

FEDERAL UNIVERSITY OF PARANÁ

HENRIQUE APARECIDO LAUREANO

MODELING THE CUMULATIVE INCIDENCE FUNCTION OF CLUSTERED
COMPETING RISKS DATA: A MULTINOMIAL GLMM APPROACH

CURITIBA

2021

HENRIQUE APARECIDO LAUREANO

MODELING THE CUMULATIVE INCIDENCE FUNCTION OF CLUSTERED
COMPETING RISKS DATA: A MULTINOMIAL GLMM APPROACH

Thesis presented to the Graduate Program of Numerical Methods in Engineering, Concentration Area in Mathematical Programming: Statistical Methods Applied in Engineering, Federal University of Paraná, as part of the requirements to the obtention of the Master's Degree in Sciences.

Supervisor: Prof. PhD Wagner Hugo Bonat
Co-supervisor: Prof. PhD Paulo Justiniano Ribeiro Jr

CURITIBA

2021

HENRIQUE APARECIDO LAUREANO

**MODELING THE CUMULATIVE INCIDENCE FUNCTION OF CLUSTERED
COMPETING RISKS DATA: A MULTINOMIAL GLMM APPROACH**

Thesis presented to the Graduate Program of Numerical Methods in Engineering, Concentration Area in Mathematical Programming: Statistical Methods Applied in Engineering, Federal University of Paraná, as part of the requirements to the obtention of the Master's Degree in Sciences.

Master thesis approved. XXX XX, 2021.

Prof. PhD Wagner Hugo Bonat
Supervisor

Prof. PhD Paulo Justiniano Ribeiro Jr
Co-supervisor

Prof. PhD ...
Internal Examinator - PPGMNE

Prof. PhD ...
Internal Examinator - PPGMNE

Prof. PhD ...
External Examiner -

CURITIBA
2021

To Celita and Olivio

ACKNOWLEDGEMENTS

Today is the shadow of tomorrow. Today is the present future of yesterday. Yesterday is the shadow of today. The darkness of the past is yesterday. And the light of the past is yesterday. The days of yesterday are all numbered in sum. In the word once. Because once upon a time there was a yesterday. Yesterday belongs to the dead. Because the dead belongs to the past. The past is yesterday. Today is the preview of tomorrow but for me. Only from a better and happier point of view. My point of view is the thought of a better or try. Reality, is today is eternity. The eternity of yesterday is dead. Yesterday is as one. The eternity of one is the eternity of the past. The past is once upon a time. Once upon a time is past. The past is yesterday- day. The past is yesterday-day. While we're searching for tomorrow.

The light of the past is the light which was. The wisdom of the past is the light of the past. The light of the future is the light which is to be. The wisdom of the future is the light of the future see. Yesterday belongs to the dead. Tomorrow belongs to the living. The past is certified as a finished product. Anything which has ended is finished. That which is perfect is finished. The perfect man is no exception to the rule. The perfect man of the past is made according to the rule of the past. The rule of the past is a law of injustice and hypocrisy. The revelation of the meaning of the law is revealed through the law itself. The wisdom of the past is the light of the past. The light which is to be the wisdom of the future. The light of the future casts the shadows of tomorrow.

As Moro said once, I'm thankful for everything and everyone.

"It's not supposed to be easy."

(Gregg Popovich)

ABSTRACT

Failure time data ...

Keywords: Competing risks.

RESUMO

Dados de tempos de falha ...

Palavras-chave: Riscos competitivos.

LIST OF FIGURES

FIGURE 1 – ILLUSTRATION OF MULTISTATE MODELS FOR A) FAILURE TIME PROCESS; B) COMPETING RISKS PROCESS; AND C) ILLNESS-DEATH MODEL, THE SIMPLEST MULTISTATE MODEL	16
FIGURE 2 – A COMPUTATIONAL GRAPH	28
FIGURE 3 – EXAMPLE OF A SIMPLE COMPUTATIONAL GRAPH	29
FIGURE 4 – R CODE FOR THE TMB IMPLEMENTATION OF A LOGISTIC MIXED MODEL	33
FIGURE 5 – R CODE FOR THE MODEL FITTING OF A LOGISTIC MIXED MODEL WRITTEN IN TMB	34
FIGURE 6 – TMB PACKAGE DESIGN	34
FIGURE 7 – ILLUSTRATION OF COEFFICIENT BEHAVIORS FOR A GIVEN CUMULATIVE INCIDENCE FUNCTION (CIF) (PROPOSED BY Cederkvist et al. (2019)), IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, AND WITH THE FOLLOWING CONFIGURATION: $\beta_2 = 0$, $u = 0$ AND $\eta = 0$; IN EACH SCENARIO ALL OTHER COEFFICIENTS ARE SET TO ZERO, WITH THE EXCEPTION OF $w_1 = 1$	39
FIGURE 8 – ILLUSTRATION OF A GIVEN CLUSTER-SPECIFIC CUMULATIVE INCIDENCE FUNCTION (CIF), PROPOSED BY Cederkvist et al. (2019), IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES AND THE FOLLOWING CONFIGURATION: $\beta_1 = -2$, $\beta_2 = -1$, $\gamma_1 = 1$, $w_1 = 3$ AND $u_2 = 0$. THE VARIATION BETWEEN FRAMES IS GIVEN BY THE LATENT EFFECTS u_1 AND η_1	40
FIGURE 9 – ILLUSTRATION OF THE PARAMETRIZATION BEHAVIOR FOR THE VARIANCE COMPONENTS, IN A), AND CORRELATION COMPONENTS, IN B)	44
FIGURE 10 – CUMULATIVE INCIDENCE FUNCTIONS (CIF) AND RESPECTIVE DERIVATIVES (dCIF) W.R.T. TIME FOR A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, LATENT EFFECTS IN ZERO, AND FIXED EFFECTS IN Equation 4.1	47

FIGURE 11 – HISTOGRAMS FOR SIMULATED PROBABILITIES WITH RESPECTIVE OUTPUT PCERCENTAGES AND FAILURE TIMES FOR A MODEL WITH TWO COMPETING CAUSES AND 50000 CLUSTERS OF SIZE TWO. THE SIMULATION FOLLOWED ALGORITHM 1 GUIDELINES WITH PARAMETER CONFIGURATIONS SPECIFIED IN Equation 4.1 AND Equation 4.2	48
FIGURE 12 – CUMULATIVE INCIDENCE FUNCTIONS (CIF) FOR A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, AND LATENT EFFECTS IN ZERO	49
FIGURE 13 – PARAMETER β_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	53
FIGURE 14 – PARAMETER β_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	53
FIGURE 15 – PARAMETER γ_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	54
FIGURE 16 – PARAMETER γ_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	54
FIGURE 17 – PARAMETER w_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	55
FIGURE 18 – PARAMETER w_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS .	55
FIGURE 19 – PARAMETER $\log(\sigma_1^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	56
FIGURE 20 – PARAMETER $\log(\sigma_2^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	56
FIGURE 21 – PARAMETER $\log(\sigma_3^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	57
FIGURE 22 – PARAMETER $\log(\sigma_4^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	57
FIGURE 23 – PARAMETER $z(\rho_{12})$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	58
FIGURE 24 – PARAMETER $z(\rho_{34})$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	58
FIGURE 25 – PARAMETERS $\{z(\rho_{13}), z(\rho_{24}), z(\rho_{14}), z(\rho_{23})\}$ BIAS WITH ± 1.96 STANDARD DEVIATIONS	59
FIGURE 26 – HIGH CUMULATIVE INCIDENCE FUNCTION (CIF) SCENARIO CURVES	60
FIGURE 27 – LOW CUMULATIVE INCIDENCE FUNCTION (CIF) SCENARIO CURVES	61
FIGURE 28 – VARIANCE PARAMETERS DENSITIES IN THE SCENARIOS OF HIGH CIF AND 60 THOUSAND DATA POINTS	63
FIGURE 29 – CORRELATION PARAMETERS DENSITIES IN THE SCENARIOS OF HIGH CIF AND 60 THOUSAND DATA POINTS	64
FIGURE 30 – COMPLETE MODEL'S PARAMETERS CORRELATION HEATMAP IN THE SCENARIO OF CLUSTER SIZE 10, HIGH CIF, AND SIXTY-THOUSAND DATA POINTS	65
FIGURE 31 – PARAMETER β_1 BIAS WITH 2.5% AND 97.5% QUANTILES . .	78
FIGURE 32 – PARAMETER β_2 BIAS WITH 2.5% AND 97.5% QUANTILES . .	78
FIGURE 33 – PARAMETER γ_1 BIAS WITH 2.5% AND 97.5% QUANTILES . .	79
FIGURE 34 – PARAMETER γ_2 BIAS WITH 2.5% AND 97.5% QUANTILES . .	79
FIGURE 35 – PARAMETER w_1 BIAS WITH 2.5% AND 97.5% QUANTILES . .	80

FIGURE 36 – PARAMETER w_2 BIAS WITH 2.5% AND 97.5% QUANTILES	80
FIGURE 37 – PARAMETER $\log(\sigma_1^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES	81
FIGURE 38 – PARAMETER $\log(\sigma_2^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES	81
FIGURE 39 – PARAMETER $\log(\sigma_3^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES	82
FIGURE 40 – PARAMETER $\log(\sigma_4^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES	82
FIGURE 41 – PARAMETER $z(\rho_{12})$ BIAS WITH 2.5% AND 97.5% QUANTILES	83
FIGURE 42 – PARAMETER $z(\rho_{34})$ BIAS WITH 2.5% AND 97.5% QUANTILES	83
FIGURE 43 – PARAMETERS $\{z(\rho_{13}), z(\rho_{24}), z(\rho_{14}), z(\rho_{23})\}$ BIAS WITH 2.5% AND 97.5% QUANTILES	84

LIST OF TABLES

LIST OF ALGORITHMS

ALGORITHM 1 SIMULATING FROM A multiGLMM FOR CLUSTERED COMPETING RISKS DATA	46
--	----

CONTENTS

1	INTRODUCTION	15
1.1	GOALS	18
1.1.1	General goals	18
1.1.2	Specific goals	18
1.2	JUSTIFICATION	19
1.3	LIMITATION	19
1.4	THESIS ORGANIZATION	20
2	GENERALIZED LINEAR MIXED MODELS: FORMULATION, OPTIMIZATION, AND IMPLEMENTATION	21
2.1	FORMULATION: OBTAINING A JOINT LIKELIHOOD FUNCTION	21
2.2	MARGINALIZATION: LAPLACE APPROXIMATION AND ALTERNATIVES	22
2.3	OPTIMIZATION: MARGINAL LIKELIHOOD FUNCTION	25
2.4	AD: AUTOMATIC DIFFERENTIATION	28
2.4.1	Forward Mode	29
2.4.2	Reverse Mode	30
2.5	TMB: TEMPLATE MODEL BUILDER	31
3	multiGLMM: A MULTINOMIAL GLMM FOR CLUSTERED COMPETING RISKS DATA	36
3.1	CLUSTER-SPECIFIC CUMULATIVE INCIDENCE FUNCTION (CIF)	36
3.2	MODEL SPECIFICATION	40
3.2.1	Parametrization	42
4	SIMULATION STUDY DATASETS	46
4.1	SIMULATING FROM THE MODEL	46
4.2	SIMULATION STUDIES DESIGN	49
5	RESULTS	51
5.1	SIMULATION STUDY	51
6	FINAL CONSIDERATIONS	66
6.1	FUTURE WORKS	66
	BIBLIOGRAPHY	67

APPENDIX	70
APPENDIX	71
APPENDIX A – ANALYTIC GRADIENT OF THE LATENT EFFECTS FOR THE JOINT LOG-LIKELIHOOD FUNCTION OF THE MULTINOMIAL GLMM FOR CLUSTERED COMPETING RISKS DATA	71
APPENDIX B – ANALYTIC HESSIAN OF THE LATENT EFFECTS FOR THE JOINT LOG-LIKELIHOOD FUNCTION OF THE MULTINOMIAL GLMM FOR CLUSTERED COMPETING RISKS DATA	72
APPENDIX C – R CODE TO SIMULATE FROM A multiGLMM WITH TWO COMPETING CAUSES AND CLUSTERS OF SIZE TWO. FOR MORE INFORMATION CHECK SECTION 4.1	76
APPENDIX D – MODEL PARAMETERS BIAS WITH 2.5% AND 97.5% QUANTILES	78

1 INTRODUCTION

Consider a cluster of random variables representing the time until the occurrence of some event. These random variables are assumed to be correlated, i.e. for some biological or environmental reason it is not adequate to assume independence between them. Also, we may be interested in the occurrence of not only one specific event, having in practice a competition of events to see which one happens first, if it happens. Such events may also be of low probability albeit severe consequences, this is the moment when the cluster correlation makes its difference: the occurrence of an event in a cluster member should affect the probability of the same happening in the others.

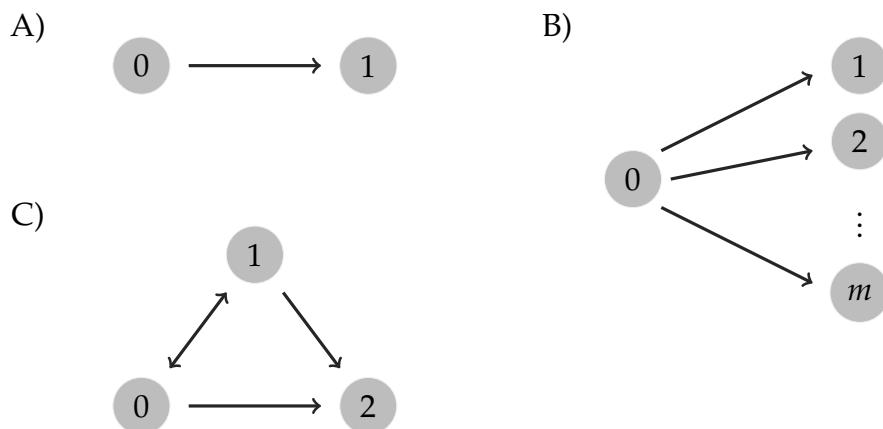
A realistic context that fits perfectly with the framework described above is the study of disease incidence in family members, where each member is indexed by a random variable and each cluster consists of a familiar structure. The inspiration to the study of these kinds of problems came from the work developed in [Cederkvist et al. \(2019\)](#), where they studied breast cancer incidence in mothers and daughters but using a complicated modeling framework. Based on that, the aim of this thesis is to propose a simpler framework taking advantage of several *state-of-art* computational libraries and see how far can we go in several scenarios. Until now we just contextualized, we still need to introduce the methodology. To this, some definitions and theoretical contexts are welcome.

When the object under study is a random variable representing the time until some event occurs, we are in the field of *failure time data* ([KALBFLEISCH; PRENTICE, 2002](#)). The occurrence of an event is generally denoted *failure*, and major areas of application are biomedical studies and industrial life testing. In this thesis, we maintain our focus on the former. As common in science, same methodologies can receive different names depending on the area. In industrial life testing is performed what is called a *reliability analysis*; in biomedical studies is performed what is called *survival analysis*. Generally, the term *survival* is applied when we are interested in the occurrence of only one event, a *failure time process*. When we are interested in the occurrence of more than one event we enter in the yard of *competing risks* and *multistate* models. A visual aid is presented on [Figure 1](#) and a comprehensive reference is [Kalbfleisch & Prentice \(2002\)](#).

Failure time and competing risks processes may be seen as particular cases of a multistate model. Besides the number of events (states) of interest, the main difference between a multistate model and its particular cases is that only in the multistate scenario we may have transient states, using a *stochastic process* language. In the particular cases, all states besides the initial state 0, are absorbents - once you reached it you do not leave.

The simplest multistate model that exemplify this behavior is the illness-death model, [Figure 1 C](#), where a patient (initially in state 0) can get sick (state 1) or die (state 2); if sick it can recover (returns to state 0) or die. We work in this thesis only with competing risks processes, and to each patient we need the time (age) until the occurrence, or not, of the event.

FIGURE 1 – ILLUSTRATION OF MULTISTATE MODELS FOR A) FAILURE TIME PROCESS; B) COMPETING RISKS PROCESS; AND C) ILNESS-DEATH MODEL, THE SIMPLEST MULTISTATE MODEL



SOURCE: The author (2021).

When for some known or unknown reason we are not able to see the occurrence of an event, we have what is denoted *censorship*. Still in the illness-death model, during the period of follow up the patient may not get sick or die, staying at state 0. This is denoted *right-censorship*; If a patient is in state 1 at the end of the study, we are *censored* to see him reaching the state 2 or returning to state 0. This is the inherent idea to censorship and must be present in the modeling framework, thus arriving in the so-called *survival models* ([KALBFLEISCH; PRENTICE, 2002](#)).

A survival model deals with the survival experience. Usually, the survival experience is modeled in the *hazard* (failure rate) scale and it can be expressed for a subject i as

$$\lambda(t | \mathbf{x}_i) = \lambda_0(t) \times c(\mathbf{x}_i \boldsymbol{\beta}) \quad \text{at time } t, \quad (1.1)$$

i.e. as the product of an arbitrary baseline hazard function $\lambda_0(\cdot)$, with a specific function form $c(\cdot)$, that will depend on the probability distribution to be chosen for the failure time and on predictors/covariates/explanatory/independent variables $\mathbf{x}_i = [x_1 \dots x_p]$, where $\boldsymbol{\beta}^\top = [\beta_1 \dots \beta_p]$ is the parameters vector.

This structure is specified for a failure time process, as in [Figure 1 A](#)). Nevertheless, the idea is easy to extend. We basically have the [Equation 1.1](#)'s model to each cause-specific (in a competing risks process) or transition (in a multistate process). A complete and extensive detailing can be, again, found in [Kalbfleisch & Prentice \(2002\)](#).

In this work we approach the case of clustered competing risks. Besides the cause-specific structure, we have to deal with the fact that the events are happening in related individuals. This configures what is denoted *family studies*, i.e. we have a cluster/group/family dependence that needs to be considered, accommodated, and modeled. This, possible, dependence is something that we do not actually measure but know (or just suppose) that exists. In the statistical modeling language this characteristic receives the name of *random* or *latent effect*. A survival model with a latent effect, association, or unobserved heterogeneity, is denoted *frailty model* (CLAYTON, 1978; VALPEL; MANTON; STALLARD, 1979). In its simplest form, a frailty is an unobserved random proportionality factor that modifies the hazard function of an individual, or of related individuals. Frailty models are extensions of [Equation 1.1](#)'s model.

In the competing risks setting, the hazard scale (focusing on the cause-specific hazard) is not the only possible scale to work on. A more attractive possibility is to work on the probability scale ([ANDERSEN et al., 2012](#)), focusing on the cause-specific cumulative incidence function (CIF). Besides the within-family dependence, in family studies there is often a strong interest in describing age at disease onset, which is directly described by the cause-specific CIF. Therefore, making the probability scale a more attractive and logical choice. Since the CIF plays a central role in this master thesis, it will be formally defined later in a place with greater emphasis. With the definitions and the theoretical context being made, let us be more specific.

To work with competing risks data on the probability scale plus a latent structure allowing for within-cluster dependence of both risk and timing, [Cederkvist et al. \(2019\)](#) proposed a pairwise composite likelihood approach based on the factorization of the cause-specific CIF as the product of a cluster-specific risk level function with a cluster-specific failure time trajectory function. A composite approach ([LINDSAY, 1988](#); [COX; REID, 2004](#); [VARIN; REID; FIRTH, 2011](#)) is a valid alternative to a full likelihood analysis in high-dimensional situations when a full approach is too computational costly or even inviable. In failure time data problems, the composite likelihood function is built from the product of marginal densities. The marginal specification implies a pairwise approach since we need to add model layers to be able to handle with the dependence structure. A clear advantage of this approach is that we do not need to care about a joint distribution specification, which generally translates also into a computational advantage. A disadvantage is the model specification, which becomes much more complicated, besides the number of small details to be workaround from the fact of being working with not an exact likelihood function.

We do not have any guarantees that a full likelihood inference procedure is not viable here, so we try to reach the same goal of [Cederkvist et al. \(2019\)](#) albeit with a simpler framework taking advantage of *state-of-art* software, something still

not so common in the statistical modeling community. This simpler framework is a generalized linear mixed model (GLMM). Instead of concentrating on failure time data and consequently having a survival/frailty model based on the hazard scale, or using a composite approach, we just build the joint/full likelihood function (a multinomial model with its link function based on the cluster-specific CIF, accounting for an appropriate latent effects structure), marginalize (integrate out the latent effects) and optimize it. A Fisherian approach per se.

To a better contextualization of our GLMM approach ([MCCULLOCH; SEARLE, 2001](#)), consider a random subject i . In a standard linear model we assume that the response variable Y_i , conditioned on the covariates x_i , follows a normal/Gaussian distribution and what we do is to model its mean, $\mu_i \equiv \mathbb{E}(Y_i | x_i)$, via a linear combination. As much well explained in [Nelder & Wedderburn \(1972\)](#), with the aid of a *link function* $g(\cdot)$, this idea is generalized to distributions of the *exponential family*. Many of its members are useful for practical modelling, such as the Poisson (for counting data), binomial (dichotomic data), gamma (continuous but positive) and Gaussian (continuous data) distributions. This extended framework received the name of generalized linear model (GLM) and a comprehensive reference is [McCullagh & Nelder \(1989\)](#).

What makes a GLM into a GLMM ([MCCULLOCH; SEARLE, 2001](#)) is the addition of a latent effect u (then, *mixed*) into the mean structure. The mean structure of a standard GLMM is defined as

$$g(\mu_i) = x_i\beta + z_iu, \quad u \sim \text{Multivariate Normal}(\mathbf{0}, \Sigma)$$

where the latent effect is assumed to follow a multivariate Gaussian distribution of zero mean and a parametrized variance-covariance matrix Σ . Its correct linkage to the mean structure is made through the i^{th} vector row of a design-matrix Z . The covariates are into x_i , the i^{th} vector row of a model-matrix X , with β being a vector of unknown parameters.

1.1 GOALS

1.1.1 General goals

Propose and study the estimability of a multinomial generalized linear mixed model (multiGLMM) to the cluster and cause-specific cumulative incidence function (CIF) of clustered competing risks data.

1.1.2 Specific goals

1. Simulate from the model, i.e. generate synthetic data to study statistical properties.

2. Write the model in the Template Model Builder (TMB) software, developed by [Kristensen et al. \(2016\)](#) and possibly the most efficient likelihood-based way of doing such task.
3. Take advantage of TMB's functionalities with special attention to the computation of gradients and Hessians via a *state-of-art* automatic differentiation (AD) implementation; and a joint likelihood marginalization via an efficient Laplace approximation routine.
4. Study the model identifiability through the proposition of different complexity level models in terms of parametric space and latent effect structures.
5. Make exact likelihood-based inference to the cluster and cause-specific CIF of clustered competing risks data.

1.2 JUSTIFICATION

In the biomedical statistical modeling literature, the study of disease occurrence in related individuals receives the name of family studies. Key points of interest are the within-family dependence and determining the role of different risk factors. The within-family dependence may reflect both disease heritability and the impact of shared environmental effects. The role of different risk factors arrives in the class of multivariate models, which options are limited in the statistical literature. Thus, the number of statistical models for competing risks data that accommodate the within-cluster/family dependence is even more limited. Some modeling options are briefly commented in [Cederkvist et al. \(2019\)](#), with his pairwise composite approach being proposed as a new and better option to model the cause-specific cumulative incidence function (CIF), describing age at disease onset, of clustered competing risks data on the probability scale. We propose to model the cause-specific CIF and accommodate the within-family dependence in the same fashion (via a latent structure that allows the absolute risk and the failure time distribution to vary between families) but with an easier framework, based on a multinomial generalized linear mixed model approach.

1.3 LIMITATION

This work restraint to the proposition and model identifiability study of a multinomial model for the cause-specific cumulative incidence function (CIF) of competing risks data, with a latent effect structure to accommodate within-family dependence with regard to both risk and timing. Given its considerable model complexity, hypothesis tests; residual analysis; and good-of-fit measures are not contemplated.

1.4 THESIS ORGANIZATION

This master thesis contains 6 chapters including this introduction. [Chapter 2](#) presents a systematic review of the main aspects involved in the formulation, optimization, and implementation of a generalized linear mixed model (GLMM). Given the modeling framework overview, [Chapter 3](#) presents our multinomial GLMM (multiGLMM) to model the cause-specific cumulative incidence function (CIF) of clustered competing risks data. In [Chapter 4](#) we describes the simulation procedure to generate synthetic data and present some model particularities. In [Chapter 5](#) the obtained results are presented, and in [Chapter 6](#) we discuss the contributions of this thesis and present some suggestions for future work.

2 GENERALIZED LINEAR MIXED MODELS: FORMULATION, OPTIMIZATION, AND IMPLEMENTATION

This chapter presents a systematic review of the main theoretical aspects involved in the formulation, estimation, and implementation of a generalized linear mixed model (GLMM). We start in [Section 2.1](#) with the model formulation framework, concluding with the so-called joint or full likelihood function. [Section 2.2](#) address the marginalization of that joint likelihood, performed here in terms of a Laplace approximation technique. [Section 2.3](#) discusses available alternatives for the marginal likelihood parameters optimization. [Section 2.4](#) present the automatic differentiation (AD) procedure, the most efficient routine for the computation of derivatives, and a key point for us. Last but not least, in [Section 2.5](#) we present the computational tool used to perform all the discussed methodology, a very exciting R ([R Core Team, 2021](#)) package called TMB: Template Model Builder, developed by [Kristensen et al. \(2016\)](#).

2.1 FORMULATION: OBTAINING A JOINT LIKELIHOOD FUNCTION

We model an n -vector of exponential family random variables \mathbf{Y} , in terms of its conditional expected value $\boldsymbol{\mu} \equiv \mathbb{E}(\mathbf{Y} | \mathbf{X}, \mathbf{u})$, via a linear combination called of linear predictor and generally expressed by

$$g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u}, \quad \mathbf{u} \sim \text{Multivariate Normal}(\mathbf{0}, \boldsymbol{\Sigma}). \quad (2.1)$$

In other words, a GLMM ([McCULLOCH; SEARLE, 2001](#)) is a generalized linear model (GLM) in which the linear predictor depends on some Gaussian latent effects, \mathbf{u} , times a latent effects design-matrix \mathbf{Z} . Since we do not observe the latent component, an exemplification of the idea embedded in matrix \mathbf{Z} is welcome. Suppose e.g., three individuals (or clusters) and that each one has two measures. This configures a repeated measures context, the most common latent structure in family studies. Also, it is reasonable to admit that each individual has its particular latent effect value. Consequently, we have

$$\mathbf{Z}\mathbf{u} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_1 \\ u_2 \\ u_2 \\ u_3 \\ u_3 \end{bmatrix} = \begin{bmatrix} u_1 \\ u_1 \\ u_2 \\ u_2 \\ u_3 \\ u_3 \end{bmatrix},$$

where $\mathbf{u}^\top = [u_1 \ u_2 \ u_3]$ and \mathbf{Z} has the role of projecting the values of \mathbf{u} to match the number of measures.

In a mixed model the mean structure is approached as a combination of probability distributions. It is a combination since we have to assume probabilistic structures for the observed and non-observed/latent data. To each observed variable y_{ij} we have a probability distribution of the exponential family, denoted by $f(y_{ij} | \boldsymbol{u}_i, \boldsymbol{\theta})$. To the latent effect we have, generally, a (multivariate) Gaussian distribution, denoted by $f(\boldsymbol{u}_i | \boldsymbol{\Sigma})$. To each individual or unity under study i , and to each measure j , we have the product of these probability densities, a likelihood contribution.

Our goal is to estimate the parameter vector $\boldsymbol{\theta} = [\boldsymbol{\beta} \ \boldsymbol{\Sigma}]^\top$ of a mean structure, as in [Equation 2.1](#). Besides the role of emphasizing the fact that $\boldsymbol{\mu}$ is a function of $\boldsymbol{\theta}$, and that we want to estimate $\boldsymbol{\theta}$, the likelihood function ties the probability densities i.e., the likelihood is the product of the product of probability densities, to each subject i . Since Y_i are mutually independent, the likelihood for $\boldsymbol{\theta}$ can be written as

$$L(\boldsymbol{\theta} | \mathbf{y}, \mathbf{u}) = \prod_{i=1}^I \prod_{j=1}^{n_i} f(y_{ij} | \boldsymbol{u}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}) f(\boldsymbol{u}_i | \boldsymbol{\Sigma}). \quad (2.2)$$

From standard probability theory is easy to see that in the right-hand side (r.h.s.) we have a joint density, consequently, [Equation 2.2](#) represents what is called a full or a joint likelihood function. What makes problematic working with this joint likelihood is that we do not have all the necessary information to just maximize it and get the desired parameter estimates. The latent effect \boldsymbol{u} is *latent* i.e., we do not observe it. To handle this we have basically two available paths.

2.2 MARGINALIZATION: LAPLACE APPROXIMATION AND ALTERNATIVES

To deal with a joint likelihood function as in [Equation 2.2](#) we have a choice to make. Be or not to be Bayesian. Each choice has its own difficulties, advantages, and characteristics.

The Bayesian path assumes that all $\boldsymbol{\theta}$ components are random variables. With all parameters being treated as random variables, and since we do not observe them, what the Bayesian framework does is try to compute the mode of each “parameter” marginal distribution, generally, via a sampling algorithm called MCMC: Markov chain Monte Carlo ([GELFAND; SMITH, 1990](#); [DIACONIS, 2009](#)).

The advantage of being Bayesian is that we can reach an MCMC algorithm to basically any statistical model, the disadvantage is that this approach is very time consuming and we have to propose prior distributions to each “parameter”. These prior proposals are not always easy to make, and the resulting marginal distributions can be very depending of it. A Bayesian approach can be applied in basically any context, without guarantees that will work - obtain convergence to all parameters is not

a straightforward task. However, in complex scenarios they can be the only available method to “maximize” the likelihood function. This is not the case here.

We have a joint density where one of the random variables is not observed, but we are not interested in it, only in the variance parameters inherent in it. Again, from standard probability theory, if we have a joint density we can just integrate out the undesired variable resulting in

$$\begin{aligned} L(\boldsymbol{\theta} \mid \mathbf{y}) &= \prod_{i=1}^I \int_{\mathcal{R}^{u_i}} \left[\prod_{j=1}^{n_i} f(y_{ij} \mid \mathbf{u}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}) f(\mathbf{u}_i \mid \boldsymbol{\Sigma}) \right] d\mathbf{u}_i \\ &= \prod_{i=1}^I \int_{\mathcal{R}^{u_i}} f(\mathbf{y}_i, \mathbf{u}_i \mid \boldsymbol{\theta}) d\mathbf{u}_i, \end{aligned} \quad (2.3)$$

a marginal density that keeps the parameters $\boldsymbol{\Sigma}$ of the integrated variable.

When the response distribution of a mixed model is Gaussian, is analytically tractable to integrate \mathbf{u} out of the joint density. Consequently, it is possible to evaluate the marginal likelihood exactly. This is the case of the linear mixed models (LMMs) and one of the main differences to the GLMMs. When the response distribution is not Gaussian, generally, it is not anymore analytically tractable to integrate out the latent effect. So what do we do? Well, we have basically two options again.

We can avoid the integrals in [Equation 2.3](#), replacing it by integrals that are more analytically tractable. This can be performed via an algorithm called Expectation-Maximization (EM), proposed by [Dempster, Laird & Rubin \(1977\)](#). This approach is considered a little bit naive and generally is not recommended if you have a better option. The other option consists of performing a numerical integration i.e., approximating the integral. The most common way of doing that in the statistical modeling literature is via an importance sampling version of the Gaussian quadrature rule, denoted adaptive Gaussian quadrature (AGQ) ([PINHEIRO; CHAO, 2006](#)). In general, adaptive Gaussian quadratures are not so simple to use (computationally expensive; we have to choose how many integration points will be used; and we also have to choose an importance distribution to approximate the integrand).

To us, the better option consists in take advantage of the exponential family structure together with the fact that we are dealing with Gaussian latent effects. These ideas converge to an adaptive Gaussian quadrature with one integration point, also called as *Laplace approximation* ([MOLENBERGHS; VERBEKE, 2005; SHUN; MCCULLAGH, 1995; TIERNEY; KADANE, 1986; WOOD, 2015](#)).

With an integral that is analytically intractable, we may approximate it to obtain a tractable closed-form expression allowing then the numerical maximization of the resulting marginal likelihood function ([BONAT; RIBEIRO, 2016](#)). The Laplace

approximation has been designed to approximate integrals in the form

$$\int_{\mathcal{R}^{\mathbf{u}_i}} \exp\{Q(\mathbf{u}_i)\} d\mathbf{u}_i \approx (2\pi)^{n_u/2} |Q''(\hat{\mathbf{u}}_i)|^{-1/2} \exp\{Q(\hat{\mathbf{u}}_i)\}, \quad (2.4)$$

where $Q(\mathbf{u}_i)$ is a known, unimodal bounded function, and $\hat{\mathbf{u}}_i$ is the value for which $Q(\mathbf{u}_i)$ is maximized. As Wood (2015) shows, a Laplace approximation consists of a second order Taylor expansion of $\log f(\mathbf{y}_i, \mathbf{u}_i | \boldsymbol{\theta})$, about $\hat{\mathbf{u}}_i$, that gives

$$\log f(\mathbf{y}_i, \mathbf{u}_i | \boldsymbol{\theta}) \approx \log f(\mathbf{y}_i, \hat{\mathbf{u}}_i | \boldsymbol{\theta}) - \frac{1}{2} (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{H} (\mathbf{u}_i - \hat{\mathbf{u}}_i),$$

where $\mathbf{H} = -\nabla_{\mathbf{u}}^2 \log f(\mathbf{y}_i, \hat{\mathbf{u}}_i | \boldsymbol{\theta})$. Hence, we can approximate the joint by

$$f(\mathbf{y}_i, \mathbf{u}_i | \boldsymbol{\theta}) \approx f(\mathbf{y}_i, \hat{\mathbf{u}}_i | \boldsymbol{\theta}) \exp \left\{ -\frac{1}{2} (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{H} (\mathbf{u}_i - \hat{\mathbf{u}}_i) \right\}. \quad (2.5)$$

From here we start to take advantage of the points mentioned above.

First, the fact that we are dealing with Gaussian distributed latent effects. In Equation 2.5 we have the core of a Gaussian density, that complete is

$$\int_{\mathcal{R}^{\mathbf{u}_i}} \frac{1}{(2\pi)^{n_u/2} |\mathbf{H}^{-1}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{H} (\mathbf{u}_i - \hat{\mathbf{u}}_i) \right\} d\mathbf{u}_i = 1$$

i.e., integrates to 1. Integrating Equation 2.5 follows that

$$\begin{aligned} \int_{\mathcal{R}^{\mathbf{u}_i}} f(\mathbf{y}_i, \mathbf{u}_i | \boldsymbol{\theta}) d\mathbf{u}_i &\approx f(\mathbf{y}_i, \hat{\mathbf{u}}_i | \boldsymbol{\theta}) \int_{\mathcal{R}^{\mathbf{u}_i}} \exp \left\{ -\frac{1}{2} (\mathbf{u}_i - \hat{\mathbf{u}}_i)^\top \mathbf{H} (\mathbf{u}_i - \hat{\mathbf{u}}_i) \right\} d\mathbf{u}_i \\ &= (2\pi)^{n_u/2} |\mathbf{H}|^{-1/2} f(\mathbf{y}_i, \hat{\mathbf{u}}_i | \boldsymbol{\theta}) \end{aligned}$$

i.e., we get Equation 2.4, a first order Laplace approximation to the integral. Careful accounting of the approximation error shows it to generally be $\mathcal{O}(n^{-1})$, where n is the sample size, and assuming a fixed length for \mathbf{u}_i (WOOD, 2015).

The second advantage of a Laplace approximation approach in a GLMM is the exponential family structure. In a usual GLMM the response follows a one-parameter exponential family distribution that can be written as

$$f(\mathbf{y}_i | \mathbf{u}_i, \boldsymbol{\theta}) = \exp \left\{ \mathbf{y}_i^\top (\mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{u}_i) - \mathbf{1}_i^\top b(\mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{u}_i) + \mathbf{1}_i^\top c(\mathbf{y}_i) \right\},$$

where $b(\cdot)$ and $c(\cdot)$ are known functions.

This general and easy to compute expression, together with a (multivariate) Gaussian distribution, highlights the convenience of the Laplace method. The $Q(\mathbf{u}_i)$ function to be maximized can be expressed as

$$\begin{aligned} Q(\mathbf{u}_i) &= \mathbf{y}_i^\top (\mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{u}_i) - \mathbf{1}_i^\top b(\mathbf{x}_i \boldsymbol{\beta} + \mathbf{z}_i \mathbf{u}_i) + \mathbf{1}_i^\top c(\mathbf{y}_i) \\ &\quad - \frac{n_u}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \mathbf{u}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{u}_i. \end{aligned} \quad (2.6)$$

The approximation in [Equation 2.4](#) requires the maximum $\hat{\boldsymbol{u}}_i$ of the function $Q(\boldsymbol{u}_i)$. As we assume a Gaussian distribution with a known mean for the latent effects, we have the perfect initial guess for a gradient-based maximization method, as the Newton-Raphson (NR) algorithm.

The NR method consists of an iterative scheme as follows:

$$\boldsymbol{u}_i^{(k+1)} = \boldsymbol{u}_i^{(k)} - Q''(\boldsymbol{u}_i^{(k)})^{-1} Q'(\boldsymbol{u}_i^{(k)}), \quad k = 0, 1, \dots$$

until convergence, which gives $\hat{\boldsymbol{u}}_i$. At this stage, all parameters $\boldsymbol{\theta}$ are considered known. [Bonat & Ribeiro \(2016\)](#) presents the generic expressions for the derivatives required by the NR method, given by the following:

$$\begin{aligned} Q'(\boldsymbol{u}_i^{(k)}) &= \{\boldsymbol{y}_i - b'(\boldsymbol{x}_i\boldsymbol{\beta} + \boldsymbol{z}_i\boldsymbol{u}_i^{(k)})\}^\top - \boldsymbol{u}_i^{(k)\top} \boldsymbol{\Sigma}^{-1}, \\ Q''(\boldsymbol{u}_i^{(k)}) &= -\text{diag}\{b''(\boldsymbol{x}_i\boldsymbol{\beta} + \boldsymbol{z}_i\boldsymbol{u}_i^{(k)})\} - \boldsymbol{\Sigma}^{-1}. \end{aligned}$$

We have the initial guesses at $k = 0$.

Finally, the marginal log-likelihood function returned by the Laplace approximation, to each individual or unit under study i , is as follows:

$$\begin{aligned} l(\boldsymbol{\theta} | \boldsymbol{y}_i) = \log L(\boldsymbol{\theta} | \boldsymbol{y}_i) &= \frac{n}{2} \log(2\pi) - \frac{1}{2} \log \left| \text{diag}\{b''(\boldsymbol{x}_i\boldsymbol{\beta} + \boldsymbol{z}_i\hat{\boldsymbol{u}}_i)\} + \boldsymbol{\Sigma}^{-1} \right| \\ &\quad + \boldsymbol{y}_i^\top (\boldsymbol{x}_i\boldsymbol{\beta} + \boldsymbol{z}_i\hat{\boldsymbol{u}}_i) - \boldsymbol{1}_i^\top b(\boldsymbol{x}_i\boldsymbol{\beta} + \boldsymbol{z}_i\hat{\boldsymbol{u}}_i) + \boldsymbol{1}_i^\top c(\boldsymbol{y}_i) \\ &\quad - \frac{n_u}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \hat{\boldsymbol{u}}_i^\top \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{u}}_i, \end{aligned}$$

that can now be numerically maximized over the model parameters $\boldsymbol{\theta} = [\boldsymbol{\beta} \ \boldsymbol{\Sigma}]^\top$.

2.3 OPTIMIZATION: MARGINAL LIKELIHOOD FUNCTION

At this point it is already clear that we have two optimizations to be performed, an “inside” and an “outside” optimization. The inside one is made into the Laplace approximation layer via a Newton-Raphson algorithm, a Newton’s method. The outside optimization is made with the Laplace approximation outputs i.e., the maximization of [Equation 2.3](#)’s marginal log-likelihood over its parameters $\boldsymbol{\theta}$. This task is usually performed via a quasi-Newton method, we focus on two of the most traditional ones: the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm and the PORT routines.

The inside optimization is the numerical maximization of the joint log-likelihood with respect to (w.r.t.) its latent effects. This is kind of a simple task since all model parameters are considered as fixed, and we “know” that the latent effects are distributed with zero mean i.e., we have the perfect initial guess. In this context, the use of a Newton’s method is straightforward. When we talk about the outside optimization

it is a completely different scenario, it is not straightforward to find a good initial guess or reach convergence. Thus, more robust methods are a good choice.

In optimization, Newton methods are algorithms for finding local maxima and minima of functions i.e., the search for the zeroes of the gradient of that function. Newton methods are characterized by the use of a symmetric matrix of function's second derivatives, the Hessian matrix. Quasi-Newton methods are based on Newton's method and are seen as an alternative to it. They can be used if the Hessian is unavailable or if is too expensive to compute it at every iteration.

As shown in [Nocedal & Wright \(2006\)](#), major advantages of quasi-Newton methods over Newton's method are that the Hessian matrix does not need to be computed, it is approximated; and it also does not need to be inverted. Newton's method requires the Hessian to be inverted, typically by solving a system of linear equations - often quite costly. In contrast, quasi-Newton methods usually generate an estimate of it directly. As in Newton's method, they use a second-order approximation to find the minimum of a function $f(\mathbf{x})$. The Taylor series of $f(\mathbf{x})$ around an iterate is

$$f(\mathbf{x}_k + \Delta\mathbf{x}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^\top \mathbf{B} \Delta\mathbf{x},$$

where $\nabla f(\cdot)$ is the gradient, and \mathbf{B} an approximation to the Hessian matrix. The gradient of this approximation w.r.t. $\Delta\mathbf{x}$ is

$$\nabla f(\mathbf{x}_k + \Delta\mathbf{x}) \approx \nabla f(\mathbf{x}_k) + \mathbf{B} \Delta\mathbf{x},$$

setting this gradient to zero provides the Newton step:

$$\Delta\mathbf{x} = -\mathbf{B}^{-1} \nabla f(\mathbf{x}_k).$$

The Hessian approximation \mathbf{B} is chosen to satisfy

$$\nabla f(\mathbf{x}_k + \Delta\mathbf{x}) = \nabla f(\mathbf{x}_k) + \mathbf{B} \Delta\mathbf{x},$$

which is called the *secant* equation i.e., the Taylor series of the gradient itself. Solving for \mathbf{B} and applying the Newton's step with the updated value is equivalent to the *secant* method. Quasi-Newton methods are a generalization of the secant method to find the root of the first derivative for multidimensional problems. The various quasi-Newton methods differ in their choice of the solution to the secant equation.

In a general quasi-Newton method, the unknown \mathbf{x}_k is updated applying the Newton's step calculated using the current approximate Hessian matrix \mathbf{B}_k in the following fashion:

- $\Delta\mathbf{x}_k = -\alpha_k \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k)$, with α chosen to satisfy the so called Wolfe conditions ([NOCEDAL; WRIGHT, 2006](#), p. 34);

- $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta\mathbf{x}_k$;
- The gradient computed at the new point $\nabla f(\mathbf{x}_{k+1})$, and $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$ is used to update the approximate Hessian \mathbf{B}_{k+1} , or directly its inverse $\mathbf{H}_{k+1} = \mathbf{B}_{k+1}^{-1}$.

The most popular quasi-Newton method is the BFGS algorithm, named for its discoverers Broyden, Fletcher, Goldfarb, and Shanno. It has the following update formula

$$\begin{aligned}\mathbf{B}_{k+1} &= \mathbf{B}_k + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \Delta\mathbf{x}_k} - \frac{\mathbf{B}_k \Delta\mathbf{x}_k (\mathbf{B}_k \Delta\mathbf{x}_k)^\top}{\Delta\mathbf{x}_k^\top \mathbf{B}_k \Delta\mathbf{x}_k}, \\ \mathbf{H}_{k+1} = \mathbf{B}_{k+1}^{-1} &= \left(\mathbf{I} - \frac{\Delta\mathbf{x}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \Delta\mathbf{x}_k} \right) \mathbf{H}_k \left(\mathbf{I} - \frac{\mathbf{y}_k \Delta\mathbf{x}_k^\top}{\mathbf{y}_k^\top \Delta\mathbf{x}_k} \right) + \frac{\Delta\mathbf{x}_k \Delta\mathbf{x}_k^\top}{\mathbf{y}_k^\top \Delta\mathbf{x}_k}.\end{aligned}$$

Another quasi-Newton method popular in the statistical modeling literature, is the one based on the PORT routines (<http://www.netlib.org/port/>). It is a Fortran mathematical subroutine library designed to be *portable* over different types of computers, developed by David Gay in the Bell Labs (GAY, 1990). Is a quasi-Newton adaptive nonlinear least-squares algorithm (DENNIS; GAY; WELSCH, 1981) with the following update formula

$$\begin{aligned}\mathbf{B}_{k+1} &= \mathbf{B}_k \\ &+ \frac{(\mathbf{y}_k - \mathbf{B}_k \Delta\mathbf{x}_k) \Delta\mathbf{x}_k^\top \mathbf{B}_k + \mathbf{B}_k \Delta\mathbf{x}_k (\mathbf{y}_k - \mathbf{B}_k \Delta\mathbf{x}_k)^\top}{\Delta\mathbf{x}_k^\top \mathbf{B}_k \Delta\mathbf{x}_k} \\ &- \frac{\Delta\mathbf{x}_k^\top (\mathbf{y}_k - \mathbf{B}_k \Delta\mathbf{x}_k) \mathbf{B}_k \Delta\mathbf{x}_k \Delta\mathbf{x}_k^\top \mathbf{B}_k}{(\Delta\mathbf{x}_k^\top \mathbf{B}_k \Delta\mathbf{x}_k)^\top \Delta\mathbf{x}_k^\top \mathbf{B}_k \Delta\mathbf{x}_k}.\end{aligned}$$

As Nocedal & Wright (2006) points out, each quasi-Newton method iteration can be performed at a cost of $\mathcal{O}(n^2)$ arithmetic operations (plus the cost of function and gradient evaluations); there are no $\mathcal{O}(n^3)$ operations such as linear system solves or matrix-matrix operations. In the BFGS algorithm is known that the rate of convergence is superlinear, which is a valid assumption to any quasi-Newton method and is fast enough for most practical purposes. Even though Newton's method converges more rapidly, quadratically, its cost per iteration usually is higher because of its need for second derivatives and solution of a linear system.

In this thesis, the used BFGS implementation is the one in the R (R Core Team, 2021) function `base::optim()`, and the PORT routine used is the one implemented in the R function `base::nlminb()`.

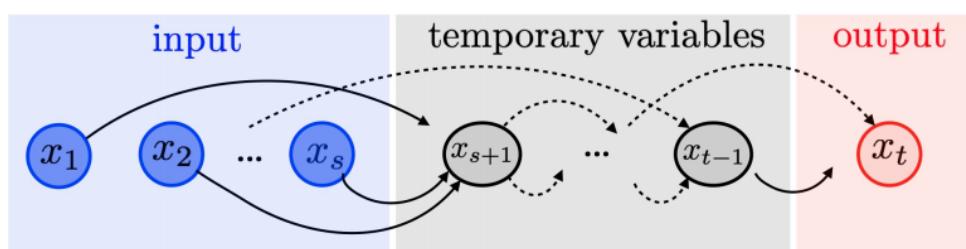
2.4 AD: AUTOMATIC DIFFERENTIATION

The computation of gradients, $\nabla f(\mathbf{x})$, are a fundamental and crucial task but also the main computational bottleneck to any Newton and quasi-Newton method. We choose to use the most efficient manner of computing gradients, and one of the best scientific computing techniques but still not so famous in the statistical modeling literature, the *automatic differentiation* (AD) procedure. AD has two modes, the so-called forward and reverse mode. We will talk a bit about both but we will use only the reverse mode. The reason can be illustrated by a simple example, given later.

Automatic differentiation, also called algorithmic differentiation or computational differentiation, is a set of techniques to numerically and recursively evaluate the derivative of a function specified by a computer program. AD techniques are based on the observation that any function, no matter how complicated, is evaluated by performing a sequence of simple elementary operations involving just one or two arguments at a time. Derivatives of arbitrary order can be computed automatically, automatized and accurately to working precision. Most of the information in this section was taken of [Peyré \(2020\)](#), but [Wood \(2015, p. 120\)](#) and [Nocedal & Wright \(2006, p. 204\)](#) are also very good references.

The most common differentiation approaches are finite differences (FD) and symbolic calculus. Considering a function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ and the goal of deriving a method to evaluate $\nabla f : \mathbb{R}^p \rightarrow \mathbb{R}^p$, the approximation of this vector field via FD would require $p + 1$ evaluations of f . The same task via reverse mode AD has in most cases a cost proportional to a single evaluation of f . AD is similar to symbolic calculus in the sense that it provides an exact gradient computation, up to machine precision. However, symbolic calculus does not takes into account the underlying algorithm which compute the function, while AD factorizes the computation of the derivative according to an efficient algorithm. The use of AD is inherent to the use of a computational graph, as exemplified in [Figure 2](#).

FIGURE 2 – A COMPUTATIONAL GRAPH



SOURCE: [Peyré \(2020, p. 31\)](#).

Assuming that f is implemented in an algorithm, the goal is to compute the

derivatives

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_k} \in \mathbb{R}^{n_t \times n_k},$$

for a numerical algorithm (succession of functions) of the form

$$\forall k = s + 1, \dots, t, \quad \mathbf{x}_k = f_k(\mathbf{x}_1, \dots, \mathbf{x}_{k-1}),$$

where f_k is a function which only depends on the previous variables. The computational graph as in [Figure 2](#), has the role of represent the linking of the variables involved in f_k to \mathbf{x}_k . The evaluation of $f(\mathbf{x})$ corresponds to a forward traversal of this graph.

Now, how we evaluate f through the graph? Via one of the AD modes.

2.4.1 Forward Mode

The forward mode correspond to the usual way of computing differentials. The method initialize with the derivative of the input nodes

$$\frac{\partial \mathbf{x}_1}{\partial \mathbf{x}_1} = \text{Id}_{n_1 \times n_1}, \quad \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1} = \mathbf{0}_{n_2 \times n_1}, \quad \frac{\partial \mathbf{x}_s}{\partial \mathbf{x}_1} = \mathbf{0}_{n_s \times n_1},$$

and then iteratively make use of the following recursion formula

$$\begin{aligned} \forall k = s + 1, \dots, t, \\ \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_1} = \sum_{l \in \text{father}(k)} \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_l} \times \frac{\partial \mathbf{x}_l}{\partial \mathbf{x}_1} = \sum_{l \in \text{father}(k)} \frac{\partial}{\partial \mathbf{x}_l} f_k(\mathbf{x}_1, \dots, \mathbf{x}_{k-1}) \times \frac{\partial \mathbf{x}_l}{\partial \mathbf{x}_1}. \end{aligned}$$

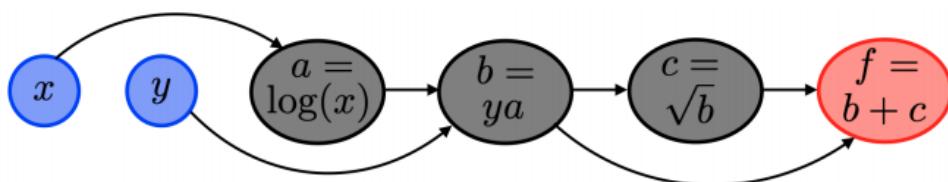
The notation “father(k)” denotes the nodes $l < k$ of the graph that are connected to k . We make use of [Peyré \(2020\)](#), p. 32’s simple example.

Example. Consider the function

$$f(x, y) = y \log(x) + \sqrt{y \log(x)}$$

with the corresponding computational graph being displayed in [Figure 3](#).

FIGURE 3 – EXAMPLE OF A SIMPLE COMPUTATIONAL GRAPH



SOURCE: [Peyré \(2020\)](#), p. 33).

The forward mode iterations to compute the derivative w.r.t. x following the computational graph, are given by

$$\begin{aligned}
 \frac{\partial x}{\partial x} &= 1, \quad \frac{\partial y}{\partial x} = 0 \\
 \frac{\partial a}{\partial x} &= \frac{\partial a}{\partial x} \frac{\partial x}{\partial x} = \frac{1}{x} \frac{\partial x}{\partial x} && \{x \mapsto a = \log(x)\} \\
 \frac{\partial b}{\partial x} &= \frac{\partial b}{\partial a} \frac{\partial a}{\partial x} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial x} = y \frac{\partial a}{\partial x} + 0 && \{(y, a) \mapsto b = ya\} \\
 \frac{\partial c}{\partial x} &= \frac{\partial c}{\partial b} \frac{\partial b}{\partial x} = \frac{1}{2\sqrt{b}} \frac{\partial b}{\partial x} && \{b \mapsto c = \sqrt{b}\} \\
 \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial x} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial x} = 1 \frac{\partial b}{\partial x} + 1 \frac{\partial c}{\partial x} && \{(b, c) \mapsto f = b + c\}
 \end{aligned}$$

To compute the derivative w.r.t. y we run another forward process

$$\begin{aligned}
 \frac{\partial x}{\partial y} &= 0, \quad \frac{\partial y}{\partial y} = 1 \\
 \frac{\partial a}{\partial y} &= \frac{\partial a}{\partial x} \frac{\partial x}{\partial y} = 0 && \{x \mapsto a = \log(x)\} \\
 \frac{\partial b}{\partial y} &= \frac{\partial b}{\partial a} \frac{\partial a}{\partial y} + \frac{\partial b}{\partial y} \frac{\partial y}{\partial y} = 0 + a \frac{\partial y}{\partial y} && \{(y, a) \mapsto b = ya\} \\
 \frac{\partial c}{\partial y} &= \frac{\partial c}{\partial b} \frac{\partial b}{\partial y} = \frac{1}{2\sqrt{b}} \frac{\partial b}{\partial y} && \{b \mapsto c = \sqrt{b}\} \\
 \frac{\partial f}{\partial y} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial y} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial y} = 1 \frac{\partial b}{\partial y} + 1 \frac{\partial c}{\partial y} && \{(b, c) \mapsto f = b + c\}
 \end{aligned}$$

2.4.2 Reverse Mode

Instead of evaluating the differentials for all the input nodes, which is problematic for a large number of nodes, the reverse mode evaluates the differentials of the output node w.r.t. all the inner nodes.

The method is based on a backward adjoint chain rule and initialize with the derivative of the final node

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_t} = \text{Id}_{n_t \times n_t},$$

and then from the last to the first node, iteratively make use of the following recursion formula

$$\begin{aligned}
 \forall k &= t-1, t-2, \dots, 1, \\
 \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} &= \sum_{m \in \text{son}(k)} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_m} \times \frac{\partial \mathbf{x}_m}{\partial \mathbf{x}_k} = \sum_{m \in \text{son}(k)} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_m} \times \frac{\partial}{\partial \mathbf{x}_k} f_m(\mathbf{x}_1, \dots, \mathbf{x}_m).
 \end{aligned}$$

The notation “ $\text{son}(k)$ ” denotes the nodes $m < k$ of the graph that are connected to k . To be clear, the same simple example.

Example. Consider again the function

$$f(x,y) = y \log(x) + \sqrt{y \log(x)}.$$

The iterations of the reverse mode are given by

$$\begin{aligned} \frac{\partial f}{\partial f} &= 1 \\ \frac{\partial f}{\partial c} &= \frac{\partial f}{\partial f} \frac{\partial f}{\partial c} = \frac{\partial f}{\partial f} 1 & \{c \mapsto f = b + c\} \\ \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial f}{\partial f} \frac{\partial f}{\partial b} = \frac{\partial f}{\partial c} \frac{1}{2\sqrt{b}} + \frac{\partial f}{\partial f} 1 & \{b \mapsto c = \sqrt{b}, b \mapsto f = b + c\} \\ \frac{\partial f}{\partial a} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial a} = \frac{\partial f}{\partial b} y & \{a \mapsto b = ya\} \\ \frac{\partial f}{\partial y} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial y} = \frac{\partial f}{\partial b} a & \{y \mapsto b = ya\} \\ \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} = \frac{\partial f}{\partial a} \frac{1}{x} & \{x \mapsto a = \log(x)\} \end{aligned}$$

This is the advantage of reverse mode over the forward mode. A single traversal over the computational graph allows to compute both derivatives w.r.t. x and y , while the forward mode necessities two processes.

An drawback of the reverse mode is the need to store the entire computational graph, which is needed for the reverse sweep. In principle, storage of this graph is not too difficult to implement. However, the main benefit of AD is higher accuracy, and in many applications the cost is not critical.

2.5 TMB: TEMPLATE MODEL BUILDER

Note that the goal of AD is not to define an efficient computational graph, it is up to the user to provide it. However, computing an efficient graph associated to a mathematical formula is a complicated combinatorial problem. Thus, since our goal is to be able to fit our desired statistical models, a computational tool able to efficiently define and implement this computational graph is make necessary. To solve this and many other tasks we have the Template Model Builder (TMB), developed by [Kristensen et al. \(2016\)](#).

TMB (<http://tmb-project.org>) is an R ([R Core Team, 2021](#)) package for fitting statistical latent variable models to data, inspired by AD Model Builder (ADMB) ([FOURNIER et al., 2012](#)). ADMB is a statistical application for fitting nonlinear statistical models and solve optimization problems, that implements AD using C++ classes and a native template language. Unlike most R packages, in TMB the model is formulated in C++. This characteristic provides great flexibility but requires some familiarity with the C/C++ programming language.

With TMB a user should be able to quickly implement complex latent effect models through simple C++ templates. As an illustrative example let us consider a simple mixed logistic regression i.e., a binomial GLMM with a logistic link function. The latent structure is in the context of repeated measures, the same subject is observed three times. Trying to keep it simple, no covariates. A hierarchical model's description is given by

$$\begin{aligned} y_{ij} \mid u_i &\sim \text{Binomial}(n, p_{ij}) \\ u_i &\sim \text{Normal}(0, 1) \\ g(p_{ij}) = \text{logit}(p_{ij}) &= \log \frac{p_{ij}}{1 - p_{ij}} = \beta + u_i \\ p_{ij} &= \frac{\exp\{\beta + u_i\}}{1 + \exp\{\beta + u_i\}}, \quad i \text{ the subject, } j \text{ the subject observation (1,2,3).} \end{aligned}$$

The TMB implementation of this model is provided in Figure 4.

To keep the coherence would be more adequate to fit here a multinomial model. However, I want to show you that in TMB we can also simulate data from the model but not all r-distributions are implemented, as is the case of the multinomial. For this reason, a binomial model was the choice. An even easier implementation of a logistic model is available in TMB, called `dbinom_robust()`, where we pass directly the $\text{logit}(\cdot)$ but we do not have an `rbinom_robust()` implementation available. Just for the sake of completeness, in Figure 5 we have the R code showing how to manipulate the C++ model template in terms of object definitions and parameters estimation and extraction for the logistic mixed model.

In this chapter we describe step-by-step all the processes involved in the formulation and parameter estimation of a GLMM. With TMB all this is put in practice in an efficient and robust fashion. The user needs to provide the negative joint log-likelihood function writing in a C++ template as exemplified in Figure 4, using specialized macros that pass the parameters, latent effects and data from R, as exemplified in Figure 5.

When the model presents latent effects, during the model template compilation the latent effects are integrated out via an efficient Laplace approximation routine with the inner optimization made by a Newton algorithm, and the negative marginal log-likelihood gradient is computed, via AD. The negative marginal log-likelihood is returned into an R object that can then be optimized using the user's favorite quasi-Newton routine, available in R. All these procedures are briefly exemplified for a logistic mixed model in Figure 4 and Figure 5.

To accomplish all that, TMB combines some state-of-art software

- CppAD, a C++ AD package (<https://coin-or.github.io/CppAD/>);
- Eigen ([GUENNEBAUD; JACOB et al., 2010](#)), a C++ templated matrix-vector library;

- CHOLMOD, C sparse matrix routines available from R, used to obtain an efficient implementation of the Laplace approximation with exact derivatives (<https://developer.nvidia.com/cholmod>);
- Parallelism through BLAS (<http://www.netlib.orgblas/>), a Fortran tuned set of Basic Linear Algebra Subprograms;

FIGURE 4 – R CODE FOR THE TMB IMPLEMENTATION OF A LOGISTIC MIXED MODEL

```

1 dll <- 'model'
2 filename <- paste0(dll, '.cpp')
3 writeLines({'// A LOGISTIC MIXED MODEL (RANDOM INTERCEPT)
4 #include <TMB.hpp>
5 template<class Type>
6 Type objective_function<Type>::operator() ()
7 {
8     // SPECIFY THE MODEL INPUTS AS DATA_
9     DATA_VECTOR(y);
10    DATA_SPARSE_MATRIX(Z);
11    DATA_SCALAR(n);
12
13    // SPECIFY THE MODEL PARAMETERS AND LATENT EFFECTS AS PARAMETER_
14    PARAMETER(beta);
15    PARAMETER(logsd);           Type sd = exp(logsd);
16    PARAMETER_VECTOR(u); vector<Type> Zu = Z*u;
17
18    // IMPLEMENT THE MODEL
19    vector<Type> risk = exp(beta+Zu);
20    vector<Type> level = 1+risk;
21    vector<Type> prob = risk/level;
22
23    // nll: NEGATIVE LOG-LIKELIHOOD
24    parallel_accumulator<Type> nll(this); // DO THE MODEL IN PARALLEL
25    nll -= dnorm(u, Type(0), sd, true).sum();
26    nll -= dbinom(y, n, prob, true).sum();
27
28    // TMB ALLOWS THE USER TO WRITE THE SIMULATION CODE AS AN INTEGRATED
29    // PART OF THE C++ MODEL TEMPLATE
30    SIMULATE {
31        y = rbinom(Type(100), prob);
32        REPORT(y);
33    }
34    // WE MODEL THE LOG STANDARD DEVIATION (logsd) BUT TMB ALLOWS US TO
35    // MAKE DIRECT INFERENCE TO PARAMETER TRANSFORMATIONS, IN THIS CASE
36    // THE STANDARD DEVIATION (sd)
37    ADREPORT(sd);
38
39    return nll;
40 }, con=filename)
41
42 library(TMB) ## install.packages('TMB')
43 TMB::compile(filename)
44 dyn.load(TMB::dynlib(dll)) ## loading the C++ model to R

```

SOURCE: The author (2021).

- Matrix ([BATES; MAECHLER, 2019](#)), a rich hierarchy sparse and dense matrix classes and methods using LAPACK (<http://www.netlib.org/lapack/>) and SuiteSparse (<https://sparse.tamu.edu/>) libraries.

FIGURE 5 – R CODE FOR THE MODEL FITTING OF A LOGISTIC MIXED MODEL WRITTEN IN TMB

```

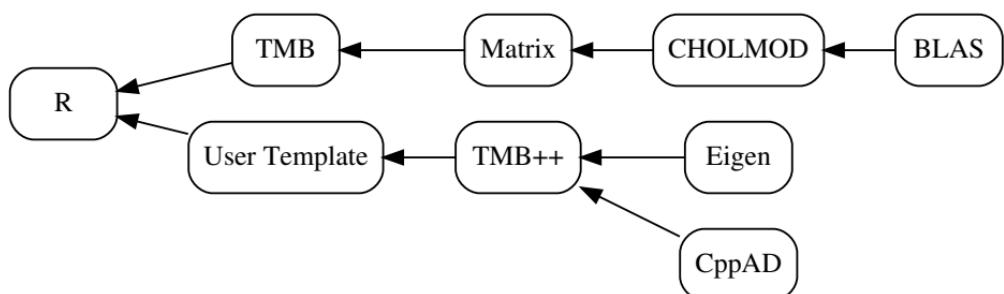
1 library(Matrix) ## install.packages('Matrix')
2 beta <- 2
3 sd <- 1
4 cs <- 3 ## cluster size
5 nc <- 50 ## number of cluster
6 n <- nc * cs
7 Z <- Matrix:::bdiag(replicate(nc, rep(1, cs), simplify=FALSE))
8 u <- rnorm(nc, mean=0, sd=sd)
9 u0 <- numeric(nc) ## empty vector (initial guess)
10 risk <- exp(beta + Z %*% u)
11 prob <- risk/(1 + risk)
12 y <- numeric(n)
13 ## base:::rbinom() is not vectorized, do a raw loop is an option
14 for (i in seq(y)) y[i] <- rbinom(n=1, size=n, prob=prob[i])
15
16 ## building objective functions with derivatives based on the compiled
17 ## C++ template
18 obj <- TMB::MakeADFun(data      =list(y=y, Z=Z, n=n),
19                         parameters=list(beta=beta, logsd=log(sd), u=u0),
20                         DLL       =dll,
21                         random    ='u')
22 set.seed(1) ## optional
23 obj$simulate() ## generating a simulation
24 ## obj$simulate(complete=TRUE)
25 (opt <- nlminb(obj$par, obj$fn, obj$gr)) ## parameters estimation
26 sdr <- TMB:::sdreport(obj) ## standard deviations
27 summary(sdr, select='fixed') ## extracting model parameters
28 summary(sdr, select='report') ## ... reported variables
29 cbind(u, summary(sdr, select='random')) ## ... random effects

```

SOURCE: The author (2021).

An overview of the package design is shown in [Figure 6](#).

FIGURE 6 – TMB PACKAGE DESIGN



SOURCE: [Kristensen et al. \(2016\)](#).

Reinforcing, some key characteristics are

- TMB employs AD to calculate first and second order derivatives of the log-likelihood function or of any objective function written in C++;
- The objective function, and its derivatives, can be called from R. Hence, parameter estimation via `base::optim()` or `base::nlminb()` is easy to be performed;
- Standard deviations of any parameter, or derived parameter, can be obtained via the *delta method* ([Ver HOEF, 2012](#)) implemented in `TMB::sdreport()`.

Here we focus on GLMMs, but basically any statistical model with a latent structure (or not), linear (or not), can be fitted with TMB. In times of *big data* and with the TMB's authors having a professional preference for state-space and spatial models, TMB has also automatic sparseness detection and some other nice built tools. Pre and post-processing of data should be done in R.

A TMB Users' mailing list exists, and it is extremely helpful for taking doubts and questions (<https://groups.google.com/g/tmb-users>). Also, a very didactic and comprehensive documentation with several examples is available online (<https://kaskr.github.io/adcomp/book/Tutorial.html>).

3 multiGLMM: A MULTINOMIAL GLMM FOR CLUSTERED COMPETING RISKS DATA

The clustered competing risks setting is a specific survival data structure. Although, we are using a general statistical modeling framework, a generalized linear mixed model (GLMM). Consequently, the data structure characteristics have to be properly accommodated into the modeling construction.

To model competing risks data we need a multivariate model and we have to choose in which scale to work on. We may work on the hazard scale and deal with the cause-specific hazard function or on the probability scale and deal with the cause-specific cumulative incidence function (CIF). By choosing the correct link function, we are able to construct an appropriate multivariate GLMM to work on the probability scale.

Our goal is to be able to deal with complex family studies, where there is generally a strong interest in describing age at disease onset in the scenarios of within-cluster dependence. The distribution of age at disease onset is directly described by the cause-specific CIF. To build a multivariate GLMM for this type of data we need to accommodate the cause-specific CIFs and the censorings. Assuming the conditional distribution for our model response as multinomial we already deal with both left-truncation and right-censoring, avoiding the specification of a censoring distribution. The cause-specific CIFs can be modeled via the link function of our, then, multinomial GLMM (multiGLMM). The multinomial distribution also guarantees that the CIFs of all causes are modeled.

Our choice for a general framework tries to make the inference of this complex model, easier. Besides, taking advantage of all the computational procedures mentioned in the previous chapter. This chapter presents our multiGLMM for clustered competing risks data and is divided into two sections. In [Section 3.1](#) we discuss in detail the cluster-specific cumulative incidence function (CIF) and in [Section 3.2](#) we present the complete modeling framework.

3.1 CLUSTER-SPECIFIC CUMULATIVE INCIDENCE FUNCTION (CIF)

Consider that the observed follow-up time of an individual is given by $T = \min(T^*, C)$, where T^* denote the failure time and C denote the censoring time. Given the possible covariates x , for a cause-specific of failure k , the cumulative incidence

function (CIF) is defined as

$$\begin{aligned} F_k(t | \mathbf{x}) &= \mathbb{P}[T \leq t, K = k | \mathbf{x}] = \int_0^t f_k(z | \mathbf{x}) dz \\ &= \int_0^t \lambda_k(z | \mathbf{x}) S(z | \mathbf{x}) dz, \quad t > 0, \quad k = 1, \dots, K, \end{aligned}$$

where $f_k(t | \mathbf{x})$ is the (sub)density for the time to a type k failure. This is the general definition of a CIF, and to define it we need to define the functions that compose the subdensity. The first is the cause-specific hazard function or process

$$\lambda_k(t | \mathbf{x}) = \lim_{h \rightarrow 0} \frac{1}{h} \mathbb{P}[t \leq T < t + h, K = k | T \geq t, \mathbf{x}], \quad t > 0, \quad k = 1, \dots, K.$$

In words, the cause-specific hazard function $\lambda_k(t | \mathbf{x})$, represents the instantaneous rate for failures of type k at time t given \mathbf{x} and all other failure types (competing causes). If we sum up all cause-specific hazard functions we get the overall hazard function,

$$\lambda(t | \mathbf{x}) = \sum_{k=1}^K \lambda_k(t | \mathbf{x}).$$

From the overall hazard function we arrive in the overall survival function,

$$S(t | \mathbf{x}) = \mathbb{P}[T > t | \mathbf{x}] = \exp \left\{ - \int_0^t \lambda(z | \mathbf{x}) dz \right\},$$

the second function that compose the subdensity $f_k(t | \mathbf{x})$. A comprehensive reference for all these definitions is the book of [Kalbfleisch & Prentice \(2002\)](#).

Until this point, we were talking about a general CIF's definition. We need now a precise framework telling us how to take into consideration our clustered/family structure. We use the same CIF specification of [Cederkvist et al. \(2019\)](#) i.e., the approach that motivated this thesis.

For two competing causes of failure, the cause-specific CIFs are specified in the following manner

$$F_k(t | \mathbf{x}, u_1, u_2, \eta_k) = \underbrace{\pi_k(\mathbf{x}, u_1, u_2)}_{\text{cluster-specific risk level}} \times \underbrace{\Phi[w_k g(t) - \mathbf{x}\gamma_k - \eta_k]}_{\text{cluster-specific failure time trajectory}}, \quad t > 0, \quad k = 1, 2, \quad (3.1)$$

i.e., as the product of a cluster-specific risk level with a cluster-specific failure time trajectory, resulting in a cluster-specific CIF. What makes these components cluster-specific are $\mathbf{u} = \{u_1, u_2\}$ and $\boldsymbol{\eta} = \{\eta_1, \eta_2\}$, Gaussian distributed latent effects with zero mean and potentially correlated i.e.,

$$\begin{bmatrix} u_1 \\ u_2 \\ \eta_1 \\ \eta_2 \end{bmatrix} \sim \text{Multivariate Normal} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{u_1}^2 & \text{cov}(u_1, u_2) & \text{cov}(u_1, \eta_1) & \text{cov}(u_1, \eta_2) \\ \sigma_{u_2}^2 & \text{cov}(u_2, \eta_1) & \text{cov}(u_2, \eta_2) \\ \sigma_{\eta_1}^2 & \text{cov}(\eta_1, \eta_2) \\ \sigma_{\eta_2}^2 \end{bmatrix} \right).$$

The cluster-specific survival function is given by $S(t | \mathbf{x}, \mathbf{u}, \boldsymbol{\eta}) = 1 - F_1(t | \mathbf{x}, \mathbf{u}, \eta_1) - F_2(t | \mathbf{x}, \mathbf{u}, \eta_2)$. Since we use the same CIF specification of [Cederkvist et al. \(2019\)](#), the following details are essentially the same encountered in the paper.

Focusing first on the second component of [Equation 3.1](#), the cluster-specific failure time trajectory

$$\Phi[w_k g(t) - \mathbf{x}\boldsymbol{\gamma}_k - \eta_k], \quad t > 0, \quad k = 1, 2,$$

where $\Phi(\cdot)$ is the cumulative distribution function of a standard Gaussian distribution. Instead of $w_k g(t)$, in [Cederkvist et al. \(2019\)](#) is specified $\alpha_k(g(t))$ with $\alpha_k(\cdot)$ being a monotonically increasing function known up to a finite-dimensional parameter vector, w_k . Examples are monotonically increasing B-splines or piecewise linear functions. However, to simplify the model structure we consider just the finite-dimensional parameter vector. The bottom line is that the authors do the same approach in their applications. With regard to the function $g(t)$, it plays a crucial role since the CIF separation in [Equation 3.1](#) is only possible with it. It is used a time t transformation given by

$$g(t) = \operatorname{arctanh} \left(\frac{t - \delta/2}{\delta/2} \right), \quad t \in (0, \delta), \quad g(t) \in (-\infty, \infty),$$

where δ depends on the data and cannot exceed the maximum observed follow-up time τ i.e., $\delta \leq \tau$. With this Fisher-based transformation the value of the cluster-specific failure time trajectory is equal 1, at time δ . Consequently, $F_k(\delta | \mathbf{x}, \mathbf{u}, \eta_k) = \pi_k(\mathbf{x} | \mathbf{u})$ and we can interpret $\pi_1(\mathbf{x} | \mathbf{u})$ and $\pi_2(\mathbf{x} | \mathbf{u})$ as the cause-specific cluster-specific risk levels, at time δ .

The cluster-specific risk levels are modeled by a multinomial logistic regression model with latent effects i.e.,

$$\pi_k(\mathbf{x}, \mathbf{u}) = \frac{\exp\{\mathbf{x}\boldsymbol{\beta}_k + \mathbf{u}_k\}}{1 + \exp\{\mathbf{x}\boldsymbol{\beta}_1 + \mathbf{u}_1\} + \exp\{\mathbf{x}\boldsymbol{\beta}_2 + \mathbf{u}_2\}}, \quad k = 1, 2, \quad (3.2)$$

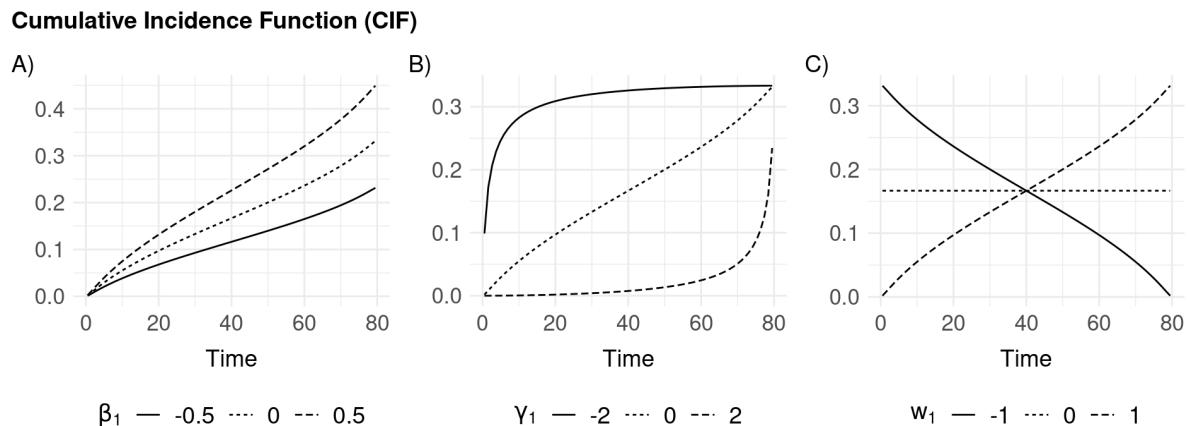
where the $\boldsymbol{\beta}_k$'s are the coefficients responsible for quantifying the impact of the covariates in the cause-specific risk levels. For individuals from the same cluster/family, at the same time point, the $\boldsymbol{\beta}_k$'s have the well-known odds ratio interpretation.

A direct understanding of all coefficients/parameters of [Equation 3.1](#) can be reached via the illustrations in [Figure 7](#). To really understand what is going on, we simplify the model. We still consider just two competing causes but without covariates and we plot just the cluster-specific CIF of one failure cause. In [Figure 7 A\)](#) we see that the β 's are also related with the curve's maximum value i.e., bigger the β , highest the CIF will be.

The γ_k 's are the coefficients responsible for quantifying the impact of the covariates in the cause-specific failure time trajectories i.e., the shape of the cumulative

incidence. In Figure 7 B) we see that the γ 's are also related with an idea of midpoint and consequently, growth speed. The fact that γ_k enters negatively in the cluster-specific failure time trajectory makes that a negative value causes an advance towards the curve, whereas a positive value causes a delay. Last but not least, the w 's in Figure 7 C). With negative values, we have a decreasing curve and with positive values an increasing curve i.e., we are interested only on the positive side.

FIGURE 7 – ILLUSTRATION OF COEFFICIENT BEHAVIORS FOR A GIVEN CUMULATIVE INCIDENCE FUNCTION (CIF) (PROPOSED BY [Cederkvist et al. \(2019\)](#)), IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, AND WITH THE FOLLOWING CONFIGURATION: $\beta_2 = 0$, $u = 0$ AND $\eta = 0$; IN EACH SCENARIO ALL OTHER COEFFICIENTS ARE SET TO ZERO, WITH THE EXCEPTION OF $w_1 = 1$



SOURCE: The author (2021).

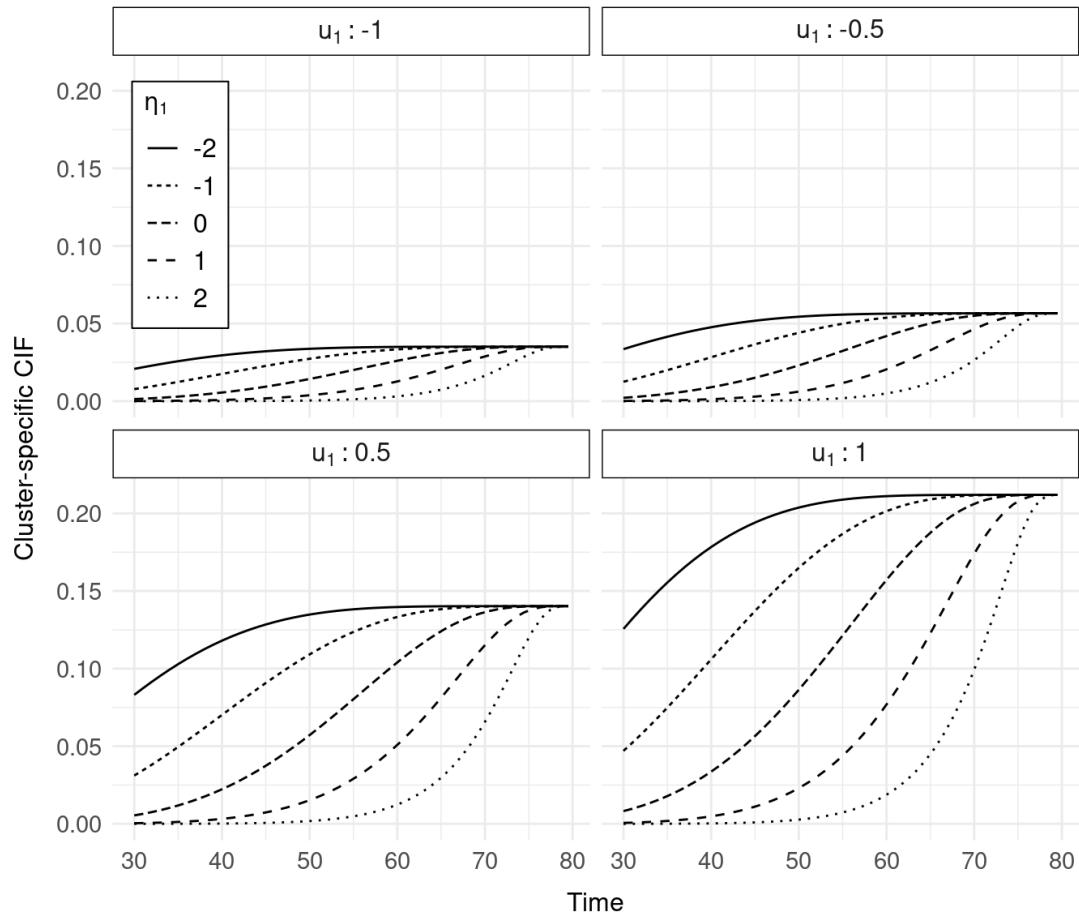
Remains to talk about the within-cluster dependence induced by the latent effects in u and η . Unfortunately, they do not have an easy interpretation. To help in the discussion, Figure 8 illustrates the cluster-specific CIF for a given failure cause in a model without covariates, let us call it failure cause 1 (in total we have two).

The latent effects u_1 and u_2 always appear together in the cluster-specific risk level, as consequence they have a joint effect on the cumulative incidence of both causes. As we can see in Figure 8, an increase in u_k will increase the risk of failure from cause k . The interpretation of $\text{cov}(\eta_1, \eta_2)$ and $\text{cov}(u_1, u_2)$ is straightforward, and those values are in most of the cases positive, as said in Cederkvist et al. (2019). With regard to $\text{cov}(u_k, \eta_k)$, negative values are the common situation. A negative correlation between η_k and u_k imply that when η_k decreases, u_k increases and conversely when η_K increases, u_k decreases. In other words, an increased risk level is reached quickly and a decreased risk level is reached later, respectively.

Practical situations with a positive within-cause correlation are hard to find i.e., where an increased risk level is associated with a late onset and vice versa. However, a positive cross-cause correlation between η and u sounds much more realistic i.e., where

late onset of one failure cause is associated with a high absolute risk of another failure cause.

FIGURE 8 – ILLUSTRATION OF A GIVEN CLUSTER-SPECIFIC CUMULATIVE INCIDENCE FUNCTION (CIF), PROPOSED BY [Cederkvist et al. \(2019\)](#), IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES AND THE FOLLOWING CONFIGURATION: $\beta_1 = -2$, $\beta_2 = -1$, $\gamma_1 = 1$, $w_1 = 3$ AND $u_2 = 0$. THE VARIATION BETWEEN FRAMES IS GIVEN BY THE LATENT EFFECTS u_1 AND η_1



SOURCE: The author (2021).

The latent effects $\{u_k, \eta_k\}$ are assumed independent across clusters and shared by individuals within the same cluster/family.

3.2 MODEL SPECIFICATION

The multiGLMM for clustered competing risks data is specified in the following hierarchical fashion. By simplicity, we focus on two competing causes of failure but an extension is straightforward.

For two competing causes of failure, a subject i , in the cluster/family j , in time

t , we have

$$y_{ijt} \mid \{u_{1j}, u_{2j}, \eta_{1j}, \eta_{2j}\} \sim \text{Multinomial}(p_{1ijt}, p_{2ijt}, p_{3ijt})$$

$$\begin{bmatrix} u_1 \\ u_2 \\ \eta_1 \\ \eta_2 \end{bmatrix} \sim \text{MN} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_{u_1}^2 & \text{cov}(u_1, u_2) & \text{cov}(u_1, \eta_1) & \text{cov}(u_1, \eta_2) \\ \sigma_{u_2}^2 & \text{cov}(u_2, \eta_1) & \text{cov}(u_2, \eta_2) \\ \sigma_{\eta_1}^2 & \text{cov}(\eta_1, \eta_2) \\ \sigma_{\eta_2}^2 \end{bmatrix} \right)$$

$$\begin{aligned} p_{kijt} &= \frac{\partial}{\partial t} F_k(t \mid \mathbf{x}, u_1, u_2, \eta_k) \\ &= \frac{\exp\{\mathbf{x}_{kij}\beta_k + u_{kj}\}}{1 + \sum_{m=1}^{K-1} \exp\{\mathbf{x}_{mij}\beta_m + u_{mj}\}} \\ &\times w_k \frac{\delta}{2\delta t - 2t^2} \phi\left(w_k \text{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_k - \eta_{kj}\right), \end{aligned} \quad (3.3)$$

$$k = 1, 2.$$

The probabilities are given by the derivative w.r.t. time t of the cluster-specific CIF. The choice of a multinomial logistic regression model ensures that the sum of the predicted cause-specific CIFs does not exceed 1.

Considering two competing causes of failure, we have a multinomial with three classes. The third class exists to handle the censorship and its probability is given by the complementary to reach 1. This framework in [Equation 3.3](#) results in what we call multiGLMM, a multinomial GLMM to handle the CIF of clustered competing risks data. For a random sample, the corresponding marginal likelihood functions in given by

$$\begin{aligned} L(\boldsymbol{\theta}; y) &= \prod_{j=1}^J \int_{\Re^4} \pi(y_j \mid \mathbf{r}_j) \times \pi(\mathbf{r}_j) d\mathbf{r}_j \\ &= \prod_{j=1}^J \int_{\Re^4} \underbrace{\left\{ \prod_{i=1}^{n_j} \prod_{t=1}^{n_{ij}} \left(\frac{(\sum_{k=1}^K y_{kijt})!}{y_{1ijt}! y_{2ijt}! y_{3ijt}!} \prod_{k=1}^K p_{kijt}^{y_{kijt}} \right) \right\}}_{\text{fixed effect component}} \times \\ &\quad \underbrace{(2\pi)^{-2} |\Sigma|^{-1/2} \exp\left\{ -\frac{1}{2} \mathbf{r}_j^\top \Sigma^{-1} \mathbf{r}_j \right\}}_{\text{latent effect component}} d\mathbf{r}_j \\ &= \prod_{j=1}^J \int_{\Re^4} \underbrace{\left\{ \prod_{i=1}^{n_j} \prod_{t=1}^{n_{ij}} \prod_{k=1}^K p_{kijt}^{y_{kijt}} \right\}}_{\text{fixed effect}} \underbrace{(2\pi)^{-2} |\Sigma|^{-1/2} \exp\left\{ -\frac{1}{2} \mathbf{r}_j^\top \Sigma^{-1} \mathbf{r}_j \right\}}_{\text{latent effect component}} d\mathbf{r}_j, \quad (3.4) \end{aligned}$$

where $\boldsymbol{\theta} = [\boldsymbol{\beta} \ \boldsymbol{\gamma} \ \mathbf{w} \ \boldsymbol{\sigma}^2 \ \boldsymbol{\rho}]^\top$ is the parameters vector to be maximized. In our framework, a subject can fail from just one competing cause or get censor, at a given time. Thus, the

fraction of factorials in the fixed effect component is made only by 0's and 1's. Finally, returning the value 1. The matrix Σ is the variance-covariance matrix, which parameters are given by σ^2 and ρ .

Now, [Equation 3.4](#) in words. To each cluster/family j we have a product of two components. The fixed effect component, given by a multinomial distribution with its probabilities specified through the cluster-specific CIF ([Equation 3.1](#)) and, the latent effect component, given by a multivariate Gaussian distribution.

To each subject i that composes a cluster j we have its specific fixed effects contribution. The likelihood in [Equation 3.4](#) is the most general as possible, allowing for repeated measures to each subject. Since all subjects of a given cluster shares the same latent effect, we have just one latent effect contribution multiplying the product of fixed effect contributions. As we do not observe the latent effect variables, r_j , we integrate out in it. With two competing causes of failure, we have four latent effects (a multivariate Gaussian distribution in four dimensions). Consequently, for each cluster, we approximate an integral in four dimensions. The product of these approximated integrals results in the called marginal likelihood, to be maximized in θ .

Independent of the parameters value choice, the probabilities for the failure causes will be small, providing always a bigger probability mass to the censorship. More will be talked about this model characteristic in the next chapter.

3.2.1 Parametrization

We have to choose in which terms we parameterize the variance-covariance matrix Σ . Besides the latent effects variances $\{\sigma^2\}$, we have to choose if we will estimate its covariances or correlations. By the name *variance-covariance* matrix, it is natural to think on covariance terms. However, this option is not very attractive since its interpretation is not clear. A more attractive choice is in terms of correlation.

The covariance between two terms is defined as a triple product: the two terms standard deviations times the correlation ρ . Still thinking in two competing causes of failure, we have then an Σ matrix with six correlations

$$\Sigma = \begin{bmatrix} \sigma_{u_1}^2 & \rho_{u_1, u_2} \sigma_{u_1} \sigma_{u_2} & \rho_{u_1, \eta_1} \sigma_{u_1} \sigma_{\eta_1} & \rho_{u_1, \eta_2} \sigma_{u_1} \sigma_{\eta_2} \\ \rho_{u_1, u_2} \sigma_{u_1} \sigma_{u_2} & \sigma_{u_2}^2 & \rho_{u_2, \eta_1} \sigma_{u_2} \sigma_{\eta_1} & \rho_{u_2, \eta_2} \sigma_{u_2} \sigma_{\eta_2} \\ \rho_{u_1, \eta_1} \sigma_{u_1} \sigma_{\eta_1} & \rho_{u_2, \eta_1} \sigma_{u_2} \sigma_{\eta_1} & \rho_{\eta_1, \eta_2} \sigma_{\eta_1} \sigma_{\eta_2} & \\ \rho_{u_1, \eta_2} \sigma_{u_1} \sigma_{\eta_2} & \rho_{u_2, \eta_2} \sigma_{u_2} \sigma_{\eta_2} & \sigma_{\eta_1}^2 & \\ & & & \sigma_{\eta_2}^2 \end{bmatrix}.$$

With the matrix parametrization being chosen, we have that the parameters to be estimated are the components of the vector $\theta = [\beta \gamma w \sigma^2 \rho]^\top$. There we have the fixed effects or mean components $\{\beta \gamma w\}$, the easiest to estimate in a statistical modeling framework; we have variance components $\{\sigma^2\}$, the intermediate ones; and the correla-

tion components $\{\rho\}$, the hardest ones. This idea of easy or hard to estimate may be justified by three, connected, arguments.

The first comes from the fact that we are modeling the mean of a probability distribution in a hierarchical and structured fashion, consequently, the easiest parameters to estimate will be the mean components. We may make the analogy that to estimate the mean parameters we need data (resources); to estimate the variance parameters we need more data (more resources), and to estimate the correlation parameters we need much more data (even more resources). The second argument comes to also explain the first one via the parametric space constraints.

Generally, the fixed effect components do not present constraints i.e., they can vary in all \mathbb{R} . The same can not be said from the variance components, constrained by definition into the \mathbb{R}_*^+ . Finally, we have the correlation components, constrained to the interval $[-1, 1]$. These parametric space constraints drive us again to the first argument since we need more data/resources/information to be able to estimate coefficients constrained to some interval. Nevertheless, this may not be enough. Without providing some extra information in terms of an e.g., constrained algorithm, it is very reasonable to expect that during the optimization procedure some unrealistic areas of the parametric space could be visited and jeopardize the stability or even the whole optimization procedure. To overcome these possible difficulties, parameter reparametrizations are more than welcome.

The variance and correlation parameters are modeled in terms of the matrix Σ . This matrix is symmetric and more important, positive semi-definite. This last characteristic is also the third argument to justify why is so difficult to estimate these parameters. Since the estimates should lead to a positive semi-definite matrix, the employment of a parametrization is welcome to enforces this condition.

In the subject of choosing the components parametrization for a positive-definite matrix Σ , we have basically two big options available in the statistical modeling literature. One of them consists of just transform the scale. By practical reasons, let us think in a 2×2 matrix

$$\Sigma = \begin{bmatrix} \exp\{\log \sigma_1^2\} & z^{-1}(z(\rho_{1,2})) \sqrt{\exp\{\log \sigma_1^2\}} \sqrt{\exp\{\log \sigma_1^2\}} \\ & \exp\{\log \sigma_2^2\} \end{bmatrix}$$

i.e., in the main diagonal we may now estimate the log-variances and in the off-diagonal we may estimate Fisher z-transformed correlations.

The estimation of the log-variances has two big advantages

- Since the natural logarithm is a real-valued function, we overcome the parametric space constraint problem;

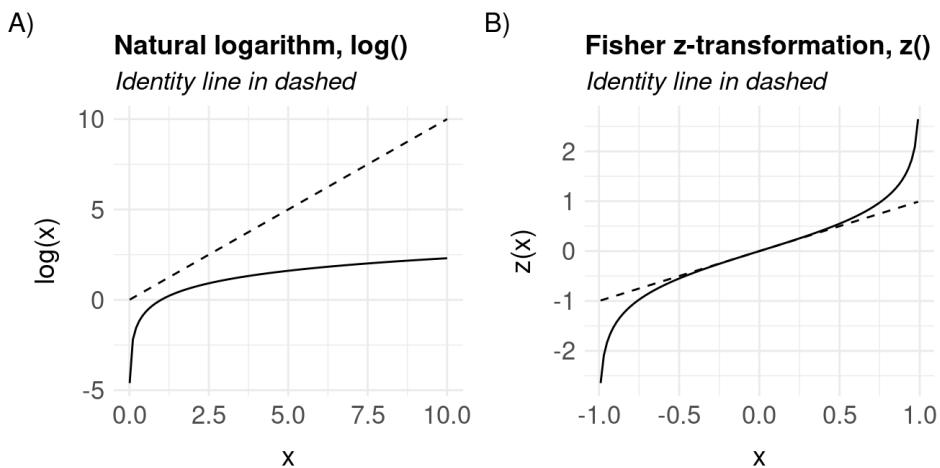
- High variances are problematic for many reasons but in the context of seeing them as the diagonal components of a restricted matrix, being able to control its magnitudes is a crucial task to the stability of any optimization routine. With the natural logarithm transformation we shrink the parametric space as illustrated in [Figure 9 A\)](#), avoiding some eventual numerical cumbersome.

With the correlation components we proceed with the estimation of its Fisher z-transformation. This transformation, and its inverse, are defined as

$$z(\rho) = \frac{1}{2} \log \left(\frac{1+\rho}{1-\rho} \right) = \operatorname{arctanh}(\rho), \quad z^{-1}(\rho) = \frac{\exp\{2\rho\} - 1}{\exp\{2\rho\} + 1} = \tanh(\rho).$$

The Fisher z-transformation plays the role of stretching the small correlation parametric space but doing this in a smooth fashion, as illustrated in [Figure 9 B\)](#).

FIGURE 9 – ILLUSTRATION OF THE PARAMETRIZATION BEHAVIOR FOR THE VARIANCE COMPONENTS, IN A), AND CORRELATION COMPONENTS, IN B)



SOURCE: The author (2021).

The other parametrization option consist in estimate the elements of a factorization or decomposition of the positive-definite matrix Σ . The most common is the Cholesky factorization or decomposition ([PINHEIRO; BATES, 1996](#)). For two competing causes of failure, a standard Cholesky decomposition of Σ may be expressed as

$$\Sigma = \begin{bmatrix} c_1 & 0 & 0 & 0 \\ c_2 & c_3 & 0 & 0 \\ c_4 & c_5 & c_6 & 0 \\ c_7 & c_8 & c_9 & c_{10} \end{bmatrix} \begin{bmatrix} c_1 & c_2 & c_4 & c_7 \\ 0 & c_3 & c_5 & c_8 \\ 0 & 0 & c_6 & c_9 \\ 0 & 0 & 0 & c_{10} \end{bmatrix} = LL^\top,$$

where $\{c_i\}_{i=1}^{10}$ are the unconstrained coefficients to be estimated.

A disadvantage in the use of a decomposition as the Cholesky is the lack of a straightforward interpretation to the elements $\{c_i\}_{i=1}^{10}$. However, with the application of the delta method, already implemented in TMB's ([KRISTENSEN et al., 2016](#))

TMB::sdreport(), it is straightforward to get back the Σ elements together with its respective standard errors. The main advantage of this parametrization apart from the fact that it ensures positive definiteness, is that it is computationally simple and stable.

Just to mention another viable possibilities, we could use a modified Cholesky decomposition (POURAHMADI, 2007) providing a better statistical interpretation of the decomposition elements or, we could parametrize the precision matrix, $Q = \Sigma^{-1}$. Since we use Σ^{-1} in the marginal likelihood of [Equation 3.4](#), parametrizing directly its inverse save us some computations.

Besides the popularity of the Cholesky method, there is another factorization scheme available and efficiently implemented in TMB. It is a factorization based on a vector scale transformation of an unstructured correlation matrix. For two competing causes of failure the decomposition is specified in the following fashion

$$\Sigma = VD^{-1/2}LL^\top D^{-1/2}V^\top,$$

where

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ c_1 & 1 & 0 & 0 \\ c_2 & c_3 & 1 & 0 \\ c_4 & c_5 & c_6 & 1 \end{bmatrix}, \quad D = \text{diag}(LL^\top) \quad \text{and} \quad W = \text{diag}(\{\sigma_i\}_{i=1}^4).$$

This scheme is based initially on the factorization of a correlation matrix defined as $D^{-1/2}LL^\top D^{-1/2}$. The elements $\{c_i\}_{i=1}^6$ to be estimated has the advantage of being unconstrained and guarantees the matrix symmetry and positive definiteness. The variances are scaled via the diagonal matrix V , its elements $\{\sigma_i\}_{i=1}^4$ are then the standard deviations to be estimated.

4 SIMULATION STUDY DATASETS

This chapter describes how to simulate from our multiGLMM and describes the performed simulation studies. The general simulation procedure is addressed in Section 4.1. In Section 4.2 the simulation studies are presented in detail.

4.1 SIMULATING FROM THE MODEL

Being able to simulate data from a model is a key task, fundamental to assess the finite-sample properties and the estimation procedure liability of a given statistical model. The step-by-step describing the simulation procedure of our multiGLMM is presented on Algorithm 1, following the model hierarchical structure stipulated in Equation 3.3.

ALGORITHM 1 SIMULATING FROM A multiGLMM FOR CLUSTERED COMPETING RISKS DATA

- 1: Set J , the number of clusters
- 2: Set n_j , the number of cluster elements ▷ can be of different sizes
- 3: Set $K - 1$, the number of competing causes of failure
- 4: Set the model parameter values $\theta = [\beta \gamma w \sigma^2 \varrho]^\top$
- 5: Sample J latent effect vectors from a Multivariate Normal $_{(K-1) \times (K-1)}(\mathbf{0}, \Sigma(\sigma^2, \varrho))$
- 6: Set δ ▷ maximum follow-up time
- 7: Set the failure times t_{ij}
- 8: Compute the competing risks probabilities

$$p_{kij} = \frac{\exp\{x_{kij}\beta_{kj} + u_{kj}\}}{1 + \sum_{m=1}^{K-1} \exp\{x_{mij}\beta_{mj} + u_{mj}\}}$$

$$\times w_k \frac{\delta}{2\delta t_{ij} - 2t_{ij}^2} \phi\left(w_k \text{arctanh}\left(\frac{t_{ij} - \delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}\right),$$

$$\text{Censorship : } p_{Kij} = 1 - \sum_{k=1}^{K-1} p_{kij}, \quad k = 1, 2, \dots, K - 1$$

- 9: Sample $J \times n_j$ vectors from a Multinomial($p_{1ij}, p_{2ij}, \dots, p_{Kij}$)
 - 10: If $t_{ij} = \delta$, moves to class K ▷ any failure at time δ is a censorship
 - 11: **return** multinomial vectors and their respective failure/censoring times
-

SOURCE: The author (2021).

The model described in Equation 3.3 is in a general form, allowing for varying coefficients between clusters. However, we focus on a simpler structure with just fixed

intercepts. Fixing the latent effects in its distribution mean, zero, and using the following fixed effects configuration for two competing causes of failure

$$\begin{aligned}\beta &= [-2 \ 1.5]^\top \\ \gamma &= [1.2 \ 1]^\top \\ w &= [3 \ 5]^\top,\end{aligned}\tag{4.1}$$

we get the CIF's and failure probabilities (CIF derivatives w.r.t. time t , dCIF) presented respectively in [Figure 10](#).

FIGURE 10 – CUMULATIVE INCIDENCE FUNCTIONS (CIF) AND RESPECTIVE DERIVATIVES (dCIF) W.R.T. TIME FOR A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, LATENT EFFECTS IN ZERO, AND FIXED EFFECTS IN [Equation 4.1](#)



SOURCE: The author (2021).

By adding a complete latent structure,

$$\begin{bmatrix} u_1 \\ u_2 \\ \eta_1 \\ \eta_2 \end{bmatrix} \sim \text{Multivariate Normal} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.4 & -0.5 & 0.4 \\ 0.4 & 1 & 0.4 & -0.3 \\ -0.5 & 0.4 & 1 & 0.4 \\ 0.4 & -0.3 & 0.4 & 1 \end{bmatrix} \right),\tag{4.2}$$

we are able to apply [Algorithm 1](#) and generate a complete model sample with 50000 clusters of size two (pairs), summarized in [Figure 11](#). The R function written to simulate the data is available in [Appendix C](#).

Varying the parameters configuration we are able to build basically any CIF's format. However, its dCIF will be always small i.e., the generated probabilities for the failure causes will always be (really) small, passing all the probability mass to the censorship class. Low probabilities imply few failures, making the multiGLMM even harder to fit. All these behaviors are seen in [Figure 11](#).

FIGURE 11 – HISTOGRAMS FOR SIMULATED PROBABILITIES WITH RESPECTIVE OUTPUT PCERCENTAGES AND FAILURE TIMES FOR A MODEL WITH TWO COMPETING CAUSES AND 50000 CLUSTERS OF SIZE TWO. THE SIMULATION FOLLOWED ALGORITHM 1 GUIDELINES WITH PARAMETER CONFIGURATIONS SPECIFIED IN [Equation 4.1](#) AND [Equation 4.2](#)



SOURCE: The author (2021).

Thinking of epidemiological or public health problems, the reasonable CIF behaviors are the ones presented in [Figure 10](#). In the simulation routine, the failure-/censorship times are based on the random sampling of values between 30 and 80-time units. Another approach could be performing this random sampling by age group, to have something similar to a realistic population pyramid. However, by performing some tests we saw that having a realistic CIF is enough. Even with a uniform time distribution, the failure times will not be uniform, as the CIF curves impose their form. We can see this happening in the bottom graphs of [Figure 11](#).

[Cederkvist et al. \(2019\)](#) does something different through the sampling of the censorship times from a $U(0, \delta)$, the sampling of $\zeta \sim U(0, 1)$, and the computation of the cause-specific failure times by solving

$$\zeta = \Phi \left(w_k \operatorname{arctanh} \left(\frac{t_{ij} - \delta/2}{\delta/2} \right) - x_{kij} \gamma_k - \eta_{kj} \right) \quad \text{for } t_{ij},$$

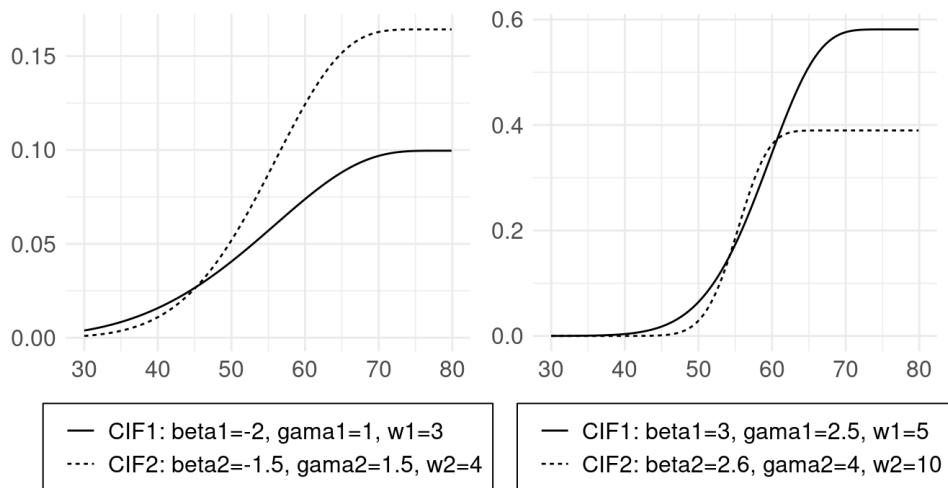
with i being the subject, j the cluster, and k the failure cause. This approach implies a parametric form for the failure times, which we do not know if holds in the real world.

4.2 SIMULATION STUDIES DESIGN

To stress the model and check the parameters identifiability, we propose seventy-two scenarios. All scenarios with two competing causes. Thus, a three classes multinomial distribution.

We consider two CIF scenarios, in summary, a low and a high CIF scenario. The parameter configurations are presented together with their curves in [Figure 12](#). As mentioned before, independent of the configuration the censorship level is basically the same.

FIGURE 12 – CUMULATIVE INCIDENCE FUNCTIONS (CIF) FOR A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, AND LATENT EFFECTS IN ZERO



SOURCE: The author (2021).

We use four multiGLMMs, that can be discriminated by their latent effect structures

$$\sigma^2 = [\sigma_{u_1}^2 = 1 \quad \sigma_{u_2}^2 = 0.6 \quad \sigma_{\eta_1}^2 = 0.7 \quad \sigma_{\eta_2}^2 = 0.9];$$

$$\rho = [\rho_{u_1, u_2} = 0.1 \quad \rho_{u_1, \eta_1} = -0.5 \quad \rho_{u_1, \eta_2} = 0.3 \quad \rho_{u_2, \eta_1} = 0.3 \quad \rho_{u_2, \eta_2} = -0.4 \quad \rho_{\eta_1, \eta_2} = 0.2].$$

In

$$\Sigma = \begin{bmatrix} R & C \\ C^\top & T \end{bmatrix} \quad \text{with}$$

$$R = \begin{bmatrix} \sigma_{u_1}^2 & \text{cov}(\rho_{u_1, u_2}) \\ & \sigma_{u_2}^2 \end{bmatrix}, \quad T = \begin{bmatrix} \sigma_{\eta_1}^2 & \text{cov}(\rho_{\eta_1, \eta_2}) \\ & \sigma_{\eta_2}^2 \end{bmatrix}, \quad C = \begin{bmatrix} \text{cov}(\rho_{u_1, \eta_1}) & \text{cov}(\rho_{u_1, \eta_2}) \\ \text{cov}(\rho_{u_2, \eta_1}) & \text{cov}(\rho_{u_2, \eta_2}) \end{bmatrix}.$$

They are

risk model A model with latent effects only on the risk level i.e., $\Sigma_{2 \times 2} = R$;

time model A model with latent effects only on the trajectory time level i.e., $\Sigma_{2 \times 2} = T$;

block-diag model A model with latent effects on the risk and trajectory time level only i.e., $\Sigma_{4 \times 4} = \text{diag}(R, T)$;

complete model A model with a *complete* latent effects structure i.e., $\Sigma_{4 \times 4} = \Sigma$.

All models are based on some decomposition of the symmetric matrix

$$\Sigma = \begin{bmatrix} R & C \\ C^\top & T \end{bmatrix} = \begin{bmatrix} 1.0 & 0.0775 & -0.4183 & 0.2846 \\ & 0.6 & 0.1944 & -0.2939 \\ & & 0.7 & 0.1587 \\ & & & 0.9 \end{bmatrix}.$$

Propose a 4×4 matrix like this is not so simple since we should not only look at the values coherence. The matrix should also be positive-definite, and more than that, the submatrices R , T , and $\text{diag}(R, T)$, should also be positive-definite since they are also used as variance-covariances matrices.

Three sample sizes are used, 5000; 30000; and 60000 data points. They are combined with three cluster sizes, size 2; size 5; and size 10 clusters. Thus,

- | <i>5000 data points,</i> | <i>30000 data points,</i> | <i>60000 data points,</i> |
|---|---|--|
| <ul style="list-style-type: none"> • 2500 clusters of size 2 • 1000 clusters of size 5 • 500 clusters of size 10 | <ul style="list-style-type: none"> • 15000 clusters of size 2 • 6000 clusters of size 5 • 3000 clusters of size 10 | <ul style="list-style-type: none"> • 30000 clusters of size 2 • 12000 clusters of size 5 • 6000 clusters of size 10 |

Besides the sample size per se, by increasing it we also increase the number of Laplace approximations to be performed since each cluster implies a Laplace approximation. This affects hugely the computational time.

In summary, we have

Two CIF configurations, four latent effect structures, three sample sizes, and three cluster sizes; $2 \times 4 \times 3 \times 3 = 72$ scenarios. For each scenario, we simulate 300 samples and fitted the corresponding model.

5 RESULTS

This chapter presents the simulation study results. We have seventy-two simulation scenarios, as detailed in [Chapter 4](#). For each scenario we simulate 300 samples. In total, we fit 21600 models.

5.1 SIMULATION STUDY

Let us just recap the parameter values used

High CIF configuration : $\{\beta_1 = -2, \beta_2 = -1.5, \gamma_1 = 1, \gamma_2 = 1.5, w_1 = 3, w_2 = 4\}$;

Low CIF configuration : $\{\beta_1 = 3, \beta_2 = 2.5, \gamma_1 = 2.6, \gamma_2 = 4, w_1 = 5, w_2 = 10\}$.

$$\begin{array}{l} \sigma_{u_1}^2 = 1 \\ \sigma_{u_2}^2 = 0.7, \\ \sigma_{\eta_1}^2 = 0.6 \\ \sigma_{\eta_2}^2 = 0.9 \end{array} \quad \text{Correlation structure} = \begin{pmatrix} u_1 & u_2 & \eta_1 & \eta_2 \\ 1 & 0.1 & -0.5 & 0.3 \\ & 1 & 0.3 & -0.4 \\ & & 1 & 0.2 \\ & & & 1 \end{pmatrix} \begin{matrix} u_1 \\ u_2 \\ \eta_1 \\ \eta_2 \end{matrix} .$$

The parameter values per se are not important. What is important is to keep in mind the behaviors implied by them, and see if the proposed model is able to learn the true values in several different scenarios and measure the quality of this learning.

The take-home message for the fixed-effect parameters, is to show that we can construct different level CIF scenarios. The β s are responsible for the curve maximum point or plateau, the γ s and w s are responsible for basically the curve shape. Its interpretation is presented in detail in [Chapter 3](#). About the latent-effects, the chosen covariance structure is considerably high but still acceptable. The underlying idea was to try to build a realistic covariance scenario and consequently be able to check how the model performs in such conditions.

In the following pages we have several graphs summarizing the parameters bias. In each figure, we have the parameter bias and its uncertainty described by a Wald-based interval i.e., ± 1.96 the bias standard deviation. This is a good uncertainty representation choice since it is symmetric and robust to outliers. In the [Appendix D](#), we have the same parameters bias but with the corresponding 2.5 and 97.5% quantiles.

The seventy-two scenarios are accommodated. We have up to four blocks of bars, each block representing a model. In each block we have eighteen bars, each bar representing the 300 fits in each of the eighteen scenarios, $4 \times 18 \times 300 = 21600$.

Each scenario name consists of a combination of three strings

- The cluster size (cs), 2, 5, and 10;
- The CIF configuration, high and low;
- The sample size, 5, 30, and 60 thousand.

We have tried to fit a total of 21600 models but not all converged. Besides that inconvenience, from the ones that converged do not all converged in the strict sense of the word. To show these two characteristics respectively, we control the bar widths and colors. The ideal convergence is the traditional one i.e., based on the gradient. However, in complicated optimization scenarios as the one we are here, we can reach a convergence based on a metric different from the gradient one.

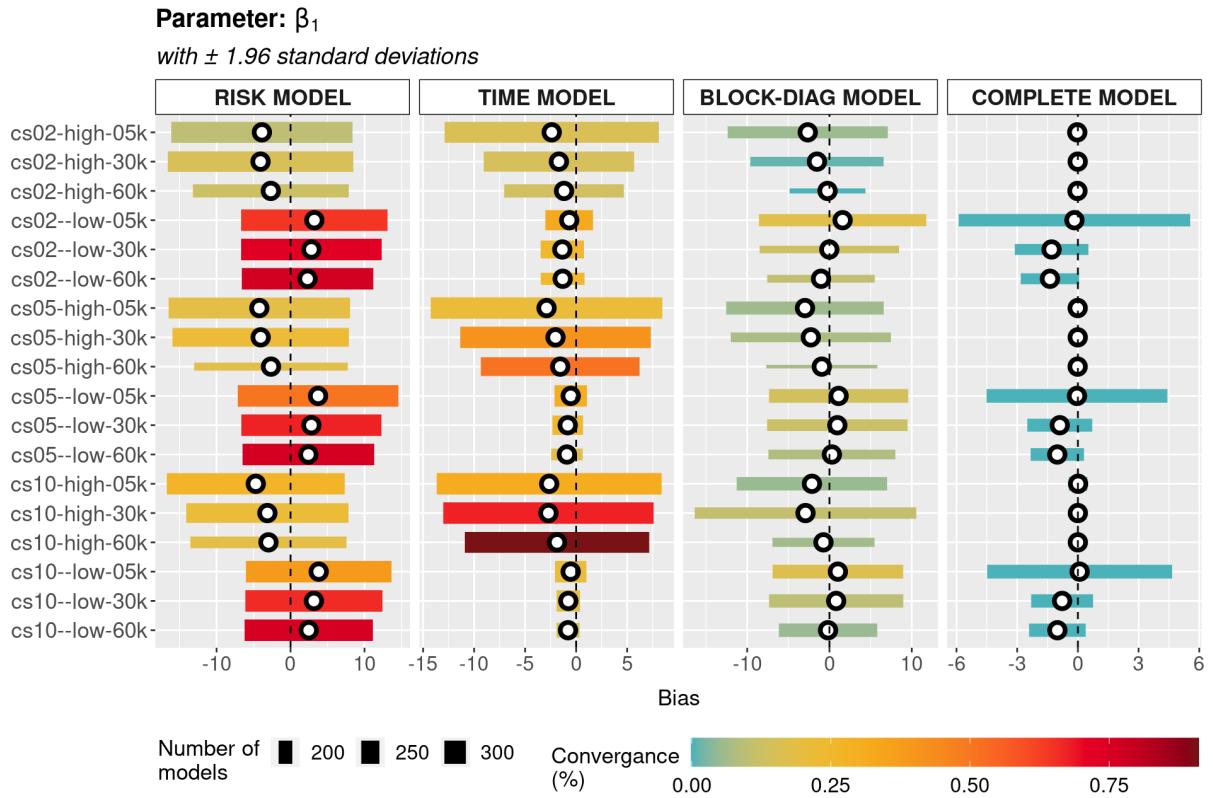
In the base::`nlmimb()` R's PORT implementation, when the default gradient-based convergence is reached we get a 0 flag message. When the optimization converges but not gradient-based, we get a 1 flag message. Looking at the technical report of [Gay \(1990\)](#), we see that this different convergence means that the last algorithm iterate x appears to be at a certain scaled distance of the locally optimal point x^* . This scaled distance, $p(x, x^*)$, is defined by

$$p(x, y) = \max_{1 \leq i \leq p} \{d_i, |x_i - y_i|\} / \max_{1 \leq j \leq p} \{d_j, |x_j| + |y_j|\},$$

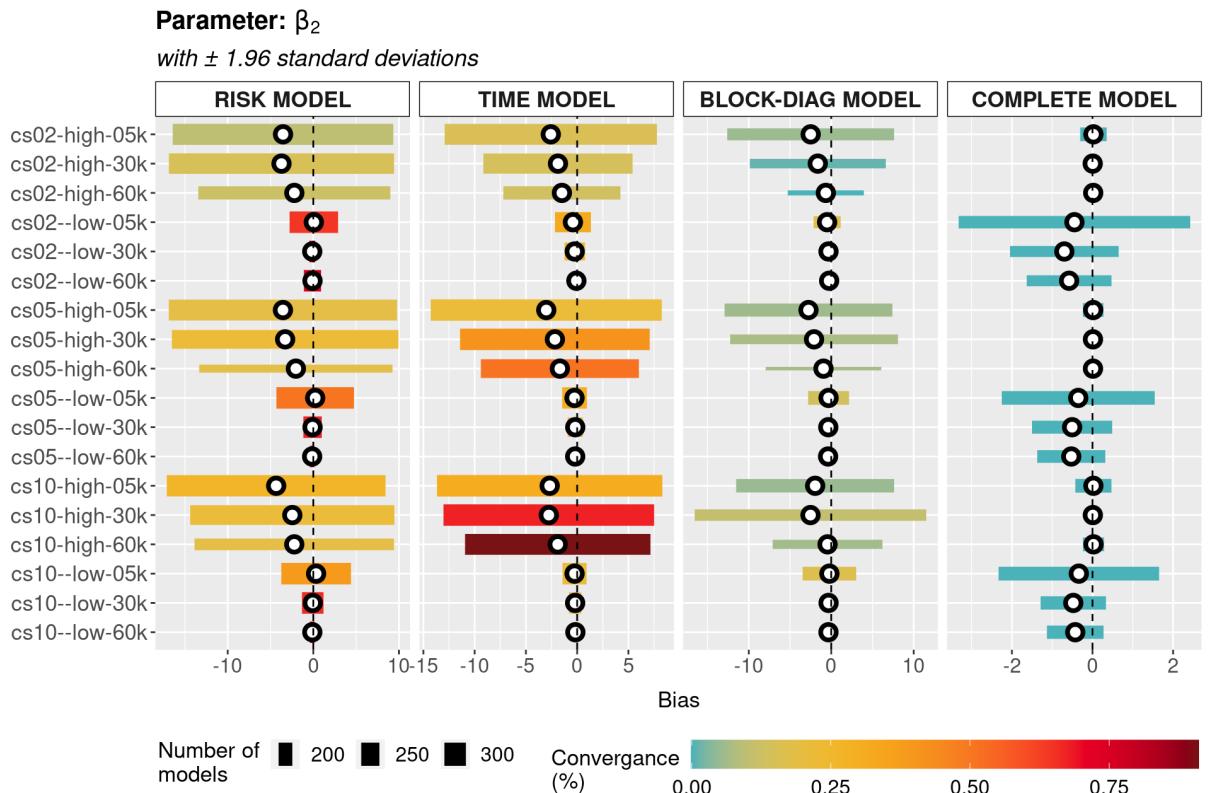
where d is a certain scale vector. In all figures, the presented convergence percentage is from the gradient-based one.

Something specific can be said about each parameter but let us keep the focus on the general remarks. Starting from the fixed-effect parameters in [Figure 13](#), [Figure 14](#), [Figure 15](#), [Figure 16](#), [Figure 17](#), and [Figure 18](#), we have some very nice results that already show a strong inclination towards the complete model's choice.

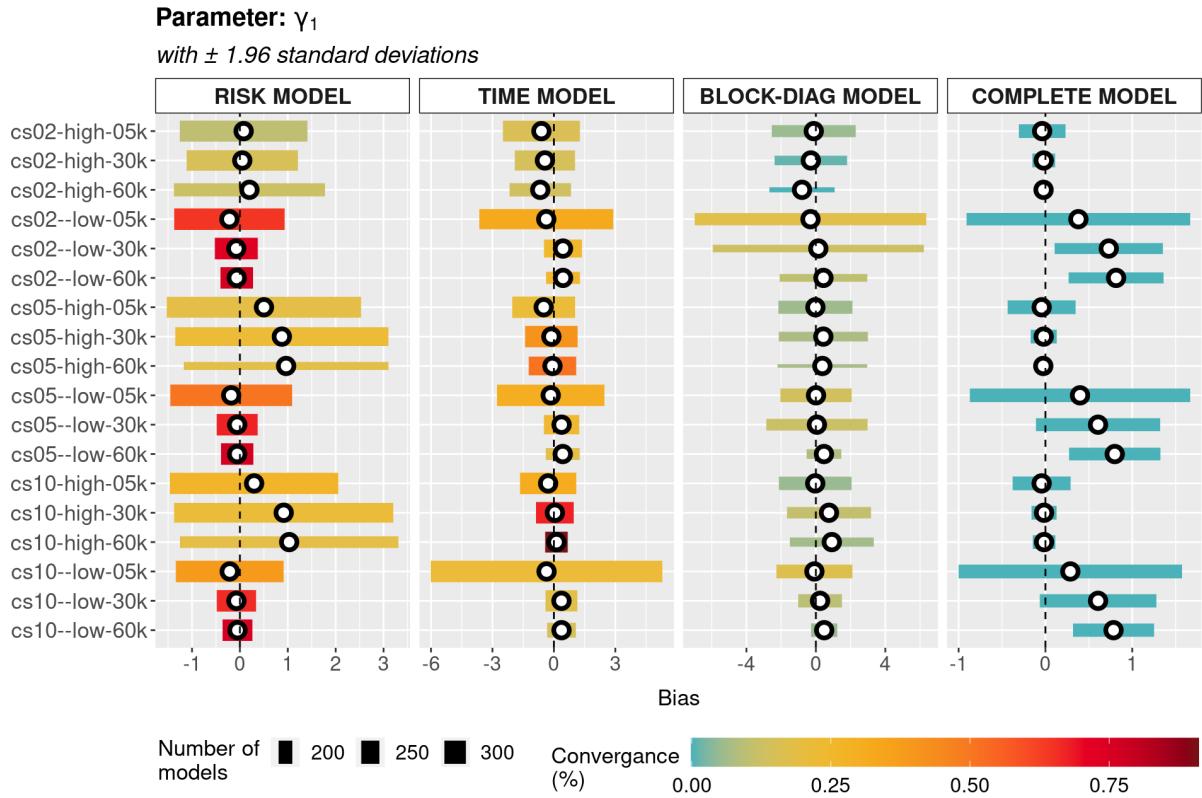
With a latent structure only in the risk level or in the trajectory time level, the low CIF scenarios are the ones with a much smaller bias-variance. In general, the mean-bias is small but the variances are high. When we have a latent structure on both levels but assuming the cross-correlations as zero (block-diag model), the results get a little bit better. Nevertheless, when we assume a non-zero cross-correlation structure (complete model), basically everything changes. The mean biases get even closer to zero, the standard deviations decrease 50% or more, and mainly, now the high CIF scenarios are the ones with a much smaller bias-variance. All this is accomplished through the consideration of the cross-correlations.

FIGURE 13 – PARAMETER β_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS

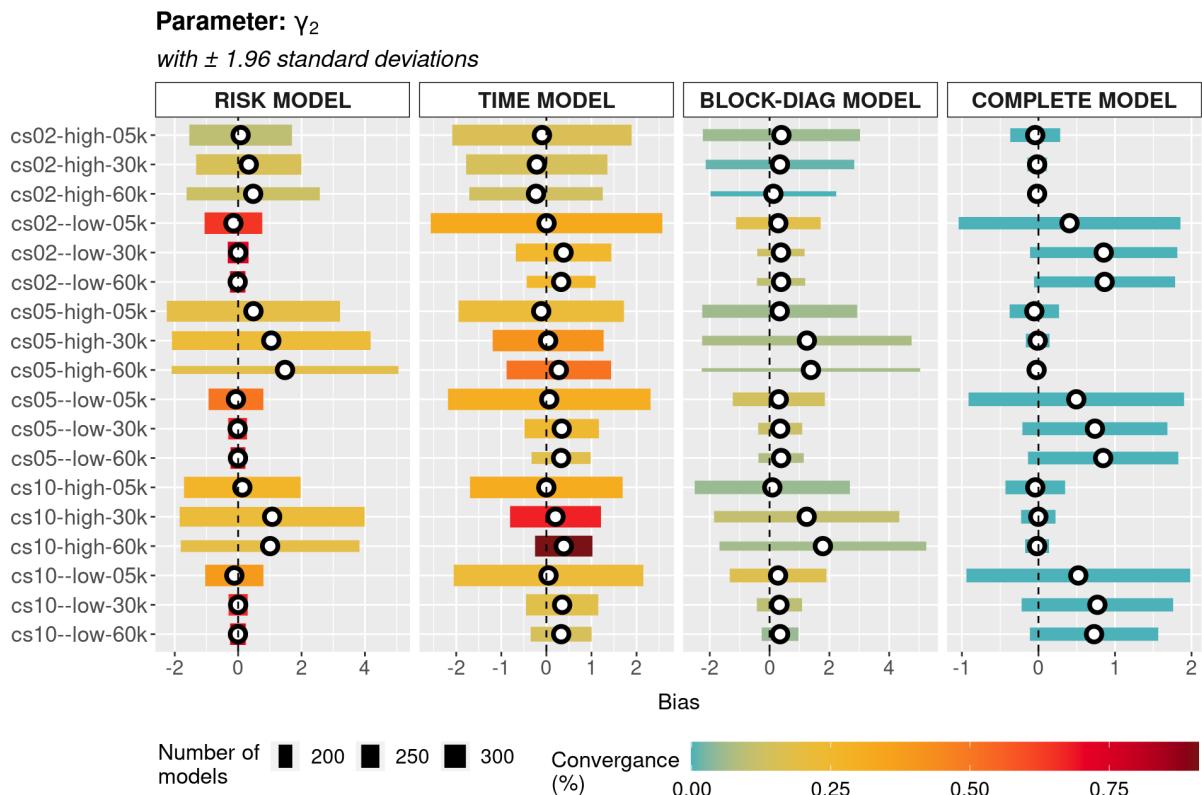
SOURCE: The author (2021).

FIGURE 14 – PARAMETER β_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS

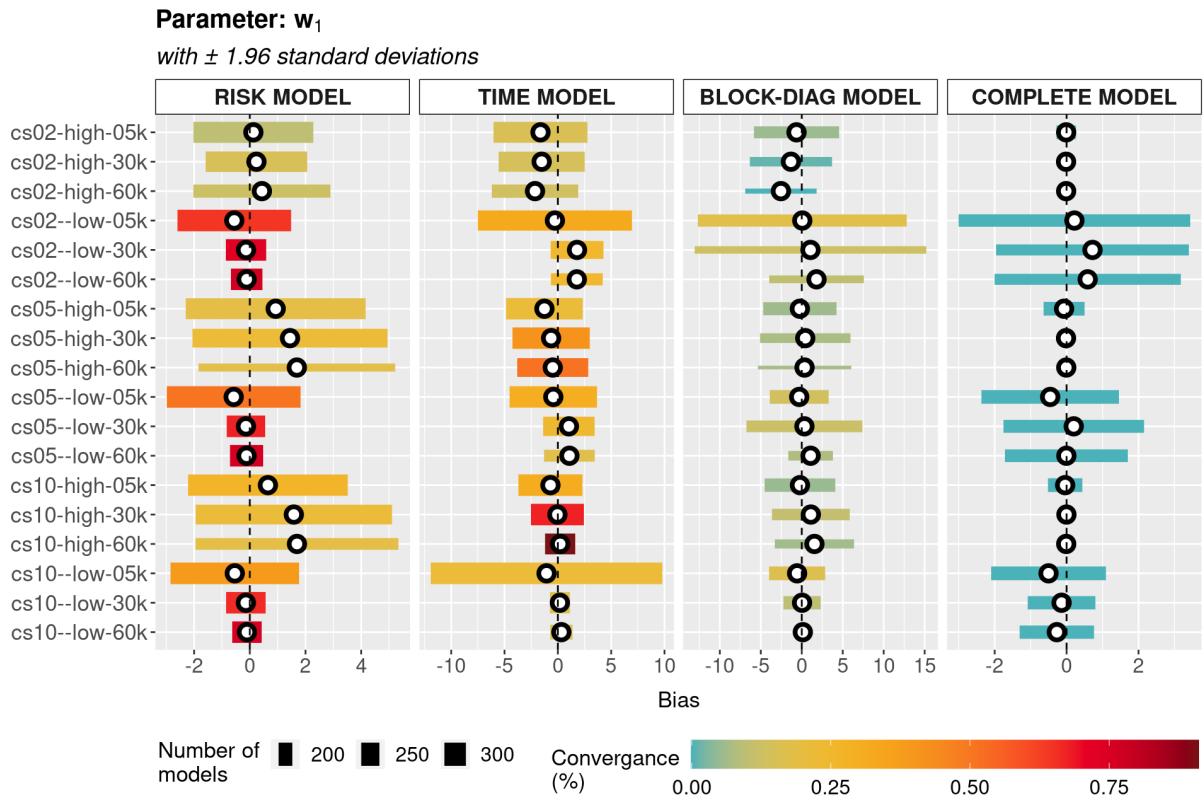
SOURCE: The author (2021).

FIGURE 15 – PARAMETER γ_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS

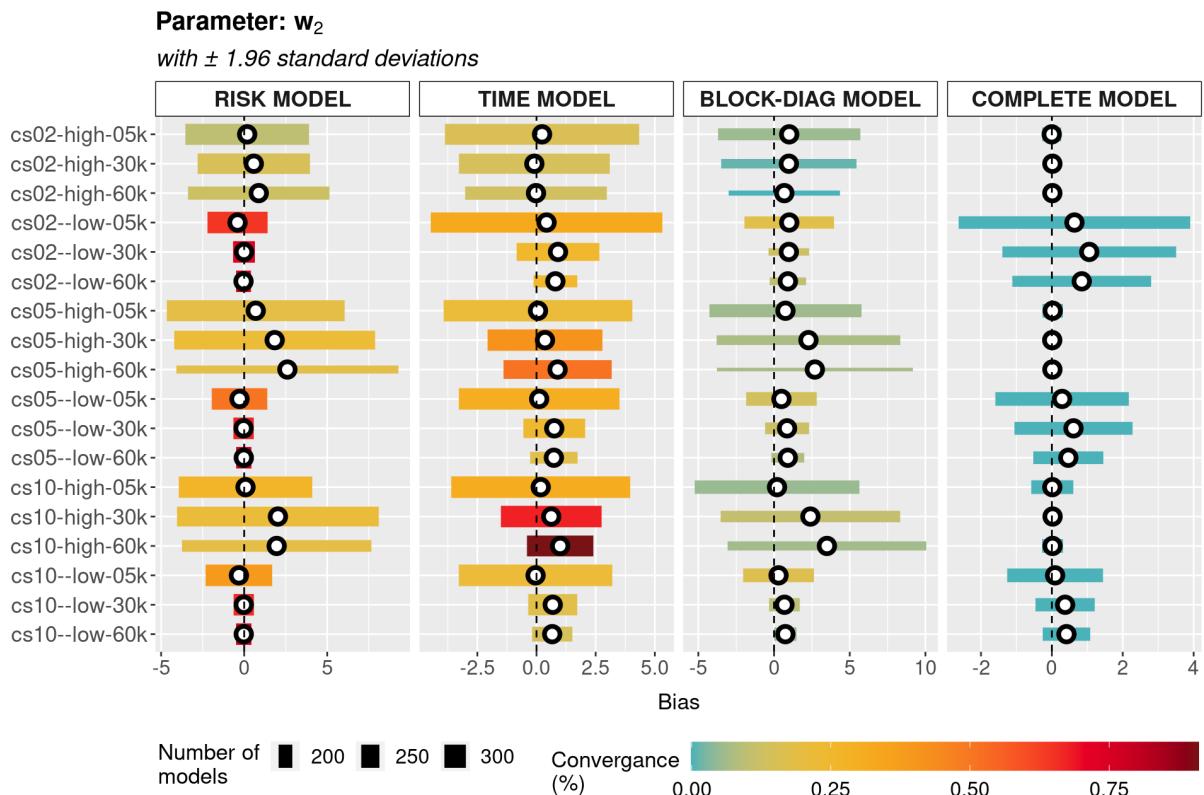
SOURCE: The author (2021).

FIGURE 16 – PARAMETER γ_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS

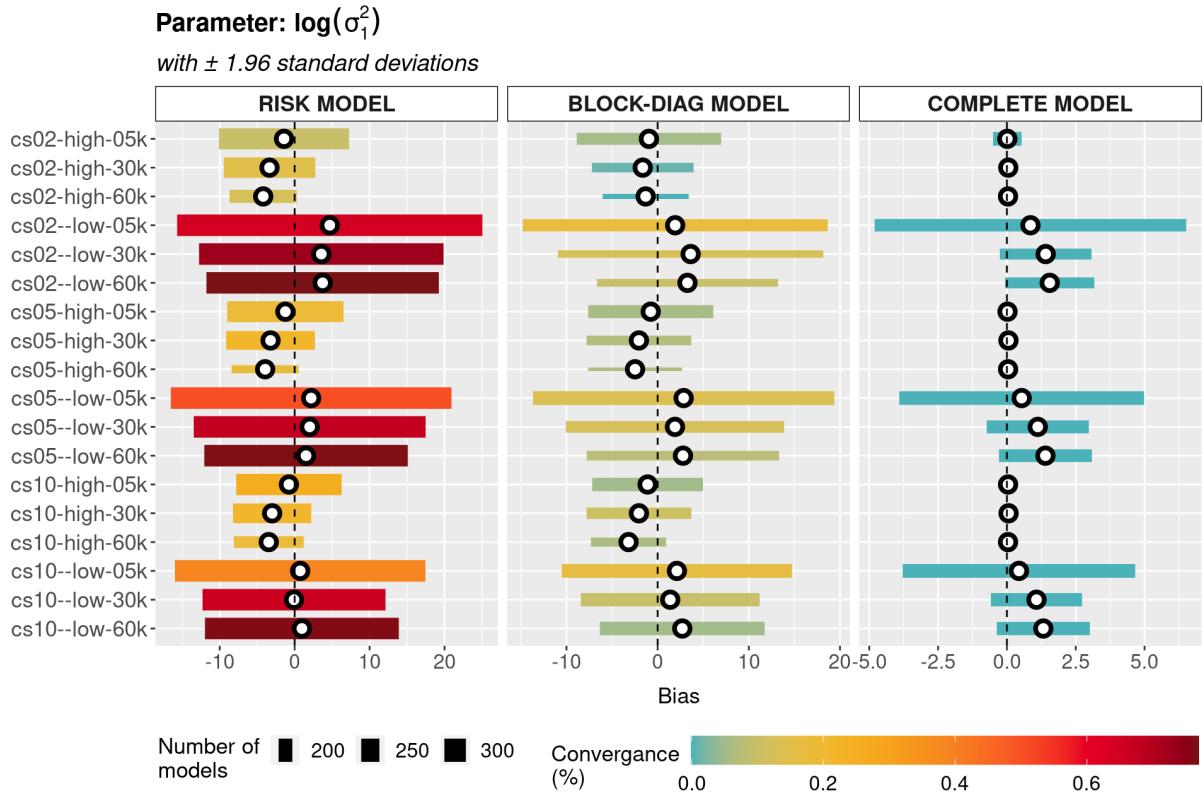
SOURCE: The author (2021).

FIGURE 17 – PARAMETER w_1 BIAS WITH ± 1.96 STANDARD DEVIATIONS

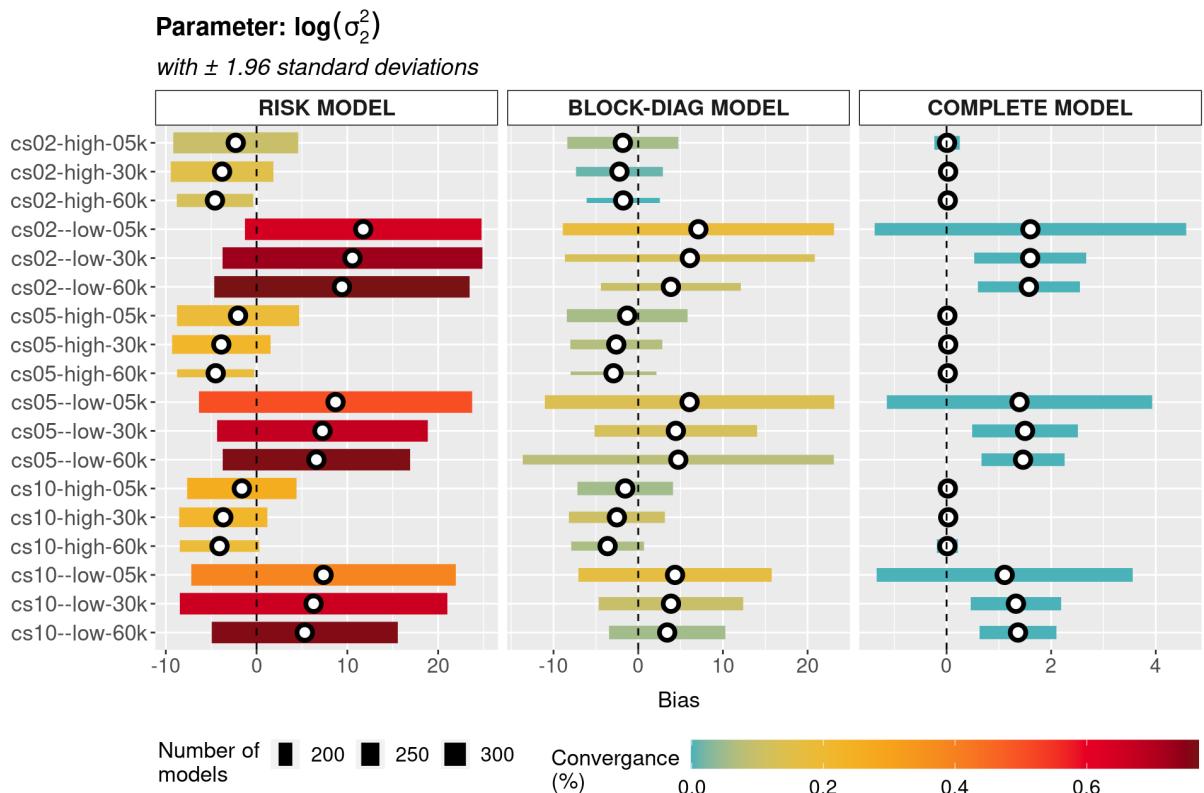
SOURCE: The author (2021).

FIGURE 18 – PARAMETER w_2 BIAS WITH ± 1.96 STANDARD DEVIATIONS

SOURCE: The author (2021).

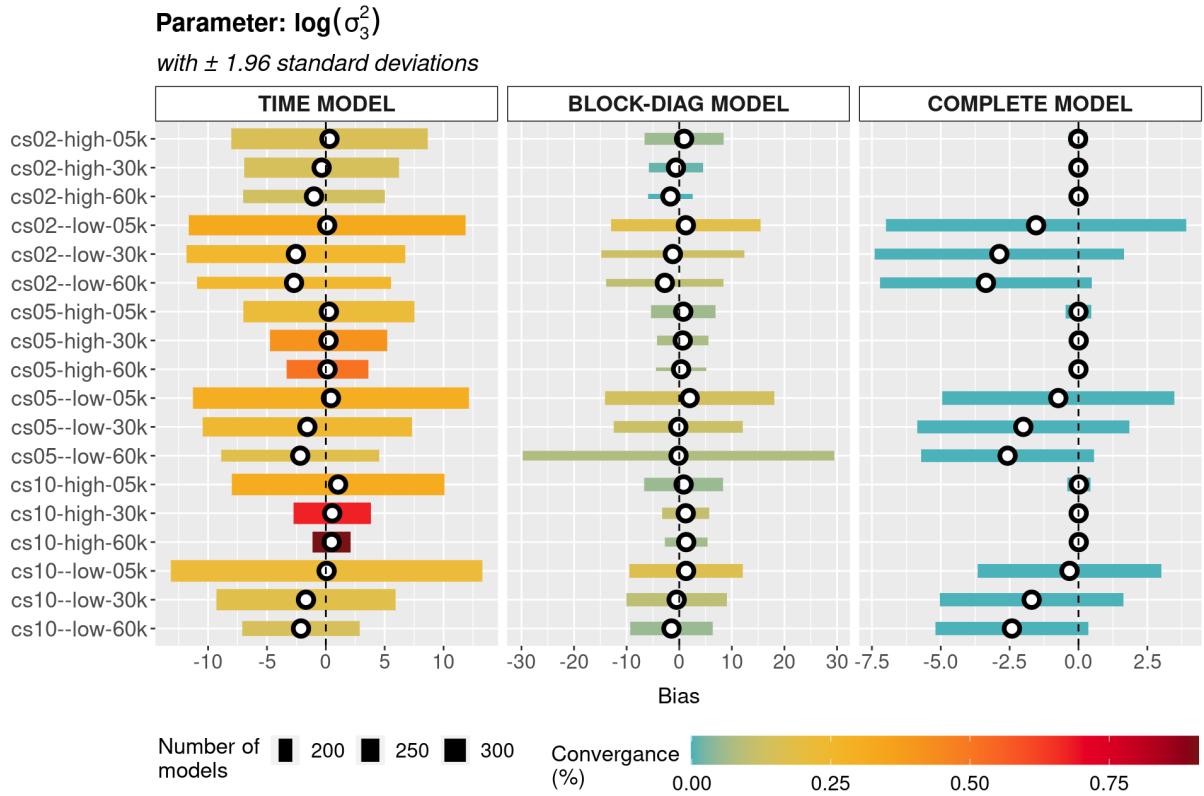
FIGURE 19 – PARAMETER $\log(\sigma_1^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS

SOURCE: The author (2021).

FIGURE 20 – PARAMETER $\log(\sigma_2^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS

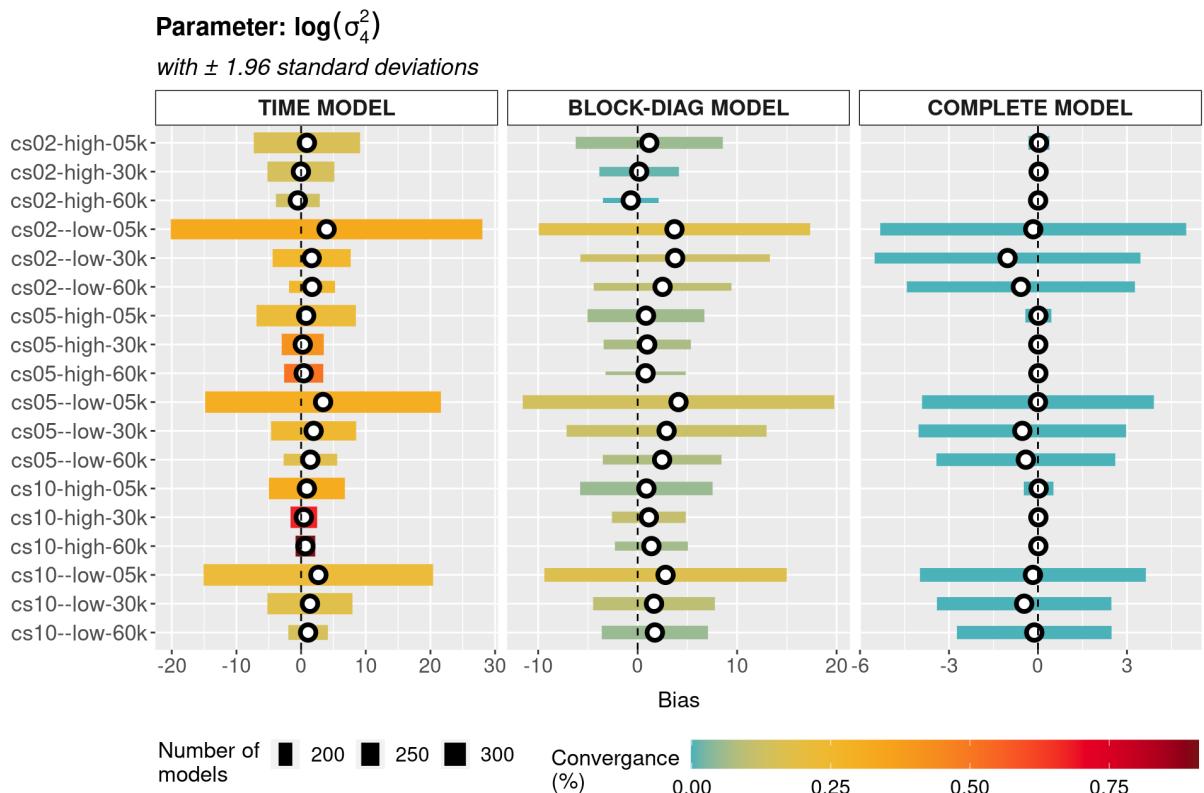
SOURCE: The author (2021).

FIGURE 21 – PARAMETER $\log(\sigma_3^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS

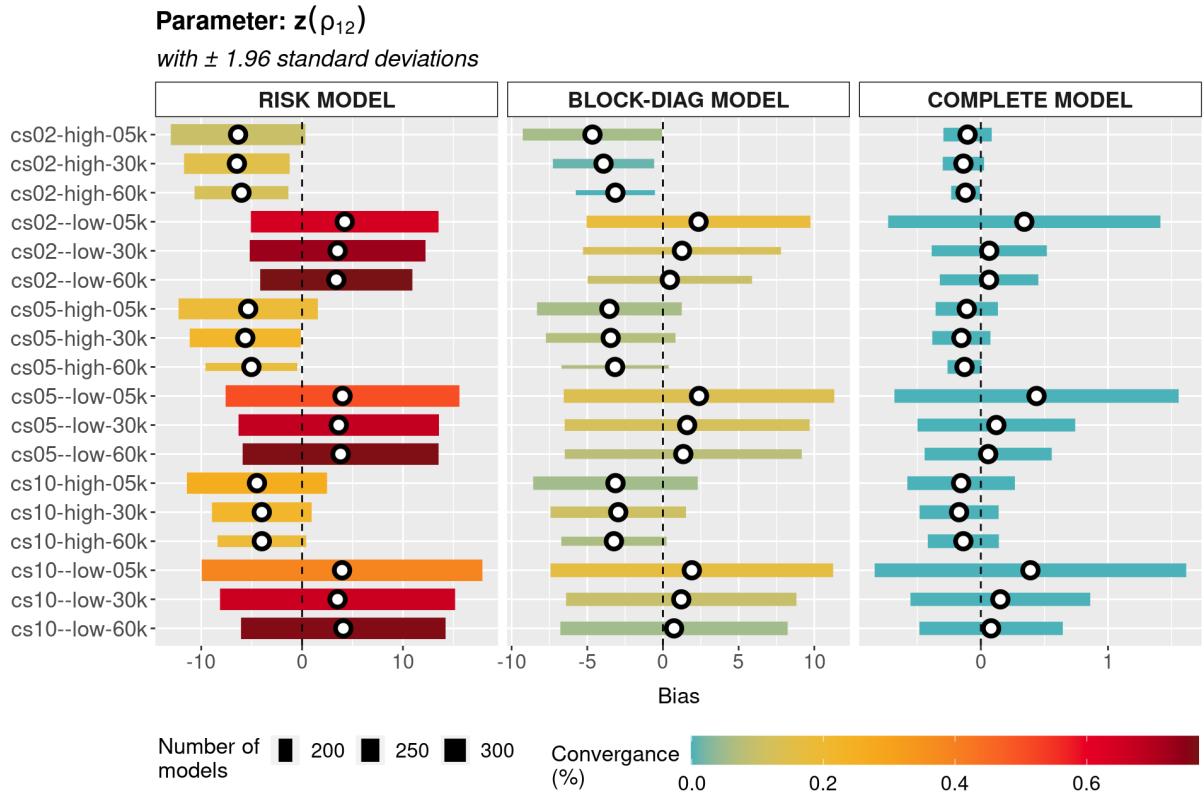


SOURCE: The author (2021).

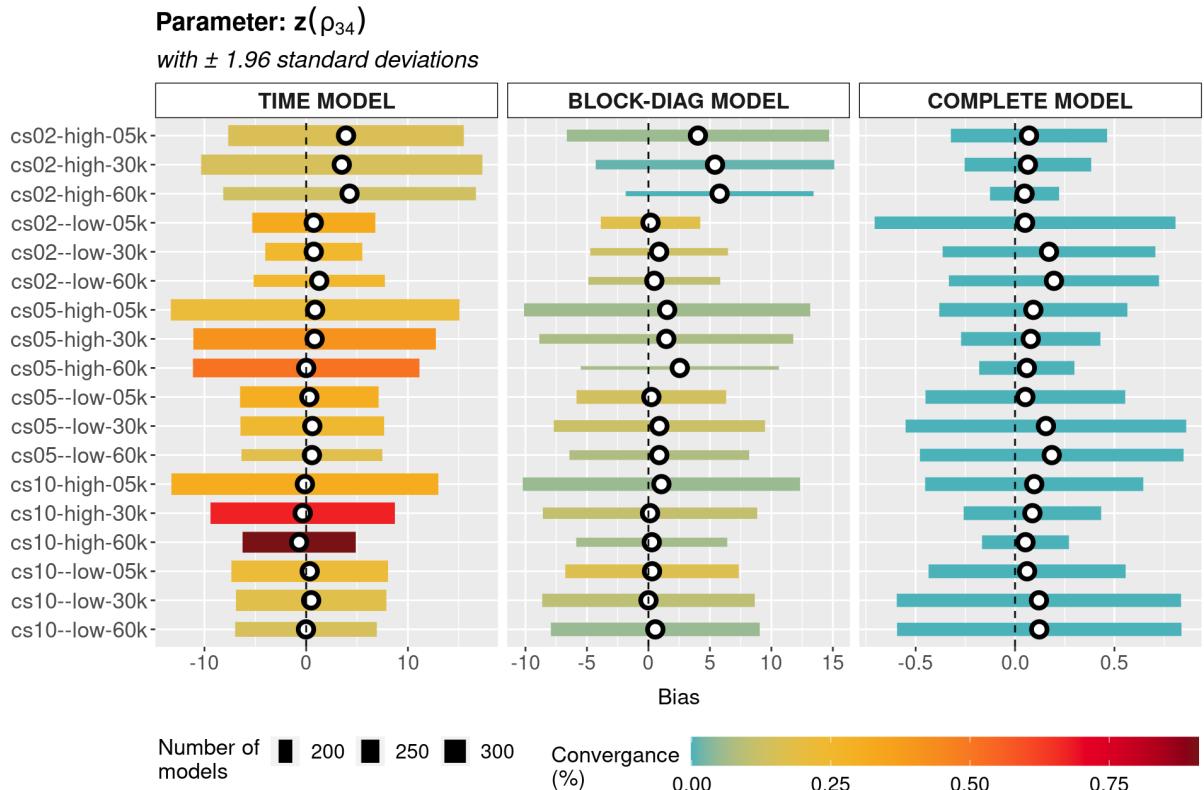
FIGURE 22 – PARAMETER $\log(\sigma_4^2)$ BIAS WITH ± 1.96 STANDARD DEVIATIONS



SOURCE: The author (2021).

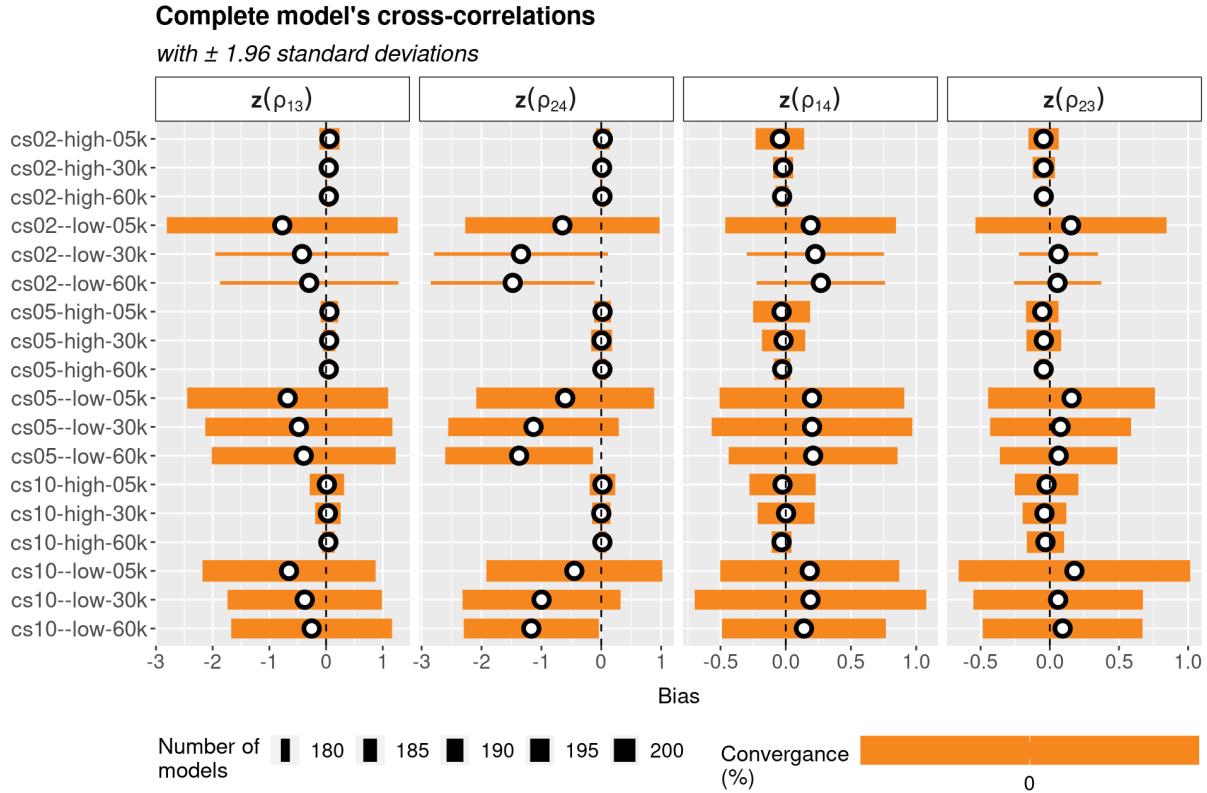
FIGURE 23 – PARAMETER $z(\rho_{12})$ BIAS WITH ± 1.96 STANDARD DEVIATIONS

SOURCE: The author (2021).

FIGURE 24 – PARAMETER $z(\rho_{34})$ BIAS WITH ± 1.96 STANDARD DEVIATIONS

SOURCE: The author (2021).

FIGURE 25 – PARAMETERS $\{z(\rho_{13}), z(\rho_{24}), z(\rho_{14}), z(\rho_{23})\}$ BIAS WITH ± 1.96 STANDARD DEVIATIONS



SOURCE: The author (2021).

In the *simpler* models with a latent structure just in one level, is hard to see some significant difference between the cluster and sample sizes. In the complete model, in the other hand, the difference is clear: as we increase the cluster and the sample size, the bias-variance decreases. The mean-bias is basically always the same. In the risk and block-diag models is hard to point-out a scenario as the best or worst. For the time model, in two of the parameters, with the scenario of cluster size 10; low CiF; and 5 thousand data points, we get a much bigger standard deviation.

With the log-variances in [Figure 19](#), [Figure 20](#), [Figure 21](#), and [Figure 22](#), we have instead a similar behavior through the models. For all models, the high CIF scenarios are the ones with a smaller mean and bias-variance. From the risk/time model to the block-diag model, we do not see a significant improvement in terms of bias reduction. Such improvement, however, is clear when we look at the complete model. Again, the magick of considering the cross-correlations.

The same said about the log-variances, can be applied to the risk correlation in [Figure 23](#), with one addendum: the bias reduction is even bigger. With the time correlation in [Figure 24](#), we get the same behavior observed with the fixed-effect parameters i.e., in the simpler models, the smaller biases are observed with the low CIF scenarios. However, with the complete model, we get the opposite. With the cross-correlations in

[Figure 25](#), the mean and bias-variances are much smaller in the high CIF scenarios.

The biggest bias-variances obtained are with the log-variances. A final remark to be made is about convergences. With the simpler models, not all of them work, having in some scenarios (generally the ones with sixty thousand points) a 50~60% convergence rate. However, from this 50~60%, most reach the ideal gradient-based convergence. The complete model is the complete opposite. Basically, all fits reach a convergence (~100% performance) but never the desired gradient-based convergence.

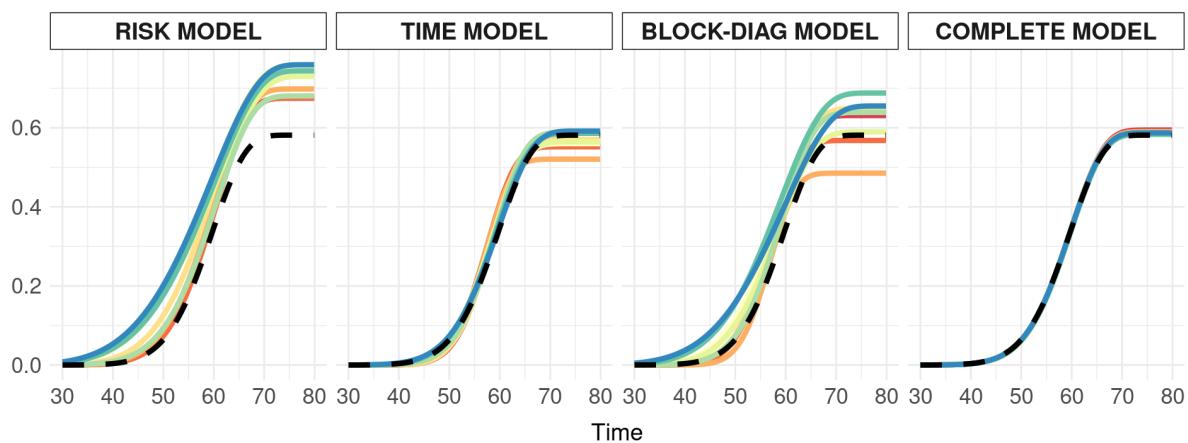
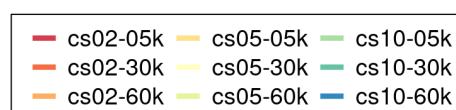
After looking at the parameter biases, let us take a look at the implied mean-CIF curves. To nicely accommodate all seventy-two scenarios we split the curves by level-CIF. In [Figure 26](#) we have the high CIF scenario curves and in [Figure 27](#) the low CIF scenario curves.

FIGURE 26 – HIGH CUMULATIVE INCIDENCE FUNCTION (CIF) SCENARIO CURVES

HIGH CIF SCENARIO

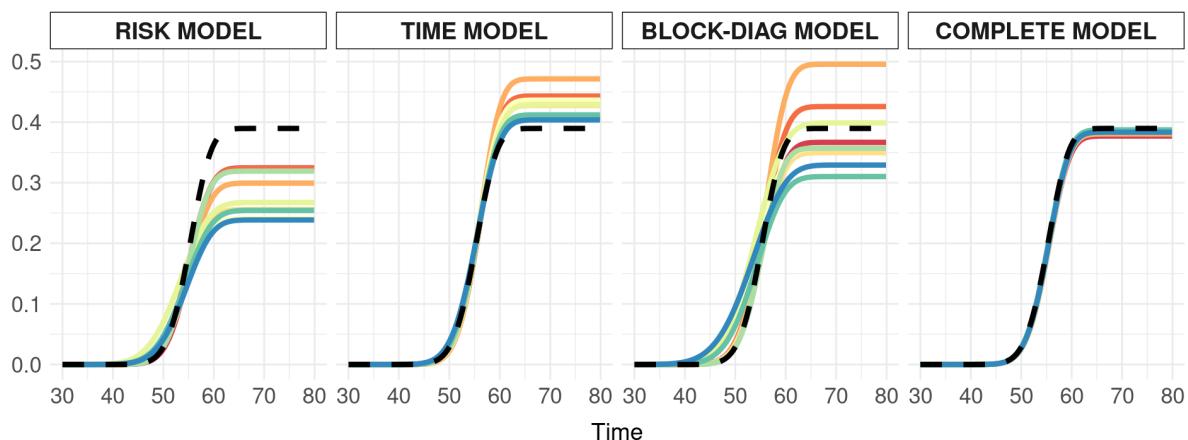
CIF of failure cause 1

True curve in dashed black



CIF of failure cause 2

True curve in dashed black



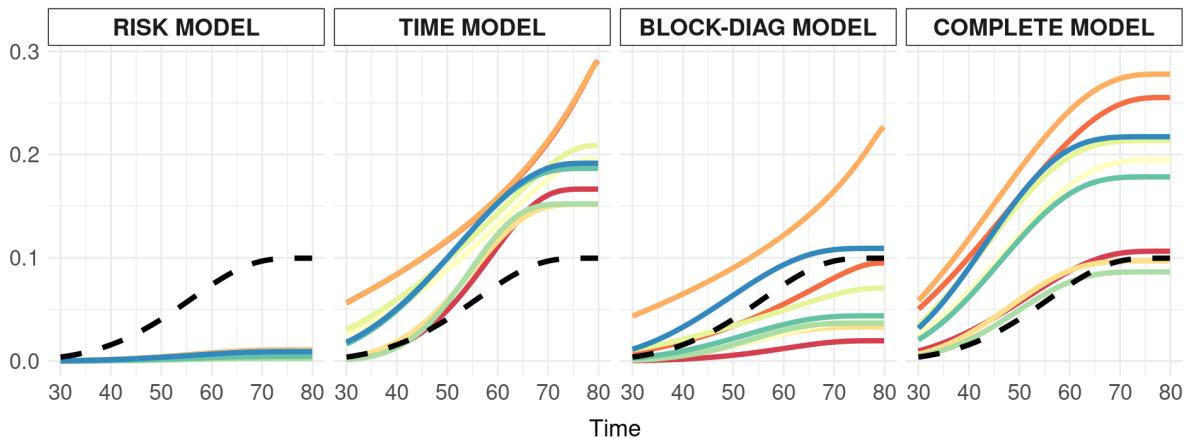
SOURCE: The author (2021).

FIGURE 27 – LOW CUMULATIVE INCIDENCE FUNCTION (CIF) SCENARIO CURVES

LOW CIF SCENARIO

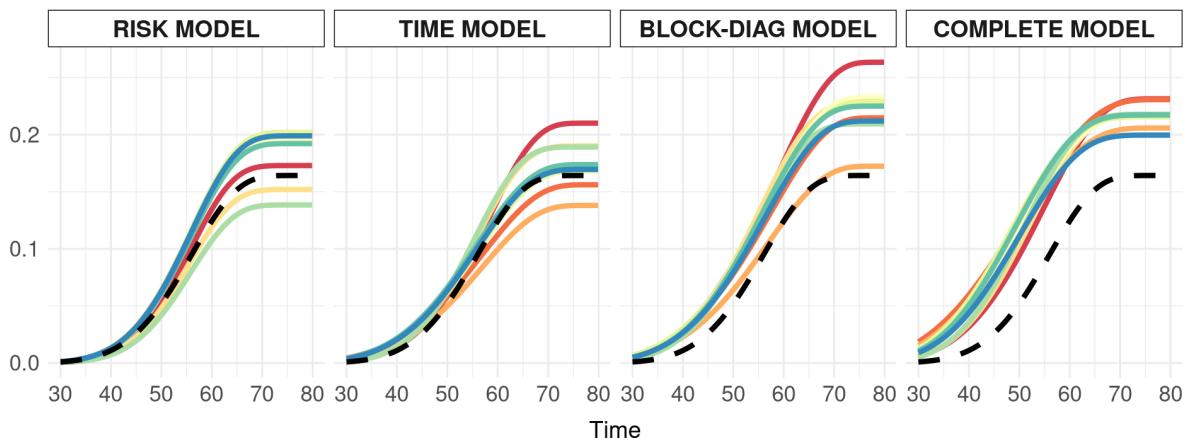
CIF of failure cause 1

True curve in dashed black



CIF of failure cause 2

True curve in dashed black



SOURCE: The author (2021).

Since in all models we have a latent structure for the within-cluster dependency, the inherent idea is that this will also affect the fixed-effect parameter estimates. By taking its average in each of the seventy-two scenarios, we are able to construct the mean CIF curves.

In [Figure 26](#) we have all the thirty-six curves obtained in the high CIF scenarios. It is clear that with the complete model we get a perfect fit in all nine scenarios. The risk and time models estimate well the curve shape parameters but they fail to learn the max incidence. A compensation between curves is clear. In the risk model, there is a super estimation of β_1 in all scenarios. For failure cause 2, there is a sub estimation. With the time model, we observe the opposite compensation but on a smaller scale.

With the time model, we get much better curves than with the risk model. The block-diag model results are a middle term between them. For the time model, the scenario with cluster size 10 and 60 thousand points is a highlight. For the block-diag model, the highlight is the scenario with cluster size 5 and 30 thousand points.

In the low CIF scenarios in [Figure 27](#), the estimation is clearly more difficult. The overall fits are bad, being impossible to select a scenario with overall good results. For one of the failure causes, the estimation quality is not so bad. The problem is when we look to the other. An interesting scenario is the one with cluster size 2 and 60 thousand data points. In this scenario we see the worst fits for failure cause 1, with a negative highlight in the block-diag configuration. However, in this same model, for failure cause 2, this is the only scenario, by far, where we learn the true curve. An interesting compensation phenomena. The best joint fit is still with the complete model.

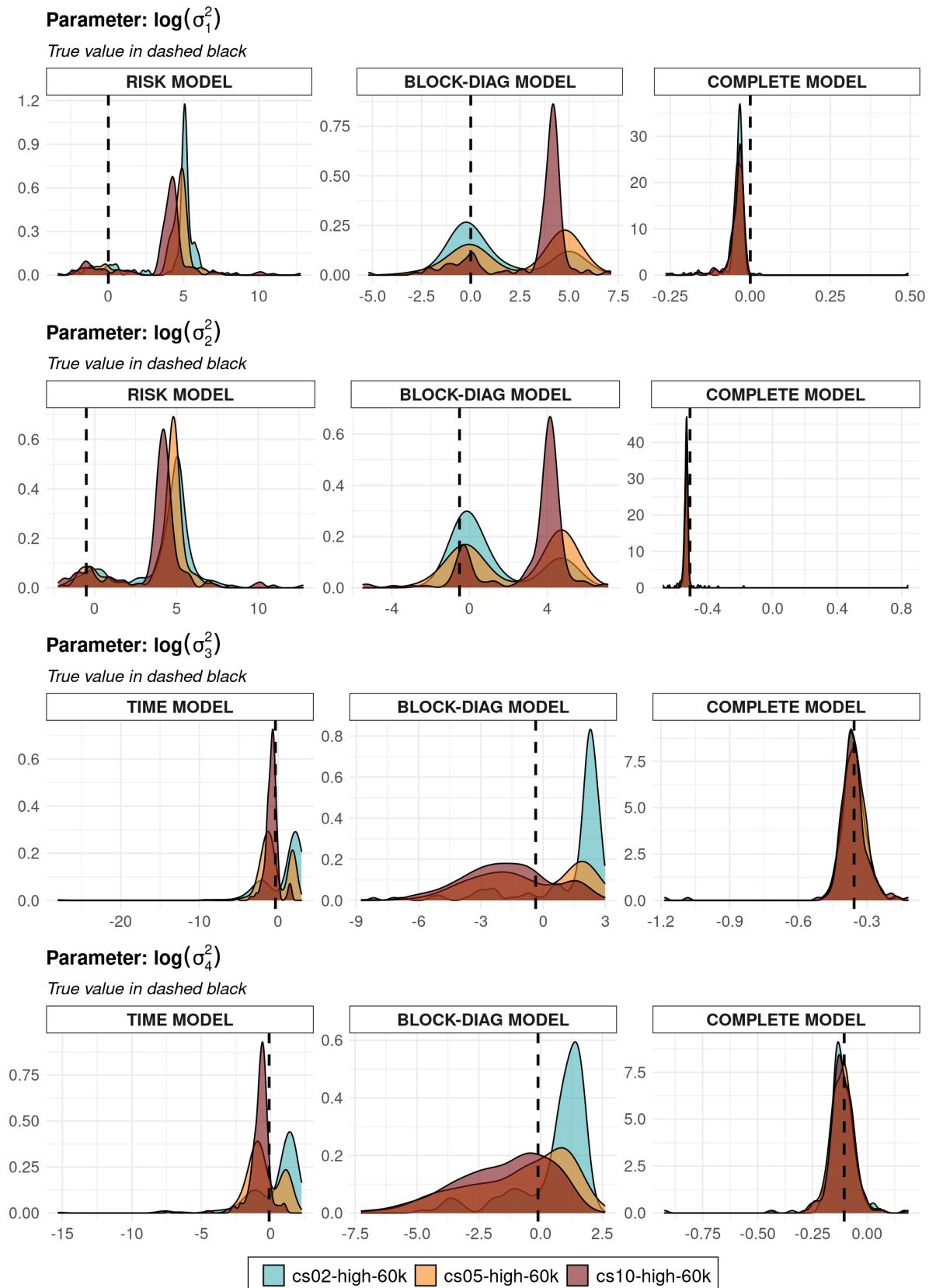
Now we look at how the latent-effect parameter estimates distribute themselves. Given the huge number of scenarios and the fact that it is harder to estimate covariance parameters, we chose to plot the parameters just in the scenarios with better performances. By the metrics of small bias and CIF shape learning, the scenarios with better results are the ones with high CIF and bigger sample sizes. We have the densities for the variance parameters, in each of these scenarios, presented in [Figure 28](#). In [Figure 29](#) we have the same for the correlation parameters.

An interesting result is the clear difference between risk and time models' covariance parameters. With the risk model, we have an evident super estimation and bigger variances. With the time model we get much better results, but still with high variances. The block-diag model generally performs better than the risk model and worse than the time model, showing again to be a compromise between them. Besides the bias itself, we should also pay attention to the values. We model the variances in the log-scale, so a value 5, in reality, implies a variance of $\exp(4) = 148$. Terrible. This kind of problem do not sound to appear with the complete model.

All correlations are quite well estimated, in all three scenarios, with the complete model. Not only the correlations but the variances also. The lack of any considerable difference between the covariance densities, indicates no quality divergences in the results for different cluster sizes. The densities in [Figure 28](#) and [Figure 29](#) are the final corroboration indicating the good estimability of the complete model. Between the four tested models, the complete model was the one with the smallest biases, better CIF shape learning, and precisest covariance parameter estimates.

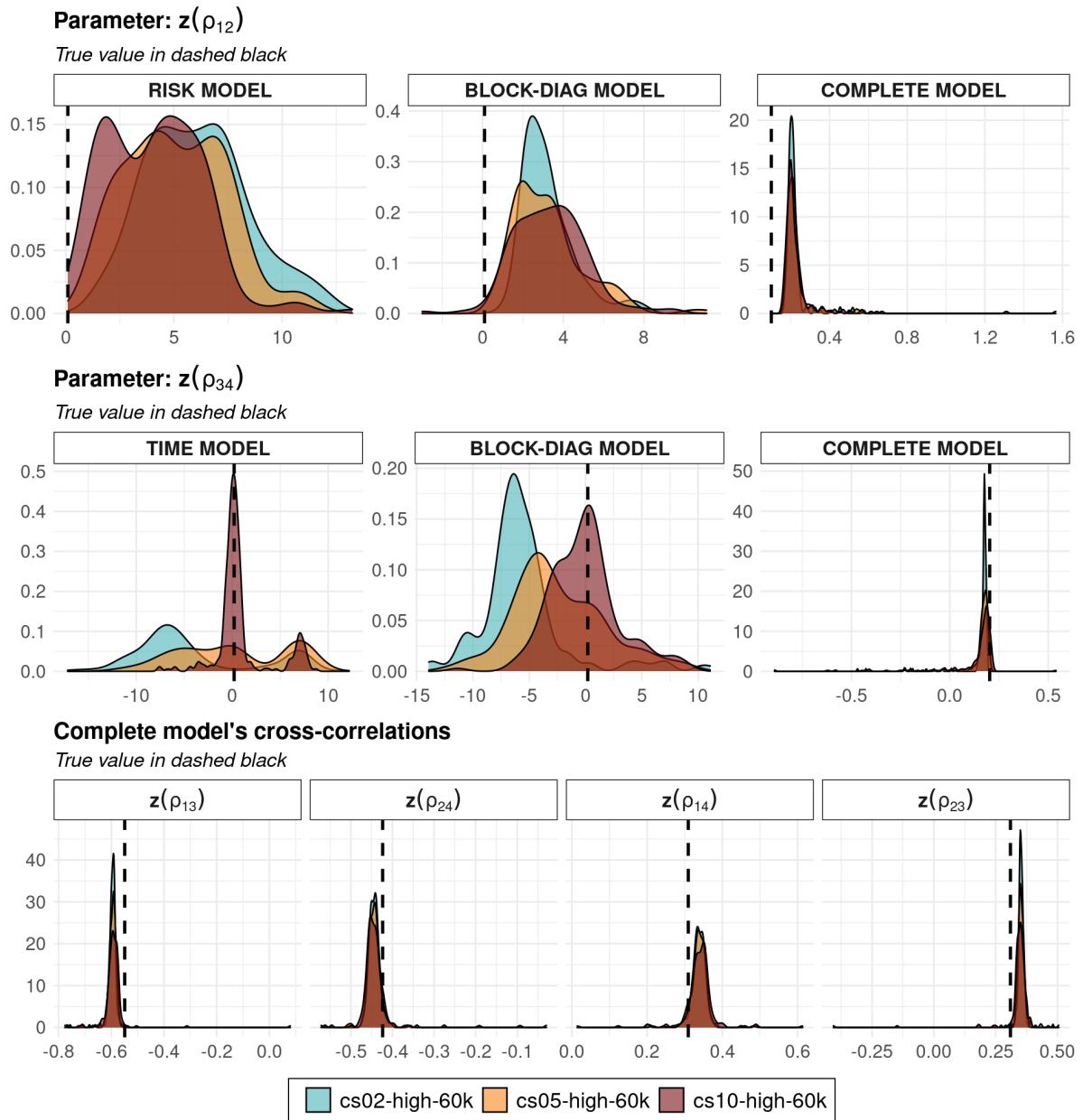
In [Figure 30](#) we have a heat-map of the correlations between parameter estimates of the complete model in the scenario with clusters of size 10, high CIF, and 60 thousand data points.

FIGURE 28 – VARIANCE PARAMETERS DENSITIES IN THE SCENARIOS OF HIGH CIF AND 60 THOUSAND DATA POINTS



SOURCE: The author (2021).

FIGURE 29 – CORRELATION PARAMETERS DENSITIES IN THE SCENARIOS OF HIGH CIF AND 60 THOUSAND DATA POINTS



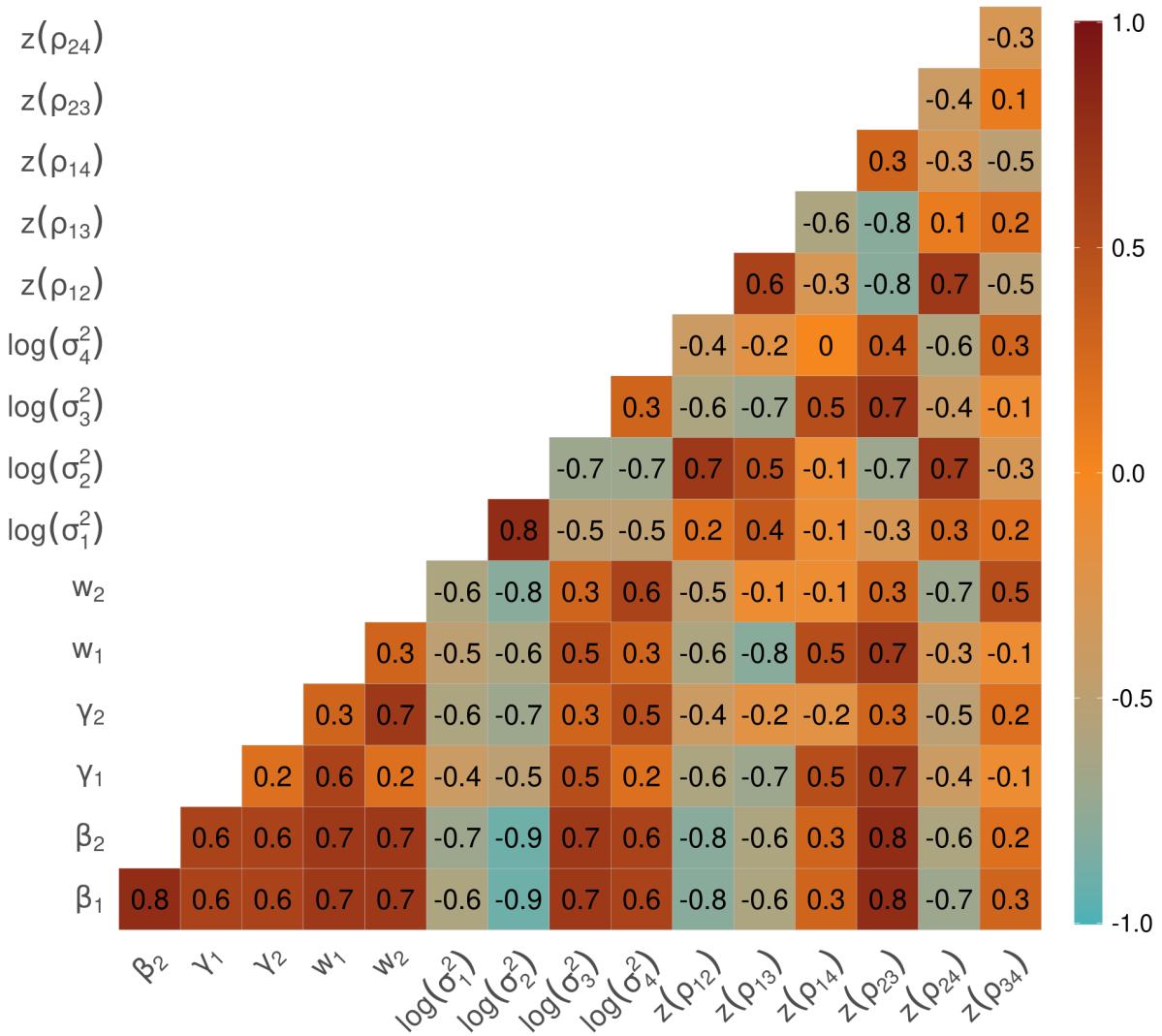
SOURCE: The author (2021).

We have a little bit of everything in Figure 30. Some correlations are very close to zero, but we also have strong positive and negative correlations. We can mention some curiosities, but nothing pathological appears to happen, at least nothing clear.

All fixed-effect parameters are positive correlated, with an emphasis on the correlation between β_1 and β_2 , and the one of the β s with the ws . Another interesting observation is the negative correlation between the β s and the risk log-variances, and also the positive correlation between the β s and the trajectory time log-variances.

FIGURE 30 – COMPLETE MODEL'S PARAMETERS CORRELATION HEAT-MAP IN THE SCENARIO OF CLUSTER SIZE 10, HIGH CIF, AND SIXTY-THOUSAND DATA POINTS

cs10-high-60k, complete model



SOURCE: The author (2021).

6 FINAL CONSIDERATIONS

6.1 FUTURE WORKS

BIBLIOGRAPHY

- ANDERSEN, P. K.; GESKUS, R. B.; WITTE, T. de; PUTTER, H. Competing risks in epidemiology: possibilities and pitfalls. *International Journal of Epidemiology*, v. 31, n. 1, p. 861–870, 2012. Cited on page [17](#).
- BATES, D.; MAECHLER, M. *Matrix: Sparse and Dense Matrix Classes and Methods*. R Foundation for Statistical Computing. Vienna, Austria, 2019. R package version 1.2-18 (<https://CRAN.R-project.org/package=Matrix>). Cited on page [34](#).
- BONAT, W. H.; RIBEIRO, P. J. Practical likelihood analysis for spatial generalized linear mixed models. *Environmetrics*, v. 27, n. 1, p. 83–89, 2016. Cited 2 times on pages [23](#) and [25](#).
- CEDERKVIST, L.; HOLST, K. K.; ANDERSEN, K. K.; SCHEIKE, T. H. Modeling the cumulative incidence function of multivariate competing risks data allowing for within-cluster dependence of risk and timing. *Biostatistics*, v. 20, n. 2, p. 199–217, 2019. Cited 9 times on pages [8](#), [15](#), [17](#), [19](#), [37](#), [38](#), [39](#), [40](#), and [48](#).
- CLAYTON, D. G. A model for association in bivariate life tables and its application in epidemiological studies of familial tendency in chronic disease incidence. *Biometrika*, v. 65, n. 1, p. 141–151, 1978. Cited on page [17](#).
- COX, D. R.; REID, N. A note on pseudolikelihood constructed from marginal densities. *Biometrika*, v. 91, n. 3, p. 729–737, 2004. Cited on page [17](#).
- DEMPSTER, A. P.; LAIRD, N. M.; RUBIN, D. B. Maximum likelihood estimation from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society, Series B (Methodological)*, v. 39, n. 1, p. 1–38, 1977. Cited on page [23](#).
- DENNIS, J. E.; GAY, D. M.; WELSCH, R. E. An Adaptive Nonlinear Least-Squares Algorithm. *ACM Transactions on Mathematical Software*, v. 7, n. 3, p. 348–368, 1981. Cited on page [27](#).
- DIACONIS, P. The Markov chain Monte Carlo revolution. *Bulletin (New Series) of the American Mathematical Society*, v. 46, n. 2, p. 179–205, 2009. Cited on page [22](#).
- FOURNIER, D. A.; SKAUG, H. J.; ANCHETA, J.; IANELLI, J.; MAGNUSSON, A.; MAUNDER, M. N.; NIELSEN, A.; SIBERT, J. AD Model Builder: using automatic differentiation for statistical inference of highly parameterized complex nonlinear models. *Optimization Methods and Software*, v. 27, n. 2, p. 233–249, 2012. Cited on page [31](#).
- GAY, D. M. *Usage summary for selected optimization routines*. Computing Science Technical Report 153, AT&T Bell Laboratories. Murray Hill, NJ, 1990. Cited 2 times on pages [27](#) and [52](#).
- GELFAND, A. E.; SMITH, A. F. M. Sampling-Based Approaches to Calculating Marginal Densities. *Journal of the American Statistical Association*, v. 85, n. 410, p. 398–409, 1990. Cited on page [22](#).
- GUENNEBAUD, G.; JACOB, B. et al. *Eigen v3*. 2010. (<http://eigen.tuxfamily.org>). Cited on page [32](#).

KALBFLEISCH, J. D.; PRENTICE, R. L. *The Statistical Analysis of Failure Time Data*. Second Edition. Hoboken, New Jersey: John Wiley & Sons, Inc., 2002. Cited 3 times on pages 15, 16, and 37.

KRISTENSEN, K.; NIELSEN, A.; BERG, C. W.; SKAUG, H. J.; BELL, B. M. TMB: Automatic Differentiation and Laplace Approximation. *Journal of Statistical Software*, v. 70, n. 5, p. 1–21, 2016. Cited 5 times on pages 19, 21, 31, 34, and 44.

LINDSAY, B. G. Composite likelihood methods. *Contemporary Mathematics*, v. 80, n. 1, p. 221–239, 1988. Cited on page 17.

MCCULLAGH, P.; NELDER, J. A. *Generalized linear models*. Second edition. London: Chapman & Hall, 1989. Cited on page 18.

MCCULLOCH, C. E.; SEARLE, S. R. *Generalized, Linear, and Mixed Models*. New York: John Wiley & Sons, Inc., 2001. Cited 2 times on pages 18 and 21.

MOLENBERGHS, G.; VERBEKE, G. *Models for Discrete Longitudinal Data*. New York: Springer, 2005. Cited on page 23.

NELDER, J. A.; WEDDERBURN, R. W. M. Generalized linear models. *Journal of the Royal Statistical Society, Series A*, v. 135, n. 3, p. 370–384, 1972. Cited on page 18.

NOCEDAL, J.; WRIGHT, S. J. *Numerical Optimization*. Second Edition. New York: Springer, 2006. (Springer Series in Operations Research and Financial Engineering). Cited 3 times on pages 26, 27, and 28.

PEYRÉ, G. *Course notes on Optimization for Machine Learning*. 2020. May 10, <https://mathematical-tours.github.io/book-sources/optim-ml/OptimML.pdf>. CNRS & DMA, École Normale Supérieure. Cited 2 times on pages 28 and 29.

PINHEIRO, J. C.; BATES, D. M. Unconstrained parametrizations for variance-covariance matrices. *Statistics and Computing*, v. 6, n. 3, p. 289–296, 1996. Cited on page 44.

PINHEIRO, J. C.; CHAO, E. C. Efficient Laplacian and Adaptive Gaussian Quadrature Algorithms for Multilevel Generalized Linear Mixed Models. *Journal of Computational and Graphical Statistics*, v. 15, n. 1, p. 58–81, 2006. Cited on page 23.

POURAHMADI, M. Cholesky decompositions and estimation of a covariance matrix: orthogonality of variance-correlation parameters. *Biometrika*, v. 94, n. 4, p. 1006–1013, 2007. Cited on page 45.

R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2021. <https://www.R-project.org/>. Cited 3 times on pages 21, 27, and 31.

SHUN, Z.; MCCULLAGH, P. Laplace approximation of high dimensional integrals. *Journal of the Royal Statistical Society, Series B (Methodological)*, v. 57, n. 4, p. 749–760, 1995. Cited on page 23.

TIERNEY, L.; KADANE, J. Accurate approximations for posterior moments and marginal densities. *Journal of the American Statistical Association*, v. 81, n. 393, p. 82–86, 1986. Cited on page 23.

VALPEL, J. W.; MANTON, K. G.; STALLARD, E. The impact of heterogeneity in Individual Frailty on the Dynamics of Mortality. *Demography*, v. 16, n. 1, p. 439–454, 1979. Cited on page [17](#).

VARIN, C.; REID, N.; FIRTH, D. An overview of composite likelihood methods. *Statistica Sinica*, v. 21, n. 1, p. 5–42, 2011. Cited on page [17](#).

Ver HOEF, J. M. Who Invented the Delta Method? *The American Statistician*, v. 66, n. 2, p. 124–127, 2012. Cited on page [35](#).

WOOD, S. N. *Core Statistics*. IMS: Institute of Mathematical Statistics, Textbooks, 2015. Cited 3 times on pages [23](#), [24](#), and [28](#).

Appendix

APPENDIX A – ANALYTIC GRADIENT OF THE LATENT EFFECTS FOR THE JOINT
LOG-LIKELIHOOD FUNCTION OF THE MULTINOMIAL GLMM FOR CLUSTERED
COMPETING RISKS DATA

The following gradient components are computed by cluster, to be used e.g., in a Newton optimization. Subject i at cluster j and for competing cause k

$$\begin{aligned} \frac{\partial}{\partial u_{kj}} \log L(\boldsymbol{\theta} | \mathbf{y}_j, \mathbf{r}_j) = & \\ & y_{kij} \frac{1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} - \left(\sum_{m \neq k}^{K-1} y_{mij} \right) \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} - \\ & y_{Kij} \frac{1}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \left(\frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} (1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\})}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \right. \\ & \left. \frac{w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}]}{1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]} - \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \right. \\ & \left. \sum_{m \neq k}^{K-1} w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{mij}\gamma_{mj} - \eta_{mj}] \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]} \right) - \\ & e_k^\top \mathbf{Q} \mathbf{r}_j, \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \eta_{kj}} \log L(\boldsymbol{\theta} | \mathbf{y}_j, \mathbf{r}_j) = & \\ & y_{kij} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}) - \\ & y_{Kij} \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \\ & w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] - \\ & 1 - \sum_{n=1}^{K-1} \frac{\exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}] \\ & e_k^\top \mathbf{Q} \mathbf{r}_j, \end{aligned}$$

with e_k^\top begin a vector with 1 at the k -th position and zero elsewhere.

APPENDIX B – ANALYTIC HESSIAN OF THE LATENT EFFECTS FOR THE JOINT
LOG-LIKELIHOOD FUNCTION OF THE MULTINOMIAL GLMM FOR CLUSTERED
COMPETING RISKS DATA

The following hessian components are computed by cluster, to be used e.g., in a Newton optimization. Subject i at cluster j and for competing cause k

$$\begin{aligned}
 & \frac{\partial^2}{\partial u_{kj}^2} \log L(\boldsymbol{\theta} | \mathbf{y}_j, \mathbf{r}_j) = \\
 & - \frac{\left(\sum_{k=1}^{K-1} y_{kij} \right) \exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \left(1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} \right)}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2} + \\
 & \frac{y_{Kij} \exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \\
 & \frac{\sum_{m \neq k}^{K-1} w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{mij}\gamma_{mj} - \eta_{mj}] \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}])} - \\
 & \frac{y_{Kij} w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}]}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \\
 & \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \left(1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} \right)}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}])} - \\
 & \frac{y_{Kij} \exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2} \left(\right. \\
 & \frac{\sum_{m \neq k}^{K-1} w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{mij}\gamma_{mj} - \eta_{mj}] \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right)^2} - \\
 & \frac{w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \left(1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} \right)}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right)^2} \left. \right) \\
 & \times \left(\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right) + \right. \\
 & \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right) \times \right. \\
 & \left. \left(1 - w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \right) - \mathbf{e}_k^\top \mathbf{Q}, \right)
 \end{aligned}$$

$$\begin{aligned}
& \frac{\partial^2}{\partial \eta_{kj}^2} \log L(\boldsymbol{\theta} \mid \mathbf{y}_j, \mathbf{r}_j) = \\
& -y_{kij} - y_{Kij} \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \left(\right. \\
& w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \times \\
& \left. \left(w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj} \right)^2 - 1 \right. \\
& \left. \left(1 - \sum_{n=1}^{K-1} \frac{\exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}] \right. \right. \\
& \left. \left(w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \right)^2 \right. \\
& \left. \left(1 - \sum_{n=1}^{K-1} \frac{\exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}] \right)^2 \right) \\
& \left. - \mathbf{e}_k^\top \mathbf{Q}, \right)
\end{aligned}$$

$$\begin{aligned}
& \frac{\partial^2}{\partial u_{kj} u_{mj}} \log L(\boldsymbol{\theta} \mid \mathbf{y}_j, \mathbf{r}_j) = \\
& \left(\sum_{k=1}^{K-1} y_{kij} \right) \frac{\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2} + \\
& \frac{y_{Kij} \exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \left(\right. \\
& w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{mij}\gamma_{mj} - \eta_{mj}] \\
& \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right. \right. \\
& w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \\
& \left. \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right) \right. \right. \\
& \frac{y_{Kij}}{\left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2} \left(\exp\{\mathbf{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \left(\right. \right. \\
& \sum_{m \neq k}^{K-1} w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{mij}\gamma_{mj} - \eta_{mj}] \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} \\
& \left. \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right)^2 \right. \right. \\
& w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{kij}\gamma_{kj} - \eta_{kj}] \left(1 + \sum_{m \neq k}^{K-1} \exp\{\mathbf{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} \right) \\
& \left. \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\mathbf{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t - \delta/2}{\delta/2}\right) - \mathbf{x}_{nij}\gamma_{nj} - \eta_{nj}]) \right)^2 \right) \right)
\end{aligned}$$

$$\begin{aligned}
& \left) \times \left(\exp\{\boldsymbol{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} (1+ \right. \right. \\
& \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} (1 - w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{nij}\gamma_{nj} - \eta_{nj}]) + \\
& \exp\{\boldsymbol{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\} (1 - w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{mij}\gamma_{mj} - \eta_{mj}]) (1+ \\
& \left. \left. \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right) \right) - \boldsymbol{e}_k^\top \boldsymbol{Q},
\end{aligned}$$

$$\begin{aligned}
& \frac{\partial^2}{\partial \eta_{kj} \eta_{mj}} \log L(\boldsymbol{\theta} \mid \boldsymbol{y}_j, \boldsymbol{r}_j) = \\
& - y_{Kij} \frac{\exp\{\boldsymbol{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \\
& w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{kij}\gamma_{kj} - \eta_{kj}] \times \\
& \left(1 - \sum_{n=1}^{K-1} \frac{\exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{nij}\gamma_{nj} - \eta_{nj}] \right)^2 \times \\
& \frac{\exp\{\boldsymbol{x}_{mij}\boldsymbol{\beta}_{mj} + u_{mj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_m \frac{\delta}{2\delta t - 2t^2} (w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{mij}\gamma_{mj} - \eta_{mj}) \times \\
& \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{mij}\gamma_{mj} - \eta_{mj}] - \boldsymbol{e}_k^\top \boldsymbol{Q},
\end{aligned}$$

$$\begin{aligned}
& \frac{\partial^2}{\partial \eta_{kj} u_{kj}} \log L(\boldsymbol{\theta} \mid \boldsymbol{y}_j, \boldsymbol{r}_j) = \\
& y_{Kij} \frac{\exp\{\boldsymbol{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} \times \\
& w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{kij}\gamma_{kj} - \eta_{kj}] \times \\
& \left(1 - \sum_{n=1}^{K-1} \frac{\exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{nij}\gamma_{nj} - \eta_{nj}] \right)^2 \times \\
& \left(\sum_{n \neq k}^{K-1} \frac{\exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \exp\{\boldsymbol{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\}}{\left(1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2} \times \right. \\
& w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{nij}\gamma_{nj} - \eta_{nj}] - \\
& \left. \exp\{\boldsymbol{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \left(\left(1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right) - \exp\{\boldsymbol{x}_{kij}\boldsymbol{\beta}_{kj} + u_{kj}\} \right) \times \right. \\
& \left. \left(1 + \sum_{n=1}^{K-1} \exp\{\boldsymbol{x}_{nij}\boldsymbol{\beta}_{nj} + u_{nj}\} \right)^2 \right. \\
& w_k \frac{\delta}{2\delta t - 2t^2} \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - \boldsymbol{x}_{kij}\gamma_{kj} - \eta_{kj}] \left) - \right.
\end{aligned}$$

$$\begin{aligned}
& \frac{y_{Kij} \frac{\exp\{x_{kij}\beta_{kj} + u_{kj}\} ((1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}) - \exp\{x_{kij}\beta_{kj} + u_{kj}\})}{(1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\})^2}}{1 - \sum_{n=1}^{K-1} \frac{\exp\{x_{nij}\beta_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{nij}\gamma_{nj} - \eta_{nj}] \times \\
& w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}) \times \\
& \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}] - e_k^\top Q,
\end{aligned}$$

$$\begin{aligned}
& \frac{\partial^2}{\partial \eta_{kj} u_{mj}} \log L(\theta | y_j, r_j) = \\
& y_{Kij} \frac{\exp\{x_{kij}\beta_{kj} + u_{kj}\} \exp\{x_{mij}\beta_{mj} + u_{mj}\}}{(1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\})^2} \times \\
& w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}] + \\
& \frac{1 - \sum_{n=1}^{K-1} \frac{\exp\{x_{nij}\beta_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{nij}\gamma_{nj} - \eta_{nj}]}{1 - \sum_{n=1}^{K-1} \frac{\exp\{x_{nij}\beta_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}}} \times \\
& y_{Kij} \frac{\exp\{x_{kij}\beta_{kj} + u_{kj}\}}{1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}} \times \\
& w_k \frac{\delta}{2\delta t - 2t^2} (w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}) \phi[w_k \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{kij}\gamma_{kj} - \eta_{kj}] \times \\
& \left(1 - \sum_{n=1}^{K-1} \frac{\exp\{x_{nij}\beta_{nj} + u_{nj}\}}{1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}} w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{nij}\gamma_{nj} - \eta_{nj}]\right)^2 \times \\
& \left(\sum_{n \neq m}^{K-1} \frac{\exp\{x_{nij}\beta_{nj} + u_{nj}\} \exp\{x_{mij}\beta_{mj} + u_{mj}\}}{(1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\})^2} \times \right. \\
& \left. w_n \frac{\delta}{2\delta t - 2t^2} \phi[w_n \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{nij}\gamma_{nj} - \eta_{nj}] - \right. \\
& \left. \exp\{x_{mij}\beta_{mj} + u_{mj}\} \left((1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}) - \exp\{x_{mij}\beta_{mj} + u_{mj}\} \right) \right. \\
& \left. \left(1 + \sum_{n=1}^{K-1} \exp\{x_{nij}\beta_{nj} + u_{nj}\}\right)^2 \right) \times \\
& w_m \frac{\delta}{2\delta t - 2t^2} \phi[w_m \operatorname{arctanh}\left(\frac{t-\delta/2}{\delta/2}\right) - x_{mij}\gamma_{mj} - \eta_{mj}] - e_k^\top Q,
\end{aligned}$$

with e_k^\top begin a vector with 1 at the k -th position and zero elsewhere.

APPENDIX C – R CODE TO SIMULATE FROM A multiGLMM WITH TWO COMPETING CAUSES AND CLUSTERS OF SIZE TWO. FOR MORE INFORMATION CHECK SECTION [4.1](#)

```

1 library(mvtnorm)
2 library(tidyverse)
3 library(mc2d)
4
5 datasimu <- function(J,      ## number of clusters
6                      cs,      ## clusters size (all the same size)
7                      time,    ## failure, censorship times
8                      Z,       ## latent effects design-matrix
9                      S,       ## variance-covariance matrix
10                     delta=80,
11                     beta =c( beta1=-2.0,   beta2=-1.5),
12                     gamma=c(gamma1= 1.2, gamma2= 1.0),
13                     w     =c(     w1= 3.0,       w2= 5.0),
14                     seed1,
15                     seed2)
16 {
17   out <- tibble::tibble(i=rep(seq(cs), times=J), ## cluster element
18                         j=rep(seq(J),  each=cs), ## cluster
19                         time=time,
20                         p1=NA,
21                         p2=NA,
22                         p3=NA)
23
24   K      <- dim(S)[1]/2 + 1
25   ladim <- 2*(K-1) ## latent effects dimension
26   set.seed(seed1)
27   U      <- mvtnorm::rmvnorm(J, mean=rep(0, ladim), sigma=S)
28   ZU     <- Z%*%U
29   risk1 <- exp(beta['beta1'] + ZU[, 1])
30   risk2 <- exp(beta['beta2'] + ZU[, 2])
31   level <- 1 + risk1 + risk2
32   gt     <- atanh(2*time/delta - 1)
33   dgt    <- delta/(2*time*(delta - time))
34   x1    <- w['w1']*gt - gamma['gamma1'] - ZU[, 3]
35   x2    <- w['w2']*gt - gamma['gamma2'] - ZU[, 4]
36
37   out$p1 <- risk1/level*w['w1']*dgt*dnorm(x1)
38   out$p2 <- risk2/level*w['w2']*dgt*dnorm(x2)
39
40   out    <- out %>% dplyr::mutate(p3=1-p1-p2)
41   set.seed(seed2)
42   y      <- mc2d::rmultinomial(cs*J, 1, prob=out%>%select(p1:p3))
43
44   out    <- out %>%

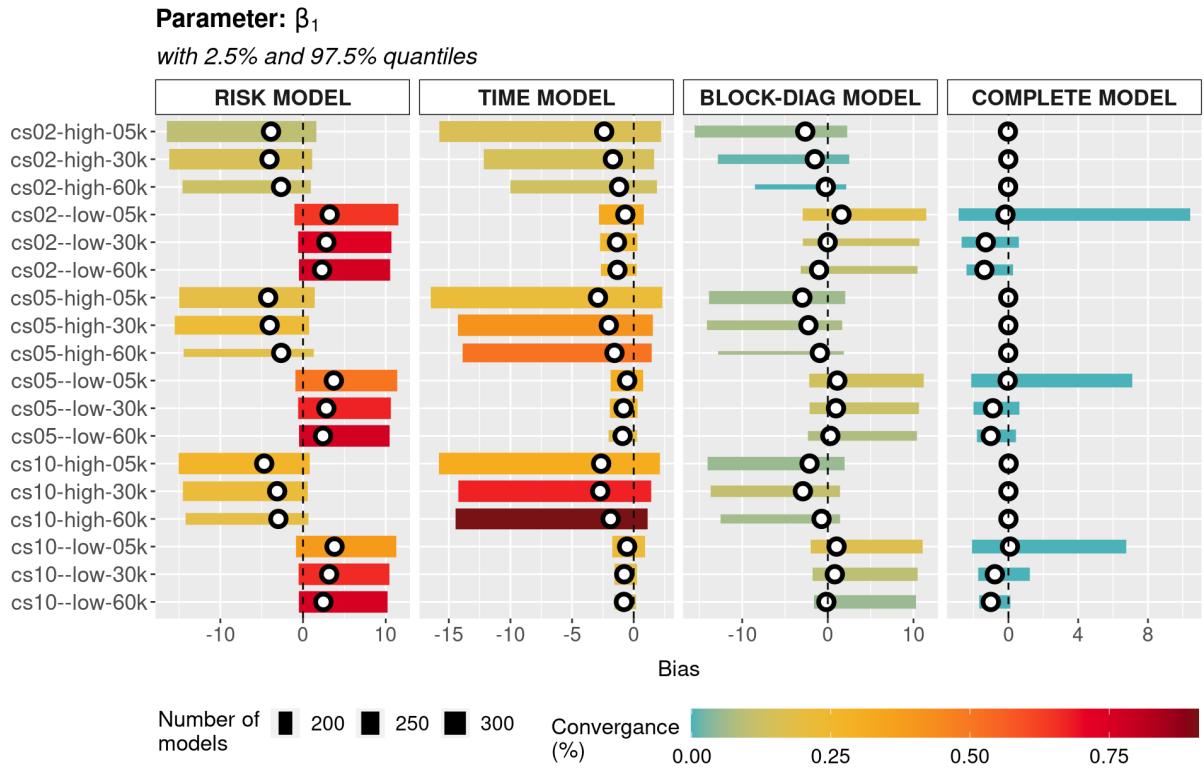
```

```
45      dplyr::bind_cols(tibble::as_tibble(y)) %>%
46      dplyr::rename(y1=V1, y2=V2, y3=V3)
47  return(out)
48 }
49 J    <- 50e3
50 cs   <- 2
51 time <- runif(n=cs*J, min=30, max=79.9)
52 Z    <- Matrix::bdiag(replicate(J, rep(1, cs), simplify=FALSE))
53 S    <- matrix(c( 1.0,  0.4, -0.1,  0.4,
54                  0.4,  1.0,  0.4, -0.1,
55                  -0.1,  0.4,  1.0,  0.4,
56                  0.4, -0.1,  0.4,  1.0), 4)
57 dat  <- datasimu(J=J, cs=cs, time=time, Z=Z, S=S, seed1=1, seed2=2)
```

SOURCE: The author (2021).

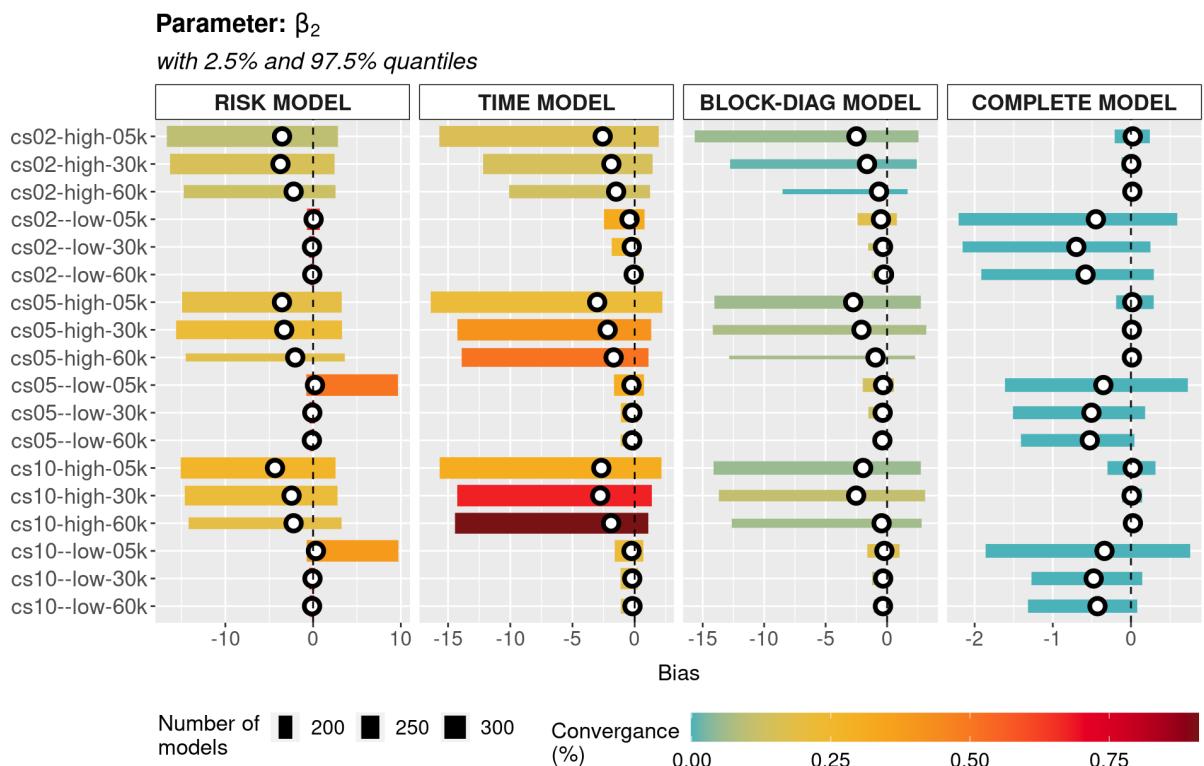
APPENDIX D – MODEL PARAMETERS BIAS WITH 2.5% AND 97.5% QUANTILES

FIGURE 31 – PARAMETER β_1 BIAS WITH 2.5% AND 97.5% QUANTILES

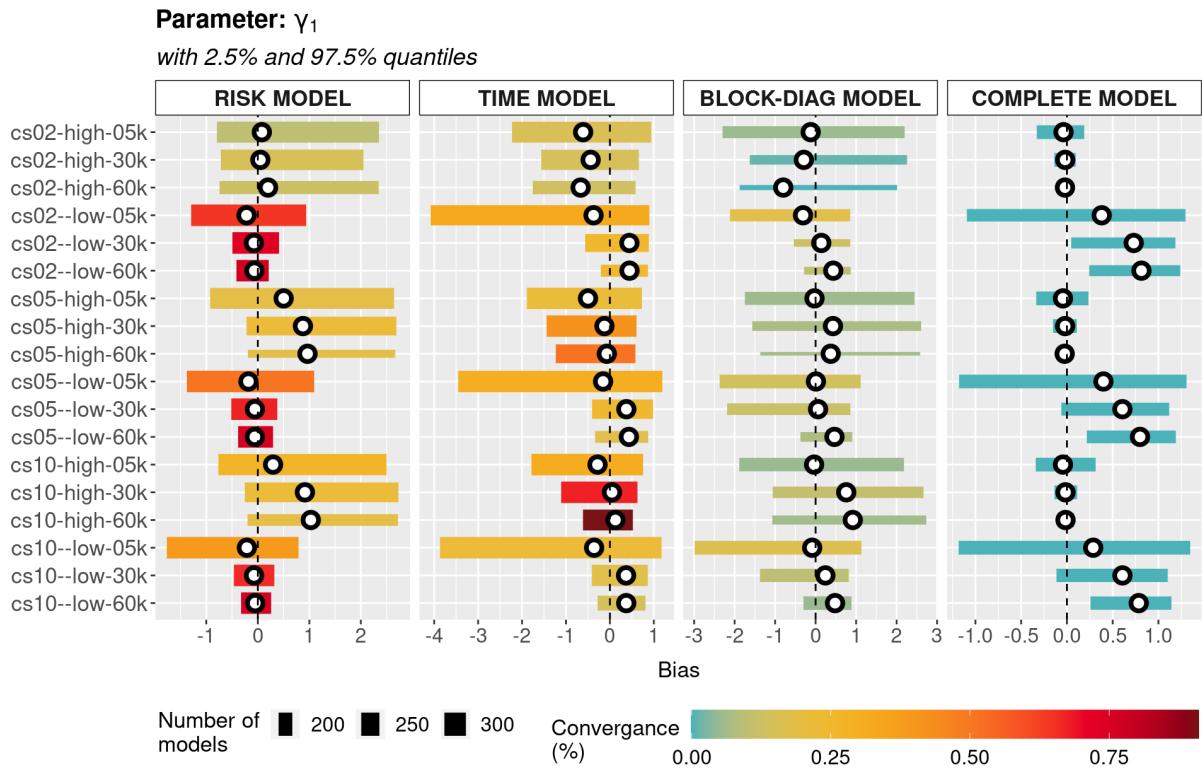


SOURCE: The author (2021).

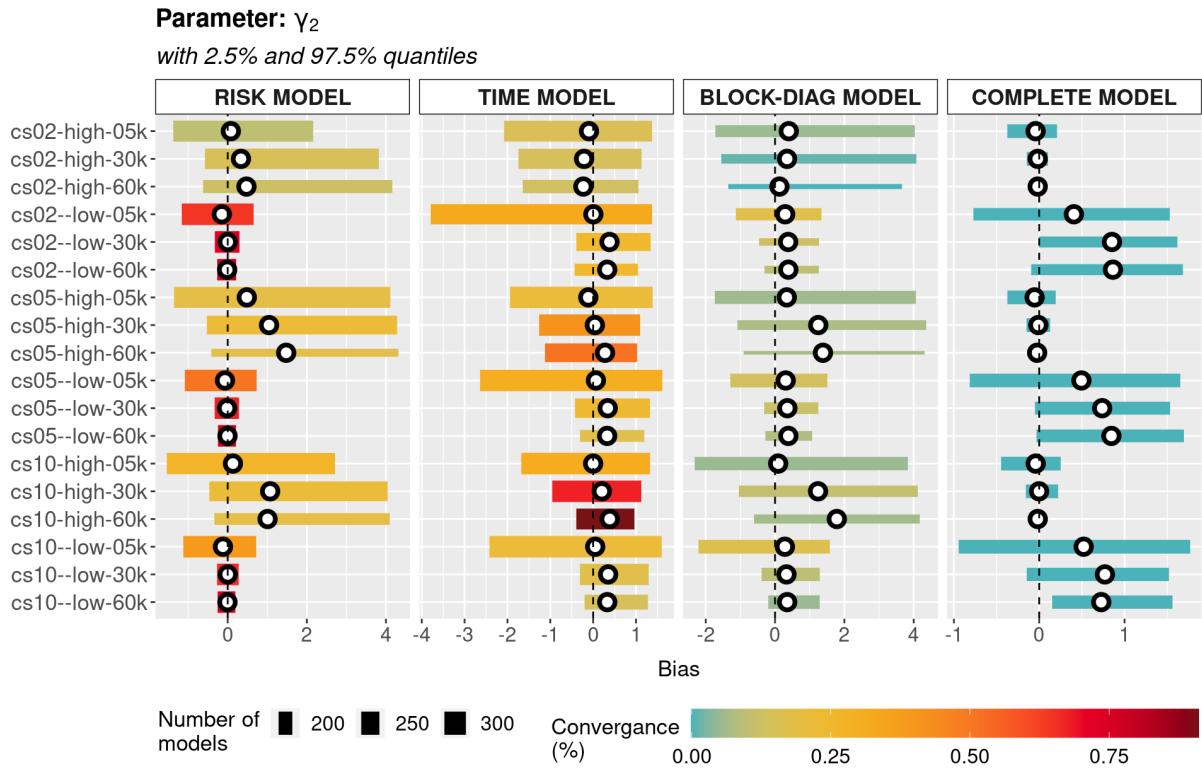
FIGURE 32 – PARAMETER β_2 BIAS WITH 2.5% AND 97.5% QUANTILES



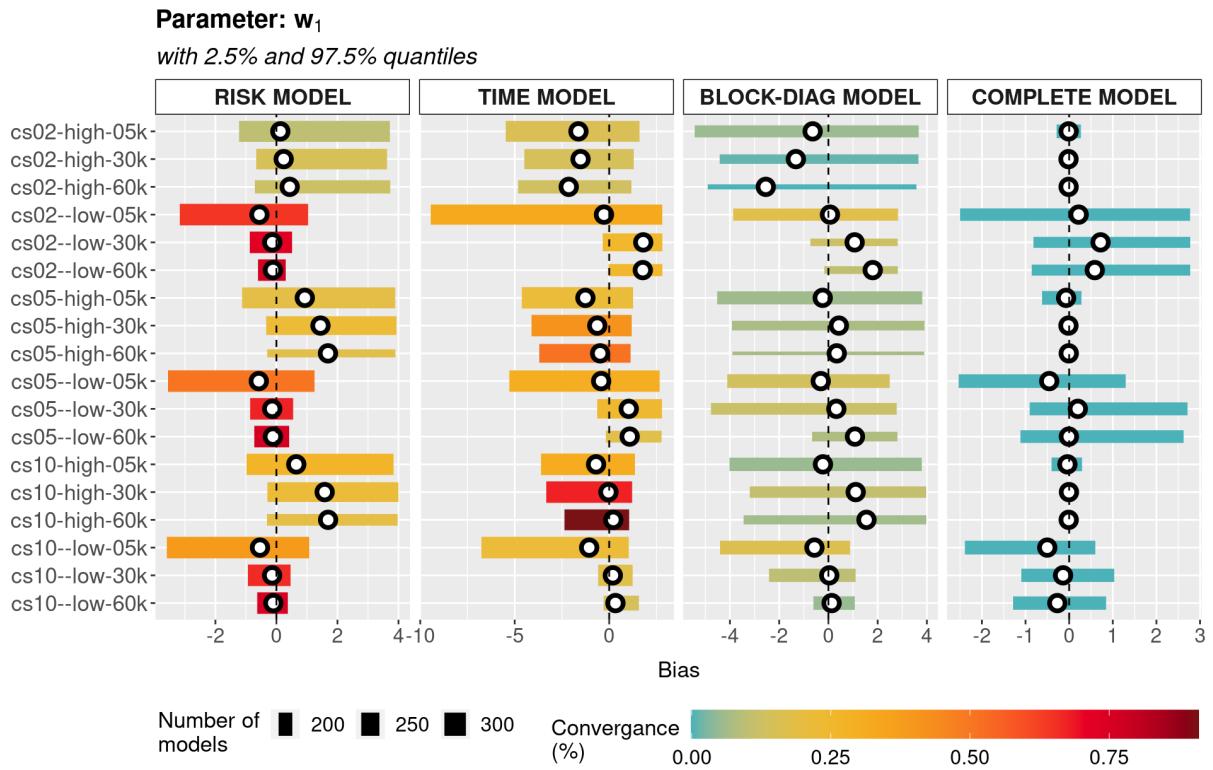
SOURCE: The author (2021).

FIGURE 33 – PARAMETER γ_1 BIAS WITH 2.5% AND 97.5% QUANTILES

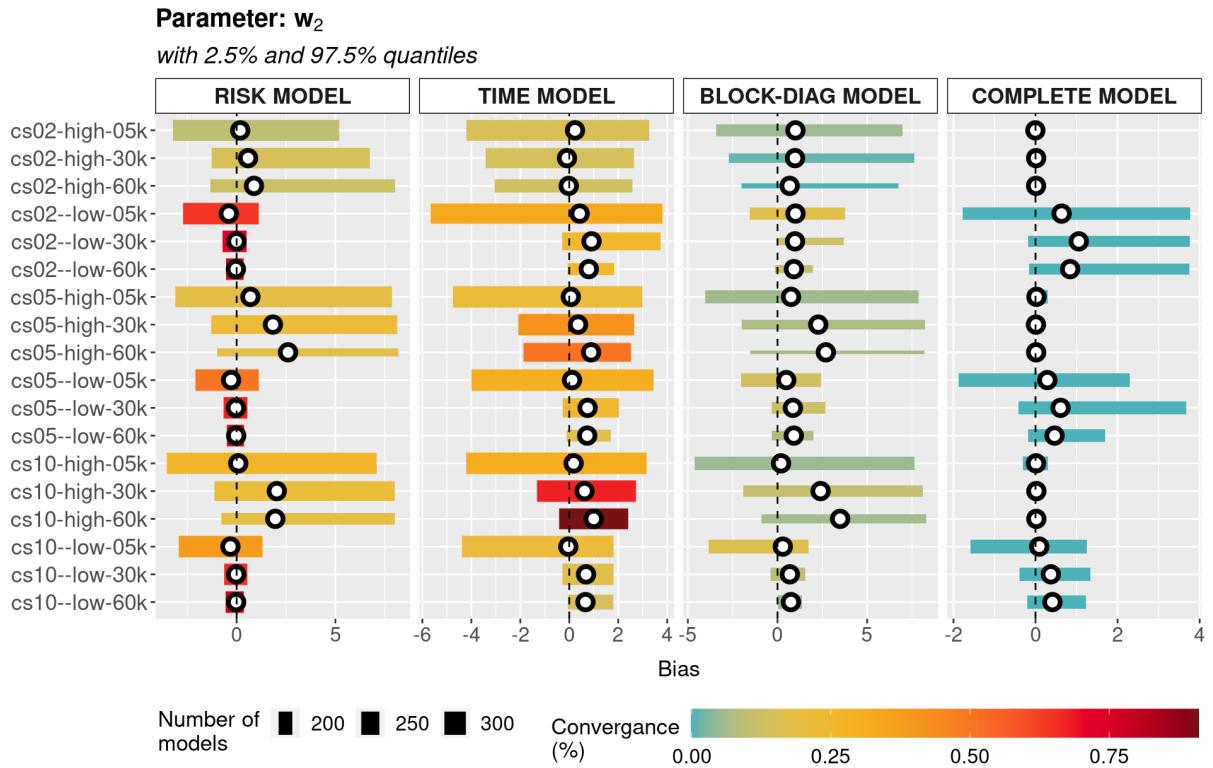
SOURCE: The author (2021).

FIGURE 34 – PARAMETER γ_2 BIAS WITH 2.5% AND 97.5% QUANTILES

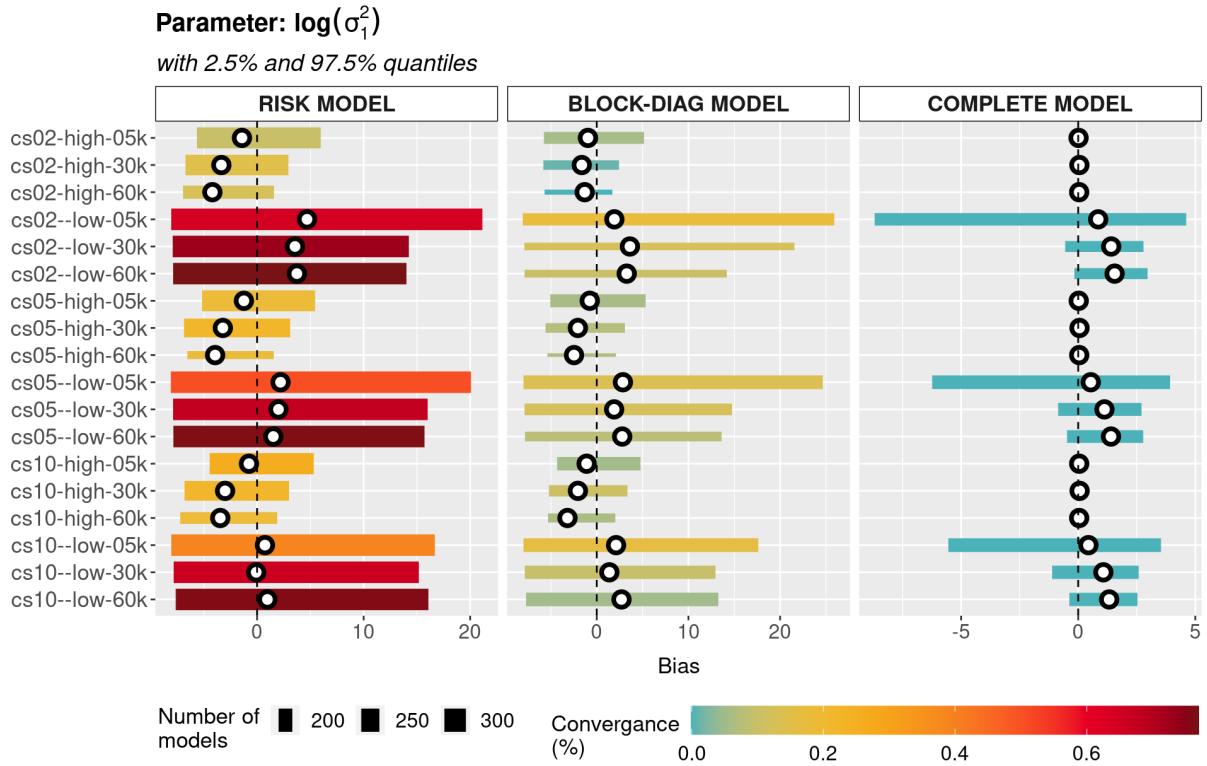
SOURCE: The author (2021).

FIGURE 35 – PARAMETER w_1 BIAS WITH 2.5% AND 97.5% QUANTILES

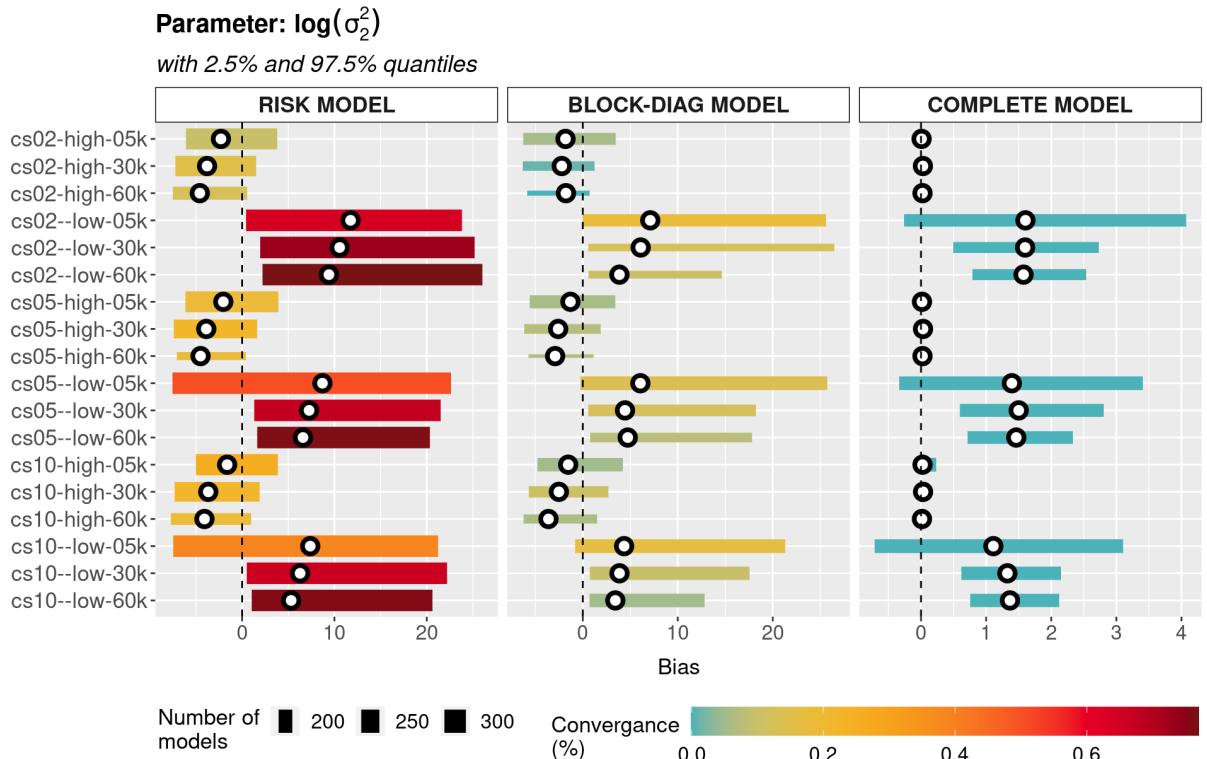
SOURCE: The author (2021).

FIGURE 36 – PARAMETER w_2 BIAS WITH 2.5% AND 97.5% QUANTILES

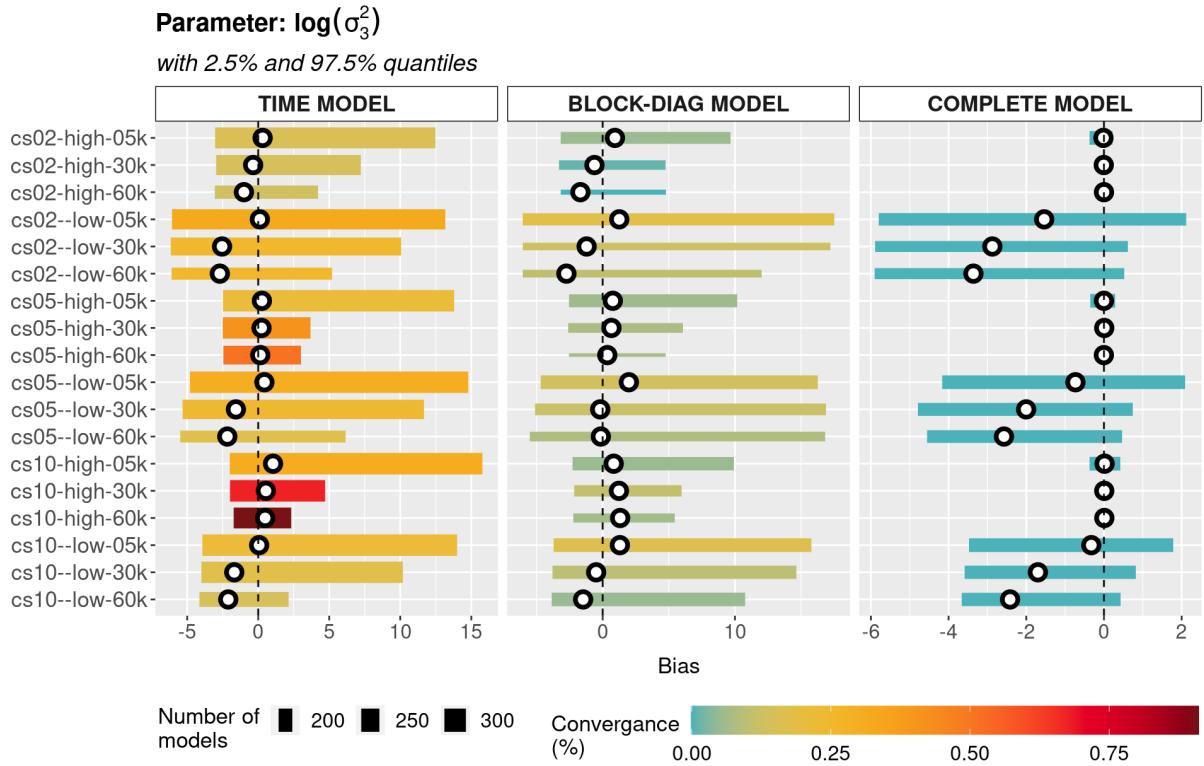
SOURCE: The author (2021).

FIGURE 37 – PARAMETER $\log(\sigma_1^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES

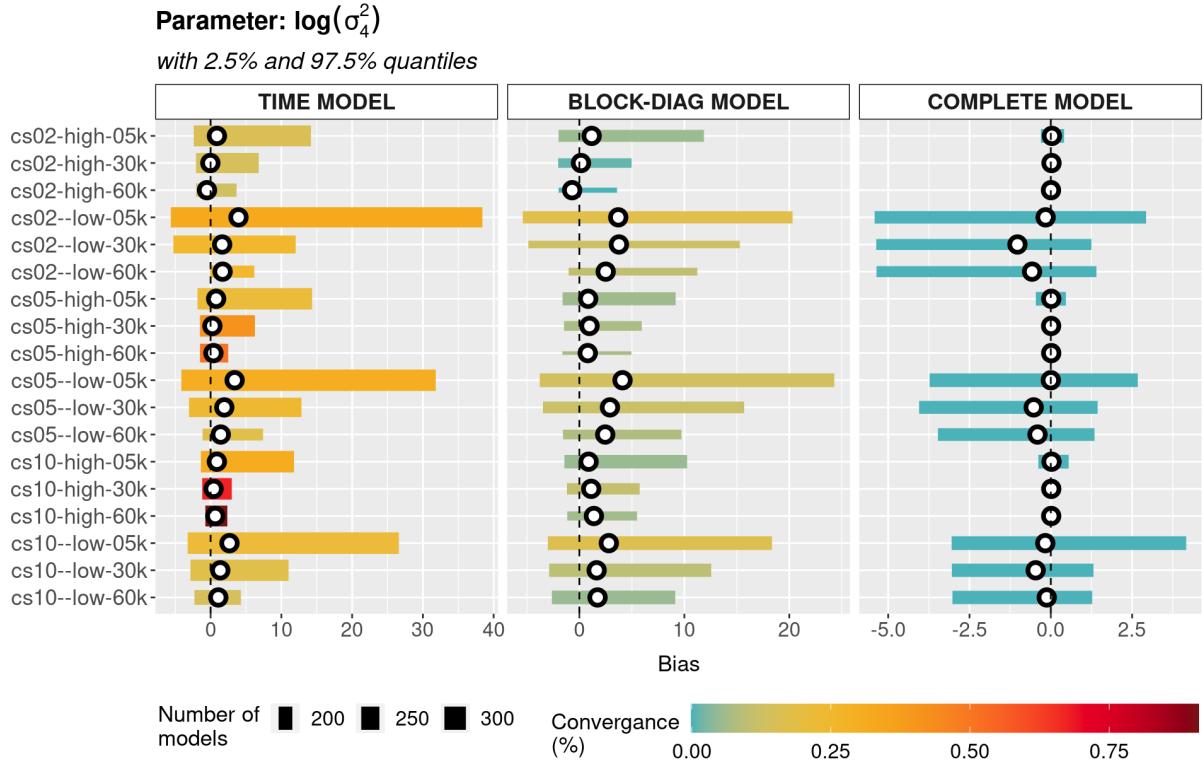
SOURCE: The author (2021).

FIGURE 38 – PARAMETER $\log(\sigma_2^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES

SOURCE: The author (2021).

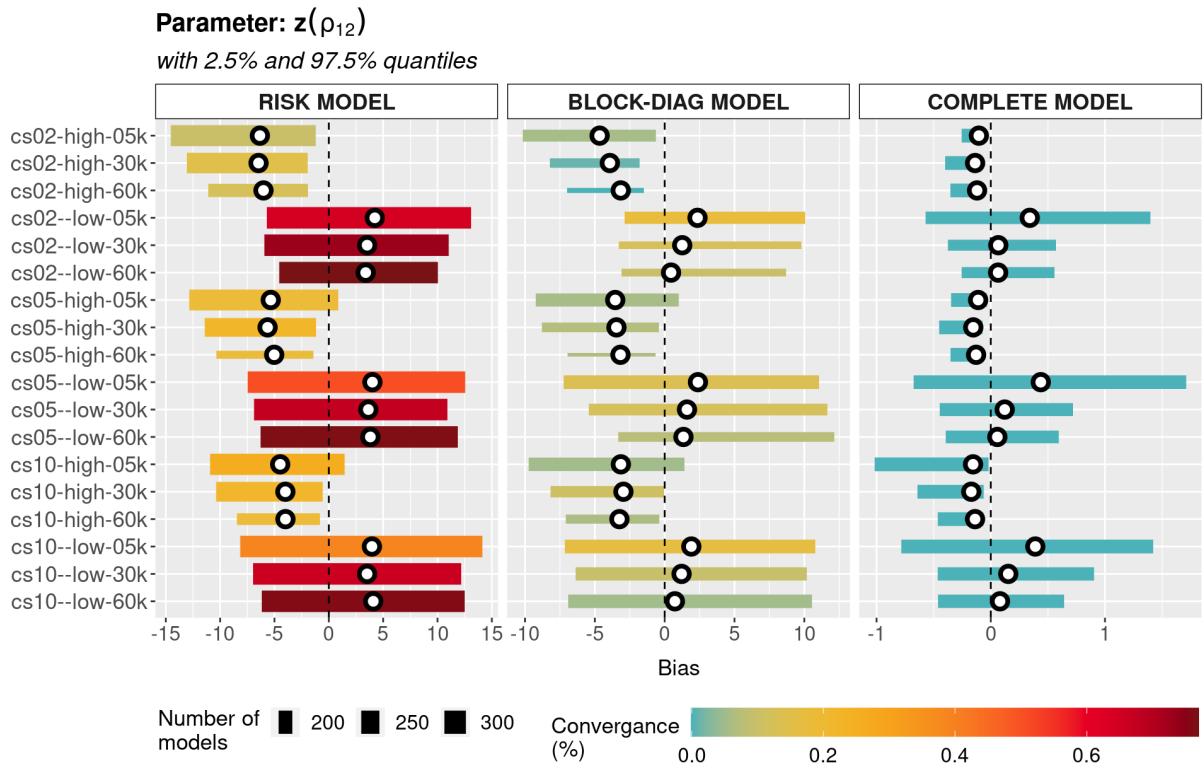
FIGURE 39 – PARAMETER $\log(\sigma_3^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES

SOURCE: The author (2021).

FIGURE 40 – PARAMETER $\log(\sigma_4^2)$ BIAS WITH 2.5% AND 97.5% QUANTILES

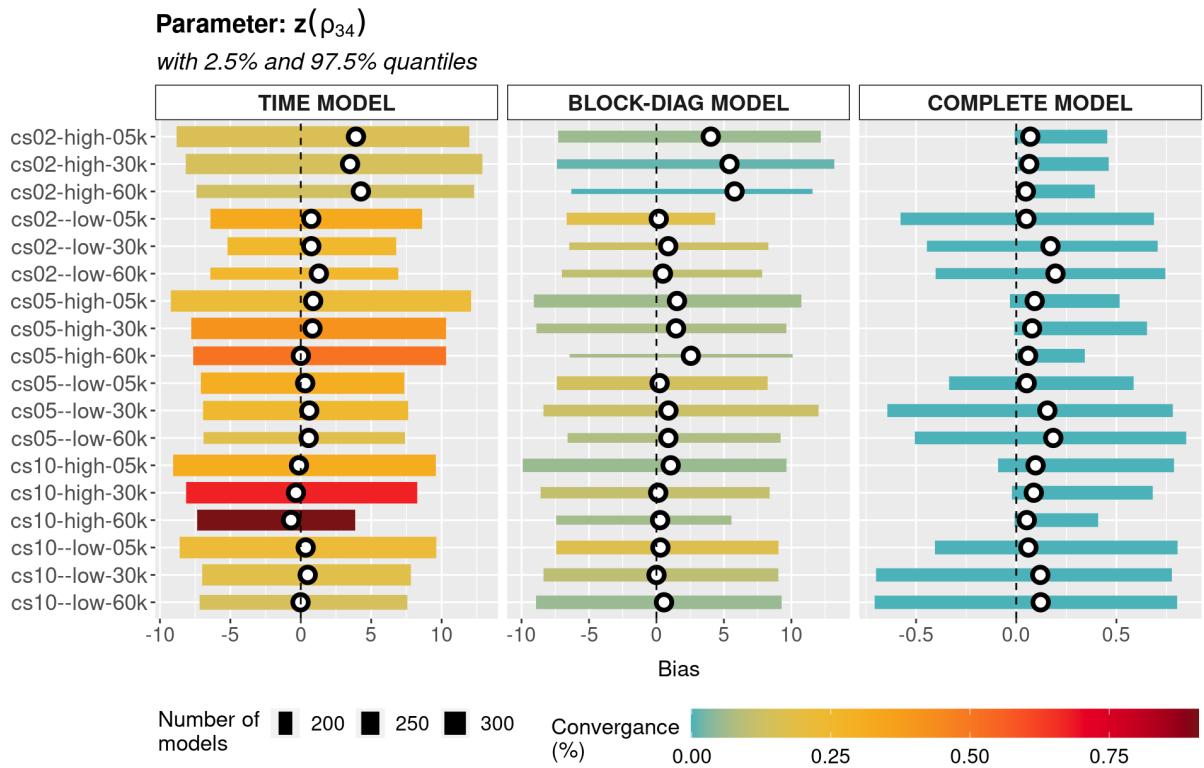
SOURCE: The author (2021).

FIGURE 41 – PARAMETER $z(\rho_{12})$ BIAS WITH 2.5% AND 97.5% QUANTILES



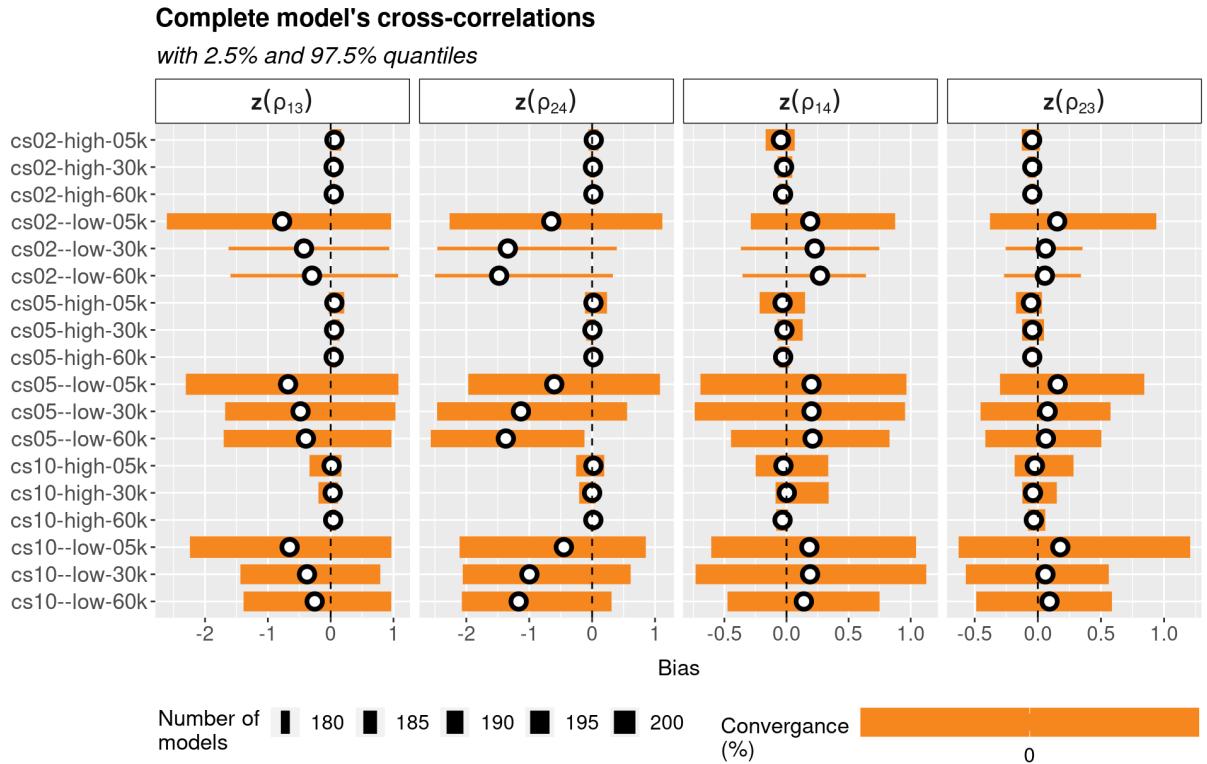
SOURCE: The author (2021).

FIGURE 42 – PARAMETER $z(\rho_{34})$ BIAS WITH 2.5% AND 97.5% QUANTILES



SOURCE: The author (2021).

FIGURE 43 – PARAMETERS $\{z(\rho_{13}), z(\rho_{24}), z(\rho_{14}), z(\rho_{23})\}$ BIAS WITH 2.5% AND 97.5% QUANTILES



SOURCE: The author (2021).