

Modeling the cumulative incidence function of clustered competing risks data: a multinomial GLMM approach

Henrique Aparecido Laureano* Wagner Hugo Bonat*

April 19, 2021

Abstract

Clustered competing risks data are a complex failure time data scheme. Its main characteristics are the cluster structure, which implies a latent within-cluster dependence between its elements, and its multiple variables competing to be the one responsible for the occurrence of an event, the failure. To handle this kind of data, we propose a full likelihood approach, based on a generalized linear mixed model instead a usual complex frailty model. We model the competing causes in the probability scale, in terms of the cumulative incidence function (CIF). A multinomial distribution is assumed for the competing causes and censorship, conditioned on the latent effects. The latent effects are accommodated via a multivariate Gaussian distribution. The CIF is specified as the product of an instantaneous risk level function with a failure time trajectory level function. The estimation procedure is performed through the R package TMB (Template Model Builder), an C++ based framework with efficient Laplace approximation and automatic differentiation routines. A large simulation study is performed, based on different latent structure formulations. The model presents to be of difficult estimation, with our results converging to a latent structure where the risk and failure time trajectory levels are correlated.

Keywords: Clustered competing risks; Within-cluster dependence; Multinomial generalized linear mixed model (GLMM); TMB: Template Model Builder; Laplace approximation; Automatic differentiation (AD).

1 Introduction

Regression models are the main statistical tool for investigating the relationship between a response variable and a set of explanatory variables. When the response is the time until the occurrence of some event, we are in the yard of failure time data with its regression models being called survival models ([Kalbfleisch and Prentice; 2002](#)), with the word *survival* meaning that the modeling focus is the survival experience. We could be interested in the occurrence of an event for multiple causes or even in systematic occurrences, having then multiple survival experiences to be modeled. These possibilities

*Laboratory of Statistics and Goeinformation, Departament of Statistics, Paraná Federal University, Curitiba, Brazil. E-mail: laureano@ufpr.br

imply different failure time data designs. They can be of three classes, as schematized in Figure 1.

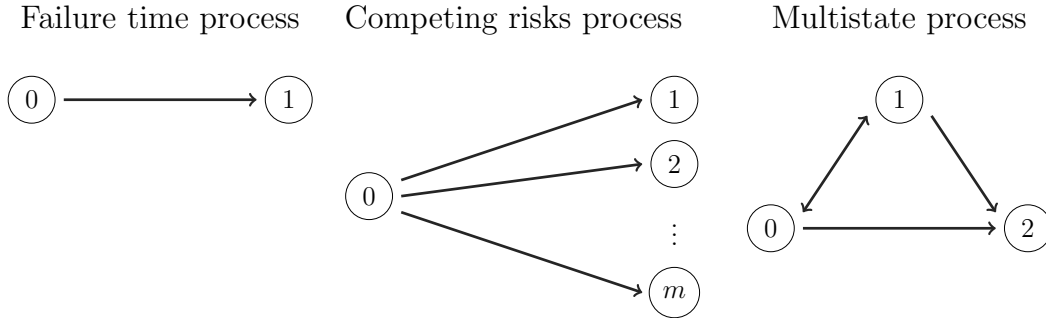


Figure 1: Illustration of failure time designs.

The first two are special cases of the last, a multistate process. The special cases are characterized by the presence of only absorbent states, besides the initial state 0. A multistate process is characterized by the presence of at least one intermediate state. In this work, our focus is on the competing risks process, more specifically, its clustered version. This means that we have groups of elements sharing some non-observed latent dependence structure.

In survival analysis, the survival experiences are usually modeled in the hazard (failure rate) scale, and when the latent within-cluster dependence is accommodated we have the so-called frailty models.

A less usual scale but convenient when dealing with competing risks process is the probability scale.

The class of generalized linear models (GLMs) (Nelder and Wedderburn; 1972) is probably the most popular statistical modelling framework. Despite its flexibility, the GLMs are not suitable for dependent data. In the case of longitudinal data, it is essential that the regression model take into account the longitudinal and/or grouped data structure. According to Diggle et al. (2002) longitudinal data are repeated measures evaluated on the same subjects over time, that are potentially correlated. Dependent data can also arise in studies with block designs, spatial and multilevel data (Verbeke and Molenberghs; 2001; Fitzmaurice et al.; 2008). For the analysis of such data several methods have been proposed over the last four decades.

Laird and Ware (1982) proposed the random effects regression models for longitudinal data analysis. Breslow and Clayton (1993) presented the generalized linear mixed models (GLMMs) for the analysis of non-Gaussian outcomes. Masarotto and Varin (2012) developed a class of marginal models for modelling dependence structures in the analysis of longitudinal data, time series and spatial based on Gaussian copula models.

The main goal of this study is to propose the . In this paper, we will investigate the as an alternative to . R (R Core Team; 2021) package TMB (Kristensen et al.; 2016).

The main contributions of this article are: (i) introducing the unit gamma distribution into the GLMMs framework; (ii) performing an extensive simulation study to check the properties of the maximum likelihood estimator to deal with longitudinal continuous bounded outcomes; (iii) applying the proposed model in two data sets from different fields of application; (iv) providing R code and C++ implementation for the unit gamma mixed models.

The work is organized as follows. Section 2, Section 3, Section 4. Finally, the main contributions of the article are discussed in Section 5.

2 multiGLMM: a multinomial GLMM for clustered competing risks data

Cumulative Incidence Function (CIF)

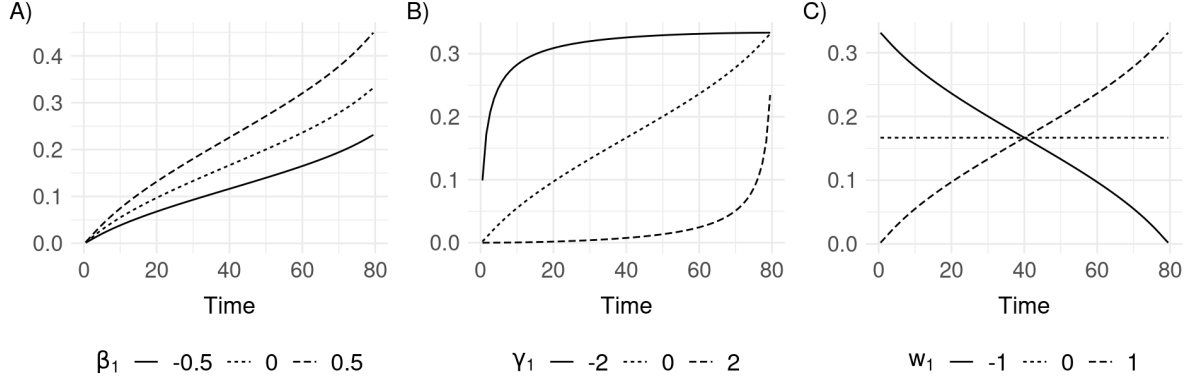


Figure 2: ILLUSTRATION OF COEFFICIENT BEHAVIORS FOR A GIVEN CUMULATIVE INCIDENCE FUNCTION (CIF) (PROPOSED BY SCHEIKE), IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES, AND WITH THE FOLLOWING CONFIGURATION: $\beta_2 = 0$, $u = 0$ AND $\eta = 0$; IN EACH SCENARIO ALL OTHER COEFFICIENTS ARE SET TO ZERO, WITH THE EXCEPTION OF $w_1 = 1$

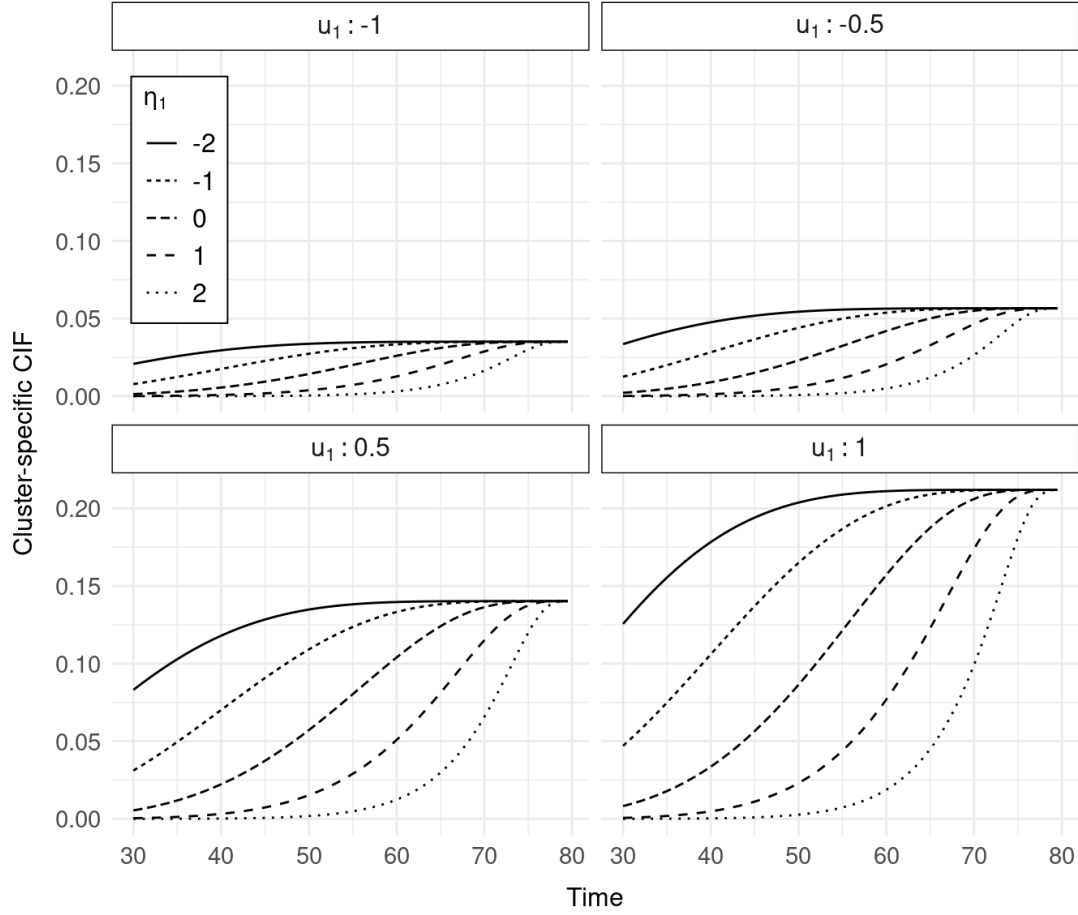


Figure 3: ILLUSTRATION OF A GIVEN CLUSTER-SPECIFIC CUMULATIVE INCIDENCE FUNCTION (CIF), PROPOSED BY SCHEIKE, IN A MODEL WITH TWO COMPETING CAUSES OF FAILURE, WITHOUT COVARIATES AND THE FOLLOWING CONFIGURATION: $\beta_1 = -2$, $\beta_2 = -1$, $\gamma_1 = 1$, $w_1 = 3$ AND $u_2 = 0$. THE VARIATION BETWEEN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

3 Estimation and inference

4 Simulation studies

5 Discussion

Supplementary material

References

- Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models, *Journal of the American statistical Association* **88**(421): 9–25.
- Diggle, P., Heagerty, P., Liang, K. Y. and Zeger, S. (2002). *Analysis of Longitudinal Data*, second Edition edn, Oxford University Press, United Kingdom.

- Fitzmaurice, G., Davidian, M., Verbeke, G. and Molenberghs, G. (2008). *Longitudinal Data Analysis*, CRC Press.
- Kalbfleisch, J. D. and Prentice, R. L. (2002). *The Statistical Analysis of Failure Time Data*, second Edition edn, John Wiley & Sons, Inc., Hoboken, New Jersey.
- Kristensen, K., Nielsen, A., Berg, C. W., Skaug, H. J. and Bell, B. M. (2016). TMB: Automatic Differentiation and Laplace Approximation, *Journal of Statistical Software* **70**(5): 1–21.
- Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data, *Biometrics* **38**(4): 963–974.
- Masarotto, G. and Varin, C. (2012). Gaussian copula marginal regression, *Electronic Journal of Statistics* **6**(1): 1517–1549.
- Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized linear models, *Journal of the Royal Statistical Society, Series A* **135**(3): 370–384.
- R Core Team (2021). *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>.
- Verbeke, G. and Molenberghs, G. (2001). *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, New York.