

**A bivariate joint frailty model with mixture framework for survival analysis of recurrent events with dependent censoring and cure fraction**

**Author**

Tawiah, Richard, McLachlan, Geoffrey J, Ng, Shu Kay

**Published**

2020

**Journal Title**

Biometrics

**Version**

Accepted Manuscript (AM)

**DOI**

[10.1111/biom.13202](https://doi.org/10.1111/biom.13202)

**Rights statement**

© 2020 John Wiley & Sons Ltd. This is the peer reviewed version of the following article: A bivariate joint frailty model with mixture framework for survival analysis of recurrent events with dependent censoring and cure fraction, *Biometrics*, 2020, 76 (3), pp. 753-766, which has been published in final form at <https://doi.org/10.1111/biom.13202>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Wiley or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Wiley's version of record on Wiley Online Library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Wiley Online Library must be prohibited.

**Downloaded from**

<http://hdl.handle.net/10072/390211>

## Funder(s)

ARC

## Grant identifier(s)

DP170100907

## Griffith Research Online

<https://research-repository.griffith.edu.au>



# A bivariate joint frailty model with mixture framework for survival analysis of recurrent events with dependent censoring and cure fraction

**Richard Tawiah**

School of Medicine and Menzies Health Institute Queensland, Griffith University, Nathan QLD 4111, Australia.

School of Psychology, University of New South Wales, Sydney NSW 2052, Australia.

and

**Geoffrey J. McLachlan**

Department of Mathematics, University of Queensland, St. Lucia QLD 4072, Australia.

and

**Shu-Kay Ng**

School of Medicine and Menzies Health Institute Queensland, Griffith University, Nathan QLD 4111, Australia.

*email:* s.ng@griffith.edu.au

**SUMMARY:** In the study of multiple failure time data with recurrent clinical endpoints, the classical independent censoring assumption in survival analysis can be violated when the evolution of the recurrent events is correlated with a censoring mechanism such as death. Moreover, in some situations, a cure fraction appears in the data because a tangible proportion of the study population benefits from treatment and becomes recurrence free and unsusceptible to death related to the disease. A bivariate joint frailty mixture cure model is proposed to allow for dependent censoring and cure fraction in recurrent event data. The latency part of the model consists of two intensity functions for the hazard rates of recurrent events and death, wherein a bivariate frailty is introduced by means of the generalized linear mixed model methodology to adjust for dependent censoring. The model allows covariates and frailties in both the incidence and the latency parts and it further accounts for the possibility of cure after each recurrence. It includes the joint frailty model and other related models as special cases. An EM-type algorithm is developed to provide residual maximum likelihood estimation of model parameters. Through simulation studies, the performance of the model is investigated under different magnitudes of dependent censoring and cure rate. The model is applied to data sets from two colorectal cancer studies to illustrate its practical value.

**KEY WORDS:** bivariate frailty, cure proportion, EM algorithm, informative censoring, joint model, mixture model, random effect, terminal event.

This paper has been submitted for consideration for publication in *Biometrics*

## 1. Introduction

Data on recurrent events such as multiple tumor relapses on the same patient are frequently observed in biomedical studies. In many clinical applications, analysis of such data are of interest, to respond to questions that provide scientific approaches for investigating the relative effect of new therapies on disease progression or examine the prognosis of the disease evolution process. Several statistical techniques, including the Andersen-Gill model, the frailty models and the marginal models are often used to analyze these data. These models are built on a classical assumption in survival analysis that stipulates that the observed recurrent disease (or clinical) events are independent of any phenomenon (e.g. death, loss to follow-up, end of study) that induces right censoring in the data. This assumption is only realistic in situations where the observed recurrent disease events have negligible or no effect on the censoring phenomenon. However, merely considering this assumption can be statistically inappropriate and too restrictive in practice, because dependent censoring (also known as informative censoring) can also arise in many situations. For example, in cardiovascular disease studies in which patients undergo repeated heart attacks, dependent censoring occurs when right censoring is completely or partly caused by a terminal event (e.g. death) that is correlated with the occurrences of heart attacks or in other words, when the occurrences of the heart attacks increase the risk of death.

Due to the advent and the use of effective medicines and advanced technologies in modern healthcare systems, there has been an increasing number of patients who get cured (or who at least survive for a long time) from many chronic diseases such as cancer (López-Cheda et al., 2017). Therefore, the possibility of a tangible fraction of cured patients is sometimes of concern in the analysis of recurrent event data from recent biomedical studies. In recent times, efforts have been made to present statistical models to handle cure fraction in recurrent event data. An example is the class of two-component mixture based cure frailty

models (Rondeau et al., 2011; Tawiah et al., 2019a) which employ the logistic model to accommodate cure fraction and a random effect survival model (i.e. standard frailty model) to model the latency distribution of failure times (e.g. gap or calendar times) in the uncured patient. Although, cure fraction and dependent censoring may possibly arise together, this class of models, however, did not incorporate these two features in tandem. Some studies (Huang and Wolfe 2002; Liu et al., 2004; Huang and Liu 2007; Chen et al., 2012) have made some amount of contribution to propose methods to handle dependent censoring in different kinds of survival data by means of joint frailty models, which employ a shared frailty term in the hazard function for the failure of interest and the censoring phenomenon, to induce their dependence. Instead, Ghosh and Lin (2003) used semiparametric joint models, wherein the marginal distribution of the failure process and dependent censoring time is formulated through scale-change models. Nonetheless, these studies did not consider cure fraction in their proposed methods. Dependent censoring has been considered in conjunction with cure fraction in some few studies published in recent times, see the works of Liu et al. (2016), Liu et al. (2017) and Bernhardt (2016), for example. Liu et al. (2017) and Bernhardt (2016) did not consider recurrent event data, but Liu et al. (2016) did and carried out estimation on an integrated likelihood, obtained by integrating out the frailties by means of the Gaussian quadrature method. The model of Liu et al. (2016) does not allow individuals who have one or more disease events to have a chance of being cured after each occurrence. In some biomedical data sets, patients with long-term censored times following one or more recurrences can be observed. A typical case is the data set from the colorectal cancer hospital readmission study (Gonzalez et al., 2005); about 9.7% of the patients had at least one recurrence followed by censored failure times that exceed 1000 days, approximately 2.7 years. In this regard, it is important to allow flexibility in how cure fraction is modelled in recurrent event data, in order to consider the possibility of cure after each recurrence. Furthermore,

their model employed a common frailty term that is shared by the hazard functions for recurrent disease events and death to account for dependent censoring in the joint intensity function of the uncured patients. The use of a shared frailty can be a limitation because it forces the frailties to be the same for the failure times of recurrent events and death, which may not fit well in some applications. Moreover, the shared frailty only permits positive associations. However, in some situations negative associations exist among paired survival times (Xue and Brookmeyer 1996). As an example, suppose patients may be admitted into a hospital several times for the same disease. The association between hospital stay and death can be negative since a long hospital stay implies a severe episode for a patient, who is more likely to die shortly after discharge.

Motivated by these concerns, our present study aims to develop a more general mixture cure model that considers a bivariate frailty, viz two possibly correlated random effects in the latency component, to jointly model the dependence between the hazard rates of recurrent events and death in the proportion of uncured patients. The model, further accounts for frailty effect in the incidence (logistic) part, that explains the patient-specific unobservable characteristics (frailties) that affect the cure probability. Moreover, the model considers all observations in the incidence and the latency components and it allows patients with one or more events to have a chance of being cured after each event. This approach increases statistical power (Rondeau et al., 2011), compared to the alternative procedure that restrict the incidence part to the first event and ignore subsequent ones.

To present an estimation procedure for our model, we use the residual maximum likelihood (REML) method (McGilchrist 1993). This method arises from the spirit of the generalized linear mixed model (GLMM) methodology, in which the fixed effect and the random effect terms are estimated alike, by maximizing a BLUP-type log-likelihood, analogous to the maximum penalized partial likelihood approach (Ripatti and Palmgren 2000), while the

variance component parameters are obtained through estimating equations formed from the first order derivative of the REML log-likelihood. For semiparametric survival models that involve cure fraction, the application of the GLMM method is not straightforward, because the underlying complete-data likelihood involves missing data and also, the unspecified baseline survival function of the uncured patients does not cancel out from the likelihood. Because the EM algorithm is very suitable for dealing with incompleteness or missing data problems (Ng 2013), we adopt it to implement the GLMM method for estimation of our model. To examine the performance of the model and the estimation procedure, simulation studies are conducted in a small sample setting. Also, we illustrate the practical importance of the model by analyzing two biomedical data sets from a colorectal cancer hospital readmission study and a metastatic colorectal cancer clinical trial.

## 2. Notation, models and likelihoods

### 2.1 Notation of data

Consider a longitudinal follow-up study that measures recurrent clinical events from  $M$  independent patients. We assume that the observation of the recurrent event process is subject to right censoring, possibly due to death, loss to follow-up or end of the study (i.e. administrative censoring). Define  $R_{jk}$  as the gap time of the  $k$ th recurrent event on the  $j$ th patient,  $D_j$  is the time to death from the last recurrent event and  $C_j$  is the censoring time due to lost to follow-up or end of study. Here, we assume that  $C_j$  is independent of both  $R_{jk}$  and  $D_j$ . We suppose the follow-up time is  $T_{jk} = \min(R_{jk}, D_j, C_j)$ , with binary indicator for recurrent events denoted by  $\delta_{jk}^R$ , which takes on a value of 1 if  $T_{jk} = R_{jk}$  and 0 if  $T_{jk} = D_j$  or  $T_{jk} = C_j$ . Let  $T_j$  denote the time between the last recurrent event and death or censoring time, which represents the last follow-up time for patient  $j$ , defined as  $T_j = \min(D_j, C_j)$ . Here, we note,  $\delta_j^D$  is the binary indicator for death, which is 1 when  $T_j = D_j$  and 0, if censored. Notice that, the data we observe on the  $j$ th patient is  $O_j =$

$\{(t_{jk}, \delta_{jk}^R, \delta_j^D, x_j), j = 1, \dots, M, k = 1, \dots, n_j\}$ , where  $t_{jk}$  consists of recurrent gap times  $t_{jk}^R$  and death time  $t_j^D$  which are both subject to censoring,  $n_j$  is the number of recurrent events experienced by the  $j$ th patient and  $x_j = (x_{j1}, \dots, x_{jp})^T$  is a  $p$  dimensional covariate vector on the  $j$ th patient, that may include some treatment or prognostic variables, e.g. chemotherapy, sex, age, BMI and comorbidity. The superscript  $T$  denotes a vector transpose and  $\sum_{j=1}^M n_j = N$  is the total number of observations in the data.

For notations relevant for the description of cure, let  $Y_{jk}$  denote a binary indicator, that is 1 for any patient who experiences the  $k$ th recurrent event (i.e. an uncured patient) and 0 otherwise (i.e. a cured patient). Denote by  $Y_j$  an indicator for death related to the disease, that is 1 when a patient dies and 0, otherwise (i.e. a long-term survivor). We assume that cured patients can neither experience recurrent events nor death due to the disease. Note that  $Y_{jk} = 1$  when  $\delta_{jk}^R = 1$  and also  $Y_j = 1$  when  $\delta_j^D = 1$ . However, when the censoring indicators  $\delta_{jk}^R$  and  $\delta_j^D$  are zeros,  $Y_{jk}$  and  $Y_j$  are not observed. To examine the survival behaviour of a study population that consist of a subgroup of uncured patients who are susceptible to disease relapses and death, and a subgroup of non-susceptible patients, we consider the joint frailty mixture cure model.

## 2.2 The bivariate joint frailty mixture cure model

Suppose in our data described in Section 2.1, follow-up has been sufficiently long and the observed right censoring times are due to death, loss to follow-up and also a heavy magnitude of administrative censoring (i.e. end of study censoring time). Furthermore, we assume that death times are possibly correlated with the gap times between recurrent events and also the presence of a heavy magnitude of administrative censoring, given a sufficiently long follow-up pre-supposes that there may be a non-negligible fraction of patients who are potentially cured and thus not susceptible to recurrent disease events and death related to recurrence of the disease. Among the patients who are not cured, that is, those who are susceptible to



recurrent events and death, their gap times and death times are not only affected by the observed covariates but also their own frailties, e.g. individual lifestyle, genetic predisposition, etc. The mixture cure model for the marginal population may be written as

$$S(t_{jk}) = 1 - \pi_k(x_j) + \pi_k(x_j) S_u(t_{jk}; x_j), \quad j = 1, \dots, M \quad \text{and} \quad k = 1, \dots, n_j, \quad (1)$$

where  $S(t_{jk})$  is a bivariate survival function for recurrent events and death for the marginal population,  $\pi_k(x_j) = P(Y_{jk} = 1)$  is the probability of experiencing the  $k$ th recurrent event after the  $(k-1)$ th event,  $1 - \pi_k(x_j)$  is the probability of being cured after the  $k$ th event and  $S_u(t_{jk}; x_j)$  is the conditional bivariate survival function for recurrent events and death in the proportion of uncured patients.

When a patient is cured the bivariate survival function for recurrent events and death degenerates to 1. The probability of being uncured  $\pi_k(x_j)$  is potentially related to the covariate vector  $x_j$  through the logistic regression model. For recurrent event data, the logistic model may be fitted to the first observations within each individual or all observations available on each individual (Rondeau et al., 2011; Tawiah et al., 2019a). In the case of the later, multiple observations within an individual may be potentially correlated. An unobservable random effect term may be useful to adjust for the correlation among observations within patients. In principle with the GLMM framework, the model may be specified as

$$\pi_k(x_j) = \exp(\xi_{jk}) / [1 + \exp(\xi_{jk})], \quad \xi_{jk} = w_j^T \alpha + u_j, \quad (2)$$

where  $\xi_{jk}$  is the linear predictor,  $\alpha$  is the fixed effect vector that measures the effects of covariates  $x_j$  on the uncured probability,  $w_j = (1 \ x_j^T)^T$ , and  $u_j$  is the frailty term affecting the uncured probability. Let  $u = (u_1, \dots, u_M)^T$  denote the realization of  $u_j$  on  $M$  patients. We assume that the realizations of  $u$  are independent and identically distributed from the normal distribution  $N(0, \theta_u^2 I_M)$ , where  $I_M$  is an  $M$  dimensional identity matrix. The variance component  $\theta_u^2$  is a measure for unobserved heterogeneity in the probability of the proportion of cured patients.

Notice that the latency function  $S_u(t_{jk}; x_j)$  is a bivariate function that depends on  $t_{jk}^R$  and  $t_j^D$ . For the proportional hazards (PH) cure model, the latency part may be modelled through a survival function (Peng and Taylor 2016) or a hazard function (Yau and Ng 2001; Lai and Yau 2008, 2009). We follow the hazard based approach to model the hazard rate for recurrent events and death in the uncured patients. We adopt two separate PH models which are correlated through unobservable patient-specific random effects  $v_j$  and  $V_j$ , which represent the frailties affecting the hazard rates of disease recurrence and death on the  $j$ th uncured patient, respectively. Given the random effects  $v_j$  and  $V_j$  and covariate vector  $x_j$ , the hazards function for recurrent events  $h_u(t_{jk}^R; x_j)$  and death  $\lambda_u(t_j^D; x_j)$  for the  $j$ th uncured patients are defined, respectively,

$$\begin{aligned} h_u(t_{jk}^R; x_j) &= h_{u0}(t_{jk}^R) \exp(\eta_{jk}); \quad \eta_{jk} = x_j^T \beta + v_j, \\ \lambda_u(t_j^D; x_j) &= \lambda_{u0}(t_j^D) \exp(\zeta_j); \quad \zeta_j = x_j^T \gamma + V_j, \end{aligned} \quad (3)$$

where  $\eta_{jk}$  and  $\zeta_j$  are the linear predictors corresponding to the two hazard models,  $h_{u0}(t_{jk}^R)$  and  $\lambda_{u0}(t_j^D)$ , are the baseline hazard functions for recurrent events and death in the uncured patients, which are specified non-parametrically,  $\beta$  is the fixed effect parameter vector that relates  $x_j$  to the hazard rate for recurrent events and  $\gamma$  is the fixed effect parameter vector that measures the effect of  $x_j$  on the hazard for death. The first model in Equation (3) is a frailty model for multivariate survival data (McGilchrist, 1993). It assumes that frailty arises from unobserved heterogeneity and intra-subject correlation of recurrent event times. The second model conceptualizes frailty in terms of unobserved covariates in univariate (independent) failure time data to allow for individual differences in mortality hazard rate (Hougaard 1995). The framework of Equations (1), (2) and (3) constitute the bivariate joint frailty mixture cure model proposed in this work. Equation (3) denotes a bivariate joint frailty model. Like the frailty terms, the covariate vector  $x_j$  may not be necessarily the same for the hazard rates of recurrent events and death, likewise the covariates of the uncured probability in the logistic part.

Let  $v = (v_1, \dots, v_M)^T$  and  $V = (V_1, \dots, V_M)^T$  denote the random vectors of  $v_j$  and  $V_j$  such that  $q = (v^T, V^T)^T$ . We assume that  $q$  follows a bivariate normal distribution  $BVN(0, \Sigma)$ , where the variance covariance matrix  $\Sigma = \Gamma \otimes I_M$ . Notice that  $\otimes$  is a Kronecker product and  $\Gamma$  is defined in the following

$$\Gamma = \begin{bmatrix} \theta_v^2 & \rho\theta_v\theta_V \\ \rho\theta_v\theta_V & \theta_V^2 \end{bmatrix},$$

where  $\rho$  is a correlation parameter that incorporates dependence between the hazard rate for recurrent events and death. If, for example,  $\rho = 0$ , the hazard rate for death is independent of the hazard rate for disease recurrences. Also,  $\theta_v^2$  and  $\theta_V^2$  measure heterogeneity in the frailty for the hazard rates of disease recurrence and death. Notice that the frailties in the model are (bivariate) log-normal since the exponential function of the random effects on the normal distribution are log-normal. Other distributions, such as the gamma distribution can also be used; see, for example, Wienke et al. (2003).

Let  $W_{(N \times (p+1))}$ ,  $X_{1(N \times p)}$ ,  $X_{2(M \times p)}$ ,  $R_{(N \times M)}$ ,  $Z_{1(N \times M)}$  and  $Z_{2(M \times M)}$  denote the design matrices corresponding to  $\alpha, \beta, \gamma, u, v$  and  $V$ , respectively. From here onwards, we write the design matrices without their dimensions. The linear predictors,  $\xi_{jk}$ ,  $\eta_{jk}$  and  $\zeta_j$  may be rewritten in terms of the design matrices and the vectors of the fixed and the random effects, as follows

$$\xi = W\alpha + Ru, \quad \eta = X_1\beta + Z_1v \quad \text{and} \quad \zeta = X_2\gamma + Z_2V.$$

To formulate a likelihood function for the model, we assume that a recurrent event and death cannot be observed at the same time point within the same patient. Considering the bivariate joint frailty cure model in Equations (1), (2) and (3), and the conditional independence of recurrent events and death within patient  $j$  given  $v_j$  and  $V_j$ , as well as  $u_j$ , the complete-data likelihood can be expressed as  $L^C = L_1(\alpha, u_j; y_{jk}, y_j) L_2(\beta, \gamma, H_{u0}, \Lambda_{u0}, v_j, V_j; y_{jk}, y_j)$ , where  $L_1(\alpha, u_j; y_{jk}, y_j)$  is the likelihood contribution on the basis of the parameters of the random effect logistic regression model and  $L_2(\beta, \gamma, H_{u0}, \Lambda_{u0}, v_j, V_j; y_{jk}, y_j)$  is that of the bivariate

joint frailty model. For brevity and conciseness, we simply write  $L_1$  and  $L_2$  to denote these two likelihood components. Notice that  $L_2$  can be written as

$$L_2 = \{h_{u0}(t_{jk}^R) \exp(\eta_{jk})\}^{\delta_{jk}y_{jk}} \exp\{-y_{jk}H_{u0}(t_{jk}^R) \exp(\eta_{jk})\} \\ \times \{\lambda_{u0}(t_j^D) \exp(\zeta_j)\}^{\delta_j y_j} \exp\{-y_j \Lambda_{u0}(t_j^D) \exp(\zeta_j)\}. \quad (4)$$

The first line of expressions in Equation (4) corresponds to the likelihood of the recurrent event gap times and the second line is the likelihood of death times. In standard frailty models, where the non-informative censoring assumption is considered, the first line of expressions are only used. If cure fraction is not assumed, then  $y_{jk} = 1$  and  $y_j = 1$ , for all  $j$  and  $k$  and the likelihood  $L_2$  reduces to that of a bivariate joint frailty model. In that case, Equation (4) is analogous to the derivation given by Huang and Wolfe (2002) for the joint-shared frailty model. The joint likelihood  $L_2$  raises concerns about issues pertaining to identifiability because in survival analysis the usual independent censoring assumption is required in order to avert an identifiability problem (Ebrahimi et al., 2003). Nevertheless, we note in passing that  $L_2$  is practically identifiable once the data provide sufficient information that distinguishes one censoring mechanism from the other, for example as discussed in the case of the data notation  $(t_{jk}, \delta_{jk}^R, \delta_j^D, x_j)$  in Section 2.1. Furthermore, with the use of random effects, the issue of identifiability on the lines of  $L_2$  is not a problem, since independent censoring assumption is not displaced in the definition of the likelihood and the bivariate joint frailty model. By this we mean that, the recurrent event times and the death censoring times are assumed to be independent, conditional on the random effects  $v$  and  $V$ , while dependent censoring is captured by means of their covariance structure  $\Sigma$ . However, identifiability issues may arise within the framework of mixture cure models, due to improper separation of the mixture (Yau and Ng 2001; Peng and Taylor 2016; Ng et al., 2019). It is because it is sometimes unclear whether a censored individual is cured or follow-up has not been pursued long enough for the event to occur in such individual. In practice,

it is necessary to consider the context where cure models are appropriate (e.g. long-term follow-up, non-zero marginal survival probability/heavy long-term censored times).

Denote by  $R_r$  and  $D_i$ , the respective recurrent event gap/censoring time and the death/censoring time which are reordered in an increasing order of magnitude, with their corresponding reordered linear predictors  $\eta_r$  and  $\zeta_i$ . Following an argument in Klein (1992) and Sy and Taylor (2000) on the lines of the profile likelihood construction for standard PH model, it can be shown that replacing  $H_{u0}(t_{jk}^R)$  and  $\Lambda_{u0}(t_j^D)$  by the Aalen-Nelson estimator (Breslow 1972), Equation (4) reduces to a joint partial likelihood of  $\beta$  and  $\gamma$  in terms of  $\eta$  and  $\zeta$ , with the random effects  $v$  and  $V$  conditionally fixed. Letting  $l_2$  to denote the logarithm of  $L_2$ , it follows that

$$l_2 = \sum_{r=1}^K \left\{ \eta_k - \log \sum_{l \in R(k)} y_l \exp(\eta_l) \right\} + \sum_{i=1}^P \left\{ \zeta_i - \log \sum_{s \in R(i)} y_s \exp(\zeta_s) \right\}, \quad (5)$$

where  $R(k)$  and  $R(i)$  are the risk sets corresponding to the gap times of the recurrent events and the death times, respectively. Also,  $t_1 < \dots < t_K$  and  $t_1 < \dots < t_P$  denote the respective distinct reordered gap times and death times. Similarly, let  $l_1 = \log(L_1)$ , then

$$l_1 = \sum_{j=1}^M \sum_{k=1}^N \{y_{jk} \log \pi_k(x_j) + (1 - y_{jk}) \log(1 - \pi_k(x_j))\} = \sum_{r=1}^N \{y_r \xi_r - \log(1 + \exp(\xi_r))\}. \quad (6)$$

Within the realm of the GLMM framework, the logarithm of the joint density function of  $u$  and  $q$ , may be expressed as

$$l_3 = -\frac{1}{2} \left\{ M \log(2\pi\theta_u^2) + \frac{1}{\theta_u^2} u^T u \right\} - \frac{1}{2} \{ M \log(2\pi|\Sigma|) + q^T \Sigma^{-1} q \}, \quad (7)$$

where  $u$  is independent of  $q$ . From the concept of the GLMM framework (McGilchrist 1993), the complete-data BLUP-type log-likelihood for the random effect logistic model and the bivariate joint frailty model can be presented as

$$l_\Phi = l_1 - \frac{1}{2} \left\{ M \log(2\pi\theta_u^2) + \frac{1}{\theta_u^2} u^T u \right\} \quad (8)$$

and

$$l_{\Omega} = l_2 - \frac{1}{2} \{M \log(2\pi|\Sigma|) + q^T \Sigma^{-1} q\}, \quad (9)$$

where  $\Phi = (\alpha, u)$  and  $\Omega = (\beta, \gamma, q)$ .

### 3. Estimation procedure

The MLE and the REML are the most commonly used estimators within the GLMM methodology. Previous research has established that the REML method performs satisfactorily well in terms of bias of fixed effect and variance component estimates, compared to the ML method (McGilchrist 1993; Yau 2001; Tawiah et al., 2019b). Therefore, we use the REML method to estimate the unknown parameters of the proposed bivariate joint frailty mixture cure model. As presented in Equations (8) and (9), the proposed model yields the complete-data log-likelihood. Accordingly, it is natural to apply the EM algorithm to achieve maximization of the REML estimation. Note that the log-likelihoods involve two binary indicators  $y_{jk}$  and  $y_j$  which are partially missing. The E-step completes the missing data on  $y_{jk}$  and  $y_j$  through the computation of the expectation of the complete-data log-likelihood conditional on the observed data and the current values of the parameters  $\Phi = (\alpha, u)$ ,  $\Omega = (\beta, \gamma, q)$ ,  $\Psi = (\theta_u^2, \theta_v^2, \theta_V^2, \rho)$ ,  $S_{u0}(t_{jk}^R)$  and  $S_{u0}(t_j^D)$ . The M-step, on the other hand, maximizes this conditional expectation of the complete-data log-likelihood over the unknown parameters. Intuitively, the E-step reduces to the conditional expectation of  $y_{jk}$  and  $y_j$ , calculated to be

$$g_{jk} = E \{y_{jk} | (\alpha, \beta, u, v, S_{u0}(t_{jk}^R))\} = \delta_{jk}^R + \frac{(1 - \delta_{jk}^R) \pi_k(x_j) S_{u0}(t_{jk}^R)^{\exp(\eta_{jk})}}{1 - \pi_k(x_j) + \pi_k(x_j) S_{u0}(t_{jk}^R)^{\exp(\eta_{jk})}} \quad (10)$$

and

$$g_j = E \{y_j | (\alpha, \gamma, u, V, S_{u0}(t_j^D))\} = \delta_j^D + \frac{(1 - \delta_j^D) \pi(x_j) S_{u0}(t_j^D)^{\exp(\zeta_j)}}{1 - \pi(x_j) + \pi(x_j) S_{u0}(t_j^D)^{\exp(\zeta_j)}}, \quad (11)$$

which respectively update  $y_{jk}$  and  $y_j$ , in  $l_{\Phi}$  and  $l_{\Omega}$ , where  $g_{jk}$  is the posterior probability of the  $j$ th patient being uncured, given the  $k$ th possible recurrent event, and  $g_j$  is the posterior

probability of the  $j$ th patient experiencing a terminal event. At convergence, the estimates of  $1 - g_{jk}$  and  $1 - g_j$  can respectively be interpreted as the cure and the longer-term survival rate for the  $j$ th patient. The terms  $S_{u0}(t_{jk}^R)$  and  $S_{u0}(t_j^D)$  are respectively defined as the conditional baseline survival function for recurrent event and death in the uncured patients. As in previous work (Sy and Taylor 2000; Lai and Yau 2008), we adopt the Breslow-type estimator to estimate  $S_{u0}(t_{jk}^R)$  and  $S_{u0}(t_j^D)$  and their tails are smoothly approximated to zero by means of the ETAIL completion method (Peng 2003).

Let  $\Phi_0$ ,  $\Omega_0$  and  $\Psi_0$  denote the vectors of the initial values of the parameter vectors  $\Phi$ ,  $\Omega$  and  $\Psi$ , respectively. Given the current estimates of  $g_{jk}$ ,  $g_j$ ,  $\Phi$ ,  $\Omega$  and  $\Psi$ , the M-step maximizes  $l_\Phi$  and  $l_\Omega$  through the Newton-Raphson iterative methods, viz

$$\begin{bmatrix} \hat{\alpha} \\ \hat{u} \end{bmatrix} = \begin{bmatrix} \alpha_0 \\ u_0 \end{bmatrix} + G^{-1} \begin{bmatrix} \partial l_\Phi / \partial \alpha \\ \partial l_\Phi / \partial u \end{bmatrix} \quad (12)$$

and

$$\begin{bmatrix} \hat{\beta} \\ \hat{\gamma} \\ \hat{v} \\ \hat{V} \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \alpha_0 \\ v_0 \\ V_0 \end{bmatrix} + H^{-1} \begin{bmatrix} \partial l_\Omega / \partial \beta \\ \partial l_\Omega / \partial \gamma \\ \partial l_\Omega / \partial v \\ \partial l_\Omega / \partial V \end{bmatrix}, \quad (13)$$

respectively, where  $G^{-1}$  and  $H^{-1}$  are the inverse of the information matrices corresponding to  $l_\Phi$  and  $l_\Omega$  and

$$\frac{\partial l_\Phi}{\partial \alpha} = W^T \frac{\partial l_1}{\partial \xi}; \quad \frac{\partial l_\Phi}{\partial u} = R^T \frac{\partial l_1}{\partial \xi} - \theta_u^{-2} u; \quad \frac{\partial l_\Omega}{\partial \beta} = X_1^T \frac{\partial l_2}{\partial \eta}; \quad \frac{\partial l_\Omega}{\partial \gamma} = X_2^T \frac{\partial l_2}{\partial \zeta};$$

$$\frac{\partial l_\Omega}{\partial v} = Z_1^T \frac{\partial l_2}{\partial \eta} - \frac{v\theta_v^2 - V\rho\theta_v\theta_V}{\theta_v^2\theta_V^2(1-\rho^2)} \quad \text{and} \quad \frac{\partial l_\Omega}{\partial V} = Z_2^T \frac{\partial l_2}{\partial \zeta} - \frac{V\theta_v^2 - v\rho\theta_v\theta_V}{\theta_v^2\theta_V^2(1-\rho^2)}.$$

For the derivation of  $\partial l_1 / \partial \xi$ ,  $\partial l_2 / \partial \eta$  and  $\partial l_2 / \partial \zeta$ , see Web Appendix A. Letting  $D_{\eta\eta} = -\partial^2 l_2 / \partial \eta \partial \eta^T$ ,  $D_{\zeta\zeta} = -\partial^2 l_2 / \partial \zeta \partial \zeta^T$  and  $D_{\eta\zeta} = \partial^2 l_2 / \partial \eta \partial \zeta^T = D_{\zeta\eta} = \partial^2 l_2 / \partial \zeta \partial \eta^T = 0$  and following the derivation given in the Web Appendix A, matrices  $G$  and  $H$  may be written

as blocks, in the form

$$G = \begin{bmatrix} G_{\beta,\beta} & G_{\beta,u} \\ G_{u,\beta} & G_{u,u} \end{bmatrix} \quad \text{where, } G^{-1} = \begin{bmatrix} A_{\alpha,\alpha} & A_{\alpha,u} \\ A_{u,\alpha} & A_{u,u} \end{bmatrix}$$

and

$$H = \begin{bmatrix} H_{\beta,\beta} & H_{\beta,\gamma} & H_{\beta,q} \\ H_{\gamma,\beta} & H_{\gamma,\gamma} & H_{\gamma,q} \\ H_{q,\beta} & H_{q,\gamma} & H_{q,q} \end{bmatrix} \quad \text{where, } H^{-1} = \begin{bmatrix} B_{\beta,\beta} & B_{\beta,\gamma} & B_{\beta,q} \\ B_{\gamma,\beta} & B_{\gamma,\gamma} & B_{\gamma,q} \\ B_{q,\beta} & B_{q,\gamma} & B_{q,q} \end{bmatrix}.$$

At convergence the Newton-Raphson iterative methods in Equations (12) and (13) provide estimates of the random effects which can be interpreted as the individual fragility of the patients. Prediction intervals of the frailties can also be obtained from the underlying information matrix using empirical Bayes (EB) or the conditional mean squared error of prediction (CMSEP) method to estimate the variances of the random effects (Tawiah et al., 2019b).

To this end, the simplification of the system of equations of the first order derivative of the REML log-likelihood (B.1), given in Web Appendix B, yields the respective REML estimators of  $\theta_v^2$ ,  $\theta_V^2$  and  $\rho$ , expressed in the following

$$\hat{\theta}_v^2 = \frac{1}{M} \mathfrak{S}_1, \quad \hat{\theta}_V^2 = \frac{1}{M} \mathfrak{S}_3, \quad \text{and } \hat{\rho} = \frac{1}{\sqrt{\mathfrak{S}_1 \mathfrak{S}_3}} \mathfrak{S}_2, \quad (14)$$

where  $\mathfrak{S}_1 = \text{tr} \{K_1 (B_{q,q} + qq^T)\}$ ,  $\mathfrak{S}_2 = \text{tr} \{K_2 (B_{q,q} + qq^T)\} / 2$ ,  $\mathfrak{S}_3 = \text{tr} \{K_3 (B_{q,q} + qq^T)\}$  and  $\text{tr}$  denotes the trace of a matrix. The block matrices  $K_1$ ,  $K_2$  and  $K_3$  are defined in Web Appendix B. Similarly, the REML estimator of  $\theta_u^2$  is given by

$$\hat{\theta}_u^2 = \frac{1}{M} \{ \text{tr} (A_{u,u}) + u^T u \}. \quad (15)$$

A summary of the computational iterative routine for the foregoing estimation procedure is given in Web Appendix C. The computational procedure is initialized by setting the starting values of  $\Phi$  and  $\Omega$  to zero and those of  $\Psi$  to relatively small values. The method alternates between the E-step and the M-step, then to the REML estimation of the variance components and it recycles until convergence. Once convergence is obtained, the underlying observed



complete-data information matrices  $G$  and  $H$  are calculated. The asymptotic standard errors of the fixed effect parameters  $\alpha$ ,  $\beta$  and  $\gamma$  are respectively calculated from  $A_{\alpha,\alpha}$ ,  $B_{\beta,\beta}$  and  $B_{\gamma,\gamma}$ , while those of the variance components are obtained by inverting the REML information matrices. Details for the later can be found in Web Appendix B.

#### 4. Simulation studies

In this section, we provide simulation studies to investigate the performance of the proposed model and the estimation procedure. In the simulation design, we assumed a sample of 400 patients who are followed up to a maximum of 2000 days, approximately 5.5 years from the beginning of the study. We considered a treatment variable  $x_{j1}$  as the covariate appearing in the linear predictors of the components of the proposed joint-mixture cure frailty model. The covariate  $x_{j1}$  is generated from the Bernoulli distribution with a 0.5 chance of randomization for each patient. The recurrent time-to-event data  $(t_{jk}, \delta_{jk}^R, \delta_j^D)$  characterized by dependent terminal event and cure fraction were simulated as follows:

- i. We first generate the independent censoring time  $c_j$  for each patient using the uniform distribution with a maximum of 2000 days.
- ii. Next, we generate  $u_j$  from  $N(0, \theta_u^2)$  and by means of the Bernoulli distribution, the binary indicators  $y_{jk}$  and  $y_j$  are generated for each patient, with the probability of being uncured and that of death, given by the random effect logistic regression model in Equation (2).
- iii. When a patient is cured ( $y_{jk} = 0$ ), we assign the time of recurrence  $t_{jk}^R$  to  $\infty$ . Likewise, if a patient has no terminal event ( $y_j = 0$ ), we assign the death time  $t_j^D$  to  $\infty$ . Here, whenever  $y_{jk} = 0$  we restrict  $y_j = 0$ .
- iv. We generate  $q_j = (v_j^T, V_j^T)^T$  from  $BVN(0, \Sigma_j)$ . When a patient is not cured ( $y_{jk} = 1$ ), we generate the time of recurrence  $t_{jk}^R$  repeatedly from the frailty model  $h_u(t_{jk}^R; x_{j1}) = h_{u0}(t_{jk}^R) \exp(\eta_{jk})$ ;  $\eta_{jk} = x_{j1}^T \beta + v_j$  until the corresponding uncured patient is either cured, dies or until  $\sum_{k=1}^{n_j} t_{jk}^R \geq c_j$ . Similarly, when a patient dies, the death time  $t_j^D$  is

thus generated from the frailty model  $\lambda_u(t_j^D; x_{j1}) = \lambda_{u0}(t_j^D) \exp(\zeta_j)$ ;  $\zeta_j = x_{j1}^T \gamma + V_j$ . In both settings, we assumed that the functions of the baseline hazards follow the Weibull distribution  $\mu \tau t^{\tau-1}$ , with  $\mu = 0.0005$ ,  $\tau = 1.8$  for recurrent event gap times and  $\mu = 0.0003$ ,  $\tau = 1.5$  for the death times.

- v. The observed failure time  $t_{jk}$  is then generated from  $\min(t_{jk}^R, t_j^D, c_j)$ , and  $\delta_{jk}^R$  and  $\delta_j^D$  are obtained according to their definitions given in Section 2.1.

The settings of the true parameter values were considered as follows:  $\alpha_0 = 1.0, 0.3$ ;  $\alpha_1 = -0.5, -0.8$ ;  $\beta_1 = -0.8, -0.6$ ;  $\gamma_1 = -0.5, -0.3$ ;  $\theta_u^2 = 0.5$ ;  $\theta_v^2 = 0.3, 0.5$ ;  $\theta_V^2 = 0.5$  and  $\rho = 0.2, 0.4, 0.6, 0.8, -0.2, -0.8$ . The first set of true parameter values for  $\alpha_0$  and  $\alpha_1$  were chosen to achieve a moderately low cure rate of 37.8%, while the second set gives a moderately high cure rate of 62.2% in the treatment group. Moreover, the varying settings of  $\rho$  were used to validate the model in different magnitudes of dependent censoring, that is, relatively low, moderate and high positive correlations as well as negative correlations. Several simulation scenarios were considered, in each case 500 data sets were generated and the EM REML procedure detailed in Section 3 was applied for estimation. We present in Tables 1 and S1 in Web Appendix D, a summary of the estimation results given in terms of average bias (Ave. bias), average of the standard error estimates (SEE), the standard error (SE) of the parameter estimates and the coverage probability (CP) of 95% confidence interval (CI) based on the normal approximation. For the simulation scenarios with low cure rate (Table 1), the fixed effect parameter  $\alpha_1$  is empirically unbiased under all the settings of dependent censoring, while  $\alpha_0$ ,  $\beta_1$  and  $\gamma_1$  have small biases. Similarly, no substantial bias is observed for  $\alpha_0$ ,  $\alpha_1$ ,  $\beta_1$  and  $\gamma_1$  in the simulation with high cure rate (Table S1) considering all the magnitudes of dependent censoring. However, when cure rate is high and the magnitude of  $\rho$  and  $\theta_v^2$  increases simultaneously,  $\beta_1$  has slightly positive bias, though the magnitude of the bias of  $\gamma_1$  reduces (Simulation 6 in Table S1). Comparing all the simulation sets, it is seen that the size of bias of  $\gamma_1$  decreases sufficiently when the correlation moves from a low level to a higher

order. This is particularly apparent when the cure rate decreases. The fixed effect parameters and the variance components have acceptable biases when  $\rho$  goes to zero in both the low and high cure rate. The accuracy of the standard error estimates of the fixed effect parameters is acceptable, comparing SEE and SE. Under the positive correlation schemes, the CPs of  $\alpha_0$ ,  $\alpha_1$  and  $\beta_1$  are close to the nominal level, while those of  $\gamma_1$  are slightly below the nominal level for both low and high cure rate. However, it is noticed that the CP of  $\gamma_1$  has the tendency of drawing closer to the nominal level when it decreases in bias. The CP of  $\alpha_0$  shrinks slightly from the nominal level when dependent censoring is of higher magnitude in the negative direction. This is possibly due to the increment in bias in the estimates of  $\alpha_0$  under strongly negative correlation scheme, as observed in both low and high cure rates. For the variance component parameters, we found that the bias is reasonably small. Nevertheless, the bias of  $\theta_V^2$  increases towards null when the correlation goes to a higher order in the positive direction for both the low and high cure rates. The standard errors of  $\theta_v^2$ ,  $\theta_V^2$  and  $\rho$  are quite well estimated, comparing SEE and SE. Nonetheless, in the case of  $\theta_u^2$  the estimates of SEE are noticeably larger than those of SE and its CP is above the nominal level in all the simulation settings. In this regard, the estimator (B.5) (see Web Appendix B) may overestimate the standard errors of this variance parameter. As a remark, it is not advisable to directly apply the standard error estimates of the variance components to conduct a test of statistical significance, even in the case where the standard errors are accurately estimated (Tawiah et al., 2019a). The reason is that the use of the normal approximation of the null test statistic is inappropriate because the null hypothesis of the variance components lies on the boundary of parameter space (Vaida and Xu 2000; Lai and Yau 2009). Some formal tests such as the score test of homogeneity, the likelihood ratio test with corrected null distribution or random effect inference based on EB prediction interval may be used to highlight the significance of random effects (Vaida and Xu 2000; Tawiah et al., 2019b). In this paper, we consider the EB

prediction interval approach, which tests, for example,  $H_0 : v_j = 0$  rather than  $H_0 : \theta_v^2 = 0$ . For the correlation coefficient, we provided a test for  $H_0 : \rho = 0$  versus  $H_0 : \rho \neq 0$  based on the conventional approach using standard error estimates since the null hypothesis is not affected by the boundary condition related to the restricted parameter space. Another issue that worth noting is that, the use of the normal approximation for the CIs of the variance component parameters could possibly affect the accuracy of their CPs because the estimates of these parameters have right-skewed distributions. Further simulation studies are outlined to examine the robustness of the proposed model to misspecification of the normal assumption underlying the distribution of the random effects; see Table S2 and Web Appendix D for results and discussion.

[Table 1 about here.]

## 5. Application

### 5.1 Hospital readmissions and death in a colorectal cancer study

We analyze the data set from a colorectal cancer hospital readmission study (Gonzalez et al., 2005), which involved 403 patients who received surgery to remove tumors after being diagnosed with colorectal cancer. During follow-up, some patients encountered several hospital readmissions related to the colorectal cancer. The first admission time was defined as the time between the date of surgery and the first hospital admission date. The subsequent readmission times were considered as the difference between the last date of discharge and the current date of hospitalization. In total, 861 readmissions were observed, about 200 patients had no recurrence at all and 112 patients died during follow-up. The maximum number of readmission was 22, with mean 2.25 and median 1.0. The times to hospital readmission and death are both important outcome measures. The study provides data on some variables which may potentially affect the readmission and the death times. These

include chemotherapy, gender, tumor stage measured by Dukes classification (i.e. A-B, C and D) and comorbidity measured by Charlson index (i.e. 0, 1-2,  $\geq 3$ ).

Figure 1(a)-(c) present the Kaplan-Meier (KM) survival curves of the gap times for the first three successive readmissions, classified by treatment arms. It appears that the recurrence free rate is higher in the chemotherapy arm than the control arm for the first readmission, while in the case of the second and the third readmissions the recurrence free rate are statistically equivalent for both treatment arms. The right tails of the curves for both arms are quite stable. While the control arm stays away from zero survival probability for the three successive readmissions, the arm for chemotherapy drops to zero. The zero-survival probability in the tails of the chemotherapy arm is due to the fact that the largest gap times in this arm are uncensored for the first, second and the third readmissions. Therefore, it would be unrealistic to rule out the possibility of a cure fraction in the data. In Figure 1(d) the KM curve stratified by treatment arm is illustrated for the death times following the last readmission. We see that the tails of the curves for both arms level-off above 0.6. Empirically, this feature suggests that a fraction of long-term survivors exists in the data. This could be a reflection of patients who may have been cured and therefore not susceptible to death related to the disease. Also, dependent censoring may appear in these data because hospital readmissions and death observed on the same patient are likely to be correlated. This dependence feature may be due to patient specific unobserved frailties. Hence, a frailty model accommodating cure fraction and dependent terminal event (i.e. death) in the presence of recurrent events (i.e. readmissions) should be considered for modelling these data.

[Figure 1 about here.]

We apply the proposed model to these data. Again we consider two reduced models, that is, a joint frailty model that do not incorporates a cure fraction and a standard frailty model to allow comparison of results. Table 2 gives a summary of the results. For all the three

models, the reference group consists of male patients who did not receive chemotherapy, had Dukes stage A-B and Charlson comorbidity index 0. From the proposed model, it is seen that chemotherapy has negative coefficient estimates in the logistic and the hazard components. The corresponding  $p$ -values suggest that chemotherapy minimizes the hazard rate of readmission in the uncured patients, but it leads to a small, insignificant in curing the disease as well as minimizing the hazard rate of death in the uncured patients. Female gender has significant negative estimates in all the components of the model. From the logistic component, the results clearly depict that female patients are substantially more likely to be cured than males. On the hazard part, the uncured female patients have significantly longer times to be readmitted and significantly longer survival rate (i.e. time to death) after readmission compared to the male patients who are uncured. Dukes stage C and D decrease the probability of being cured and also increase the hazard rate of readmission and death. Higher comorbidity (i.e. Charlson index  $\geq 3$ ) leads to a significant lower cure probability and significant increasing hazard rate for hospital readmission and death.

[Table 2 about here.]

Figure 2 displays the point estimates of the random effects  $u$ ,  $v$  and  $V$  and their 95% EB prediction interval. The figure provides inferences concerning heterogeneity of the individual frailties among the 403 patients, where patients are arranged in ascending order of the number of events encountered. For the frailty terms  $u$  and  $V$ , no significant heterogeneity is to be expected since their CIs do not depart from zero. Patients 397, 402 and 403 would have higher risk of recurrence because the CIs of their individual frailties  $v_j$ 's are significantly higher than zero, though in general, most of the interval lengths of the CIs of  $v$  involves zero, thus also suggesting a small heterogeneity in the hazard rate for hospital readmission. As confirmed by the results from the model (Table 2), the estimates of the variance components in the hazard part,  $\theta_u^2$  and  $\theta_v^2$ , and that in the logistic part  $\theta_u^2$  are mild. The correlation  $\rho$

is significant and moderate in magnitude, indicating that there is a meaningful dependence between the hazard rate of readmission and death. However, when the correlation is ignored (i.e.  $\hat{\rho} = 0$ ), the estimates of the regression coefficients, the variance components and their standard errors do not change noticeably (see Table S3 in Web Appendix E). In the analysis of the second data example on metastatic colorectal cancer (Ducreux, et al., 2011) presented in Web Appendix E,  $\rho$  is significant and  $\hat{\rho} = 0.963$  (Table S4). Neglecting this correlation yields different regression estimates (see Web Appendix E and Table S5 for further discussion).

From Table 2, the results from the reduced models are comparable with those from the proposed model, although some slight differences can also be seen. In particular, in the proposed model chemotherapy has a negative coefficient estimate ( $\gamma_1 = -0.009$ ) on the hazard rate for death, but it is positive ( $\gamma_1 = 0.264$ ) in the joint frailty model, although the estimate is not significant in both models. The reduced models depict that Dukes stage D has significant increasing effect on the hazard rate of readmission, but it is insignificant in the proposed model. The estimate of the frailty variance parameters  $\theta_v^2$ ,  $\theta_V^2$  and  $\rho$  in the joint frailty model are considerably larger than those in the proposed. These differences could be due to the inclusion of cure fraction in the proposed model because the adjustment for this feature may explain some aspects of unobserved heterogeneity across patients. The proposed joint frailty mixture cure model and the joint frailty model were fitted to the data set using author written R codes and the standard frailty model was fitted in R using the `coxme` package (Therneau 2018). The implementation of the package is based on the penalized partial likelihood estimation (Ripatti and Palmgren 2000) for log-normal frailty model. The current version of the package does not provide the standard errors of the frailty variance parameters. Estimation of the joint frailty model is based on the REML method (McGilchrist 1993; Yau 2001), without EM implementation.

[Figure 2 about here.]

## 6. Discussion

In this paper, a semiparametric joint frailty mixture cure model has been developed to deal concurrently with dependent censoring and cure fraction in recurrent event data. Differing from the use of shared frailty, we propose to use the bivariate frailty approach, by this way the failure time for the recurrent events and the censoring mechanism are not constrained to have a common frailty. This approach provides means of dealing with situations with either positive or negative dependence structure and it overcomes the limitations of the shared frailty approach that provides one parameter to model correlation and variance. One essential advantage of the model is that, it yields clinically important information beyond those that can be observed from standard frailty models and cure models as well as the recently developed frailty cure models. These models arise as special cases of the proposed model.

Our proposed EM-based REML estimation contrasts with the Markov chain Monte Carlo (MCMC) EM procedure (e.g. see, Huang and Wolfe 2002; Liu et al., 2004) that is used when the integrals in the E-step are not available in closed form. The MCMC EM method is, however, computationally intensive and it leads to issues related to the assessment of the convergence of the MCMC algorithm (Abrahantes and Burzykowski 2005). A potential alternative is the numerical integration method, but this approach is not feasible when the dimension of the random effect goes to a higher order (Vaida and Xu 2000). The relative advantage of the proposed EM REML estimation is that it circumvents these complications. This method allows for flexibility in modeling additional correlation structures, for instance, the correlation of  $u_j$  with  $v_j$  or  $V_j$  can be considered as a topic for further research. Another advantage of the proposed model is its ability to predict individual fragility of patients. This information is valuable to identify high-risk individuals with poorer outcomes. The analysis



of the data example shows that the presence of cure fraction can explain some aspect of between-patient heterogeneity in the data, consistent with the results of Liu et al. (2016).

The probability of experiencing the  $k$ th recurrent event ( $\pi_k(x_k)$ ) can depend on the frequency of previous recurrent events by means of a time-varying covariate. However, the absence of such covariate highlights the necessity of random patient effect in the logistic part. Other practical problems where mixture modeling of dependent censoring can be useful include the two-component survival mixture model proposed by Ng et al. (2004).

#### ACKNOWLEDGEMENTS

The authors are grateful to the reviewers and an associate editor for their helpful comments. Richard Tawiah gratefully acknowledges the financial support from Griffith University International Postgraduate Research Scholarship and Griffith University Postgraduate Research Scholarship (HTH). This work was supported by the Australian Research Council.

#### REFERENCES

- Abrahantes, J. C., and Burzykowski, T. (2005). A version of the EM algorithm for proportional hazard model with random effects. *Biometrical Journal* **47**, 847–862.
- Bernhardt, P. W. (2016). A flexible cure rate model with dependent censoring and a known cure threshold. *Statistics in Medicine* **35**, 4607–4623.
- Breslow, N. E. (1972). Contribution to the discussion of D. R. Cox (1972). *Journal of the Royal Statistical Society, Series B* **34**, 216–217.
- Chen, C-M., Lu, T-F. C., Chen, M-H., and Hsu, C-M. (2012). Semiparametric transformation models for current status data with informative censoring. *Biometrical Journal* **54**, 641–656.
- Ducreux, M., Malka, D., Mendiboure, J., Etienne, P. L., Texereau, P., Auby, D., et al., (2011). Sequential versus combination chemotherapy for the treatment of advanced colorectal

- cancer (FFCD 2000-05): an open-label, randomised, phase 3 trial. *The Lancet Oncology* **12**, 1032–1034.
- Ebrahimi, N., Molefe, D. and Ying, Z. (2003). Identifiability and censored data. *Biometrika* **90**, 724–727.
- Ghosh, D., and Lin, D. Y. (2003). Semiparametric analysis of recurrent events data in the presence of dependent censoring. *Biometrics* **59**, 877–885.
- Gonzalez, J. R., Fernandez, E., Moreno, V., Ribes, J., Peris, M., Navarro, M., et al., (2005). Sex differences in hospital readmission among colorectal cancer patients. *Journal of Epidemiology and Community Health* **59**, 506–511.
- Hougaard, P. (1995). Frailty models for survival data. *Lifetime Data Analysis* **1**, 255–273.
- Huang, X., and Wolfe, R. A. (2002). A frailty model for informative censoring. *Biometrics* **58**, 510–520.
- Huang, X., and Liu, L. (2007). A joint frailty model for survival and gap times between recurrent events. *Biometrics* **63**, 389–397.
- Klein, J. P. (1992). Semiparametric estimation of random effects using the Cox model based on the EM algorithm. *Biometrics* **48**, 795–806.
- Lai, X., and Yau, K. K. W. (2009). Multilevel mixture cure models with random effects. *Biometrical Journal* **51**, 456–466.
- Lai, X. and Yau, K. K. W. (2008). Long-term survivor model with bivariate random effects: applications to bone marrow transplant and carcinoma study data. *Statistics in Medicine* **27**, 5692–5708.
- Liu, L., Wolfe, R. A., and Huang, X. (2004). Shared frailty models for recurrent events and a terminal event. *Biometrics* **60**, 747–756.
- Liu, L., Huang, X., Yaroshinsky, A., and Cormier, J. N. (2016). Joint frailty models for zero-inflated recurrent events in the presence of a terminal event. *Biometrics* **72**, 204–214.

- Liu, Y., Hu T., and Sun, J. (2017). Regression analysis of current status data in the presence of a cured subgroup and dependent censoring. *Lifetime Data Analysis* **23**, 626–650.
- Lopóz-Cheda, A., Cao, R., Jácome, M. A., and Keilegom, I. V. (2017). Nonparametric incidence estimation and bootstrap bandwidth selection in mixture cure models. *Computational Statistics and Data Analysis* **105**, 144–165.
- McGilchrist, C. A. (1993). REML estimation for survival models with frailty. *Biometrics* **49**, 221–225.
- Ng, S-K. (2013). Recent developments in expectation-maximization methods for analyzing complex data. *Wiley Interdisciplinary Reviews: Computational Statistics* **5**, 415–431.
- Ng, S.K., McLachlan, G.J., Yau, K.K.W. and Lee, A.H. (2004). Modelling the distribution of ischaemic stroke-specific survival time using an EM-based mixture approach with random effect adjustment. *Statistics in Medicine* **23**, 2729–2744.
- Ng, S.K., Xiang, L. and Yau, K.K.W. (2019). *Mixture Modelling for Medical and Health Sciences*. Boca Raton, FL: Chapman & Hall/CRC Press.
- Peng, Y., and Taylor, J. M. G. *Cure models*. In: Klein, J. P., Van Houwelingen, H. C., Ibrahim, J. G., and Scheike, T. H., eds. *Handbook of survival analysis*. CRC Press; 2016.
- Peng, Y. (2003). Estimating baseline distribution in proportional hazards cure models. *Computational Statistics & Data Analysis* **42**, 187–201.
- Ripatti, S. and Palmgren, J. (2000). Estimation of multivariate frailty models using penalized partial likelihood. *Biometrics* **56**, 1016–1022.
- Rondeau, V., Schaffner, E., Corbière, F., Gonzalez, J. R., and Mathoulin-Pélissier, S. (2011). Cure frailty models for survival data: application to recurrences for breast cancer and to hospital readmissions for colorectal cancer. *Statistical Methods in Medical Research* **22**, 243–260.
- Sy, J. P., and Taylor, J. M. G., (2000). Estimation in a Cox proportional hazards cure model.

*Biometrics* **56**, 227–236.

Tawiah, R., McLachlan, G. J., and Ng, S. K. (2019a). Mixture cure models with time-varying and multilevel frailties for recurrent event data. *Statistical Methods in Medical Research*

**DOI: 10.1177/0962280219859377.**

Tawiah, R., Yau, K. K. W., McLachlan, G. J., Chambers, S. K., and Ng, S-K. (2019b). Multilevel model with random effects for clustered survival data with multiple failure outcomes. *Statistics in Medicine* **38**, 1036–1055.

Therneau, R. (2018). Mixed effects Cox models. <https://cran.r-project.org/web/packages/coxme/vignettes/coxme.pdf> (accessed April 9, 2019).

Vaida, F., and Xu, R. (2000). Proportional hazards model with random effects. *Statistics in Medicine* **19**, 3309–3324.

Wienke, A., Lichtenstein, P. and Yashin A.I. (2003). A Bivariate Frailty Model with a Cure Fraction for Modeling Familial Correlations in Diseases. *Biometrics* **59**, 1178–1183.

Xue, X., and Brookmeyer, R. (1996). Bivariate frailty model for the analysis of multivariate survival time. *Life Time Data Analysis* **2**, 277–289.

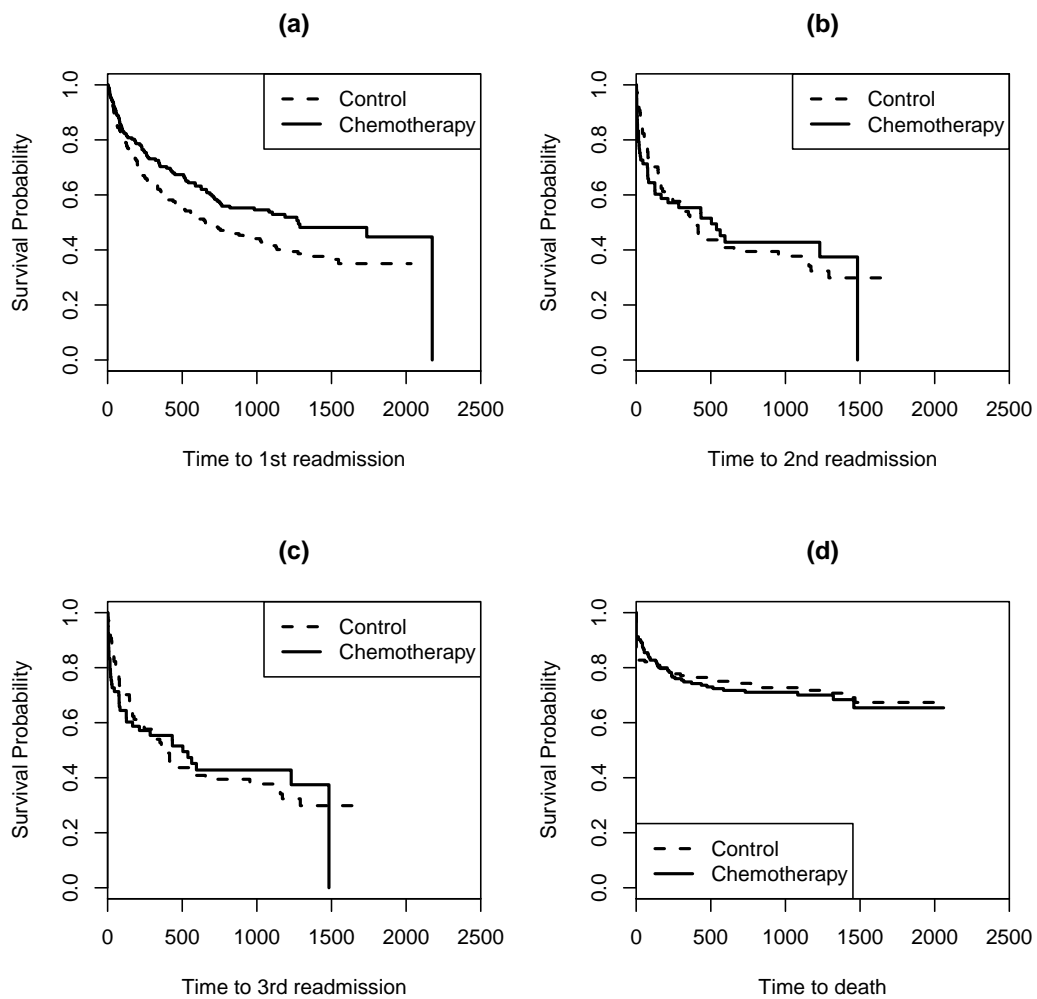
Yau, K. K. W. (2001). Multilevel models for survival analysis with random effects. *Biometrics* **57**, 96–102.

Yau, K. K. W., and Ng, A. S. K. (2001). Long-term survivor mixture model with random effects: application to a multi-centre clinical trial of carcinoma study data. *Statistics in Medicine* **20**, 1591–1607.

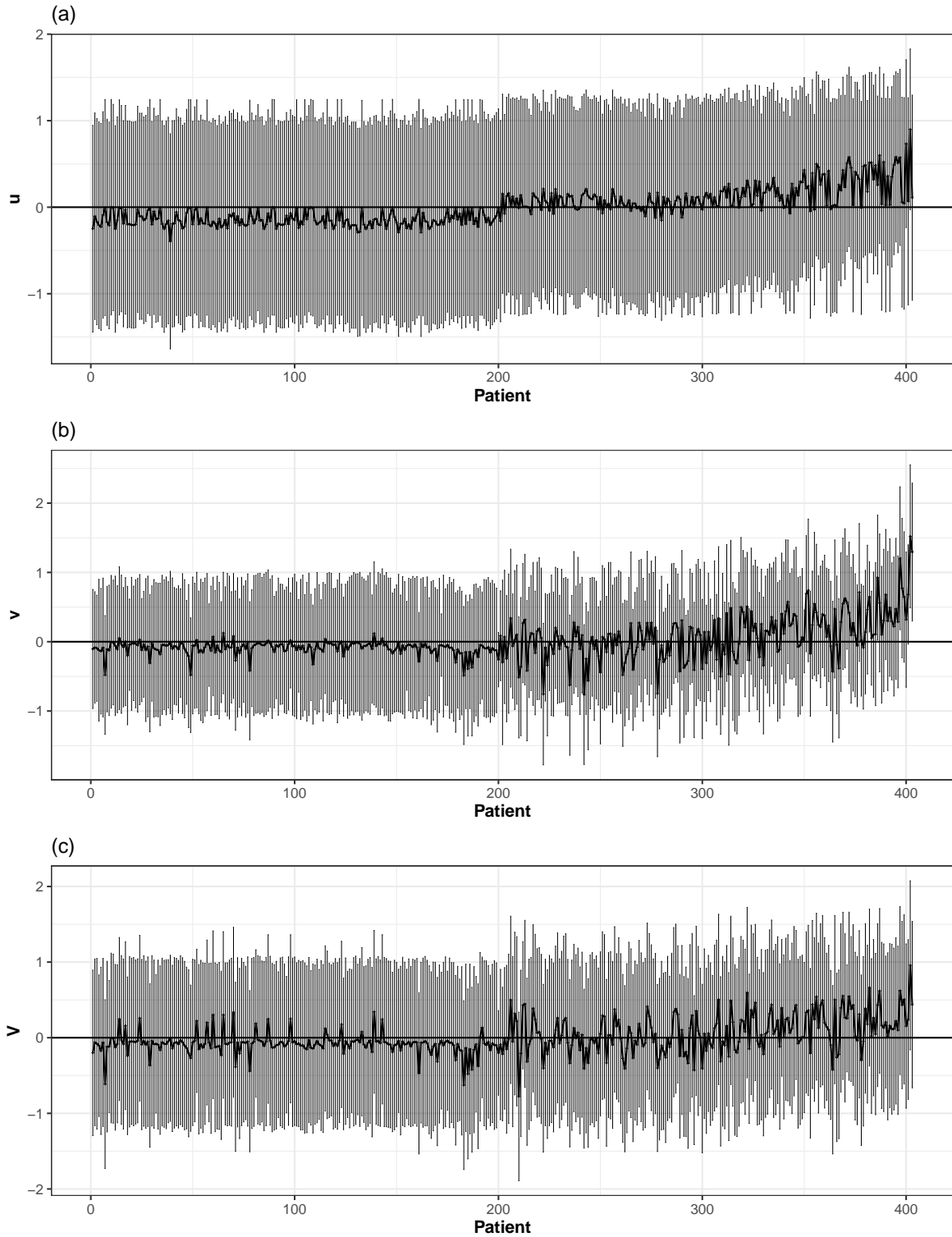
#### SUPPORTING INFORMATION

Web Appendices A-E, Tables S1-S5, and Figure S1 referenced in Sections 3-5 are available with this paper at the Biometrics website on Wiley Online Library. The R code and an example data for implementing the proposed model can be found on the library.

*Received October 2007. Revised February 2008. Accepted March 2008.*



**Figure 1:** Kaplan-Meier curves of the gap times for colorectal cancer hospital readmissions and death, stratified by treatment arm. Times to the first three readmissions are shown.



**Figure 2:** Patient-specific frailties and their 95% empirical Bayes prediction interval for heterogeneity in the (a) cure probability and the hazard rates for (b) hospital readmission and (c) death. Patients are sorted in increasing order of number of hospitalizations encountered.

Table 1: Bias, standard error and coverage probability of the REML estimators for the bivariate joint frailty cure model with 37.8% cure rate in treatment group based on 500 replications of simulated data.

Parameter	True value	Ave. bias	SEE	SE	CP	Parameter	True value	Ave. bias	SEE	SE	CP
Simulation 1						Simulation 2					
$\alpha_0$	1.0	0.050	0.115	0.118	0.922	$\alpha_0$	1.0	0.045	0.116	0.119	0.928
$\alpha_1$	-0.5	0.000	0.163	0.167	0.954	$\alpha_1$	-0.5	0.001	0.164	0.166	0.950
$\beta_1$	-0.8	0.053	0.110	0.102	0.932	$\beta_1$	-0.8	0.058	0.107	0.100	0.922
$\gamma_1$	-0.5	-0.068	0.188	0.230	0.884	$\gamma_1$	-0.5	-0.058	0.185	0.228	0.892
$\theta_u^2$	0.5	0.013	0.174	0.035	1.000	$\theta_u^2$	0.5	0.028	0.176	0.049	1.000
$\theta_v^2$	0.3	-0.044	0.052	0.028	0.956	$\theta_v^2$	0.3	-0.045	0.049	0.029	0.968
$\theta_V^2$	0.5	0.026	0.046	0.032	0.992	$\theta_V^2$	0.5	-0.032	0.042	0.037	0.968
$\rho$	0.2	-0.021	0.063	0.062	0.976	$\rho$	0.4	-0.027	0.062	0.069	0.994
Simulation 3						Simulation 4					
$\alpha_0$	1.0	0.035	0.116	0.119	0.938	$\alpha_0$	1.0	0.032	0.116	0.118	0.940
$\alpha_1$	-0.5	0.002	0.164	0.166	0.948	$\alpha_1$	-0.5	0.002	0.164	0.166	0.948
$\beta_1$	-0.8	0.060	0.106	0.099	0.916	$\beta_1$	-0.8	0.060	0.105	0.099	0.920
$\gamma_1$	-0.5	-0.037	0.175	0.215	0.896	$\gamma_1$	-0.5	-0.027	0.171	0.210	0.902
$\theta_u^2$	0.5	0.041	0.177	0.049	1.000	$\theta_u^2$	0.5	0.043	0.177	0.049	1.000
$\theta_v^2$	0.3	-0.038	0.040	0.026	0.980	$\theta_v^2$	0.3	-0.044	0.037	0.024	0.962
$\theta_V^2$	0.5	-0.066	0.034	0.026	0.961	$\theta_V^2$	0.5	-0.087	0.034	0.024	0.965
$\rho$	0.6	0.026	0.045	0.038	0.983	$\rho$	0.8	-0.017	0.030	0.022	0.977
Simulation 5						Simulation 6					
$\alpha_0$	1.0	0.063	0.116	0.120	0.903	$\alpha_0$	1.0	0.048	0.115	0.120	0.926
$\alpha_1$	-0.5	-0.002	0.163	0.170	0.944	$\alpha_1$	-0.5	-0.001	0.163	0.167	0.950
$\beta_1$	-0.8	0.065	0.109	0.101	0.931	$\beta_1$	-0.8	0.067	0.106	0.101	0.917
$\gamma_1$	-0.5	-0.068	0.188	0.237	0.876	$\gamma_1$	-0.5	-0.028	0.171	0.216	0.900
$\theta_u^2$	0.5	0.008	0.175	0.035	1.000	$\theta_u^2$	0.5	0.019	0.175	0.030	1.000
$\theta_v^2$	0.5	-0.047	0.052	0.027	0.944	$\theta_v^2$	0.5	-0.051	0.038	0.021	0.960
$\theta_V^2$	0.5	0.024	0.046	0.031	0.995	$\theta_V^2$	0.5	-0.093	0.036	0.020	0.988
$\rho$	0.2	-0.024	0.063	0.060	0.980	$\rho$	0.8	0.033	0.023	0.014	0.979
Simulation 7						Simulation 8					
$\alpha_0$	1.0	0.066	0.116	0.119	0.902	$\alpha_0$	1.0	0.078	0.116	0.120	0.886
$\alpha_1$	-0.5	-0.004	0.164	0.167	0.956	$\alpha_1$	-0.5	-0.007	0.164	0.168	0.948
$\beta_1$	-0.8	0.065	0.113	0.103	0.932	$\beta_1$	-0.8	0.072	0.110	0.104	0.918
$\gamma_1$	-0.5	-0.082	0.190	0.230	0.892	$\gamma_1$	-0.5	-0.051	0.173	0.211	0.882
$\theta_u^2$	0.5	0.015	0.175	0.030	1.000	$\theta_u^2$	0.5	0.016	0.176	0.029	1.000
$\theta_v^2$	0.3	-0.018	0.063	0.029	0.986	$\theta_v^2$	0.3	-0.041	0.051	0.027	0.976
$\theta_V^2$	0.5	0.034	0.052	0.033	0.984	$\theta_V^2$	0.5	0.047	0.048	0.026	0.991
$\rho$	-0.2	-0.011	0.063	0.042	0.993	$\rho$	-0.8	0.061	0.022	0.012	0.964

Abbreviations: Ave. bias, average bias; SE, standard error; SEE, standard error estimates; CP, coverage probability.

Table 2: Estimates of the proposed bivariate frailty mixture cure model, a joint frailty model and a standard frailty model based on the colorectal cancer hospital readmission data.

	Joint frailty cure model			Joint frailty model			Frailty model		
	Estimate	OR/HR	SE	Estimate	HR	SE	Estimate	HR	SE
Logistic component									
Intercept	0.251		0.202						
Chemotherapy	-0.067	0.935	0.200						
Sex (female)	-0.476 <sup>a</sup>	0.621	0.188						
Dukes stage									
C	0.487 <sup>a</sup>	1.627	0.205						
D	3.767 <sup>a</sup>	43.250	0.587						
Charlson index									
1-2	0.503	1.654	0.395						
≥ 3	0.486 <sup>b</sup>	1.626	0.252						
Patient frailty									
$\theta_u^2$	0.398		0.198						
Hazard component									
<b>Readmission</b>									
Chemotherapy	-0.230 <sup>b</sup>	0.795	0.137	-0.138	0.871	0.143	-0.177	0.838	0.139
Sex (female)	-0.376 <sup>a</sup>	0.687	0.135	-0.457 <sup>a</sup>	0.633	0.139	-0.457 <sup>a</sup>	0.633	0.135
Dukes stage									
C	0.116	1.123	0.156	0.327	1.390	0.161	0.285 <sup>b</sup>	1.329	0.156
D	0.223	1.250	0.178	1.144 <sup>a</sup>	3.140	0.191	1.052 <sup>a</sup>	2.863	0.187
Charlson index									
1-2	0.291	1.338	0.253	0.286	1.330	0.257	0.407	1.502	0.252
≥ 3	0.303 <sup>a</sup>	1.354	0.131	0.371 <sup>a</sup>	1.450	0.134	0.320 <sup>a</sup>	1.378	0.133
Patient frailty									
$\theta_v^2$	0.354		0.041	0.610		0.071	0.494		
<b>Death</b>									
Chemotherapy	-0.009	0.991	0.237	0.264	1.300	0.265			
Sex (female)	-0.468 <sup>a</sup>	0.626	0.227	-0.418 <sup>b</sup>	0.658	0.252			
Dukes stage									
C	1.015 <sup>a</sup>	2.759	0.346	0.928 <sup>a</sup>	2.530	0.355			
D	2.327 <sup>a</sup>	10.247	0.366	2.621 <sup>a</sup>	13.700	0.376			
Charlson index									
1-2	0.529	1.697	0.641	0.403	1.500	0.682			
≥ 3	1.405 <sup>a</sup>	4.076	0.258	1.373 <sup>a</sup>	3.950	0.280			
Patient frailty									
$\theta_v^2$	0.429		0.036	1.235		0.141			
$\rho$	0.489 <sup>a</sup>		0.047	0.949 <sup>a</sup>		0.007			

Abbreviations: OR, odds ratio; HR, hazard ratio; SE, standard error.  
 Male, reference category for gender; A-B, reference category for Dukes's stage;  
 0, reference category for Charlson index.

<sup>a</sup> p-value < 0.05

<sup>b</sup> p-value < 0.10

Author Manuscript