	11	776	573 (97,7)	571 (100)
	0,50	17	778	577 (97,2)	574 (99,5)
60 %	0	1	490	283 (86,3)	287 (100)
	0,10	1	489	282 (83,6)	286 (100)
	0,20	3	490	283 (89,1)	287 (100)
	0,30	6	494	289 (86,3)	293 (100)
	0,40	11	494	287 (88,6)	291 (100)
	0,50	16	497	294 (89,1)	298 (100)
90 %	0	1	189	20 (71,3)	37 (100)
	0,10	1	192	22 (65,0)	39 (100)
	0,20	2	194	23 (64,5)	41 (100)
	0,30	3	194	24 (70,4)	41 (100)
	0,40	4	191	23 (64,5)	39 (100)
	0,50	6	196	26 (66,8)	44 (100)

(b)  $\sigma^2 = 1$  et second regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,00	1	831	810 (99,5)	809 (74,5)
	0,10	2	792	770 (99,5)	769 (89,5)
	0,20	4	817	795 (99,1)	795 (99,1)
	0,30	8	799	777 (99,5)	777 (99,5)
	0,40	14	831	812 (99,5)	812 (100)
	0,50	21	825	803 (99,1)	803 (100)
30 %	0,00	1	617	595 (99,5)	595 (79,5)
	0,10	1	621	599 (99,5)	599 (96,8)
	0,20	4	645	623 (99,1)	623 (99,1)
	0,30	7	627	605 (99,5)	605 (100)
	0,40	14	622	600 (99,5)	600 (100)
	0,50	19	657	635 (98,6)	635 (100)
60 %	0,00	1	361	338 (98,6)	338 (83,6)
	0,10	1	387	365 (99,5)	365 (90,4)
	0,20	3	354	332 (99,1)	332 (100)
	0,30	6	386	365 (98,6)	365 (100)
	0,40	12	379	357 (99,5)	357 (100)
	0,50	16	387	364 (100)	364 (100)
90 %	0,00	1	97	75 (94,5)	75 (95,0)
	0,10	1	95	73 (95,4)	73 (98,1)
	0,20	1	98	76 (95,9)	76 (98,6)
	0,30	3	99	77 (93,6)	77 (100)
	0,40	4	98	76 (94,5)	76 (100)
	0,50	6	103	81 (96,3)	81 (100)

(c)  $\sigma^2 = 4$  et premier regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,00	1	1 772	2 260 (51,2)	2 251 (0)
	0,10	1	1 953	2 271 (50,5)	2 262 (0)
	0,20	1	1 996	2 277 (52,9)	2 268 (0)
	0,30	2	1 951	2 256 (49,8)	2 247 (0)
	0,40	2	2 000	2 273 (53,8)	2 264 (0)
	0,50	3	2 035	2 263 (52,0)	2 254 (0)
30 %	0,00	1	1 116	1 574 (60,8)	1 570 (0,6)
	0,10	1	1 537	1 584 (61,8)	1 579 (1,0)
	0,20	1	1 638	1 571 (60,6)	1 566 (2,4)
	0,30	2	1 410	1 587 (64,0)	1 582 (5,8)
	0,40	3	1 606	1 597 (63,4)	1 593 (13,2)
	0,50	4	1 357	1 597 (66,4)	1 592 (22,0)
60 %	0,00	1	1 013	831 (78,1)	840 (77,2)
	0,10	1	854,5	845 (77,6)	854 (78,4)
	0,20	1	989,6	850 (72,8)	859 (73,2)
	0,30	2	1 052	853 (71,2)	862 (77,6)
	0,40	3	1 001	862 (69,6)	871 (77,6)
	0,50	5	1 020	858 (72,4)	866 (80,8)
90 %	0,00	1	303	110 (40,4)	154 (5,9)
	0,10	1	304	108 (52,2)	152 (10,9)
	0,20	1	308	111 (47,7)	155 (56,3)
	0,30	2	315	116 (44,5)	162 (96,8)
	0,40	3	311	114 (44,5)	159 (100)
	0,50	5	311	115 (50,4)	159 (99,0)

(d)  $\sigma^2 = 4$  et second regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,00	1	1 777	2 005 (97,5)	2 005 (6,0)
	0,10	1	1 627	1 897 (95,5)	1 896 (8,0)
	0,20	2	637	1 964 (97,0)	1 964 (12,5)
	0,30	3	1 792	1 967 (94,5)	1 967 (19,0)
	0,40	5	1 801	1 916 (97,0)	1 916 (24,0)
	0,50	7	1 914	1 944 (98,5)	1 944 (34,5)
30 %	0,00	1	1 477	1 564 (100)	1 564 (15,5)
	0,10	1	1 436	1 565 (100)	1 565 (18,0)
	0,20	2	1 388	1 524 (99,0)	1 523 (23,5)
	0,30	3	1 536	1 551 (97,0)	1 550 (31,0)
	0,40	4	1 279	1 588 (99,5)	1 588 (37,0)
	0,50	6	1 478	1 543 (97,0)	1 542 (45,5)
60 %	0,00	1	885	905 (97,5)	905 (25,0)
	0,10	1	894	912 (97,5)	912 (27,0)
	0,20	2	929	916 (96,0)	916 (46,0)
	0,30	3	908	915 (97,5)	915 (53,5)
	0,40	4	905	900 (96,5)	900 (72,5)
	0,50	7	890	898 (98,5)	898 (85,5)
90 %	0,00	1	211	191 (84,0)	193 (67,5)
	0,10	1	218	197 (77,2)	198 (67,7)
	0,20	2	188	200 (81,3)	201 (85,0)
	0,30	2	210	196 (80,4)	197 (90,9)
	0,40	3	160	204 (84,5)	206 (92,7)
	0,50	5	218	198 (83,1)	199 (96,8)

TAB. E.8 – Estimation des paramètres et puissance du test du maximum de vraisemblance – premier regroupement,  $\beta_0 = 0$  et fragilité gamma de variance  $\sigma^2 = 4$ .

Censure	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne							
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**
0 %	0.00	0.000	0.063	8,1	0,7	0.000	0.071	4,4	0,5	0.060	8,8	2,9	0.006	0.070	2,9	0,7	-3,051	4,377	0.000	0.070	2,2	0,7	10 <sup>-17</sup>	16,504
	0.10	0.014	0.063	7,8	2,3	0.093	0.071	27,3	14,8	0.059	8,6	3,9	0.095	0.070	29,7	12,5	-3,078	4,412	0.094	0.070	28,9	11,7	10 <sup>-18</sup>	16,683
	0.20	0.046	0.063	11,6	0,7	0.201	0.071	84,5	60,4	0.060	12,4	3,8	0.201	0.070	86,0	67,4	-3,093	4,466	0.199	0.070	84,5	66,6	10 <sup>-18</sup>	16,691
	0.30	0.066	0.063	18,4	3,8	0.307	0.071	91,5	90,0	0.060	20,0	9,2	0.304	0.070	96,9	95,3	-3,043	4,375	0.302	0.070	96,9	95,3	10 <sup>-17</sup>	16,325
	0.40	0.080	0.063	20,1	8,0	0.406	0.071	97,5	96,7	0.060	25,8	11,3	0.406	0.070	100	100	-3,077	4,403	0.403	0.070	100	100	10 <sup>-18</sup>	16,617
0.50	0.105	0.063	38,7	18,3	0.502	0.072	98,6	97,1	0.060	40,8	25,3	0.504	0.071	100	100	-3,045	4,375	0.500	0.071	100	100	10 <sup>-18</sup>	16,475	
30 %	0.00	-0.007	0.075	2,4	0,6	-0.008	0.084	6,1	3,0	0.074	3,7	0,6	-0.008	0.083	4,9	2,4	-2,503	4,337	-0.008	0.083	5,5	1,8	10 <sup>-18</sup>	8,420
	0.10	0.024	0.075	5,6	0,6	0.097	0.084	24,8	9,3	0.075	5,6	0,6	0.093	0.083	26,7	9,3	-2,536	4,384	0.092	0.083	25,4	8,6	10 <sup>-18</sup>	8,564
	0.20	0.054	0.075	7,8	3,2	0.213	0.084	68,1	47,4	0.075	7,1	3,9	0.203	0.083	66,2	43,5	-2,499	4,369	0.202	0.083	65,5	42,8	10 <sup>-18</sup>	8,387
	0.30	0.076	0.075	19,4	5,8	0.310	0.084	94,8	89,6	0.075	16,8	5,8	0.296	0.083	96,1	89,6	-2,533	4,420	0.294	0.083	96,1	88,9	10 <sup>-18</sup>	8,569
	0.40	0.107	0.075	24,8	9,9	0.422	0.085	98,1	97,5	0.074	27,3	10,5	0.403	0.083	100	98,1	-2,558	4,413	0.400	0.083	100	98,1	10 <sup>-18</sup>	8,663
0.50	0.132	0.075	45,3	18,0	0.512	0.085	98,7	98,1	0.075	46,5	22,3	0.488	0.083	100	100	-2,557	4,409	0.485	0.084	100	100	10 <sup>-18</sup>	8,673	
60 %	0.00	-0.002	0.100	1,9	0,9	0.000	0.109	9,0	3,3	0.100	2,8	0,0	-0.001	0.108	7,1	1,4	-1,893	4,309	-0.001	0.108	5,7	1,4	10 <sup>-18</sup>	4,390
	0.10	0.048	0.100	9,0	1,6	0.113	0.110	15,9	7,9	0.099	9,5	2,1	0.105	0.110	12,7	4,7	-1,948	4,457	0.104	0.109	13,8	4,7	10 <sup>-19</sup>	4,528
	0.20	0.073	0.100	11,6	5,0	0.202	0.110	48,4	31,3	0.099	13,6	4,5	0.185	0.110	45,4	25,7	-1,962	4,474	0.183	0.109	44,9	24,2	10 <sup>-18</sup>	4,567
	0.30	0.127	0.100	25,6	14,7	0.326	0.110	76,3	58,1	0.098	24,6	14,3	0.305	0.108	74,4	53,2	-1,973	4,488	0.302	0.110	73,9	52,2	10 <sup>-19</sup>	4,577
	0.40	0.172	0.100	38,4	18,7	0.427	0.111	95,5	84,7	0.099	41,3	19,7	0.401	0.109	94,1	83,7	-1,989	4,545	0.397	0.109	93,6	83,2	10 <sup>-17</sup>	4,651
0.50	0.211	0.100	61,1	38,4	0.530	0.111	99,0	96,7	0.100	61,1	37,4	0.495	0.110	99,0	96,6	-1,960	4,474	0.492	0.110	98,5	96,7	10 <sup>-18</sup>	4,601	
90 %	0.00	0.020	0.201	2,8	0,0	0.010	0.213	0,6	0,0	0.197	5,2	0,4	0.013	0.211	3,6	0,0	-1,050	4,533	-0.013	0.210	2,4	0,0	10 <sup>-18</sup>	2,015
	0.10	0.093	0.201	9,2	2,4	0.125	0.214	8,4	1,6	0.197	8,8	4,4	0.119	0.211	7,2	0,6	-1,026	4,424	0.115	0.210	7,2	0,8	10 <sup>-18</sup>	1,982
	0.20	0.148	0.201	12,0	2,8	0.211	0.215	20,4	7,6	0.195	12,4	5,6	0.194	0.212	20,0	4,8	-1,055	4,561	0.190	0.211	18,8	4,4	10 <sup>-18</sup>	2,020
	0.30	0.227	0.202	22,8	6,0	0.302	0.216	35,2	16,4	0.196	21,6	7,6	0.282	0.213	29,6	14,4	-1,087	4,691	0.279	0.212	28,0	14,4	10 <sup>-19</sup>	2,070
	0.40	0.296	0.203	37,5	14,9	0.396	0.216	56,4	29,0	0.199	37,5	15,8	0.370	0.213	50,0	27,4	-1,067	4,592	0.363	0.212	48,8	26,6	10 <sup>-18</sup>	2,044
0.50	0.430	0.205	50,8	31,6	0.542	0.219	70,0	46,4	0.202	54,0	31,2	0.514	0.215	65,6	40,8	-1,052	4,515	0.507	0.215	64,4	39,2	10 <sup>-18</sup>	2,020	

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.9 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0$  et fragilité gamma de variance  $\sigma^2 = 4$ .

Censure	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne									
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Exp.**	Var.**	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Exp.**	Var.**	
0 %	0.00	-0.002	0.063	5,2	0,4	-0.010	0.064	8,5	2,0	0,048	19,5	7,7	-0.009	0.064	8,1	2,0	-2,990	4,325			-0.009	0.064	7,7	2,0	10 <sup>-18</sup>	15,916
	0.10	0.028	0.063	7,7	3,2	0.100	0.064	26,6	14,3	0.049	21,3	9,4	0.100	0.064	29,9	14,7	-2,731	4,014			0.100	0.064	29,9	14,7	10 <sup>-17</sup>	14,090
	0.20	0.005	0.063	17,9	7,5	0.200	0.064	83,3	62,5	0.050	39,6	23,3	0.193	0.064	85,4	64,5	-2,931	4,281			0.193	0.064	85,4	64,6	10 <sup>-17</sup>	15,673
	0.30	0.094	0.063	30,3	15,7	0.291	0.064	99,6	95,8	0.051	44,0	31,1	0.297	0.064	99,6	97,9	-2,908	4,127			0.296	0.064	99,6	97,9	10 <sup>-17</sup>	15,417
	0.40	0.132	0.063	49,6	28,6	0.392	0.065	100	99,1	0.052	63,5	45,5	0.399	0.065	100	100	-2,753	4,051			0.399	0.065	100	100	10 <sup>-18</sup>	13,907
0.50	0.162	0.063	69,2	47,1	0.496	0.065	100	100	0.054	76,6	61,8	0.499	0.065	100	100	-2,779	4,059			0.499	0.065	100	100	10 <sup>-18</sup>	14,379	
30 %	0.00	0.003	0.075	4,8	0,8	-0.007	0.076	7,2	1,2	0.066	10,8	4,8	-0.006	0.076	7,2	1,2	-2,707	4,349			-0.006	0.076	6,8	1,2	10 <sup>-17</sup>	11,645
	0.10	0.036	0.075	6,1	1,2	0.092	0.076	27,0	13,9	0.066	13,1	5,7	0.094	0.076	28,2	13,5	-2,709	4,432			0.094	0.076	28,2	13,5	10 <sup>-18</sup>	11,650
	0.20	0.068	0.075	15,4	6,1	0.193	0.076	75,6	52,8	0.068	23,5	13,0	0.195	0.076	78,8	54,0	-2,551	4,066			0.195	0.076	78,8	53,6	10 <sup>-18</sup>	10,808
	0.30	0.105	0.075	25,2	14,8	0.291	0.077	95,4	88,8	0.071	32,6	17,7	0.292	0.077	97,1	91,3	-2,594	4,171			0.292	0.077	97,1	91,3	10 <sup>-18</sup>	11,051
	0.40	0.124	0.075	39,1	25,0	0.390	0.077	100	98,7	0.072	46,2	31,2	0.397	0.077	100	100	-2,750	4,384			0.397	0.077	100	100	10 <sup>-19</sup>	11,606
0.50	0.164	0.075	52,8	30,6	0.491	0.078	100	99,5	0.075	53,6	36,2	0.493	0.077	100	100	-2,653	4,239			0.494	0.077	100	100	10 <sup>-18</sup>	11,345	
60 %	0.00	0.000	0.100	9,2	2,8	-0.000	0.101	7,2	0,8	0.087	13,2	5,6	-0.000	0.101	6,8	0,8	-2,334	4,584			-0.000	0.101	6,8	0,8	10 <sup>-18</sup>	7,342
	0.10	0.052	0.100	8,1	1,6	0.096	0.101	13,4	4,4	0.088	12,2	6,5	0.095	0.101	13,8	4,0	-2,318	4,319			0.095	0.101	13,8	4,4	10 <sup>-18</sup>	7,150
	0.20	0.105	0.100	16,2	3,2	0.209	0.101	49,8	25,9	0.092	22,6	11,7	0.209	0.101	49,8	25,9	-2,364	4,467			0.209	0.101	50,2	25,9	10 <sup>-18</sup>	7,379
	0.30	0.134	0.100	32,0	16,6	0.281	0.101	84,5	67,0	0.091	38,3	25,4	0.282	0.101	84,5	67,5	-2,360	4,437			0.282	0.101	84,5	67,5	10 <sup>-19</sup>	7,529
	0.40	0.190	0.100	39,8	25,2	0.408	0.102	97,5	90,6	0.101	41,4	25,6	0.407	0.102	97,1	91,8	-2,283	4,241			0.407	0.102	97,5	91,8	10 <sup>-18</sup>	7,105
0.50	0.250	0.100	62,6	42,9	0.493	0.102	99,6	98,4	0.102	60,2	40,5	0.494	0.102	100	99,2	-2,277	4,225			0.494	0.102	100	99,1	10 <sup>-17</sup>	7,060	
90 %	0.00	-0.014	0.201	6,4	1,3	-0.013	0.201	6,0	2,1	0.159	16,3	8,1	-0.013	0.202	6,0	1,7	-1,020	4,102			-0.013	0.202	5,6	2,1	10 <sup>-18</sup>	3,692
	0.10	0.082	0.201	9,5	1,6	0.110	0.202	8,7	1,2	0.157	23,7	16,6	0.109	0.202	8,7	0,8	-1,635	4,124			0.109	0.202	8,7	0,8	10 <sup>-18</sup>	3,690
	0.20	0.184	0.202	14,1	2,1	0.215	0.202	18,3	5,1	0.161	31,6	20,0	0.214	0.203	17,5	4,7	-1,080	4,291			0.214	0.203	17,5	4,7	10 <sup>-18</sup>	3,816
	0.30	0.226	0.202	25,2	11,3	0.292	0.203	32,7	16,3	0.164	43,2	27,6	0.290	0.204	32,7	15,5	-1,642	4,274			0.290	0.204	32,7	15,5	10 <sup>-18</sup>	3,724
	0.40	0.337	0.203	39,3	19,6	0.403	0.204	51,4	29,2	0.157	55,2	37,6	0.402	0.205	51,0	28,8	-1,746	4,701			0.402	0.205	51,0	28,8	10 <sup>-19</sup>	3,969
0.50	0.439	0.205	51,5	28,3	0.516	0.206	68,6	42,4	0.162	64,8	50,6	0.514	0.207	67,8	42,0	-1,650	4,284			0.514	0.207	68,2	41,6	10 <sup>-18</sup>	3,779	

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.10 – Estimation des paramètres et puissance du test du maximum de vraisemblance – premier regroupement,  $\beta_0 = 0, 7$  et fragilité gamma de variance  $\sigma^2 = 1$ .

Censure	$\beta$	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne						
		$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**
0 %	0,70	0,353	0,064	100	99,5	0,772	0,073	14,0	5,2	0,061	100	99,5	0,717	0,071	4,0	0,0	-0,629	1,151	0,711	0,071	3,2	0,0	10 <sup>-19</sup>	1,875
	0,80	0,397	0,064	99,6	98,4	0,844	0,074	61,6	36,4	0,069	99,6	96,8	0,786	0,072	25,6	9,2	-0,615	1,129	0,778	0,072	22,0	8,8	10 <sup>-18</sup>	1,850
	0,90	0,440	0,064	96,3	92,3	0,973	0,075	93,5	83,1	0,068	96,0	90,0	0,900	0,073	78,1	51,4	-0,623	1,141	0,894	0,073	75,1	48,2	10 <sup>-18</sup>	1,862
	1,00	0,472	0,064	88,0	73,6	1,064	0,076	99,6	99,2	0,070	85,6	66,4	0,988	0,073	98,4	92,0	-0,631	1,165	0,981	0,073	98,0	92,0	10 <sup>-18</sup>	1,901
	1,10	0,526	0,064	73,1	54,2	1,183	0,077	100	100	0,071	67,4	46,9	1,097	0,074	100	99,6	-0,634	1,158	1,089	0,074	100	99,2	10 <sup>-18</sup>	1,925
	1,20	0,578	0,064	48,8	28,4	1,283	0,078	100	100	0,073	38,0	19,6	1,190	0,075	100	100	-0,628	1,149	1,183	0,075	100	100	10 <sup>-18</sup>	1,900
30 %	0,70	0,432	0,076	92,3	81,5	0,762	0,085	92,3	4,0	0,085	91,9	79,9	0,704	0,083	4,4	1,6	-0,575	1,161	0,703	0,083	4,0	1,2	10 <sup>-18</sup>	1,468
	0,80	0,497	0,076	79,5	55,4	0,872	0,086	48,6	24,9	0,078	78,7	54,6	0,808	0,084	20,4	10,4	-0,583	1,174	0,806	0,084	20,4	10,0	10 <sup>-18</sup>	1,495
	0,90	0,535	0,076	48,0	24,8	0,956	0,087	85,6	70,4	0,079	45,2	23,2	0,884	0,084	64,0	40,4	-0,575	1,160	0,883	0,085	64,8	38,0	10 <sup>-18</sup>	1,478
	1,00	0,613	0,076	22,1	9,6	1,081	0,088	99,2	95,1	0,080	19,3	7,6	1,000	0,085	94,3	82,6	-0,582	1,169	0,999	0,086	94,3	82,2	10 <sup>-18</sup>	1,496
	1,10	0,666	0,077	6,0	2,4	1,187	0,089	100	100	0,082	4,4	1,2	1,096	0,086	100	98,4	-0,582	1,190	1,095	0,086	100	98,4	10 <sup>-18</sup>	1,498
	1,20	0,736	0,077	10,0	3,2	1,293	0,090	100	100	0,081	8,8	3,2	1,194	0,087	100	100	-0,564	1,147	1,192	0,087	100	100	10 <sup>-18</sup>	1,455
60 %	0,70	0,541	0,102	40,4	17,6	0,764	0,111	10,4	4,0	0,102	40,4	18,4	0,708	0,110	6,0	0,8	-0,489	1,195	0,708	0,109	6,4	0,8	10 <sup>-18</sup>	1,146
	0,80	0,618	0,103	13,2	5,6	0,862	0,112	32,0	14,8	0,103	14,8	5,2	0,799	0,110	14,8	4,0	-0,482	1,181	0,799	0,110	15,6	4,0	10 <sup>-19</sup>	1,131
	0,90	0,699	0,103	6,4	0,8	0,973	0,113	62,0	42,0	0,104	6,8	1,2	0,903	0,110	42,8	18,0	-0,483	1,182	0,902	0,111	43,2	18,4	10 <sup>-18</sup>	1,136
	1,00	0,769	0,104	9,6	3,2	1,081	0,114	88,4	74,0	0,105	10,0	3,2	1,002	0,111	76,4	52,4	-0,484	1,173	1,001	0,112	75,2	53,2	10 <sup>-18</sup>	1,145
	1,10	0,847	0,105	30,4	12,8	1,178	0,116	98,0	94,4	0,108	29,6	12,4	1,093	0,113	95,6	84,8	-0,472	1,148	1,092	0,113	95,6	84,8	10 <sup>-18</sup>	1,123
	1,20	0,924	0,106	56,8	30,4	1,292	0,117	100	99,6	0,108	54,4	30,8	1,195	0,114	99,6	96,0	-0,482	1,163	1,196	0,114	99,6	96,4	10 <sup>-18</sup>	1,146
90 %	0,70	0,679	0,211	4,4	0,8	0,764	0,223	4,4	0,8	0,209	4,0	1,2	0,719	0,217	3,2	0,4	-0,235	1,143	0,714	0,216	3,6	0,8	10 <sup>-19</sup>	0,723
	0,80	0,776	0,215	3,2	0,8	0,875	0,227	8,0	1,6	0,211	3,6	0,4	0,823	0,220	4,0	0,4	-0,244	1,197	0,816	0,220	4,0	1,2	10 <sup>-19</sup>	0,745
	0,90	0,849	0,218	9,6	2,8	0,958	0,229	22,0	6,4	0,215	9,2	2,8	0,902	0,223	16,4	3,6	-0,242	1,168	0,895	0,222	14,4	3,2	10 <sup>-19</sup>	0,742
	1,00	0,960	0,223	21,2	6,8	1,081	0,235	39,2	18,0	0,216	22,4	7,6	1,017	0,228	28,4	10,4	-0,249	1,192	1,010	0,228	26,4	10,0	10 <sup>-18</sup>	0,759
	1,10	1,078	0,228	31,6	12,4	1,204	0,240	48,4	25,6	0,225	32,0	14,4	1,137	0,233	39,2	17,6	-0,239	1,156	1,128	0,233	37,2	16,8	10 <sup>-18</sup>	0,736
	1,20	1,150	0,232	48,0	24,0	1,306	0,244	67,2	41,6	0,231	47,6	27,2	1,223	0,238	59,2	32,8	-0,248	1,185	1,212	0,237	58,4	32,4	10 <sup>-20</sup>	0,760

\* : estimation de l'écart-type de  $\hat{\beta}$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.11 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0, 7$  et fragilité gamma de variance  $\sigma^2 = 1$ .

Censure	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne							
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**
0 %	0,70	0,404	0,064	95,2	89,6	0,702	0,067	6,8	1,2	0,064	96,4	90,4	0,698	0,066	5,6	1,2	-0,619	1,225	0,697	0,067	5,6	1,2	10 <sup>-18</sup>	1,713
	0,80	0,463	0,064	79,6	68,8	0,806	0,067	38,0	15,6	0,068	76,0	65,6	0,802	0,067	36,4	13,6	-0,603	1,179	0,800	0,067	35,6	13,2	10 <sup>-18</sup>	1,721
	0,90	0,519	0,065	69,2	60,8	0,904	0,069	87,2	68,0	0,072	66,0	51,2	0,899	0,068	85,6	65,6	-0,619	1,214	0,898	0,068	84,8	65,2	10 <sup>-19</sup>	1,716
	1,00	0,576	0,065	53,2	40,8	1,008	0,069	100	96,4	0,076	48,0	32,8	1,002	0,069	99,6	95,6	-0,604	1,179	1,000	0,069	99,6	95,6	10 <sup>-19</sup>	1,726
	1,10	0,617	0,065	43,6	28,8	1,105	0,070	100	100	0,082	34,8	21,6	1,098	0,070	100	100	-0,621	1,235	1,096	0,069	100	100	10 <sup>-18</sup>	1,835
	1,20	0,677	0,066	40,8	25,6	1,198	0,072	100	100	0,086	33,2	21,6	1,191	0,071	100	100	-0,590	1,172	1,189	0,071	100	100	10 <sup>-18</sup>	1,680
30 %	0,70	0,439	0,076	79,2	66,8	0,701	0,078	4,4	1,6	0,087	70,4	56,8	0,696	0,078	4,4	1,6	-0,608	1,171	0,696	0,078	4,4	1,6	10 <sup>-18</sup>	1,667
	0,80	0,530	0,076	64,8	50,8	0,807	0,080	23,6	8,8	0,096	54,4	34,8	0,803	0,080	20,4	7,6	-0,577	1,162	0,802	0,080	20,4	7,6	10 <sup>-19</sup>	1,556
	0,90	0,580	0,076	46,0	27,6	0,907	0,090	71,2	49,6	0,100	26,8	16,4	0,901	0,079	68,4	48,8	-0,594	1,185	0,901	0,079	68,8	48,0	10 <sup>-18</sup>	1,651
	1,00	0,629	0,077	33,6	19,6	1,001	0,080	95,6	92,4	0,105	20,8	11,2	0,995	0,080	95,2	91,2	-0,631	1,251	0,994	0,080	95,2	90,8	10 <sup>-18</sup>	1,781
	1,10	0,688	0,077	29,2	14,4	1,113	0,081	99,6	99,2	0,110	12,4	3,6	1,106	0,081	99,6	99,2	-0,636	1,240	1,106	0,081	99,6	99,2	10 <sup>-18</sup>	1,744
	1,20	0,773	0,077	39,6	24,4	1,198	0,082	100	100	0,114	22,0	10,4	1,190	0,082	100	100	-0,566	1,135	1,189	0,082	100	100	10 <sup>-18</sup>	1,533
60 %	0,70	0,546	0,102	30,8	14,4	0,700	0,103	6,0	0,8	0,103	34,0	14,0	0,696	0,105	6,0	0,8	-0,604	1,190	0,697	0,104	5,6	0,8	10 <sup>-18</sup>	1,556
	0,80	0,621	0,103	14,0	6,0	0,805	0,104	18,0	6,0	0,110	18,4	7,2	0,800	0,104	17,6	4,4	-0,612	1,209	0,800	0,104	17,6	5,2	10 <sup>-18</sup>	1,551
	0,90	0,703	0,104	11,6	4,8	0,907	0,105	55,6	27,2	0,111	10,4	5,6	0,902	0,106	53,6	26,0	-0,593	1,163	0,902	0,106	54,4	26,0	10 <sup>-18</sup>	1,501
	1,00	0,795	0,105	20,4	11,2	1,000	0,107	83,2	62,8	0,113	19,6	11,6	0,995	0,107	82,0	61,2	-0,563	1,122	0,995	0,107	82,4	61,2	10 <sup>-18</sup>	1,442
	1,10	0,858	0,105	42,0	21,2	1,107	0,108	96,8	90,4	0,121	36,0	20,0	1,101	0,108	96,8	89,2	-0,624	1,220	1,101	0,108	96,8	89,2	10 <sup>-19</sup>	1,600
	1,20	0,950	0,107	63,6	46,0	1,214	0,109	99,6	98,0	0,122	54,8	37,6	1,206	0,109	99,6	97,6	-0,569	1,127	1,207	0,109	99,6	97,6	10 <sup>-19</sup>	1,447
90 %	0,70	0,684	0,211	5,6	0,4	0,724	0,212	5,6	0,8	0,189	16,0	6,8	0,721	0,212	6,0	0,8	-0,499	1,118	0,721	0,212	6,0	0,8	10 <sup>-18</sup>	1,146
	0,80	0,762	0,214	6,0	1,2	0,804	0,215	7,2	2,8	0,215	12,0	4,4	0,801	0,215	7,2	2,4	-0,508	1,140	0,801	0,215	7,2	2,4	10 <sup>-18</sup>	1,168
	0,90	0,880	0,219	12,0	2,4	0,926	0,220	16,0	5,2	0,202	20,0	11,6	0,922	0,220	15,6	5,2	-0,482	1,074	0,923	0,220	15,6	5,2	10 <sup>-19</sup>	1,098
	1,00	0,967	0,223	20,0	7,2	1,024	0,224	28,4	11,2	0,192	32,0	18,8	1,019	0,224	27,2	10,4	-0,492	1,098	1,019	0,224	27,6	10,8	10 <sup>-18</sup>	1,128
	1,10	1,078	0,228	40,0	15,2	1,129	0,229	52,4	24,4	0,202	50,0	31,6	1,125	0,229	51,2	23,6	-0,484	1,078	1,125	0,229	51,2	23,6	10 <sup>-18</sup>	1,098
	1,20	1,131	0,231	47,2	20,0	1,199	0,232	58,4	30,8	0,200	56,0	36,0	1,193	0,232	56,8	27,6	-0,506	1,124	1,193	0,232	57,2	28,4	10 <sup>-18</sup>	1,153

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>c</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; <sup>d</sup> : espérance de la fragilité; <sup>e</sup> : variance de la fragilité.

ANNEXE E. RÉSULTATS COMPLETS DES SIMULATIONS

TAB. E.12 – AIC et (entre parenthèses) taux de recouvrement de l'intervalle de confiance à 95 % de  $\hat{\sigma}^2$  pour  $\beta_0 = 0, 7$ .

(a)  $\sigma^2 = 1$  et premier regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,70	31	1 074	872 (100)	867 (100)
	0,80	40	1 072	869 (100)	864 (100)
	0,90	48	1 085	883 (98,1)	878 (98,1)
	1,00	55	1 100	908 (99,1)	902 (99,5)
	1,10	67	1 130	928 (99,5)	923 (99,5)
	1,20	80	1 140	939 (100)	934 (100)
30 %	0,70	33	808	602 (97,6)	599 (99,6)
	0,80	44	824	621 (97,2)	618 (100)
	0,90	50	821	624 (96,8)	621 (99,5)
	1,00	65	845	644 (97,7)	640 (100)
	1,10	77	860	655 (96,8)	651 (100)
	1,20	96	859	658 (97,7)	654 (100)
60 %	0,70	29	523	316 (84,3)	320 (100)
	0,80	38	527	320 (88,5)	324 (100)
	0,90	48	538	331 (85,7)	334 (100)
	1,00	58	549	342 (88,1)	345 (100)
	1,10	69	554	347 (88,1)	350 (100)
	1,20	81	570	364 (88,5)	367 (100)
90 %	0,70	11	200	30 (70,4)	46 (100)
	0,80	15	205	34 (68,1)	51 (100)
	0,90	17	207	37 (69,1)	54 (100)
	1,00	21	212	42 (65,0)	59 (100)
	1,10	26	215	46 (66,3)	62 (100)
	1,20	29	220	50 (68,6)	66 (100)

(b)  $\sigma^2 = 1$  et second regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,70	41	841	820 (98,6)	819 (100)
	0,80	53	857	835 (99,5)	835 (100)
	0,90	65	872	850 (99,1)	850 (100)
	1,00	80	890	869 (99,5)	868 (100)
	1,10	91	919	897 (98,6)	897 (100)
	1,20	108	913	891 (99,5)	891 (100)
30 %	0,70	35	670	648 (99,5)	648 (100)
	0,80	51	646	624 (99,5)	623 (100)
	0,90	60	680	658 (99,5)	658 (100)
	1,00	71	715	693 (98,6)	693 (100)
	1,10	84	741	719 (98,6)	719 (100)
	1,20	105	704	682 (100)	682 (100)
60 %	0,70	30	401	379 (99,1)	379 (100)
	0,80	39	416	393 (98,6)	393 (100)
	0,90	49	419	399 (99,5)	399 (100)
	1,00	62	416	394 (99,5)	394 (100)
	1,10	71	452	430 (98,1)	430 (100)
	1,20	86	443	421 (100)	421 (100)
90 %	0,70	12	107	858 (95,0)	84 (100)
	0,80	14	112	90 (95,0)	90 (100)
	0,90	18	112	90 (95,9)	90 (100)
	1,00	22	117	95 (97,2)	268 (100)
	1,10	26	121	99 (95,9)	99 (100)
	1,20	28	121	104 (93,3)	104 (100)

(c)  $\sigma^2 = 4$  et premier regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,70	5	1 833	2 296 (55,1)	2 287 (0,1)
	0,80	7	1 818	2 305 (57,1)	2 296 (0,3)
	0,90	9	2 119	2 295 (55,6)	2 287 (1,4)
	1,00	11	1 929	2 293 (58,5)	2 285 (1,1)
	1,10	13	2 045	2 310 (57,8)	2 302 (2,0)
	1,20	14	1 866	2 299 (61,4)	2 290 (1,4)
30 %	0,70	7	1 200	1 601 (69,8)	1 597 (42,1)
	0,80	9	1 384	1 614 (66,5)	1 610 (50,1)
	0,90	11	1 652	1 620 (71,1)	1 615 (56,7)
	1,00	13	1 513	1 622 (70,1)	1 618 (65,6)
	1,10	16	1 482	1 626 (74,7)	1 621 (73,6)
	1,20	18	1 517	1 625 (70,9)	1 621 (70,3)
60 %	0,70	9	1 014	857 (73,6)	865 (80,0)
	0,80	13	1 011	865 (81,6)	873 (88,0)
	0,90	16	995	867 (76,8)	875 (83,2)
	1,00	19	1 063	872 (80,0)	880 (87,2)
	1,10	24	1 061	874 (80,0)	882 (85,6)
	1,20	28	1 057	874 (84,0)	881 (88,0)
90 %	0,70	8	319	120 (42,7)	165 (99,5)
	0,80	11	322	126 (45,4)	171 (99,5)
	0,90	12	324	125 (43,1)	170 (100)
	1,00	15	324	126 (48,1)	171 (99,5)
	1,10	18	326	130 (45,9)	174 (100)
	1,20	23	329	132 (44,5)	176 (99,5)

(d)  $\sigma^2 = 4$  et second regroupement

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,70	14	1 940	1 983 (95,9)	1 983 (45,9)
	0,80	17	1 949	1 967 (97,2)	1 967 (54,5)
	0,90	21	1 793	1 980 (95,4)	1 979 (64,1)
	1,00	25	1 884	1 993 (95,4)	1 993 (66,8)
	1,10	29	1 989	2 000 (94,5)	2 000 (70,9)
	1,20	37	1 943	1 966 (95,4)	1 965 (87,7)
30 %	0,70	10	1 526	1 587 (98,1)	1 587 (64,1)
	0,80	12	1 437	1 608 (99,5)	1 608 (70,0)
	0,90	16	1 308	1 578 (98,1)	1 578 (80,9)
	1,00	21	1 478	1 629 (96,3)	1 628 (82,2)
	1,10	26	1 515	1 575 (98,6)	1 575 (92,7)
	1,20	27	1 414	1 664 (97,7)	1 664 (90,9)
60 %	0,70	12	933	939 (96,8)	939 (90,9)
	0,80	15	929	935 (98,6)	935 (95,0)
	0,90	20	943	930 (98,6)	930 (97,2)
	1,00	24	946	939 (98,1)	939 (98,1)
	1,10	30	946	935 (98,1)	935 (98,6)
	1,20	33	966	960 (99,1)	960 (99,5)
90 %	0,70	9	174	209 (82,7)	210 (94,1)
	0,80	12	141	211 (81,8)	212 (93,1)
	0,90	15	177	207 (85,4)	208 (96,3)
	1,00	17	193	219 (79,1)	220 (93,6)
	1,10	20	123	223 (81,3)	224 (95,9)
	1,20	24	203	215 (85,9)	216 (96,8)

TAB. E.13 – Estimation des paramètres et puissance du test du maximum de vraisemblance – premier regroupement,  $\beta_0 = 0, 7$  et fragilité gamma de variance  $\sigma^2 = 4$ .

Censure	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne										
	$\beta$	$\hat{\beta}$	EC*	1 % <sup>bb</sup>	$\beta$	$\hat{\beta}$	EC*	1 % <sup>bb</sup>	EC*	5 % <sup>b</sup>	1 % <sup>bb</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>bb</sup>	EC*	5 % <sup>b</sup>	1 % <sup>bb</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>bb</sup>	EC*	5 % <sup>b</sup>	1 % <sup>bb</sup>	Var.**	
0 %	0,70	0,135	0,063	100	100	0,706	0,073	16,5	7,5	0,061	100	100	0,699	0,071	7,5	0,7	-3,116	4,435	0,694	0,071	1,5	0,7	10 <sup>-18</sup>	16,907			
	0,80	0,152	0,063	100	100	0,815	0,074	47,6	34,9	0,061	100	100	0,804	0,071	32,2	13,4	-0,312	4,455	0,799	0,071	28,8	13,4	10 <sup>-18</sup>	17,137			
	0,90	0,182	0,063	100	100	0,921	0,075	89,2	77,8	0,061	100	100	0,901	0,071	82,8	66,4	-3,091	4,415	0,896	0,071	80,7	66,4	10 <sup>-18</sup>	16,871			
	1,00	0,201	0,063	100	100	1,026	0,076	93,7	93,0	0,061	100	100	0,997	0,072	97,9	90,2	-3,070	4,400	0,990	0,072	97,2	89,6	10 <sup>-18</sup>	16,824			
	1,10	0,218	0,063	100	100	1,131	0,077	95,8	95,8	0,062	100	100	1,093	0,072	100	100	-3,100	4,425	1,085	0,072	100	100	10 <sup>-18</sup>	17,161			
	1,20	0,235	0,063	100	100	1,220	0,078	97,4	96,7	0,061	100	100	1,192	0,072	100	100	-3,065	4,387	1,184	0,072	100	100	10 <sup>-18</sup>	16,734			
30 %	0,70	0,188	0,075	100	100	0,727	0,087	13,8	6,9	0,076	100	100	0,697	0,084	6,9	1,7	-2,556	4,406	0,693	0,084	6,3	1,7	10 <sup>-19</sup>	8,675			
	0,80	0,212	0,075	100	100	0,833	0,088	42,7	27,1	0,076	100	100	0,797	0,080	25,3	10,2	-2,586	4,446	0,792	0,085	24,1	8,4	10 <sup>-18</sup>	8,828			
	0,90	0,241	0,075	100	100	0,951	0,100	79,0	62,2	0,076	100	99,4	0,900	0,085	64,6	43,1	-2,590	4,450	0,900	0,085	61,6	40,1	10 <sup>-18</sup>	8,895			
	1,00	0,268	0,076	100	100	1,051	0,090	96,0	90,4	0,077	100	100	0,997	0,086	93,2	81,9	-2,584	4,441	0,990	0,086	92,1	76,8	10 <sup>-18</sup>	8,889			
	1,10	0,295	0,076	100	100	1,174	0,091	99,4	99,4	0,078	100	100	1,108	0,086	100	98,8	-2,587	4,436	1,101	0,086	100	98,8	10 <sup>-18</sup>	8,873			
	1,20	0,317	0,076	100	98,8	1,264	0,092	99,4	99,4	0,078	99,4	97,6	1,198	0,087	100	100	-2,579	4,424	1,191	0,087	100	100	10 <sup>-18</sup>	8,826			
60 %	0,70	0,292	0,100	97,1	90,3	0,747	0,113	8,2	1,9	0,101	96,6	89,8	0,695	0,111	3,8	0,5	-1,954	4,434	0,689	0,111	4,3	0,9	10 <sup>-18</sup>	4,586			
	0,80	0,348	0,101	93,3	81,3	0,865	0,114	32,6	13,0	0,103	94,3	77,8	0,806	0,111	14,3	2,6	-1,963	4,448	0,800	0,111	11,7	2,6	10 <sup>-18</sup>	4,631			
	0,90	0,390	0,101	88,0	75,3	0,960	0,115	59,1	39,0	0,105	86,2	69,7	0,893	0,112	39,9	19,2	-1,978	4,460	0,886	0,113	38,0	17,9	10 <sup>-18</sup>	4,645			
	1,00	0,430	0,101	73,3	52,0	1,067	0,116	90,2	72,8	0,105	74,6	52,0	0,990	0,113	75,5	53,7	-1,976	4,442	0,982	0,113	73,3	51,1	10 <sup>-18</sup>	4,670			
	1,10	0,485	0,101	53,2	33,3	1,199	0,118	97,7	94,9	0,107	50,4	31,0	1,112	0,115	95,3	82,8	-1,965	4,393	1,103	0,115	94,9	81,4	10 <sup>-19</sup>	4,646			
	1,20	0,526	0,102	45,5	24,7	1,285	0,119	100	99,1	0,110	41,1	19,0	1,184	0,115	99,5	97,7	-1,959	4,366	1,184	0,115	99,5	97,3	10 <sup>-18</sup>	4,641			
90 %	0,70	0,568	0,208	11,2	3,2	0,764	0,223	6,8	3,2	0,205	13,6	4,4	0,718	0,219	6,0	1,2	-1,082	4,653	0,705	0,219	5,6	1,2	10 <sup>-18</sup>	2,075			
	0,80	0,652	0,210	5,6	2,4	0,864	0,225	12,8	3,6	0,207	7,2	2,0	0,810	0,221	5,6	1,2	-1,109	4,763	0,799	0,221	6,0	1,2	10 <sup>-20</sup>	2,111			
	0,90	0,697	0,212	3,6	1,2	0,949	0,227	20,4	9,2	0,206	5,2	1,6	0,885	0,223	14,0	4,0	-1,099	4,692	0,874	0,223	12,8	3,2	10 <sup>-19</sup>	2,097			
	1,00	0,796	0,216	7,6	0,4	1,082	0,231	39,6	18,0	0,211	11,2	1,2	1,009	0,227	26,4	11,6	-1,085	4,620	0,996	0,227	26,8	10,0	10 <sup>-19</sup>	2,083			
	1,10	0,883	0,219	12,4	2,4	1,183	0,235	52,0	29,6	0,214	12,8	4,0	1,105	0,230	40,4	18,0	-1,094	4,639	1,090	0,230	36,4	16,8	10 <sup>-18</sup>	2,094			
	1,20	1,000	0,224	23,7	52,2	1,293	0,239	67,8	46,5	0,222	26,1	8,0	1,221	0,235	58,6	32,1	-1,068	4,524	1,205	0,235	57,8	29,7	10 <sup>-18</sup>	2,068			

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>bb</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.14 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0, 7$  et fragilité gamma de variance  $\sigma^2 = 4$ .

Censure	$\beta$	Modèle naïf sans effet groupe			Modèle naïf avec effet groupe			Modèle marginal			Modèle de fragilité gamma			Modèle de fragilité gaussienne									
		$\hat{\beta}$	EC*	1 % <sup>b</sup>	$\hat{\beta}$	EC*	1 % <sup>b</sup>	EC*	5 % <sup>c</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>c</sup>	1 % <sup>b</sup>	Esp.*	Var.**							
0 %	0,70	0,227	0,064	100	0,694	0,067	8,9	4,9	0,055	100	99,6	0,701	0,061	6,1	1,6	-2,831	4,144	0,701	0,066	6,1	1,6	10 <sup>-17</sup>	15,292
	0,80	0,254	0,064	100	0,783	0,067	33,4	18,0	0,059	99,5	99,5	0,792	0,066	28,0	13,8	-2,806	4,079	0,791	0,066	28,4	13,8	10 <sup>-18</sup>	14,304
	0,90	0,283	0,064	99,5	0,889	0,068	80,6	57,5	0,060	98,3	95,7	0,898	0,067	84,4	62,6	-2,816	4,183	0,898	0,067	83,6	62,6	10 <sup>-19</sup>	15,090
	1,00	0,313	0,064	98,7	0,994	0,069	95,0	90,4	0,065	97,9	96,2	1,002	0,068	98,7	93,7	-2,787	4,095	1,001	0,068	98,7	93,7	10 <sup>-18</sup>	14,982
	1,10	0,340	0,064	98,3	1,097	0,070	97,9	97,5	0,064	95,0	92,1	1,105	0,069	100	100	-2,700	3,932	1,104	0,068	100	100	10 <sup>-17</sup>	14,514
	1,20	0,382	0,064	96,2	1,204	0,072	97,1	96,6	0,067	91,3	87,6	1,210	0,069	100	100	-2,662	3,914	1,210	0,069	100	100	10 <sup>-18</sup>	13,918
30 %	0,70	0,224	0,075	100	0,697	0,079	6,1	1,6	0,083	100	98,3	0,701	0,078	5,3	0,8	-2,725	4,342	0,701	0,078	5,3	0,8	10 <sup>-18</sup>	11,720
	0,80	0,251	0,076	100	0,792	0,080	25,9	9,8	0,091	98,3	94,2	0,799	0,079	23,0	7,4	-2,740	4,279	0,799	0,079	23,0	7,4	10 <sup>-18</sup>	12,028
	0,90	0,285	0,075	98,7	0,897	0,081	72,2	50,6	0,096	93,4	85,7	0,901	0,079	71,4	50,6	-2,634	4,173	0,901	0,079	71,4	50,2	10 <sup>-17</sup>	11,380
	1,00	0,325	0,076	94,7	1,008	0,082	94,3	86,6	0,097	87,8	78,1	1,011	0,080	96,3	89,0	-2,776	4,526	1,011	0,080	96,3	89,0	10 <sup>-18</sup>	12,011
	1,10	0,368	0,076	91,0	1,101	0,083	96,3	95,9	0,109	81,7	72,7	1,106	0,081	100	100	-2,606	4,152	1,106	0,081	100	100	10 <sup>-18</sup>	10,948
	1,20	0,372	0,076	87,3	1,191	0,084	97,9	97,1	0,110	76,7	65,7	1,206	0,082	100	99,6	-2,868	4,518	1,207	0,082	100	99,6	10 <sup>-18</sup>	12,256
60 %	0,70	0,323	0,101	88,6	0,699	0,104	6,0	0,8	0,115	82,1	66,8	0,697	0,104	5,2	0,8	-2,390	4,536	0,697	0,104	5,2	0,8	10 <sup>-18</sup>	7,568
	0,80	0,370	0,101	75,1	0,791	0,105	16,0	8,0	0,122	63,0	46,5	0,789	0,105	14,8	7,2	-2,335	4,393	0,790	0,105	14,8	7,2	10 <sup>-18</sup>	7,344
	0,90	0,430	0,101	60,6	0,894	0,106	47,8	23,7	0,126	46,5	28,5	0,890	0,106	44,9	22,5	-2,337	4,387	0,891	0,106	45,7	22,4	10 <sup>-17</sup>	7,538
	1,00	0,466	0,101	56,0	0,992	0,107	83,0	66,1	0,137	37,5	19,7	0,992	0,107	81,8	64,5	-2,382	4,418	0,993	0,107	82,2	64,5	10 <sup>-18</sup>	7,622
	1,10	0,528	0,102	45,7	1,112	0,108	93,9	87,1	0,144	30,9	17,2	1,109	0,108	94,7	85,9	-2,308	4,243	1,110	0,108	94,7	85,9	10 <sup>-18</sup>	7,330
	1,20	0,551	0,102	34,4	1,194	0,110	99,2	98,8	0,149	20,0	13,2	1,187	0,109	100	100	-2,408	4,499	1,188	0,109	100	99,2	10 <sup>-18</sup>	7,763
90 %	0,70	0,597	0,209	7,5	1,2	0,210	3,7	0,8	0,165	19,5	11,6	0,700	0,211	3,7	0,8	-1,739	4,539	0,700	0,211	3,7	0,8	10 <sup>-18</sup>	3,994
	0,80	0,701	0,212	4,5	0,8	0,214	5,4	1,6	0,172	13,2	5,4	0,832	0,214	5,4	1,6	-1,685	4,322	0,832	0,214	5,4	1,2	10 <sup>-18</sup>	3,834
	0,90	0,782	0,215	5,8	1,6	0,216	14,5	4,9	0,170	14,5	10,3	0,906	0,217	14,5	4,9	-1,633	4,157	0,905	0,217	14,5	4,9	10 <sup>-18</sup>	3,759
	1,00	0,854	0,218	11,2	2,0	0,219	26,6	10,8	0,176	27,0	16,2	1,002	0,220	25,4	10,4	-1,737	4,495	1,001	0,220	25,4	10,4	10 <sup>-18</sup>	3,981
	1,10	0,936	0,221	17,2	5,4	1,103	40,7	17,6	0,180	31,9	19,3	1,101	0,223	41,1	16,8	-1,741	4,577	1,101	0,223	41,1	17,2	10 <sup>-18</sup>	3,986
	1,20	1,024	0,226	33,2	14,9	1,190	65,5	35,2	0,180	48,5	31,9	1,189	0,228	65,9	34,4	-1,618	4,162	1,189	0,228	66,4	35,2	10 <sup>-19</sup>	3,704

\* : estimation de l'écart-type de  $\hat{\beta}$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>c</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.15 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0$  et fragilité stable positive.

Censure	Modèle naïf sans effet groupe				Modèle naïf avec effet groupe				Modèle marginal				Modèle de fragilité gamma				Modèle de fragilité gaussienne							
	$\beta$	$\hat{\beta}$	EC*	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Exp.*	Var.**	$\beta$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Exp.*	Var.**	
0 %	0.00	0.000	0.063	3.1	0.7	0.003	0.064	3.1	0.5	0.052	8.9	4.3	0.003	0.064	3.4	0.5	-1.886	2.834	0.003	0.064	3.4	0.5	10 <sup>-18</sup>	3.964
	0.10	0.005	0.063	13.4	6.1	0.106	0.064	34.1	16.1	0.055	20.1	11.6	0.108	0.064	35.1	16.9	-1.908	2.834	0.108	0.064	34.9	16.7	10 <sup>-19</sup>	3.903
	0.20	0.121	0.063	41.2	21.8	0.219	0.064	76.5	67.4	0.056	46.9	32.7	0.222	0.064	77.2	67.6	-1.863	2.823	0.222	0.064	77.2	67.4	10 <sup>-18</sup>	3.910
	0.30	0.185	0.063	68.5	53.6	0.326	0.064	84.5	83.8	0.059	70.3	57.2	0.329	0.064	84.5	84.1	-1.954	2.970	0.329	0.064	84.5	84.1	10 <sup>-18</sup>	4.145
	0.40	0.243	0.064	78.3	70.3	0.426	0.065	82.7	82.7	0.064	77.1	67.8	0.434	0.065	82.7	82.7	-1.875	2.820	0.434	0.065	82.7	82.7	10 <sup>-18</sup>	3.847
0.50	0.302	0.064	85.1	82.9	0.542	0.066	86.1	86.1	0.067	84.0	79.4	0.548	0.065	86.1	86.1	-1.785	2.765	0.547	0.065	86.1	86.1	10 <sup>-18</sup>	3.694	
30 %	0.00	-0.001	0.075	4.3	0.7	-0.001	0.076	4.1	0.5	0.068	11.8	4.5	-0.001	0.076	4.3	0.5	-1.890	2.792	-0.001	0.076	4.3	0.5	10 <sup>-18</sup>	3.901
	0.10	0.060	0.075	10.5	3.2	0.110	0.076	25.1	10.3	0.068	17.6	7.4	0.111	0.076	24.9	10.5	-1.898	2.897	0.111	0.076	25.2	10.5	10 <sup>-18</sup>	3.873
	0.20	0.117	0.075	28.7	13.8	0.214	0.076	66.3	46.9	0.070	34.5	20.5	0.217	0.076	67.2	48.5	-2.095	3.124	0.218	0.076	67.2	48.7	10 <sup>-18</sup>	4.339
	0.30	0.181	0.075	53.1	36.1	0.327	0.077	82.1	80.0	0.074	53.2	38.3	0.331	0.077	82.3	80.7	-1.798	2.722	0.331	0.077	82.3	80.5	10 <sup>-18</sup>	3.731
	0.40	0.231	0.076	66.3	54.1	0.425	0.077	83.4	83.1	0.079	65.1	49.6	0.433	0.077	83.4	83.2	-1.920	2.865	0.433	0.077	83.4	83.2	10 <sup>-18</sup>	4.043
0.50	0.309	0.076	84.0	77.1	0.546	0.078	88.0	88.0	0.085	80.5	70.0	0.552	0.077	88.0	88.0	-1.850	2.775	0.552	0.077	88.0	88.0	10 <sup>-18</sup>	3.825	
60 %	0.00	0.011	0.100	3.8	0.7	0.008	0.101	4.0	0.7	0.090	9.2	3.2	0.010	0.101	9.2	3.2	-1.866	2.850	0.008	0.101	3.8	0.5	10 <sup>-18</sup>	3.966
	0.10	0.062	0.100	9.8	2.5	0.111	0.101	16.0	6.0	0.089	14.3	16.2	0.112	0.101	15.8	5.8	-1.798	2.791	0.112	0.101	15.8	5.8	10 <sup>-18</sup>	3.693
	0.20	0.129	0.100	22.0	9.2	0.222	0.101	48.9	28.0	0.092	27.8	15.2	0.223	0.101	49.8	28.5	-1.891	2.910	0.223	0.101	50.0	28.5	10 <sup>-18</sup>	3.976
	0.30	0.182	0.100	40.1	20.0	0.327	0.101	78.3	64.1	0.097	41.2	25.4	0.331	0.101	79.8	65.4	-1.876	2.899	0.332	0.101	80.0	66.0	10 <sup>-18</sup>	3.941
	0.40	0.248	0.100	54.9	39.6	0.434	0.102	83.4	80.0	0.104	52.0	36.5	0.440	0.102	83.4	81.2	-1.813	2.801	0.440	0.102	83.6	81.2	10 <sup>-18</sup>	3.846
0.50	0.317	0.101	66.5	53.1	0.550	0.103	82.1	82.0	0.111	62.3	50.3	0.549	0.103	82.3	82.3	-1.800	2.813	0.551	0.103	82.3	82.3	10 <sup>-18</sup>	3.813	
90 %	0.00	0.011	0.201	3.2	0.1	0.005	0.201	2.7	0.9	0.158	10.7	5.4	0.006	0.202	2.7	0.9	-1.344	3.288	0.006	0.202	2.7	0.9	10 <sup>-18</sup>	2.797
	0.10	0.095	0.201	6.5	1.1	0.114	0.202	6.5	1.8	0.161	15.2	7.8	0.113	0.202	6.3	1.8	-1.379	3.375	0.113	0.202	6.3	1.8	10 <sup>-19</sup>	2.863
	0.20	0.178	0.201	9.6	2.3	0.212	0.202	11.6	4.3	0.159	20.0	10.7	0.211	0.203	11.4	4.0	-1.374	3.445	0.211	0.203	11.4	4.0	10 <sup>-18</sup>	2.878
	0.30	0.284	0.203	19.8	8.1	0.333	0.204	23.1	10.3	0.159	30.5	20.3	0.332	0.204	23.2	10.0	-1.353	3.330	0.332	0.204	23.1	10.0	10 <sup>-18</sup>	2.818
	0.40	0.373	0.204	28.9	14.1	0.442	0.205	38.5	21.8	0.158	42.5	30.5	0.439	0.205	38.1	21.6	-1.329	3.166	0.440	0.205	38.3	21.4	10 <sup>-19</sup>	2.753
0.50	0.467	0.206	40.7	24.9	0.563	0.207	51.2	38.9	0.159	50.1	39.6	0.560	0.207	51.1	37.8	-1.420	3.658	0.561	0.207	51.1	38.1	10 <sup>-18</sup>	2.938	

\* : estimation de l'écart-type de  $\beta$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

TAB. E.16 – Estimation des paramètres et puissance du test du maximum de vraisemblance – second regroupement,  $\beta_0 = 0, 7$  et fragilité stable positive.

Censure	Modèle naïf sans effet groupe					Modèle naïf avec effet groupe					Modèle marginal					Modèle de fragilité gamma					Modèle de fragilité gaussienne						
	$\beta$	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	$\hat{\beta}$	EC*	5 % <sup>b</sup>	1 % <sup>b</sup>	Esp.*	Var.**		
0 %	0,70	0,427	0,065	78,3	72,3	0,757	0,067	16,1	7,2	0,080	73,4	58,9	0,769	0,066	15,2	6,3	-1,864	2,857	0,769	0,066	15,2	6,7	10 <sup>-18</sup>	6,7	10 <sup>-18</sup>	3,835	3,835
	0,80	0,488	0,065	68,5	57,8	0,872	0,068	62,0	43,6	0,084	55,8	41,2	0,882	0,067	66,3	47,1	-1,913	2,912	0,881	0,067	66,3	46,9	10 <sup>-19</sup>	46,9	10 <sup>-19</sup>	3,943	3,943
	0,90	0,543	0,065	52,3	41,6	0,987	0,069	80,9	70,3	0,089	38,7	26,1	0,996	0,068	84,3	82,0	-1,765	2,709	0,996	0,068	84,3	81,8	10 <sup>-18</sup>	81,8	10 <sup>-18</sup>	3,698	3,698
	1,00	0,620	0,066	34,5	21,1	1,108	0,070	84,0	82,9	0,096	19,1	10,1	1,099	0,069	86,5	86,5	-1,774	2,688	1,099	0,069	86,5	86,5	10 <sup>-19</sup>	86,5	10 <sup>-19</sup>	3,634	3,634
	1,10	0,661	0,066	27,8	14,5	1,190	0,071	86,1	85,4	0,101	14,5	7,1	1,204	0,070	87,1	87,1	-1,847	2,834	1,204	0,070	87,1	87,1	10 <sup>-18</sup>	87,1	10 <sup>-18</sup>	3,808	3,808
	1,20	0,725	0,067	28,1	16,9	1,300	0,073	85,1	85,1	0,113	12,1	4,5	1,322	0,071	85,8	85,8	-1,886	2,845	1,322	0,071	85,8	85,8	10 <sup>-19</sup>	85,8	10 <sup>-19</sup>	3,819	3,819
30 %	0,70	0,406	0,076	76,7	68,3	0,756	0,079	14,3	4,5	0,099	67,4	50,0	0,766	0,078	10,9	4,0	-1,876	2,793	0,767	0,078	11,1	4,0	10 <sup>-18</sup>	4,0	10 <sup>-18</sup>	3,899	3,899
	0,80	0,485	0,076	58,5	48,1	0,871	0,080	51,8	30,3	0,080	41,6	24,9	0,883	0,079	54,1	31,4	-1,858	2,782	0,883	0,079	54,0	31,4	10 <sup>-18</sup>	31,4	10 <sup>-18</sup>	3,862	3,862
	0,90	0,541	0,076	48,3	35,4	0,977	0,081	77,4	70,3	0,115	27,4	16,0	0,989	0,080	81,6	73,1	-1,904	2,830	0,989	0,080	81,8	73,2	10 <sup>-19</sup>	73,2	10 <sup>-19</sup>	3,931	3,931
	1,00	0,591	0,076	38,7	24,0	1,087	0,082	81,8	81,1	0,123	17,8	8,7	1,101	0,081	86,1	86,1	-1,897	2,848	1,102	0,081	86,1	86,1	10 <sup>-18</sup>	86,1	10 <sup>-18</sup>	3,921	3,921
	1,10	0,65	0,077	28,1	16,0	1,186	0,083	83,8	80,9	0,132	10,1	5,6	1,207	0,081	85,8	85,8	-1,941	2,892	1,208	0,082	85,8	85,8	10 <sup>-18</sup>	85,8	10 <sup>-18</sup>	3,963	3,963
	1,20	0,719	0,077	27,1	16,5	1,304	0,084	84,0	83,8	0,142	9,4	4,3	1,317	0,082	84,9	84,9	-1,875	2,822	1,318	0,082	84,9	84,9	10 <sup>-18</sup>	84,9	10 <sup>-18</sup>	3,833	3,833
60 %	0,70	0,428	0,101	59,4	43,6	0,762	0,105	9,8	2,3	0,127	46,3	30,7	0,765	0,104	7,1	1,6	-1,909	2,939	0,766	0,105	7,2	1,8	10 <sup>-18</sup>	1,8	10 <sup>-18</sup>	3,981	3,981
	0,80	0,495	0,102	44,3	31,4	0,875	0,106	33,2	14,1	0,137	30,7	17,8	0,881	0,106	32,9	13,2	-1,881	2,870	0,883	0,106	33,2	13,4	10 <sup>-18</sup>	13,4	10 <sup>-18</sup>	3,970	3,970
	0,90	0,555	0,102	39,1	25,6	0,973	0,107	62,5	45,4	0,146	22,1	12,1	0,982	0,107	65,2	45,4	-1,928	2,954	0,984	0,107	65,2	46,5	10 <sup>-18</sup>	46,5	10 <sup>-18</sup>	4,003	4,003
	1,00	0,611	0,103	28,0	18,0	1,089	0,110	80,7	74,3	0,156	16,0	9,2	1,097	0,108	83,2	77,1	-1,928	2,915	1,100	0,108	83,2	77,2	10 <sup>-18</sup>	77,2	10 <sup>-18</sup>	4,022	4,022
	1,10	0,676	0,104	28,9	15,4	1,189	0,110	80,7	80,1	0,169	13,6	7,1	1,207	0,110	86,0	85,6	-1,948	2,944	1,210	0,110	86,0	85,6	10 <sup>-18</sup>	85,6	10 <sup>-18</sup>	4,005	4,005
	1,20	0,739	0,104	30,3	21,2	1,310	0,112	83,8	83,4	0,181	15,8	10,1	1,314	0,111	85,8	85,8	-1,894	2,867	1,317	0,112	85,8	85,8	10 <sup>-18</sup>	85,8	10 <sup>-18</sup>	3,905	3,905
90 %	0,70	0,654	0,211	4,7	1,4	0,774	0,212	3,6	0,7	0,166	16,1	9,4	0,771	0,213	3,4	0,7	-1,379	3,351	0,772	0,213	3,4	0,7	10 <sup>-18</sup>	0,7	10 <sup>-18</sup>	2,860	2,860
	0,80	0,749	0,214	4,7	0,5	0,890	0,215	10,0	1,6	0,168	14,7	8,0	0,886	0,216	9,8	1,6	-1,347	3,269	0,887	0,216	9,8	1,6	10 <sup>-18</sup>	1,6	10 <sup>-18</sup>	2,793	2,793
	0,90	0,849	0,218	8,3	1,8	0,987	0,219	16,7	5,4	0,179	18,5	5,4	0,984	0,220	15,6	5,4	-1,302	3,086	0,981	0,219	16,1	5,6	10 <sup>-18</sup>	5,6	10 <sup>-18</sup>	2,684	2,684
	1,00	0,946	0,222	13,6	6,0	1,096	0,224	28,3	14,0	0,180	23,8	14,1	1,093	0,224	29,1	12,7	-1,342	3,194	1,089	0,222	28,9	13,6	10 <sup>-18</sup>	13,6	10 <sup>-18</sup>	2,796	2,796
	1,10	1,028	0,226	20,7	9,8	1,228	0,229	43,4	26,9	0,186	31,1	21,1	1,227	0,229	43,6	26,9	-1,454	3,637	1,220	0,229	43,1	27,2	10 <sup>-18</sup>	27,2	10 <sup>-18</sup>	3,028	3,028
	1,20	1,114	0,231	28,9	14,5	1,331	0,233	52,7	39,1	0,185	40,0	27,6	1,328	0,233	52,5	39,1	-1,415	3,485	1,329	0,233	52,5	39,1	10 <sup>-18</sup>	39,1	10 <sup>-18</sup>	2,945	2,945

\* : estimation de l'écart-type de  $\hat{\beta}$ ; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 5 %; <sup>b</sup> : puissance du test du maximum de vraisemblance au seuil de 1 %; \* : espérance de la fragilité; \*\* : variance de la fragilité.

ANNEXE E. RÉSULTATS COMPLETS DES SIMULATIONS

TAB. E.17 – Fragilité stable positive – second regroupement : AIC et (entre parenthèses) taux de recouvrement de l'intervalle de confiance à 95 % de  $\hat{\sigma}^2$ .

(a)  $\beta = 0$

(b)  $\beta = 0,7$

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,00	1	1 039	1 127 (80,7)	1 127 (30,1)
	0,10	2	1 017	1 126 (84,1)	1 126 (61,1)
	0,20	4	820	1 112 (83,4)	1 112 (75,8)
	0,30	9	956	1 143 (84,1)	1 143 (79,1)
	0,40	15	975	1 116 (82,5)	1 116 (81,8)
	0,50	23	997	1 124 (85,6)	1 124 (85,8)
	30 %	0,00	1	775	1 003 (84,3)
0,10		1	880	1 002 (82,0)	1 002 (49,8)
0,20		3	976	1 054 (80,7)	1 054 (67,2)
0,30		6	899	1 002 (82,3)	1 002 (78,7)
0,40		10	884	1 038 (83,1)	1 038 (80,7)
0,50		18	902	1 007 (87,8)	1 007 (86,7)
60 %		0,00	1	670	753 (82,9)
	0,10	1	674	739 (80,3)	740 (52,0)
	0,20	2	704	749 (83,6)	749 (71,1)
	0,30	4	652	768 (86,1)	768 (78,7)
	0,40	7	700	750 (83,2)	750 (79,8)
	0,50	11	736	757 (81,8)	757 (80,7)
	90 %	0,00	1	164	187 (59,1)
0,10		1	179	193 (62,5)	194 (39,2)
0,20		2	177	189 (60,1)	191 (40,3)
0,30		3	152	190 (60,7)	191 (60,5)
0,40		4	166	191 (58,7)	192 (65,8)
0,50		6	164	199 (58,5)	201 (64,9)

Censure	$\beta$	Modèle naïf		Modèle de fragilité	
		Sans effet groupe	Avec effet groupe	Gamma	Gaussienne
0 %	0,70	45	1 044	1 145 (82,0)	1 145 (82,5)
	0,80	57	972	1 179 (85,2)	1 179 (85,4)
	0,90	70	1 047	1 182 (84,5)	1 182 (85,1)
	1,00	89	1 054	1 181 (86,1)	1 181 (86,5)
	1,10	100	1 094	1 221 (86,1)	1 221 (87,1)
	1,20	118	836	1 250 (84,9)	1 250 (85,8)
	30 %	0,70	30	991	1 055 (83,2)
0,80		42	988	1 047 (82,3)	1 047 (82,3)
0,90		53	968	1 066 (85,6)	1 066 (86,1)
1,00		63	880	1 079 (85,6)	1 079 (86,1)
1,10		76	1 021	1 095 (85,4)	1 095 (85,8)
1,20		91	940	1 090 (84,3)	1 090 (84,9)
60 %		0,70	20	721	780 (84,7)
	0,80	26	756	786 (82,3)	786 (82,5)
	0,90	33	731	797 (86,5)	797 (87,2)
	1,00	39	756	802 (87,1)	802 (87,2)
	1,10	47	771	815 (85,8)	815 (86,1)
	1,20	56	761	811 (85,2)	811 (85,8)
	90 %	0,70	11	135	203 (61,8)
0,80		14	143	203 (63,2)	204 (70,9)
0,90		17	203	202 (59,4)	203 (69,1)
1,00		21	120	208 (63,8)	209 (71,2)
1,10		24	4	223 (62,3)	224 (68,0)
1,20		28	157	222 (61,8)	224 (67,2)

## Annexe F

# COMPLÉMENTS À L'ÉTUDE ÉPIDÉMIOLOGIQUE DES GENN

Nous donnons ici les résultats complémentaires concernant l'étude épidémiologique des GENN, soient :

- l'ajustement du modèle naïf de Cox pour chacun des quatre mois ;
- les ajustements naïf, marginal et mixte avec la prise en compte de possibles interactions.

TAB. F.1a – Ajustement naïf mois par mois : décembre et janvier.

Variable	Modalité	Déc.				Janv.			
		RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>	RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>
Condition de vêlage	aide facile	1,14	0,63 - 2,03	0,29	NS	1,53	0,87 - 2,69	0,28	NS
	aide difficile	1,81	0,96 - 3,41	0,32	NS	1,19	0,53 - 2,66	0,41	NS
Vaccination de la mère contre <i>E. coli</i>	non	1,17	0,75 - 1,85	0,23	NS	1,81	1,13 - 2,92	0,24	0,014
Vaccination de la mère contre d'autres agents que GENN	non	0,74	0,45 - 1,20	0,24	NS	2,94	1,43 - 6,03	0,36	0,003
Distribution de concentré	non	1,61	1,13 - 2,30	0,18	0,008	1,15	0,84 - 1,59	0,16	NS
Aliment à base d'ensilage de maïs	non	0,58	0,41 - 0,81	0,17	0,001	0,92	0,68 - 1,24	0,15	NS
Odeur d'ammoniac perceptible dans le bâtiment	oui	1,84	0,61 - 5,50	0,56	NS	1,04	0,54 - 1,99	0,33	NS
Existence d'une infirmerie pour veaux	non	1,02	0,73 - 1,42	0,17	NS	1,69	1,23 - 2,32	0,16	0,001
Densité des veaux	forte	1,12	0,86 - 1,46	0,13	NS	1,06	0,78 - 1,44	0,15	NS
Incidence de la diarrhée la saison précédente	forte	1,91	1,19 - 3,07	0,24	0,007	1,41	0,94 - 2,10	0,20	NS
Nettoyage des bâtiments avant la saison des vêlages chez les veaux	non	1,13	0,81 - 1,59	0,17	NS	0,81	0,61 - 1,08	0,14	NS
Nettoyage des bâtiments après la saison des vêlages chez les veaux	non	1,37	0,97 - 1,93	0,17	NS	1,06	0,73 - 1,53	0,18	NS
Nettoyage suite à chaque épisode diarrhéique	non	0,62	0,27 - 1,39	0,41	NS	0,30	0,14 - 0,61	0,36	< 0,001
Propreté des vaches	faible	1,24	0,91 - 1,68	0,15	NS	0,78	0,58 - 1,06	0,15	NS
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	1,53	1,08 - 2,16	0,17	0,016	0,99	0,73 - 1,34	0,15	NS

TAB. F.1b – Ajustement naïf mois par mois : février et mars.

Variable	Modalité	Fév.				Mars			
		RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>	RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>
Condition de vêlage	aide facile	1,39	0,93 - 2,09	0,20	NS	1,43	0,97 - 2,10	0,19	NS
	aide difficile	1,41	0,82 - 2,40	0,27	NS	1,19	0,74 - 1,91	0,24	NS
Vaccination de la mère contre <i>E. coli</i>	non	0,58	0,45 - 0,74	0,12	< 0,001	0,87	0,67 - 1,13	0,13	NS
Vaccination de la mère contre d'autres agents que GENN	non	2,45	1,78 - 3,35	0,16	< 0,001	1,81	1,29 - 2,53	0,17	< 0,001
Distribution de concentré	non	1,34	1,07 - 1,69	0,11	0,011	1,29	1,06 - 1,56	0,09	0,008
Aliment à base d'ensilage de maïs	non	0,49	0,38 - 0,64	0,13	< 0,001	0,76	0,61 - 0,95	0,11	0,019
Odeur d'ammoniac perceptible dans le bâtiment	oui	1,14	0,68 - 1,90	0,26	NS	1,93	1,50 - 2,48	0,12	< 0,001
Existence d'une infirmerie pour veaux	non	0,99	0,80 - 1,22	0,10	NS	0,97	0,80 - 1,18	0,10	NS
Densité des veaux	forte	1,28	1,03 - 1,59	0,11	0,023	1,86	1,54 - 2,24	0,09	< 0,001
Incidence de la diarrhée la saison précédente	forte	1,36	1,06 - 1,74	0,12	0,015	1,25	1,01 - 1,53	0,10	0,035
Nettoyage des bâtiments avant la saison des vêlages chez les veaux	non	0,78	0,65 - 0,95	0,10	0,015	1,01	0,79 - 1,29	0,12	NS
Nettoyage des bâtiments après la saison des vêlages chez les veaux	non	1,91	1,39 - 2,61	0,16	< 0,001	1,55	1,16 - 2,07	0,14	0,002
Nettoyage suite à chaque épisode diarrhéique	non	0,57	0,35 - 0,93	0,25	0,025	0,44	0,28 - 0,69	0,23	< 0,001
Propreté des vaches	faible	1,53	1,24 - 1,88	0,10	< 0,001	1,50	1,24 - 1,82	0,10	< 0,001
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	0,95	0,76 - 1,18	0,11	NS	0,53	0,43 - 0,65	0,10	< 0,001

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>b</sup> : écart-type de  $\hat{\beta}$ ; <sup>bb</sup> : degré de significativité.

TAB. F.2 – Comparaison des différents ajustements avec la prise en compte des interactions : modèle naïf et modèle marginal.

Variable	Modalité	Naïf sans effet groupe				Marginal		
		RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>	IC**	EC <sup>b</sup>	p <sup>bb</sup>
Mois de naissance	janv.	0,16	0,08 - 0,33	0,36	< 0,001	0,07 - 0,35	0,39	< 0,001
	fév.	0,68	0,40 - 1,13	0,26	NS	0,37 - 1,24	0,30	NS
	mars	1,15	0,69 - 1,92	0,26	NS	0,62 - 2,16	0,31	NS
Condition de vêlage	aide facile	1,39	1,11 - 1,74	0,11	0,004	1,08 - 1,79	0,12	0,010
	aide difficile	1,37	1,03 - 1,82	0,14	0,030	0,99 - 1,88	0,16	NS
Vaccination de la mère contre <i>E. coli</i>	non	1,01	0,68 - 1,46	0,19	NS	0,65 - 1,52	0,21	NS
	interaction avec janv.	1,85	1,03 - 3,29	0,29	0,037	0,87 - 3,89	0,38	NS
	interaction avec fév.	0,58	0,38 - 0,89	0,21	0,013	0,33 - 1,02	0,28	NS
	interaction avec mars	0,93	0,60 - 1,44	0,22	NS	0,57 - 1,50	0,24	NS
Vaccination de la mère contre d'autres agents que GENN	non	0,98	0,68 - 1,43	0,19	NS	0,58 - 1,65	0,26	NS
	interaction avec janv.	1,88	1,10 - 3,21	0,27	0,020	1,02 - 3,47	0,31	0,042
	interaction avec fév.	2,35	1,48 - 3,74	0,23	< 0,001	1,19 - 4,66	0,35	NS
	interaction avec mars	1,83	1,16 - 2,91	0,23	0,009	0,99 - 3,40	0,31	NS
Distribution de concentré	non	1,38	1,22 - 1,57	0,06	< 0,001	1,14 - 1,68	0,10	< 0,001
Aliment à base d'ensilage de maïs	non	0,58	0,44 - 0,78	0,14	< 0,001	0,36 - 0,96	0,24	0,026
	interaction avec janv.	1,49	1,01 - 2,18	0,19	0,040	0,80 - 2,78	0,31	NS
	interaction avec fév.	0,84	0,59 - 1,20	0,17	NS	0,47 - 1,51	0,30	NS
	interaction avec mars	1,37	0,97 - 1,95	0,17	NS	0,77 - 2,43	0,29	NS
Odeur d'ammoniac perceptible dans le bâtiment	oui	1,45	1,17 - 1,79	0,10	< 0,001	1,09 - 1,93	0,14	0,011
Existence d'une infirmerie pour veaux	non	1,02	0,75 - 1,39	0,15	NS	0,62 - 1,69	0,25	NS
	interaction avec janv.	1,69	1,09 - 2,61	0,22	0,018	1,01 - 2,84	0,26	0,045
	interaction avec fév.	0,96	0,66 - 1,39	0,18	NS	0,54 - 1,71	0,29	NS
	interaction avec mars	0,98	0,68 - 1,40	0,18	NS	0,54 - 1,75	0,29	NS
Densité des veaux	forte	1,14	0,90 - 1,46	0,12	NS	0,81 - 1,62	0,17	NS
	interaction avec janv.	0,92	0,63 - 1,33	0,19	NS	0,58 - 1,45	0,23	NS
	interaction avec fév.	1,20	0,88 - 1,63	0,15	NS	0,74 - 1,94	0,24	NS
	interaction avec mars	1,51	1,12 - 2,02	0,15	0,006	0,99 - 2,29	0,21	NS
Incidence de la diarrhée la saison précédente	forte	1,41	1,23 - 1,62	0,07	< 0,001	1,03 - 1,93	0,16	0,031
Nettoyage des bâtiments avant la saison des vêlages chez les veaux	non	0,88	0,77 - 0,99	0,06	0,036	0,70 - 1,10	0,11	NS
Nettoyage des bâtiments après la saison des vêlages chez les veaux	non	1,55	1,33 - 1,81	0,07	< 0,001	1,21 - 1,98	0,12	< 0,001
Nettoyage suite à chaque épisode diarrhéique	non	0,51	0,39 - 0,66	0,13	< 0,001	0,33 - 0,78	0,22	0,002
Propreté des vaches	faible	1,31	0,99 - 1,73	0,14	NS	0,74 - 2,33	0,29	NS
	interaction avec janv.	0,66	0,45 - 0,97	0,19	0,035	0,33 - 1,33	0,35	NS
	interaction avec fév.	1,17	0,84 - 1,63	0,16	NS	0,59 - 2,33	0,35	NS
	interaction avec mars	1,12	0,79 - 1,56	0,17	NS	0,59 - 2,11	0,32	NS
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	1,41	1,04 - 1,93	0,16	0,028	0,78 - 2,54	0,30	NS
	interaction avec janv.	0,72	0,47 - 1,08	0,21	NS	0,37 - 1,40	0,33	NS
	interaction avec fév.	0,68	0,47 - 0,99	0,19	0,044	0,33 - 1,41	0,37	NS
	interaction avec mars	0,40	0,27 - 0,57	0,19	< 0,001	0,21 - 0,76	0,33	0,005

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>b</sup> : écart-type de  $\hat{\beta}$ ; <sup>bb</sup> : degré de significativité.

TAB. F.3 – Ajustement par le modèle de fragilité de Cox avec la prise en compte des interactions.

Variable	Modalité	Naïf sans effet groupe			
		RR*	IC**	EC <sup>b</sup>	p <sup>bb</sup>
Mois de naissance	janv.	0,17	0,08 - 0,35	0,35	< 0,001
	fév.	0,57	0,33 - 1,01	0,28	NS
	mars	1,05	0,59 - 1,86	0,29	NS
Condition de vêlage	aide facile	1,18	0,92 - 1,51	0,12	NS
	aide difficile	1,32	0,98 - 1,78	0,15	NS
Vaccination de la mère contre <i>E. coli</i>	non	1,01	0,80 - 1,25	0,11	NS
Vaccination de la mère contre d'autres agents que GENN	non	0,80	0,50 - 1,30	0,25	NS
	interaction avec janv.	3,24	1,76 - 5,96	0,31	< 0,001
	interaction avec fév.	2,40	1,43 - 4,01	0,26	< 0,001
	interaction avec mars	2,53	1,50 - 4,27	0,26	< 0,001
Distribution de concentré	non	1,51	1,15 - 1,97	0,14	0,003
Aliment à base d'ensilage de maïs	non	0,54	0,36 - 0,80	0,20	0,002
	interaction avec janv.	1,36	0,87 - 2,10	0,22	NS
	interaction avec fév.	0,98	0,65 - 1,45	0,20	NS
	interaction avec mars	1,48	0,97 - 2,26	0,21	NS
Odeur d'ammoniac perceptible dans le bâtiment	oui	1,38	0,80 - 2,38	0,28	NS
Existence d'une infirmerie pour veaux	non	0,78	0,51 - 1,18	0,21	NS
	interaction avec janv.	1,81	1,10 - 2,98	0,25	0,018
	interaction avec fév.	1,35	0,87 - 2,07	0,22	NS
	interaction avec mars	1,22	0,79 - 1,89	0,22	NS
Densité des veaux	forte	1,16	0,80 - 1,68	0,19	NS
	interaction avec janv.	0,85	0,56 - 1,30	0,21	NS
	interaction avec fév.	1,42	0,99 - 2,03	0,18	NS
	interaction avec mars	1,60	1,09 - 2,36	0,19	0,015
Incidence de la diarrhée la saison précédente	forte	1,51	1,15 - 1,99	0,14	0,003
Nettoyage des bâtiments avant la saison des vêlages chez les veaux	non	0,88	0,677 - 1,16	0,14	NS
Nettoyage des bâtiments après la saison des vêlages chez les veaux	non	1,36	1,01 - 1,84	0,15	0,041
Nettoyage suite à chaque épisode diarrhéique	non	0,42	0,23 - 0,74	0,29	0,002
Propreté des vaches	faible	1,32	0,90 - 1,93	0,19	NS
	interaction avec janv.	0,72	0,47 - 1,10	0,21	NS
	interaction avec fév.	1,05	0,72 - 1,56	0,19	NS
	interaction avec mars	1,04	0,70 - 1,56	0,20	NS
Supplémentation en vitamines et minéraux aux veaux à leur naissance	oui	1,41	0,93 - 2,14	0,21	NS
	interaction avec janv.	0,57	0,36 - 0,92	0,24	0,022
	interaction avec fév.	0,67	0,44 - 1,01	0,21	NS
	interaction avec mars	0,35	0,22 - 0,55	0,22	< 0,001

\* : risque relatif, soit  $\exp(\hat{\beta})$ ; \*\* : intervalle de confiance à 5 % du risque relatif; <sup>b</sup> : écart-type de  $\hat{\beta}$ ; <sup>bb</sup> : degré de significativité.

## Annexe G

# GLOSSAIRE ÉPIDÉMIOLOGIQUE

**Censure** Absence d'information (relative à l'événement étudié) à la date de point.

**Cohorte** Groupe d'individus ayant vécu une même expérience (exposition à un facteur de risque par exemple) et suivis dans le temps depuis la date de cette expérience (qui peut éventuellement être différente d'un sujet à l'autre).

**Date de dernières nouvelles** Date la plus récente à laquelle une information concernant le sujet est recueillie – elle est différente pour tous les sujets.

**Date d'origine** Date de randomisation ou de début de traitement – elle est différente pour tous les sujets.

**Date de point** Date à partir de laquelle il n'est plus tenu compte de l'information – elle est commune à tous les sujets.

**Durée de surveillance** Délai écoulé entre la date d'origine et la date de dernières nouvelles.

**Enquête cas-témoins** Sélection d'un groupe de sujets atteints de la maladie étudiée (les *cas*), supposé être représentatif – pour l'exposition au facteur de risque – de l'ensemble des malades ayant cette pathologie, et d'un (ou plusieurs) groupe(s) de sujets indemnes de cette maladie (les *témoins*), supposés être représentatif(s) de la population dont sont issus les cas. Par suite, recherche d'informations concernant l'exposition aux facteurs de risque (et aux facteurs de confusion éventuels) dans le passé des sujets de l'étude. Il s'agit d'une enquête rétrospective.

**Enquête exposés-non exposés** Définition d'une période d'observation, au début de laquelle tous les sujets retenus sont indemnes de la maladie étudiée. Possibilité dans certains cas de distinguer *a priori*, au sein de l'ensemble des sujets de l'étude, un groupe exposé au facteur de risque et un groupe non exposé ; dans d'autres cas, la distinction se fait *a posteriori*, c'est-à-dire à la fin de la période d'étude, grâce aux données recueillies individuellement sur les sujets. Il s'agit d'une enquête prospective.

**Enquête prospective** Enquête au sein de laquelle les sujets sont systématiquement suivis dans le temps, afin de mesurer d'éventuelles modifications d'exposition, et surtout d'observer l'incidence de la maladie étudiée.

**Enquête rétrospective** Enquête au sein de laquelle on recherche, dans le passé des sujets de l'étude, des informations concernant leur exposition au facteur de risque étudié.

**Exclu-vivant** Censure dont la date de dernière nouvelle est postérieure à la date de point.

**Facteur de confusion (ou concomitant)** Lors de l'observation d'une liaison entre un facteur dont on étudie le rôle et la maladie, il peut exister un autre facteur, dit « facteur de confusion », lui-même lié à la fois au facteur étudié et à la maladie, qui peut produire l'association observée initialement – même si le facteur dont on étudie le rôle et la maladie sont indépendants.

**Incidence** Nombre de nouveaux cas pendant une période donnée.

**Létalité** Phénomène de la mort spécifique à une cause, considéré du point de vue du nombre.

**Morbidité** Phénomène de la maladie, considéré du point de vue du nombre.

**Mortalité** Phénomène de la mort, considéré du point de vue du nombre.

**Odds-ratio** Voir *Risque relatif estimé*.

**Perdu de vue** Censure dont la date de dernière nouvelle est antérieure à la date de point.

**Personne-années** Somme des durées, cumulées sur l'ensemble de la population étudiée, pendant laquelle les sujets sont susceptibles d'être enregistrés comme de nouveaux cas. Mesure la « taille » d'une population. L'année peut évidemment être remplacée par une autre division du temps.

**Randomisation** Dans les essais thérapeutiques, la méthode utilisée pour contrôler les facteurs de confusion au moment de la constitution de l'échantillon est le tirage au sort : les sujets sont répartis *au hasard* (*random* en anglais) entre deux groupes, l'un recevant le premier traitement et l'autre le second.

**Recul** Délai écoulé entre la date d'origine et la date de point.

**Risque relatif** Rapport de l'incidence dans le groupe exposé et de l'incidence dans le groupe non exposé, il ne peut être calculé que dans ce cadre (enquête exposés-non exposés). C'est une mesure du rôle étiologique d'un facteur de risque.

**Risque relatif estimé (odds-ratio)** Estimation du risque relatif dans le cadre d'une enquête cas-témoins, où l'incidence ne peut (par définition) être connue.

**Temps de participation** Délai écoulé entre la date d'origine et la date de dernière nouvelle si celle-ci est antérieure à la date de point ; sinon c'est le délai écoulé entre la date d'origine et la date de point.



# Index

- A**
- Adéquation ..... 35
- Akaike ..... voir Critère
- Algorithme
- E.M. .... 68
- de Newton-Raphson ..... 33, 69
- B**
- BLP ..... 75
- BLUP ..... 76
- BP ..... 74
- Breslow
- estimateur ..... 31
- test ..... voir Test
- vraisemblance partielle ..... 43
- C**
- Cadlag ..... 11, 22, 152, 155
- Censure
- aléatoire ..... 16
- définition** ..... 15, 197
- à droite ..... 15
- à gauche ..... 15
- informative ..... 15
- par intervalle ..... 15
- non-aléatoire ..... 16
- non-informative ..... 15
- de type I ..... 16
- de type II ..... 16
- Chaîne de Markov par technique de
- Monte-Carlo (MCMC) ..... 73
- Cluster* ..... 47
- Coefficient
- de corrélation intraclasse ..... 59
- de dispersion ..... 59
- Cohorte ..... 111, **197**
- Compensateur ..... 152
- Compound symmetric* ..... 59
- Confusion ..... voir Facteur
- Coordonnée d'un processus ..... 151
- Copule ..... 48
- Cox
- modèle ..... 30
- vraisemblance partielle ..... 31
- Critère d'Information d'Akaike ..... 97
- D**
- Date
- de dernières nouvelles ..... 197
- d'origine ..... 197
- de point ..... 197
- Design effect* ..... 59
- Données
- aberrantes ..... 40
- Durée
- de surveillance ..... 197
- de vie ..... 14
- E**
- E.M., algorithme ..... 68
- EBLUP ..... 76
- Efron ..... voir Vraisemblance
- Enquête
- cas-témoins ..... 197
- exposés-non exposés ..... 197
- prospective ..... 198
- rétrospective ..... 198
- Équations
- d'estimation généralisées ..... 48
- du modèle mixte ..... 76
- Équicorrélation ..... 59
- Estimateur
- de Greenwood ..... 22
- jackknife* ..... 42, 84
- de Kaplan-Meier ..... 21
- Ex-æquo ..... 25, 42
- Exchangeable* ..... 59
- Exclu-vivant ..... 198
- F**
- Facteur de confusion ..... 198
- Filtration ..... **151**, 152
- Fisher ..... voir Matrice
- Fonction
- aléatoire ..... 151
- de risque
- cumulé ..... 14
- instantané ..... 14
- de survie ..... 14
- Fragilité
- définition ..... **62**
- partagée ..... 71
- G**
- Gamma ..... voir Loi
- GEE ..... 48
- Gehan ..... voir Test
- GENN ..... 110
- Glivenko-Cantelli ..... voir Théorème
- Greenwood ..... voir Estimateur
- Groupe ..... 47
- H**
- Hiérarchie ..... 47
- I**
- Ignorable* ..... voir Censure non-informative

- Incidence.....198  
 Intraclasse ..... voir Coefficient de corrélation  
 Inverse gaussienne ..... voir Loi
- J**
- Jackknife* ..... voir Estimateur
- K**
- Kaplan-Meier ..... voir Estimateur
- L**
- Létalité ..... 198  
 Log-rank.....voir Test  
 Loi  
   gamma ..... 64  
   inverse gaussienne ..... 64  
   stable positive.....65, 94
- M**
- Mantel-Haenszel ..... voir Test  
 Martingale ..... 151  
 Matrice d'information de Fisher ..... 33  
 Maximum de vraisemblance..... voir Test  
 Meilleur prédicteur ..... 74  
 Modèle  
   bayésien ..... 72  
   de Cox  
     avec fragilité ..... 9, 49, **62**  
     marginal ..... 9, 48, 52  
   étendu de régression du risque  
     instantané.....43  
     hiérarchique ..... 47  
     linéaire mixte ..... 73  
 Morbidité ..... 110, **198**  
 Mortalité..... 110, **198**
- N**
- Newton-Raphson ..... 33
- O**
- Odds-ratio ..... 198  
*Outliers* ..... 40
- P**
- Partagée (fragilité) ..... 71  
 Perdu de vue.....198  
 Prédicteur..... 74  
 Prentice ..... voir Test  
 Processus  
   aléatoire..... 151  
   de comptage..... 17, **152**, 153  
   coordonnée ..... 151  
   d'intensité..... 18, 152  
   d'intensité cumulée ..... 18, 152  
   prévisible croissant ..... 153  
 Produit-intégrale (ou infini) ..... 155  
 Pseudo-modèle mixte..... 49
- R**
- Randomisation..... 198  
 Rebolledo.....voir théorème  
 Recul ..... 198  
 Résidu  
   de déviance..... 40  
   martingale.....35  
   partiel de Schoenfeld ..... 41  
   du score ..... 38  
 Risque  
   cumulé ..... 14  
   instantané..... 14  
   marginal.....52  
   relatif  
   **définition** ..... 198  
   estimé..... 198  
 RR ..... voir Risque relatif
- S**
- Schoenfeld ..... voir Résidu  
 Stable positive..... voir Loi  
 Symétrie par composition ..... 59
- T**
- Tarone..... voir Test  
 Temps de participation..... 199  
 Test  
   de Breslow.....27  
   de Gehan ..... 27  
   du log-rank ..... 27  
   de Mantel-Haenszel ..... 27  
   du maximum de vraisemblance.....34  
   de Prentice ..... 28  
   du rapport de vraisemblance.....**34**, 54  
   du score.....**34**, 54  
   de Tarone et Ware.....28  
   de Wald.....**34**, 54  
 Théorème  
   de Duhamel ..... 155  
   de Glivenko-Cantelli ..... 160  
   de l'innovation.....153  
   de Jacod.....155  
   de Rebolledo ..... 154  
 Trajectoire..... 151  
 Troncature ..... 15
- V**
- Vraisemblance partielle

de Breslow .....	43
de Cox .....	31
d'Efron .....	43
exacte .....	43
pénalisée.....	86

**W**

Wald .....	voir Test
Ware .....	voir Test

## Résumé

Les méthodes et modèles statistiques destinés aux données de survie furent, dans leur grande majorité, développés sous l'hypothèse implicite d'indépendance statistique des observations individuelles. Bien que l'imposition d'une telle hypothèse soit raisonnable pour un grand nombre d'applications, il s'avère évident que dans bon nombre d'autres – et non des moins courantes –, cette hypothèse est violée. Par exemple, en science vétérinaire, de telles données corrélées apparaissent lorsque les observations individuelles sont regroupées au sein d'élevages.

Nous étudions les deux principales classes de modèles pour données de survie corrélées : les modèles de fragilité (ou conditionnels) et les modèles marginaux. Nous nous proposons une large comparaison de ces deux approches, d'une part au travers d'une étude de données vétérinaires, d'autre part au moyen de simulations.

Notre objectif est d'évaluer la sensibilité de tels modèles vis-à-vis de la structure des jeux de données qu'ils sont appelés à traiter – et plus particulièrement vis-à-vis de la taille des groupes.

**Mots-clés :** Modèle de Cox, corrélation, fragilité, modèle marginal

## Abstract

Most of the statistical models and methods for failure time data were implicitly developed under the assumption that the observations from subjects are statistically independent of each other. While sensible in many applications, this assumption is obviously violated in other situations which are not as uncommon as originally thought. For example, in veterinary science, such correlation between data occurs, specially when individuals recording single outcomes are grouped into clusters.

We study the two broad classes of models for correlated survival data : frailty (or conditional) models and marginal models. We propose a wide comparison of these two approaches ; this comparison is realized through veterinary data set and simulations.

Our goal is to assess the sensitivity of such models, and more particularly to appreciate their dependance on several parameters (for example the clusters size).

**Keywords :** Cox model, correlation, frailty, marginal model