

MODELING COUNTY LEVEL BREAST CANCER SURVIVAL DATA USING A COVARIATE-ADJUSTED FRAILTY PROPORTIONAL HAZARDS MODEL

BY HAIMING ZHOU^{*,1}, TIMOTHY HANSON^{*,1}, ALEJANDRO JARA^{†,2}
AND JIAJIA ZHANG^{*,1}

*University of South Carolina** and *Pontificia Universidad Católica de Chile†*

Understanding the factors that explain differences in survival times is an important issue for establishing policies to improve national health systems. Motivated by breast cancer data arising from the Surveillance Epidemiology and End Results program, we propose a covariate-adjusted proportional hazards frailty model for the analysis of clustered right-censored data. Rather than incorporating exchangeable frailties in the linear predictor of commonly-used survival models, we allow the frailty distribution to flexibly change with both continuous and categorical cluster-level covariates and model them using a dependent Bayesian nonparametric model. The resulting process is flexible and easy to fit using an existing R package. The application of the model to our motivating example showed that, contrary to intuition, those diagnosed during a period of time in the 1990s in more rural and less affluent Iowan counties survived breast cancer better. Additional analyses showed the opposite trend for earlier time windows. We conjecture that this anomaly has to be due to increased hormone replacement therapy treatments prescribed to more urban and affluent subpopulations.

1. Introduction. Based on data gathered for Iowa State in the Surveillance Epidemiology and End Results (SEER) program of the National Cancer Institute, we assess the effect of potential risk factors for womens' breast cancer. This involves the analysis of clustered time-to-event right-censored

Received July 2014; revised October 2014.

¹Supported by NCI Grants R03CA165110 and 5R03CA176739, and ASPIRE grant from the University of South Carolina.

²Supported by Fondecyt Grant 1141193.

Key words and phrases. Clustered time-to-event data, proportional hazards model, spatial, tailfree process.

<p>This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in <i>The Annals of Applied Statistics</i>, 2015, Vol. 9, No. 1, 43–68. This reprint differs from the original in pagination and typographic detail.</p>

data, where event times of patients from the same county of residence are expected to be associated with each other, possibly due to sharing common unobserved characteristics, such as region-specific differences in environments, treatment resources or diagnosis of the patients. As is widely known, taking into account the clustered nature of the data is a must to obtain valid statistical inferences [see, e.g., Therneau and Grambsch (2000), Chapter 8].

A standard way of modeling clustered survival data is to introduce a common random effect (frailty) into the survival model for each cluster, yielding shared frailty models. “Frailties,” termed by Vaupel, Manton and Stallard (1979), were originally introduced to deal with possible heterogeneity due to unobserved covariates and are regarded as unobserved common characteristics for each cluster able to account for the dependence among event times. In the context of the proportional hazards (PH) model, as conventionally implemented, frailties are incorporated into the linear predictor, and the median or mean of the frailty distribution is constrained to be zero to avoid identifiability problems. Conditional on the frailty, the model retains its interpretation in terms of constants of proportionality of the hazards. Survival models with frailties have been extensively used in the statistical literature, especially when the comparison of event times within cluster is of interest.

A common assumption in shared frailty survival models is the one of homogeneity, where the frailties are assumed to be independent and identically distributed (i.i.d.) random variables from a parametric or nonparametric distribution [see, e.g., Clayton and Cuzick (1985), Gustafson (1997), Qiou, Ravishanker and Dey (1999), Walker and Mallick (1997)]. Although the nonparametric approach provides flexibility in capturing a frailty distribution’s variance, skewness, shape and even modality, it essentially assumes that these frailty distributional aspects are the same across all the clusters, which may be restrictive for particular data sets [Noh, Ha and Lee (2006)]. For example, in the kidney transplantation study, Liu, Kalbfleisch and Schaubel (2011) argue that the frailty distribution may be affected by some cluster-level covariates, since “... *urban transplant facilities may exhibit more uniform practices than rural transplant hospitals, corresponding to less heterogeneity (smaller variance) for frailties of urban centers...*” Ignoring such heterogeneity can drastically affect the inference for cluster-specific effects and prediction [McCulloch and Neuhaus (2011)].

As the process generating the frailty terms is on its own right of scientific interest, different extensions of the i.i.d. frailty modeling approach have been considered. Wassell and Moeschberger (1993) studied the impact of interventions in the Framingham Heart Study by introducing a modified gamma frailty with a pairwise covariate-dependent parameter. Yashin and Iachine (1999) considered the dependence between frailty and observed covariates (BMI and smoking) in Danish twins to investigate the heritability of susceptibility to death. Noh, Ha and Lee (2006) verified frailty distribution heterogeneity in a well-known kidney infection data set by applying a dispersed

normal model. Cottone (2008) assumed either Bernoulli or normal distributions for the frailties where the frailty distribution mean or variance depends on cluster-level covariates through specified link functions. Liu, Kalbfleisch and Schaubel (2011) proposed a covariate-dependent positive stable shared frailty model with an application to kidney transplantation data from the Scientific Registry of Transplant Recipients, and demonstrated the heterogeneity in facility performance. Wang and Louis (2004) studied a related approach for binary data that has both conditional and marginal interpretation using the so-called bridge distribution instead of positive stable.

The previously described model extensions allow for particular and specific aspects of distributional shape to change with cluster-level covariates. However, a more thorough evaluation of the effect of the predictors should account for potential changes in characteristics of the frailty distribution other than just, for example, the location or scale. It is, for instance, useful to examine potential changes in the skewness, symmetry and multimodality of the frailty distribution. Therefore, a nonparametric formulation that anticipates changes in shape, skew and modality beyond simple location models is of interest.

In this paper, we propose a practicable and general framework for modeling clustered survival data as a function of covariates, based on a predictor-dependent Bayesian nonparametric model for the frailties and the Cox's PH model. Under the proposed approach the frailty distribution flexibly changes with both continuous and categorical cluster-level covariates, thus allowing for full heterogeneity across clusters. We apply this modeling approach to a subset of the SEER county-level breast cancer data consisting of 1073 women diagnosed with malignant breast cancer during 1995–1998. Important patient-level covariates include age at diagnosis, race, county of residence and the stage of the disease. Additional county-level covariates potentially associated with breast cancer survival are also available from census data, including median household income, poverty level, education and a rurality measure. These area-level socioeconomic factors have been discovered to be associated with breast cancer by many researchers [e.g., Sprague et al. (2011)]. Women living in more affluent or less rural geographic areas tend to survive breast cancer better after a diagnosis than those living in regions with indicators of low socioeconomic status. Moreover, rural counties may present more heterogeneity in access to quality care and screening for breast cancer, leading to more variability for frailties of rural counties [Zhao and Hanson (2011)]. This suggests to us that the frailty distribution could be potentially affected by these county-level socioeconomic factors. The results show that the proposed model provides better goodness of fit to the data and is predictively superior to the traditional PH spatial frailty model, as well as helping to piece together a plausible story for the data in terms of the prescribing of hormone replacement therapy.

The paper is organized as follows. In Section 2 we introduce the proposed frailty PH model, including a detailed description of the dependent Bayesian nonparametric model and the Markov chain Monte Carlo (MCMC) implementation of the posterior computations. Section 3 provides a detailed analysis of the motivating data set. Section 4 presents the results of simulation studies to evaluate the performance of the proposed model. Some concluding remarks and a final discussion are given in Section 5.

2. Covariate-adjusted frailty proportional hazards model.

2.1. *The modeling approach.* Suppose that right-censored survival data $(\mathbf{w}_{ij}, t_{ij}, \delta_{ij})$ are collected for the j th subject of the i th cluster, where $j = 1, \dots, n_i$, $i = 1, \dots, n$, \mathbf{w}_{ij} is a p -dimensional vector of exogenous covariates, t_{ij} is the recorded event time, and δ_{ij} is the censoring indicator equaling 1 if t_{ij} is an observed event time and equaling 0 if the event time is right-censored at t_{ij} . Let T_{ij} and C_{ij} be the event and censoring times, respectively, for the j th subject in the i th cluster. To take into account the within-cluster association structure, a frailty PH model is assumed for T_{ij} . The conditional PH assumption implies that the hazard function of T_{ij} is given by

$$(1) \quad \lambda(t|\mathbf{w}_{ij}, e_i) = \lambda_0(t) \exp(\mathbf{w}'_{ij}\boldsymbol{\xi} + e_i),$$

where $\mathbf{e} = (e_1, \dots, e_n)'$ is an unobserved vector of frailties, and $\lambda_0(t)$ is the baseline hazard function corresponding to the event time of a subject with covariates $\mathbf{w} = \mathbf{0}$ and $e = 0$. We additionally assume a conditionally independent censoring scheme, that is, C_{ij} and T_{ij} are independent given \mathbf{w}_{ij} and e_i . Often the frailties are assumed to be exchangeable or i.i.d. from some parametric or nonparametric distribution G . For instance, Therneau, Grambsch and Pankratz (2003) considered exchangeable Gaussian frailties and proposed an estimation procedure based on a Laplace approximation of the likelihood function leading to a penalized partial likelihood. This approach, referred to below as GF, will be compared with our method in the simulation studies.

Now consider a partition of the predictor vector $\mathbf{w}_{ij} = (\tilde{\mathbf{w}}'_{ij}, \mathbf{x}'_i)'$, where $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbf{R}^q$ is a q -dimensional vector of cluster-level covariates and $\tilde{\mathbf{w}}_{ij}$ is a $(p-q)$ -dimensional vector of subject-specific covariates, respectively, and the corresponding partition of the regression coefficient vector $\boldsymbol{\xi} = (\tilde{\boldsymbol{\xi}}', \boldsymbol{\xi}'_x)'$. On the scale of the linear predictor $\mathbf{w}'_{ij}\boldsymbol{\xi} + e_i$, the frailty e_i models the cluster-specific behavior and its distribution G is shifted by $\mathbf{x}'_i\boldsymbol{\xi}_x$. Therefore, the homogeneity assumption implies that the vector of cluster-level covariates \mathbf{x}_i modifies only the location of the distribution of cluster-specific effects but not its shape. To relax this assumption, we consider a covariate-adjusted

frailty PH model, where the frailty distribution depends on cluster-level covariates \mathbf{x}_i . That is,

$$e_i | G_{\mathbf{x}_i} \stackrel{\text{ind.}}{\sim} G_{\mathbf{x}_i},$$

where for every $\mathbf{x} \in \mathcal{X}$, $G_{\mathbf{x}}$ is a probability measure defined on \mathbb{R} ; this specifies a probability model for the entire collection of probability measures $\mathcal{G}^{\mathcal{X}} = \{G_{\mathbf{x}} : \mathbf{x} \in \mathcal{X}\}$, such that its elements are allowed to smoothly vary with the cluster-level covariates \mathbf{x} . Specifically, we consider a mixture of linear dependent tailfree processes (LDTFP) prior [Jara and Hanson (2011)] for $\mathcal{G}^{\mathcal{X}}$, denoted as

$$\mathcal{G}^{\mathcal{X}} | J, h, \theta, c, \rho \sim \text{LDTFP}(h, \Pi^{J, \theta}, \mathcal{A}^{J, c, \rho}),$$

and

$$c | Q \sim Q,$$

where $J \in \mathbb{N}$ is the level of specification of the process, $c \in \mathbb{R}_+$ is a prior precision parameter controlling the prior variability of the process, $h(\cdot) = \frac{\exp\{\cdot\}}{1 + \exp\{\cdot\}}$, $\Pi^{J, \theta}$ is a J -level sequence of binary partitions of \mathbb{R} , depending on the scale parameter $\theta \in \mathbb{R}_+$, $\mathcal{A}^{J, c, \rho} = \{2n/c\rho(1), \dots, 2n/c\rho(J)\}$ is a collection of positive numbers depending on J , c and ρ , $\rho : \mathbb{N} \rightarrow \mathbb{R}_+$ is an increasing function, and Q is a probability measure defined on \mathbb{R}_+ .

The LDTFP is specified such that for every $\mathbf{x} \in \mathcal{X}$, the process $G_{\mathbf{x}}$ is centered around an $N(0, \theta)$ distribution, that is, $E(G_{\mathbf{x}}) = N(0, \theta)$, for every $\mathbf{x} \in \mathcal{X}$. Furthermore, the process is specified such that for every $\mathbf{x} \in \mathcal{X}$, $G_{\mathbf{x}}$ is almost surely a median-zero probability measure. The latter property is important to avoid identifiability problems. The LDTFP process includes as important special cases a nonparametric exchangeable frailty model where $G_{\mathbf{x}} = G_{\mathbf{x}'}$ for $\mathbf{x}' \neq \mathbf{x}$ as well as exchangeable normal frailties $G_{\mathbf{x}} = N(0, \theta)$ for all $\mathbf{x} \in \mathcal{X}$.

As shown by Jara and Hanson (2011), dependent tailfree processes have appealing theoretical properties, such as continuity as a function of the predictors, large support on the space of conditional density functions, straightforward posterior computation relying on algorithms for fitting generalized linear models, and the process closely matches conventional Polya tree priors [see, e.g., Hanson (2006a)] at each value of the predictor, which justify its choice here. Polya trees have been extensively studied in the literature and have desirable properties in terms of support and posterior consistency. Details on the trajectories of $\text{LDTFP}(h, \Pi^{J, \theta}, \mathcal{A}^{J, c, \rho})$, useful for a complete implementation of algorithms for exploring the corresponding posterior distributions, are given in Appendix A of the supplementary material [Zhou et al. (2015)].

Other dependent processes could be considered for $\mathcal{G}^{\mathcal{X}}$, but a highly limiting requirement is that some aspect of the location, for example, mean or median, can be fixed. There are few examples where the process changes smoothly with covariates; one is the multivariate beta process of Trippa, Müller and Johnson (2011). Another approach using Dirichlet process mixtures can be found in Reich, Bondell and Wang (2010), but this latter approach would have to be extended to allow the means or variances of the two mixture components to change with covariates.

2.2. Posterior computation. The conditional likelihood for $(\boldsymbol{\xi}, \lambda_0, \mathbf{e})$ is given by

$$\mathcal{L}(\boldsymbol{\xi}, \lambda_0, \mathbf{e}) = \prod_{i=1}^n \prod_{j=1}^{n_i} [\lambda_0(t_{ij}) \exp(\mathbf{w}'_{ij} \boldsymbol{\xi} + e_i)]^{\delta_{ij}} \exp\{-\Lambda_0(t_{ij}) \exp(\mathbf{w}'_{ij} \boldsymbol{\xi} + e_i)\},$$

where $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$ is the cumulative hazard function. The piecewise exponential model provides a flexible framework to deal with the baseline hazard [see, e.g., Walker and Mallick (1997)]. We partition the time period \mathbb{R}_+ into K prespecified intervals, say, $I_k = (a_{k-1}, a_k]$, $k = 1, \dots, K$, where $a_0 = 0$ and $a_K = \max\{t_{ij}\}$. The baseline hazard is assumed to be constant within each interval, that is,

$$\lambda_0(t) = \sum_{k=1}^K \lambda_k I\{t \in I_k\},$$

where $\lambda_1, \dots, \lambda_K$ are unknown hazard values and $I\{A\}$ is the usual indicator function, that is, 1 when A is true, 0 otherwise. The prior hazard is specified by the hazard values $\{\lambda_k\}_{k=1}^K$ and cut-point vector $\mathbf{a} = (a_1, \dots, a_K)$. If the prior on the λ_k 's is taken to be independent gamma distributions and $\{I_k\}_{k=1}^K$ is a reasonably fine mesh, the gamma process [Kalbfleisch (1978)] is approximated. To determine the cut-point vector \mathbf{a} , one can set a_k to be the $\frac{k}{K}$ th quantile of the empirical distribution of the t_{ij} 's, or choose them based on other considerations (see Section 3.2). Some authors have considered random cut-points [see, e.g., Sahu and Dey (2004)]. Regardless, the resulting model implies a Poisson likelihood [Laird and Olivier (1981)] as follows. Let $K(t) = \min\{k: a_k \geq t\}$, $\Delta_k(t) = \min\{a_k, t\} - a_{k-1}$, and $y_{ijk} = \delta_{ij} I\{k = K(t_{ij})\}$. Set $\mathbf{z}_{ijk} = (\boldsymbol{\nu}'_k, \mathbf{w}'_{ij})'$ and $\boldsymbol{\gamma} = (\boldsymbol{\lambda}', \boldsymbol{\xi}')'$, where $\boldsymbol{\nu}_k$ is a K -dimensional vector of zeros except the k th element is 1 and $\boldsymbol{\lambda} = (\log(\lambda_1), \dots, \log(\lambda_K))'$. Then the likelihood for $(\boldsymbol{\gamma}, \mathbf{e})$ becomes

$$\mathcal{L}(\boldsymbol{\gamma}, \mathbf{e}) = \prod_{i=1}^n \prod_{j=1}^{n_i} [\exp\{\log(\lambda_{K(t_{ij})}) + \mathbf{w}'_{ij} \boldsymbol{\xi} + e_i\}]^{\delta_{ij}}$$

$$\begin{aligned}
& \times \left[\prod_{k=1}^{K(t_{ij})} e^{-\exp\{\log(\lambda_k) + \mathbf{w}'_{ij}\xi + e_i\}} \Delta_k(t_{ij}) \right] \\
& = \prod_{i=1}^n \prod_{j=1}^{n_i} \prod_{k=1}^{K(t_{ij})} [(\exp\{\mathbf{z}'_{ijk}\boldsymbol{\gamma} + e_i\})^{y_{ijk}} e^{-\exp\{\mathbf{z}'_{ijk}\boldsymbol{\gamma} + e_i + \log(\Delta_k(t_{ij}))\}}] \\
& \propto \prod_{i=1}^n \prod_{j=1}^{n_i} \prod_{k=1}^{K(t_{ij})} p(y_{ijk} | \boldsymbol{\gamma}, e_i),
\end{aligned}$$

where $\mu_{ijk} = \exp\{\mathbf{z}'_{ijk}\boldsymbol{\gamma} + e_i + \log(\Delta_k(t_{ij}))\}$ and $p(y_{ijk} | \boldsymbol{\gamma}, e_i)$ is the probability mass function for a Poisson distribution with mean μ_{ijk} . For each $i = 1, \dots, n$, let $N_i = \sum_{j=1}^{n_i} K(t_{ij})$, $\mathbf{y}_i = (y_{ijk})$ be an $N_i \times 1$ vector with subscript ijk in lexicographical order.

Thus, the proposed covariate-adjusted frailty PH model takes the following hierarchical structure:

$$\begin{aligned}
\mathbf{y}_i | \boldsymbol{\gamma}, e_i & \stackrel{\text{ind.}}{\sim} \prod_{j=1}^{n_i} \prod_{k=1}^{K(t_{ij})} p(y_{ijk} | \boldsymbol{\gamma}, e_i), \\
\boldsymbol{\gamma} & \sim N_{K+p}(\boldsymbol{\gamma}_0, \mathbf{S}_0), \\
e_i | G_{\mathbf{x}_i} & \stackrel{\text{ind.}}{\sim} G_{\mathbf{x}_i}, \\
\mathcal{G}^X | J, h, \theta, c, \rho & \sim \text{LDTFP}(h, \Pi^{J, \theta}, \mathcal{A}^{J, c, \rho}), \\
\theta^{-2} & \sim \Gamma(\tau_1, \tau_2), \quad c \sim \Gamma(a_c, b_c),
\end{aligned}$$

which largely simplifies computations, where $N_p(\mathbf{m}, \mathbf{S})$ refers to a p -variate normal distribution with mean \mathbf{m} and covariance matrix \mathbf{S} . This forms the basis of an efficient Markov chain Monte Carlo (MCMC) scheme for obtaining posterior inference, which can be implemented using available software for generalized linear mixed models. A full description of the MCMC algorithm is given in Appendix B of the supplementary material [Zhou et al. (2015)]. Sample R code using the `LDTFPglmm` function available in `DPpackage` [Jara et al. (2011)] is provided in Appendix C of the supplementary material [Zhou et al. (2015)].

Time-dependent subject-specific covariates that are step-processes [Hanson, Johnson and Laud (2009)] are naturally accommodated by including the times where the covariate values change across all subjects into the cut-point vector \mathbf{a} . All that is changed above is $\mathbf{z}_{ijk} = (\boldsymbol{\nu}'_k, \mathbf{w}'_{ijk})'$, that is, \mathbf{w}_{ij} is replaced with its time-varying analogue \mathbf{w}_{ijk} . Similarly, time-varying regression effects can be included by replacing $\mathbf{z}'_{ijk}\boldsymbol{\gamma}$ with $\mathbf{z}'_{ijk}\boldsymbol{\gamma}_k$ in μ_{ijk} , yielding very general models. The proposed model implies exchangeable frailties for

each subgroup with a unique $\mathbf{x} \in \mathcal{X}$. Time-dependent cluster-specific covariates are therefore naturally included in the model by simply allowing \mathbf{x} to change with time. For example, in the SEER data set analyzed over a larger time window, for subjects living in the i th county, one could include into \mathbf{x}_i the median house income of that county at their particular diagnosis year. Furthermore, the frailty distribution can itself evolve in time by simply including time as a covariate in \mathbf{x} , or a time-by-cluster covariate interaction could also be entertained.

3. Analysis of SEER county-level breast cancer data.

3.1. *The Iowa SEER data.* The SEER program of the National Cancer Institute (see <http://seer.cancer.gov/>) is an authoritative source of information on cancer incidence and survival in the US, providing county-level cancer data on an annual basis for particular states for public use. We fit our proposed covariate-adjusted frailty Cox's PH model to a subset of the Iowa SEER breast cancer survival data, which consists of a cohort of 1073 women from the 99 counties of Iowa, who were diagnosed with malignant breast cancer in 1995, with enrollment and follow-up continued through the end of 1998. The observed survival time, from 1 to 48, was calculated as the number of months from diagnosis to either death or the last follow-up. In our analysis, only deaths due to metastasis of cancerous nodes in the breast were considered to be events, while the deaths from other causes were censored at the time of death. That is, we consider cause-specific survival models assuming that all other deaths are independent of breast cancer. By the end of 1998, a total of 488 patients (45.5%) had died of breast cancer, while the remaining 585 patients were censored, either because they died of other causes or survived until the last follow-up.

For each patient, the observed survival time and county of residence at diagnosis are recorded. The data set also has individual-level covariates including age at diagnosis and the stage of the breast cancer: local (confined to the breast), regional (spread beyond the breast tissue), or distant (metastasis). We create two dummy variables for regional and distant, respectively, and treat the patients in the local group as the baseline. Although several individual-level covariates that affect breast cancer survival are not available (e.g., age at first child, age at menopause and breastfeeding), we are able to obtain county-level covariates potentially associated with breast cancer survival from census data, such as median household income (small area estimates in 1993), poverty level (percentage of families in poverty in 1990), education (percentage with Bachelor's degree or higher in 1990) and rurality (Rural–Urban Continuum Codes in 1993). The Economic Research Service Rural–Urban Continuum Codes (RUCC) vary from 1 to

TABLE 1
Iowa SEER data: Summary statistics for follow-up times and both individual- and county-level covariates

Continuous variables	Minimum	Median	Maximum
Follow-up time in months	1	19	47
Age in years	26	72	103
RUCC	2	7	9
Income ($\times 1000$)	20.627	29.110	39.356
Categorical variables	Level	Count	Proportion (%)
Status	Event	488	45.5
	Censored	585	54.5
Stage	Local	510	47.5
	Regional	355	33.1
	Distant	208	19.4

9 (see www.ers.usda.gov/data-products/rural-urban-continuum-codes), distinguishing metropolitan counties by the population size of their metro area and nonmetropolitan counties by degree of urbanization and adjacency to a metro area. Higher RUCC indicates a more rural county. Other county-level covariates mentioned above are available at <http://data.iowadatacenter.org/browse/counties.html>. Since the effects of education and poverty on the survival times are not significant based on our initial model fitting by the proposed method, we exclude them in the analysis presented below. Thus, we have three-dimensional $\tilde{\mathbf{w}}_{ij}$ and two-dimensional \mathbf{x}_i . Table 1 presents several summary statistics for the data. As shown in Figure 1, median household income and RUCC are significantly, negatively correlated.

To get an initial feeling about the role that each county-level covariate is playing, Table 2 provides the distribution of each county-level covariate stratified by the individual-level stage of disease. The gamma statistic (GK), originally proposed by Goodman and Kruskal (1954), is calculated to quantify the association between each county-level covariate and the stage of disease. The GK values range from -1 (100% negative association) to 1 (100% positive association), where the value 0 indicates no association. We see that women with a distant-stage at diagnosis are much more likely than those with a local-stage to live in counties with a high degree of urbanization (GK = -0.11 ; 95% CI: from -0.20 to -0.01), while the association between stage and income is not significant (GK = 0.04 ; 95% CI: from -0.06 to 0.13). These associations roughly imply that women living in urban counties may suffer poorer survival, assuming that women in distant-stage are more likely to die than women in other stages. Next, we carefully examine

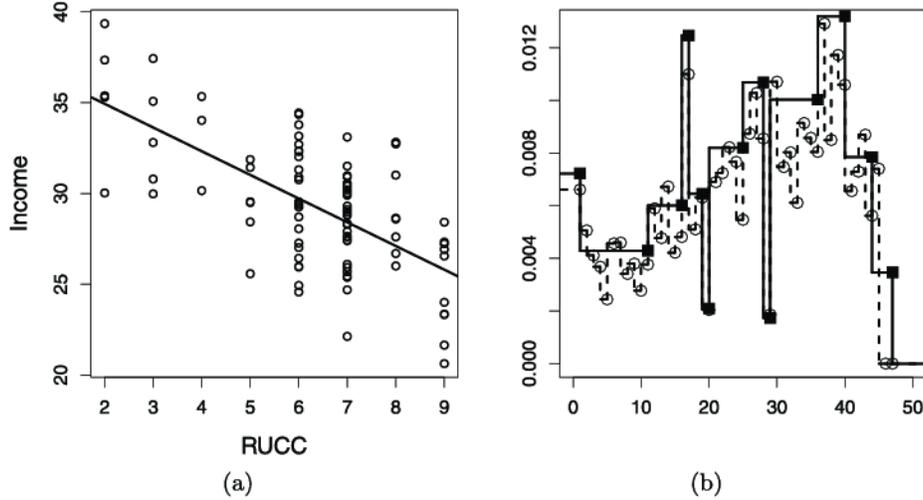


FIG. 1. Iowa SEER data: panel (a) shows the scatter plot and simple linear regression line by regressing median household income on RUCC. Panel (b) shows the baseline hazards for Model 1. The dashed line corresponds to Breslow's estimate of $\lambda_0(t)$ obtained by the GF approach, where the circles represent the hazard values at each month; the solid line is the fitted baseline hazard by our approach, where the solid squares correspond to the cut-point values $\mathbf{a} = (1, 11, 16, 17, 19, 20, 25, 28, 29, 36, 40, 44, 47)$.

both these individual-level and county-level covariates in relation to breast cancer survival, fitting the covariate-adjusted frailty proportional hazards model.

TABLE 2

Iowa SEER data: Distribution of each county-level covariate stratified by individual-level stage. The pattern of numbers is Number of women (%). Goodman and Kruskal's gamma statistics (95% confidence intervals) are -0.11 ($-0.20, -0.01$) and 0.04 ($-0.06, 0.13$) for RUCC and Income, respectively

Covariates	All women $N = 1073$	Stage		
		Local $N = 510$	Regional $N = 355$	Distant $N = 208$
RUCC				
1-3	314 (29.3)	131 (25.7)	99 (27.9)	84 (40.4)
4-7	666 (62.1)	342 (67.1)	221 (62.3)	103 (49.5)
8-9	93 (8.6)	37 (7.2)	35 (9.8)	21 (10.1)
Income ($\times 1000$)				
20-27	163 (15.2)	79 (15.5)	51 (14.4)	33 (15.9)
27-34	651 (60.7)	312 (61.2)	223 (62.8)	116 (55.8)
>34	259 (24.1)	119 (23.3)	81 (22.8)	59 (28.3)

3.2. *Models and model comparison.* We fitted the proposed covariate-adjusted frailty PH model for the Iowa SEER data with different county-level covariates, including models with RUCC only (Model 1), with median household income only (Model 2) and with both (Model 3). To see how the piecewise assumption of baseline hazard affects the predictive ability of models, we considered three specifications of cut-point vector \mathbf{a} as follows:

Case I. $\mathbf{a} = (1, 11, 16, 17, 19, 20, 25, 28, 29, 36, 40, 44, 47)$, which was determined by visually examining Breslow's estimate of $\lambda_0(t)$ using the GF approach, which is given in panel (b) of Figure 1.

Case II. $\mathbf{a} = (3, 7, 12, 16, 19, 24, 29, 34, 41, 47)$, where a_k is the $\frac{k}{10}$ th quantile of the empirical distribution of observed survival times.

Case III. $\mathbf{a} = (47)$, which yields an exponential baseline hazard.

In all cases, we set $J = 4$. We fitted all the models using the corresponding variants of the algorithm described in Appendix B of the supplementary material [Zhou et al. (2015)] and similar prior specifications suggested in the simulation study. The Markov chain mixed reasonably well despite the high dimension of our models. For each version of our model and case, we ran a single Markov chain of 1,020,000. A total number of 20,000 were discarded as burn-in period and 10,000 samples were retained for posterior inference. Moreover, we also considered another case II with 13 cut-points and cut-point specifications based on the event time quantiles from the Kaplan–Meier curve in Appendix E of the supplementary material [Zhou et al. (2015)]. The results show that carefully choosing the cut-points is more important than simply increasing the number of cut-points.

For the sake of comparison, we further fitted the exchangeable MPT frailty Cox model and the Bayesian exchangeable Gaussian frailty Cox model. We compare the models using the log pseudo marginal likelihood (LPML) developed by Geisser and Eddy (1979) and the deviance information criterion (DIC) proposed by Spiegelhalter et al. (2002). In the context of the frailty Cox model, the LPML for model M is defined as $\text{LPML} = \sum_{i=1}^n \sum_{j=1}^{n_i} \log(\text{CPO}_{ij})$, where CPO_{ij} , the ij th conditional predictive ordinate, is given by $[\lambda(t_{ij})^{\delta_{ij}} e^{-\Lambda(t_{ij})} | \mathcal{D}_{(ij)}]$ with $\mathcal{D}_{(ij)}$ denoting the remaining data after excluding the ij th data point \mathcal{D}_{ij} . One can use the simple method suggested by Gelfand and Dey (1994) to estimate the CPO statistics from MCMC output. A larger value of LPML indicates the corresponding model has better predictive ability. Furthermore, Geisser and Eddy (1979) discussed the exponentiated difference in LPML values from two models to obtain what they termed as a pseudo Bayes factor (PBF). The PBF is a surrogate for the more traditional Bayes factor and can be interpreted similarly, but is more analytically tractable, much less sensitive to prior assumptions, and does not suffer from Lindley's paradox. Set $\boldsymbol{\Omega} = (\mathbf{e}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \theta)$ as the

TABLE 3
Iowa SEER data: Deviance information criteria (DIC) and log of the pseudo marginal likelihood (LPML) for models under consideration

Model	Frailty	Case I		Case II		Case III	
		DIC	LPML	DIC	LPML	DIC	LPML
1	LDTFP	4436	-2222	4463	-2234	4495	-2247
	MPT	4441	-2225	4463	-2235	4496	-2248
	Gaussian	4444	-2225	4467	-2236	4497	-2248
2	LDTFP	4441	-2224	4465	-2235	4498	-2249
	MPT	4440	-2225	4462	-2236	4497	-2248
	Gaussian	4443	-2225	4465	-2235	4498	-2249
3	LDTFP	4438	-2223	4464	-2235	4496	-2248
	MPT	4441	-2225	4464	-2235	4498	-2249
	Gaussian	4445	-2226	4467	-2236	4498	-2248

entire collection of model parameters. The DIC for model M is defined as $DIC = \bar{D} + p_D = E_{\Omega|\mathcal{D}}\{D(\Omega)\} + p_D$, where $D(\Omega) = -2\log \mathcal{L}(\gamma, \mathbf{e})$ which is referred to as the deviance function, and $p_D = \bar{D} - D(E_{\Omega|\mathcal{D}}\{\Omega\})$ which is a measure of model complexity. Note that the DIC is also readily computed from MCMC output.

3.3. Results. Table 3 shows the DIC and LPML for all models under consideration. All models under case I provide significantly better prediction as measured by both DIC and LPML, with differences in the range of 20–55 for DIC and 10–25 for LPML, which indicates that the determination of the cut-point vector for the baseline hazard plays an important role on model prediction and fit. Comparing the frailty specifications in Model 1 across all cases, the DIC and LPML show the same trend for goodness of fit, with the proposed model based on the LDTFP frailty model outperforming both the MPT and Gaussian models, although the differences are only in the range of 1–4. Comparing between Model 2 and Model 3, the proposed model is always preferred in terms of LPML, while the MPT model is slightly better than others in term of DIC under Model 2. Comparing all the proposed models across Model 1–Model 3, the results indicate that Model 1 always performs best. Overall, allowing the frailty distribution to change with county-level covariates (especially RUCC) does improve model prediction according to LPML. In what follows, we present the results under case I.

Table 4 presents posterior medians and equal-tailed 95% credible intervals (CI) for main effects (components of ξ) under Model 1–Model 3, with covariate-adjusted frailties, and compares the individual-level covariate effects, that is, (ξ_1, ξ_2, ξ_3) , to those obtained by Zhao, Hanson and Carlin

TABLE 4
Iowa SEER data: Posterior medians (95% credible intervals) of fixed effects ξ from various models

Predictor	Model 1	Model 2	Model 3	CAR	Cox
ξ_1 (Age)	0.019 (0.013, 0.025)	0.020 (0.014, 0.026)	0.020 (0.014, 0.026)	0.018 (0.012, 0.025)	0.019 (0.013, 0.025)
ξ_2 (Regional)	0.27 (0.03, 0.49)	0.27 (0.03, 0.47)	0.27 (0.05, 0.50)	0.22 (0.01, 0.49)	0.30 (0.08, 0.52)
ξ_3 (Distant)	1.64 (1.43, 1.88)	1.67 (1.43, 1.89)	1.65 (1.43, 1.89)	1.65 (1.40, 1.93)	1.64 (1.42, 1.87)
ξ_{x_1} (RUCC)	-0.105 (-0.185, -0.041)		-0.082 (-0.179, 0.011)		
ξ_{x_2} (Income)		0.042 (0.003, 0.084)	0.011 (-0.040, 0.066)		

(2009), under the standard nonfrailty Cox model and the Cox frailty model that has a MPT prior for the baseline survival, centered at the log-logistic family, and with conditionally autoregressive (CAR) county-level spatial frailties. The best fitting Cox model reported by Zhao, Hanson and Carlin (2009) has an LPML of -2226 . Therefore, the pseudo Bayes factor for the proposed model versus the CAR model is $e^{2226-2222} \approx 55$, implying that the proposed model predicts about 55 times better than the model with CAR frailties. In addition, the proposed model offers a unique interpretation. The posterior medians and 95% CIs for all individual-level effects change little across the different versions of the proposed model, indicating that the Cox regression estimates are reasonably stable for these data, except for the estimate of “Regional stage,” for which the CAR model 95% CI is much wider than those under the considered versions of the proposed model. This may be partly due to the benefit of including county-level covariates. The best model according to LPML, Model 1, indicates that all the individual-level effects are significant at the 0.05 level. Higher age at diagnosis increases the hazard within each county. For instance, women are about $e^{0.019 \times 20} \approx 1.46$ times more likely to die from breast cancer than those twenty years younger who have the same disease stage and live in the same county. Compared with women having local stage of disease, women of the same age and living in the same county are $e^{0.27} \approx 1.31$ times more likely to die if their cancer is detected at the regional stage, and $e^{1.64} \approx 5.16$ times more likely to die if detected at the distant stage. We additionally present the fixed effects under the marginal PH model (i.e., using the R function `coxph` with option `cluster`) across Model 1–Model 3 in Appendix E of the supplementary material [Zhou et al. (2015)]. Note that the coefficient estimates under the marginal PH model have population-averaged interpretations and cannot be

directly compared with those fitted from the proposed frailty PH model due to different model structures.

Regarding the effect of county-level covariates, living in a higher median household income or urban counties is associated with poorer survival after a breast cancer diagnosis. For example, the results under Model 1 indicate that after controlling for individual covariates and county, the hazard rate of women living in urban counties (with RUCC = 2) will be $e^{0.105 \times 7} \approx 2$ times larger than that of women in rural counties (with RUCC = 9). The results under Model 2 imply that after controlling for individual covariates and frailties, women have about a 1.7 times larger hazard rate if they live in median household income counties of \$35,301 compared with median household income of \$23,354 (see also Figure 2). Under Model 3, the results indicate that when both the county-level covariates are included simultaneously, their independent effects are attenuated, partly due to the multicollinearity between them (see the middle two plots in Figure 3).

We obtain the fitted predictive frailty densities for both \mathbf{e}_i (median-zero) and $\mathbf{e}_i + \mathbf{x}'_i \boldsymbol{\xi}_x$ (full distribution) and survival curves for women with mean entry age 68.8 years and distant stage of disease who live in the counties with different levels of median household income or RUCC, under the different versions of the proposed model. The three levels are chosen from the 5%, 50% and 95% quantiles of each covariate value. The results are reported in Figures 2 and 3. Under our best fitting, Model 1 (see left three plots in Figure 2), the results indicate that higher values of RUCC increase the frailty variance and suggest a non-Gaussian shape (upper); we also see overall higher frailty after mixing over the location shift $\mathbf{x}'_i \boldsymbol{\xi}_x$ (middle) and so poorer survival (lower) in urban counties. Increasing heterogeneity as ruralness increases under Model 1 translates into increasing association among those living in more rural counties versus urban. In Appendix E of the supplementary material [Zhou et al. (2015)], Kendall's tau is computed and plotted as a function of RUCC for individuals with mean entry age 68.8 years and distant stage. Kendall's tau increases *by a factor of three* as RUCC goes from 2 to 9. Note that under a traditional gamma frailty model the association is static.

Under Model 2, the frailty densities only slightly change compared with Model 1, but we do see poorer survival in counties with higher median household income. Figure 3 demonstrates that after adjusting individual covariates and median household income (right three plots), there is little effect of RUCC on either predictive frailty densities or survival curves; while after adjusting for RUCC (left three plots), the effect of median household income is almost negligible. In Appendix E of the supplementary material [Zhou et al. (2015)], the survival curves are also compared with those obtained under the marginal PH model. Overall, the marginal PH model under-predicts survival time up to about 1 month, for example, it gives estimates of median

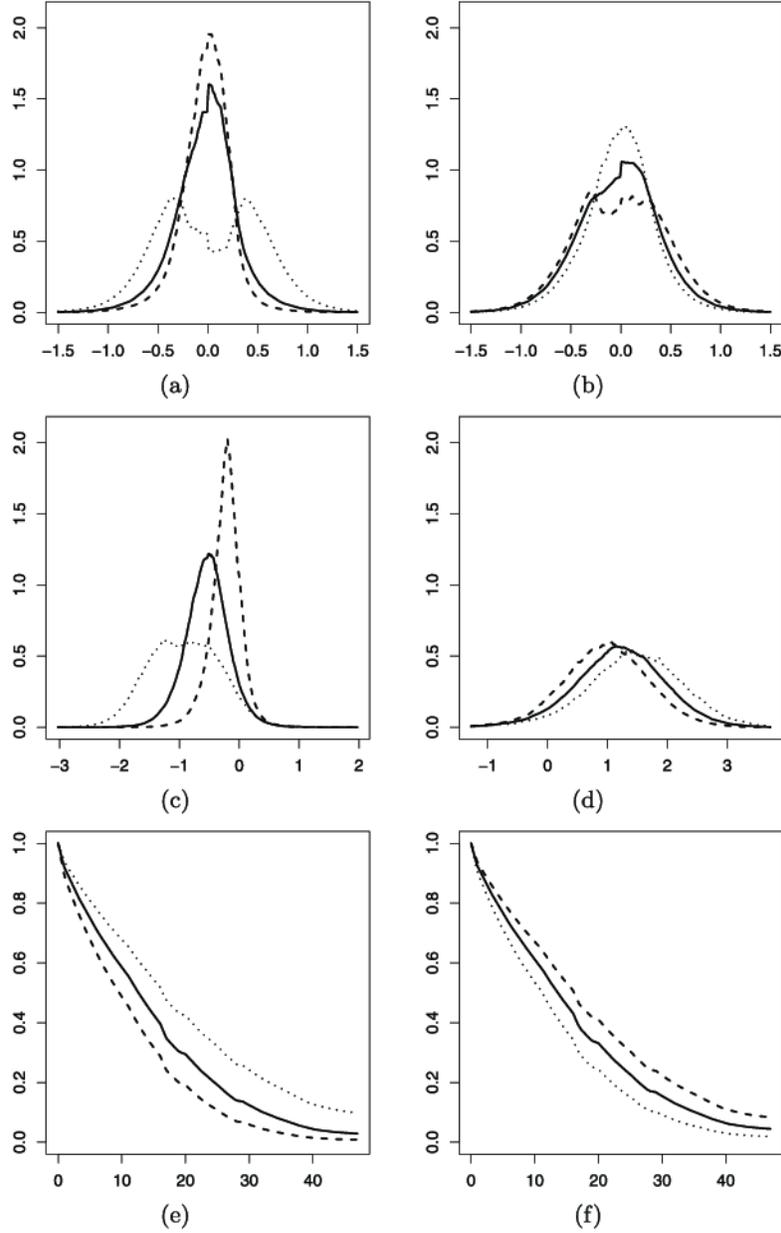


FIG. 2. Iowa SEER data: Fitted predictive frailty densities [panels (a) and (b)], frailty densities with location shifts [panels (c) and (d)] and survival curves [panels (e) and (f)] for women with mean entry age 68.8 years and distant stage of disease from different county covariate levels under Model 1 [panels (a), (c) and (e)] and Model 2 [panels (b), (d) and (f)]. In panels (a), (c) and (e), the results for RUCC = 2, 5 and 9 are displayed as dashed, continuous and dotted lines, respectively. In panels (b), (d) and (f), the results for Income = 23.354, 29.176 and 35.301 are displayed as dashed, continuous and dotted lines, respectively.

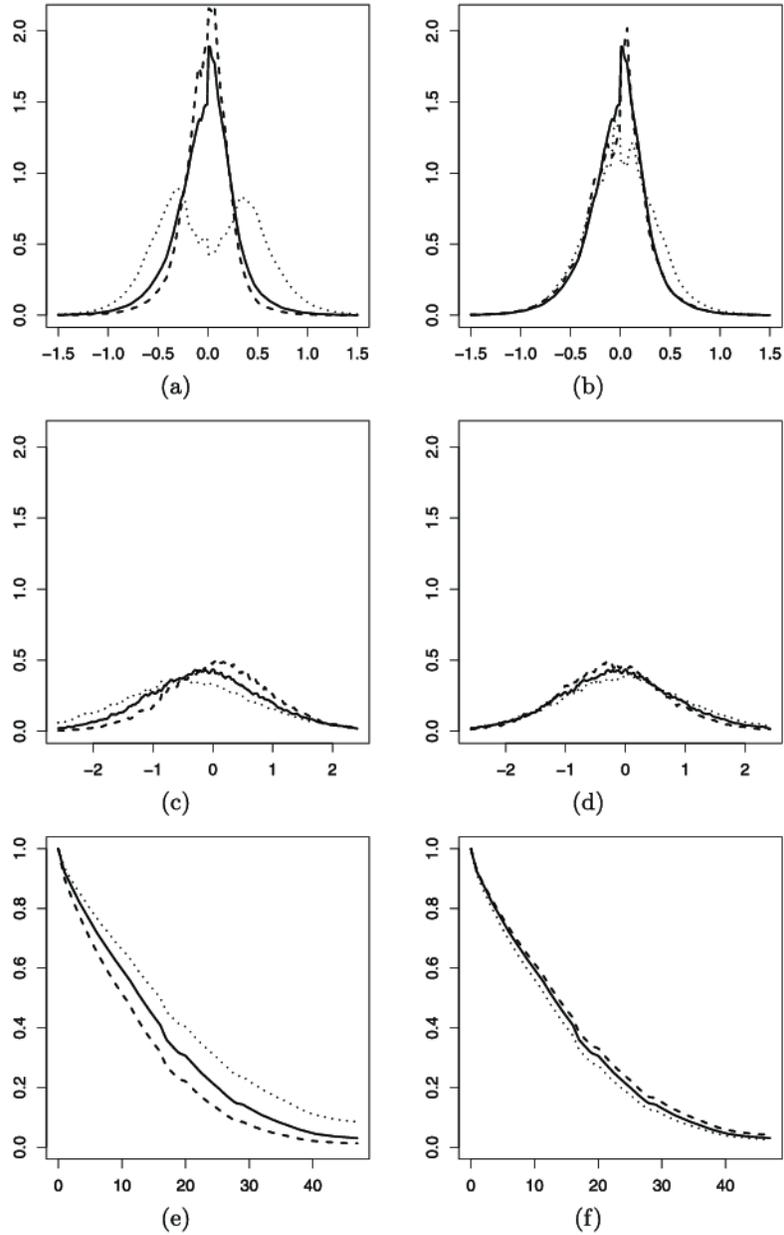


FIG. 3. Iowa SEER data: Fitted predictive frailty densities [panels (a) and (b)], frailty densities with location shifts [panels (c) and (d)] and survival curves [panels (e) and (f)] for women with mean entry age 68.8 years and distant stage of disease from different county covariate levels under Model 3. In panels (a), (c) and (e), the results for $\text{RUCC} = 2, 5$ and 9 are displayed as dashed, continuous and dotted lines, respectively. In panels (b), (d) and (f), the results for $\text{Income} = 23.354, 29.176$ and 35.301 are displayed as dashed, continuous and dotted lines, respectively.

survival a month less, compared with our proposed model for patients with mean entry age 68.8 years and distant stage of disease who live in the same county. This may be partly due to the fact that the marginal PH model averages over the changing behavior of the frailty distribution over the ruralness measure.

It is widely known that access to quality care and screening for breast cancer is more readily available to those with greater financial means and/or those living in urban areas. Therefore, our findings of increased survival for poorer and more rural counties for this cohort are initially puzzling. However, hormone replacement therapy (HRT) increased about 150% in the 1990s [Wysowski and Governale (2005)], after several observational studies linked HRT to prevention of osteoporosis and protection from heart disease. However, this increasing use of HRT abated suddenly in 2002, when the Women's Health Initiative clinical trial linked HRT to aggressively invasive breast cancer [Rossouw et al. (2002)]. In fact, overall breast cancer incidence rates peaked in 1999. Between 2001 and 2004 overall invasive breast cancer incidence declined, but fell much more drastically among women living in urban versus rural counties, and among women living in low-poverty versus high-poverty counties. Hausauer et al. (2009) attribute this discrepancy to greater use of postmenopausal estrogen/progestin hormone replacement therapy among more affluent women and/or women living in urban counties up until 2002, when the Women's Health Initiative trial was stopped prematurely on May 31, 2002, according to Rossouw et al. (2002), "*... because the test statistic for invasive breast cancer exceeded the stopping boundary for this adverse effect and the global index statistics supported risks exceeding benefits.*" It is plausible that increased risk (i.e., stochastically larger frailties) in more affluent and more urban counties has to do with a larger proportion of women being prescribed HRT in the late 1980s and 1990s. Further exploratory analyses on other cohorts of SEER Iowan breast cancer data (1975–1979, 1980–1984, 1985–1989 and 1990–1994) show a reversal of the effects of income and ruralness, agreeing with intuition. Paralleling our study, Krieger, Chen and Waterman (2010) used county-level census data on income and found rising and falling breast cancer incidence rates for the SEER data over the range 1992–2005 for caucasian women living in high-income counties, which "*mirrored the social patterning of hormone therapy use.*"

In a longer follow-up study of the Women's Health Initiative trial, Chlebowski et al. (2010) found that those on estrogen plus progestin compared to placebo had about 25% higher incidence of invasive breast cancer. Among those diagnosed with breast cancer, the two treatment arms had similar histology, but the estrogen plus progestin group were 78% more likely to have cancers that had spread to lymph nodes than placebo, and the estrogen plus progestin group were about twice as likely to die from breast cancer versus

placebo. It would appear that hormone replacement therapy fortified the virulence of breast cancer, significantly increasing both incidence and mortality. This same study showed an impressive 7% one-year drop in incidence right after the Women’s Health Initiative study was prematurely stopped and the medical community warned of a possible link between hormone replacement therapy and breast cancer.

4. Simulation studies. We performed a simulation study to assess the performance of the proposed model. The simulated data are also used to compare the proposed approach with existing models. Specifically, we consider the GF approach described in Section 2.1 and the positive stable frailty Cox model proposed by Liu, Kalbfleisch and Schaubel (2011). Under this latter model, the shape parameter is allowed to depend on cluster-level covariates. In terms of our notation, they assumed that the conditional hazard function of T_{ij} is

$$(2) \quad \lambda(t|\tilde{\mathbf{w}}_{ij}, \mathbf{x}_i, e_i) = \lambda_{0i}(t) \exp(\tilde{\mathbf{w}}_{ij}' \tilde{\boldsymbol{\xi}}_i + e_i),$$

where the baseline hazard functions $\lambda_{0i}(t)$ and regression parameters $\tilde{\boldsymbol{\xi}}_i$ are cluster-specific, and $\exp(e_i)$ follows a positive stable distribution with shape parameter $\alpha_i \in (0, 1)$, relying on the cluster-level covariates vector \mathbf{x}_i through a logit link function, denoted by $PS(\alpha_i)$. They did not deal with this model directly, but rather derived the marginal model

$$(3) \quad \lambda(t|\tilde{\mathbf{w}}_{ij}, \mathbf{x}_i) = h_0(t) \exp(\tilde{\mathbf{w}}_{ij}' \boldsymbol{\eta})$$

by imposing the restrictions $\boldsymbol{\eta} = \alpha_i \tilde{\boldsymbol{\xi}}_i$, $H_0(t) = \{\Lambda_{0i}(t)\}^{\alpha_i}$, where $H_0(t) = \int_0^t h_0(s) ds$ and $\Lambda_{0i} = \int_0^t \lambda_{0i}(s) ds$. In other words, they essentially fitted the above marginal Cox model by maximizing the pseudo partial likelihood under the working independence assumption [Wei, Lin and Weissfeld (1989)], and then utilized the imposed constraints to estimate the parameters in the frailty model. Although they considered a more flexible conditional Cox model, they made many assumptions to get the marginal model, some of which are difficult to check in practice. Moreover, they faced a nonidentifiability problem when a cluster-level covariate was included in the conditional Cox model, so cluster-level covariates had to be excluded from the marginal model as well, leading to potentially poorer prediction of the marginal survival function. Their method, referred to below as PSF, will be compared with our approach focusing on the prediction of survival functions in the simulation studies. A comparison of the two methods for the fixed effect estimates cannot be conducted, since they have different model structures. We conducted the simulation study in R. The GF and PSF approaches were implemented by using the function `coxme` and `coxph` (with the option `cluster`), respectively, included in the R packages `coxme` and `survival`.

4.1. *Simulation settings.* Two scenarios for the frailty distributions were considered. In the first case, referred to as Scenario I, a covariate-dependent family of distributions is considered, where the density shape evolves from one mode to two as the cluster-specific covariate x increases its value; this mirrors the effect of RUCC in panel (a) of Figure 3 for Model 1. In the second case, referred to as Scenario II, a covariate-dependent positive stable distribution is considered. The specific distributional forms for each setting were the following:

Scenario I. $e_i|x_i \stackrel{\text{i.i.d.}}{\sim} 0.5N(-e^{0.4x_i}, 1) + 0.5N(e^{0.4x_i}, 1)$, $x_i \stackrel{\text{i.i.d.}}{\sim} U(-3, 3)$.

Scenario II. $\exp(e_i)|x_i \stackrel{\text{i.i.d.}}{\sim} PS(\alpha_i)$, $\alpha_i = 1/(1 + e^{-0.5-0.5x_i})$, $x_i \stackrel{\text{i.i.d.}}{\sim} U(0, 2)$.

Note that the first setting is not a particular case of the proposed model; the second setting, chosen from the simulation study of Liu, Kalbfleisch and Schaubel (2011), is included to evaluate the behavior of the proposed approach when the PSF model is correct.

Given the frailties, the data under Scenario I were simulated from the conditional PH model (1) with $\lambda_0(t) = 1$, $\mathbf{w}_{ij} = (w_{1ij}, w_{2ij}, x_i)'$ and $\boldsymbol{\xi} = (\xi_1, \xi_2, \xi_x)' = (1.0, 0.5, 1.0)'$; the data under Scenario II were simulated from the PSF model (2) with $\tilde{\mathbf{w}}_{ij} = (w_{1ij}, w_{2ij})'$, $\boldsymbol{\eta} = (1, 0.5)'$ and $H_0(t) = t$. For each simulation scenario, 200 replicates of the data set were generated by assuming the following: $w_{1ij} \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$ and $w_{2ij} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(0.5)$, $i = 1, \dots, 100$, $j = 1, \dots, 10$. In each case, a noninformative censoring scheme was considered, where the censoring times were simulated from an $U(0.25, 4)$ distribution, so that the censoring rate is approximately 35% under Scenario I and 25% under Scenario II.

For each data set, the GF approach was employed, yielding point estimates of $\boldsymbol{\xi}$, $\text{var}(e_i)$ and e_i , which we denote by $\hat{\boldsymbol{\xi}}^{(0)}$, $\hat{\theta}^{2(0)}$ and $\hat{e}_i^{(0)}$, respectively. Based on these point estimates, the predictive survival function was calculated as follows:

$$(4) \quad \hat{S}_{\text{GF}}(t|\mathbf{w}) = n^{-1} \sum_{i=1}^n \exp\{-\hat{\Lambda}_0^{(0)}(t) \exp\{\mathbf{w}'\hat{\boldsymbol{\xi}}^{(0)} + \hat{e}_i^{(0)}\}\},$$

where $\hat{\Lambda}_0^{(0)}(t)$, depending on $\hat{e}_i^{(0)}$'s, denotes Breslow's estimator of $\Lambda_0(t)$ [see, e.g., Therneau, Grambsch and Pankratz (2003), Section 2]. We then fitted the proposed model, by considering $J = 4$, $K = 10$, $\tau_1 = 1.001$, $\tau_2 = 1.001\hat{\theta}^{2(0)}$, $a_c = 1$, $b_c = 1$, $\boldsymbol{\gamma}_0 = \mathbf{0}_{13}$ and $\mathbf{S}_0 = 10^3 \times \mathbf{I}_{13}$. For each data set a single Markov chain of length 55,000 was obtained by using the algorithm described in Appendix B of the supplementary material [Zhou et al. (2015)]. A burn-in period of 5000 scans was considered, and 5000 samples were retained for posterior inferences. The posterior mean of the corresponding parameters are denoted by $\hat{\boldsymbol{\xi}}$, $\hat{\theta}^2$, $\hat{g}(e|\mathbf{x})$ and $\hat{S}(t|\mathbf{w})$. Finally, the PSF approach was

considered but including the cluster-level covariates in the linear predictor, and the associated predictive survival function, based on Breslow's estimator of the underlying baseline hazard function, was obtained and is denoted by $\hat{S}_{\text{PSF}}(t|\mathbf{w})$.

The competing approaches were compared regarding the estimation of the regression coefficients and also compared by computing the weighted integrated squared error (ISE) for the estimated survival distributions, given by

$$\int_0^{\infty} \{\hat{S}_m(t|\mathbf{w}) - S(t|\mathbf{w})\}^2 f_T(t|\mathbf{w}) dt,$$

where $\hat{S}_m(t|\mathbf{w})$, $S(t|\mathbf{w})$ and $f_T(t|\mathbf{w})$ are the estimated survival function, the true survival function and density function, respectively, for a subject with covariate vector \mathbf{w} .

4.2. Simulation results. The results for the regression coefficients using the proposed model and the GF approach under Scenario I are given in Table 5, where the bias of the corresponding point estimators, the Monte Carlo mean of the posterior standard deviation/standard error (MEAN-SD), the Monte Carlo standard deviation of the point estimates (SD-MEAN) and the Monte Carlo coverage probability (CP) of the 95% credible interval/confidence intervals are presented. The results suggest that the posterior means of ξ are almost unbiased estimators and that the observed bias for ξ_x under the proposed approach is much smaller than the corresponding value obtained under the GF approach. Moreover, under the proposed model, the MEAN-SD and the SD-MEAN values are in fairly close agreement, indicating that the posterior standard deviation is an unbiased estimator of the frequentist standard error. Finally, the CPs are all around the nominal 95%.

TABLE 5

Simulation data—Scenario I: True value, bias of the point estimator, mean (across Monte Carlo simulations) of the posterior standard deviations/standard errors (MEAN-SD), standard deviation (across Monte Carlo simulations) of the point estimator (SD-MEAN) and Monte Carlo coverage probability for the 95% credible interval/confidence interval (CP) for the regression parameters. The results are presented under the proposed model and under the GF approach

Para- meters	Proposed model					GF model			
	True	BIAS	MEAN-SD	SD-MEAN	CP	BIAS	MEAN-SD	SD-MEAN	CP
ξ_1	1.0	0.011	0.052	0.054	0.930	-0.011	0.051	0.059	0.910
ξ_2	0.5	0.008	0.088	0.090	0.945	-0.003	0.088	0.091	0.950
ξ_x	1.0	-0.009	0.141	0.126	0.965	-0.052	0.083	0.142	0.775

TABLE 6

Simulated data—Scenario II: Monte Carlo mean (Monte Carlo standard deviation) for the ISE of the survival function for two different predictor values. The results for the different approaches under both simulation scenarios are presented. The numbers correspond to 10^3 times the original values

Scenario	(w_1, w_2, x)	Proposed model	GF model	PSF model
I	(2, 1, -2)	2.02 (2.48)	4.37 (3.46)	6.28 (3.49)
	(0, 1, 2)	1.94 (2.53)	10.5 (6.86)	14.3 (10.9)
II	(2, 1, 0.5)	3.17 (4.66)	3.13 (3.33)	2.19 (2.26)
	(0, 1, 1.5)	0.96 (1.18)	0.89 (1.22)	0.83 (1.10)

The same does not hold for GF, which substantially underestimates the standard error for ξ_x , leading to low coverage probabilities.

The average of the estimated frailty distributions and survival functions across simulated data sets for some specific covariate values are presented in Figure 4 for Scenario I and in Figure 5 for Scenario II. The results in Scenario I reveal that the proposed model roughly captures the modal behavior of the covariate-dependent frailty distributions. Although not perfect, the proposed model performs remarkably well given that only $n = 100$ imperfectly-observed observations were generated for each data set. The situation is much less favorable for the GF approach, which fails to correctly capture the shape of the frailty distributions, leading to poor estimated survival functions. This behavior is likely driving the underestimation of survival noted in the SEER analysis. As expected, the PSF approach also suffers from bad prediction since the underlying assumption for frailty distribution is violated. The results in Scenario II show that the proposed model is still able to capture the frailty distributional shape even when the data were truly generated from the PSF model. Regarding the estimated survival curves, the results suggest that essentially no differences among the three methods are observed; all estimated functions are close to the truth, indicating that there is little price to be paid when using the proposed model for the clustered survival data that were truly generated from the PSF model.

The results of the comparison of the estimated survival curves in terms of ISE are presented in Table 6, where the Monte Carlo mean and standard deviations for the ISE for two different predictor values are given. The results under Scenario I show a close agreement with the observed for the regression coefficients; the proposed model substantially outperforms the other two methods in terms of smaller means and standard deviations of the ISE. Even under Scenario II, the proposed model still provides almost the same results as the PSF model in terms of ISE.

In Appendix D of the supplementary material [Zhou et al. (2015)], additional simulation results are presented which show that, under Scenario I, for

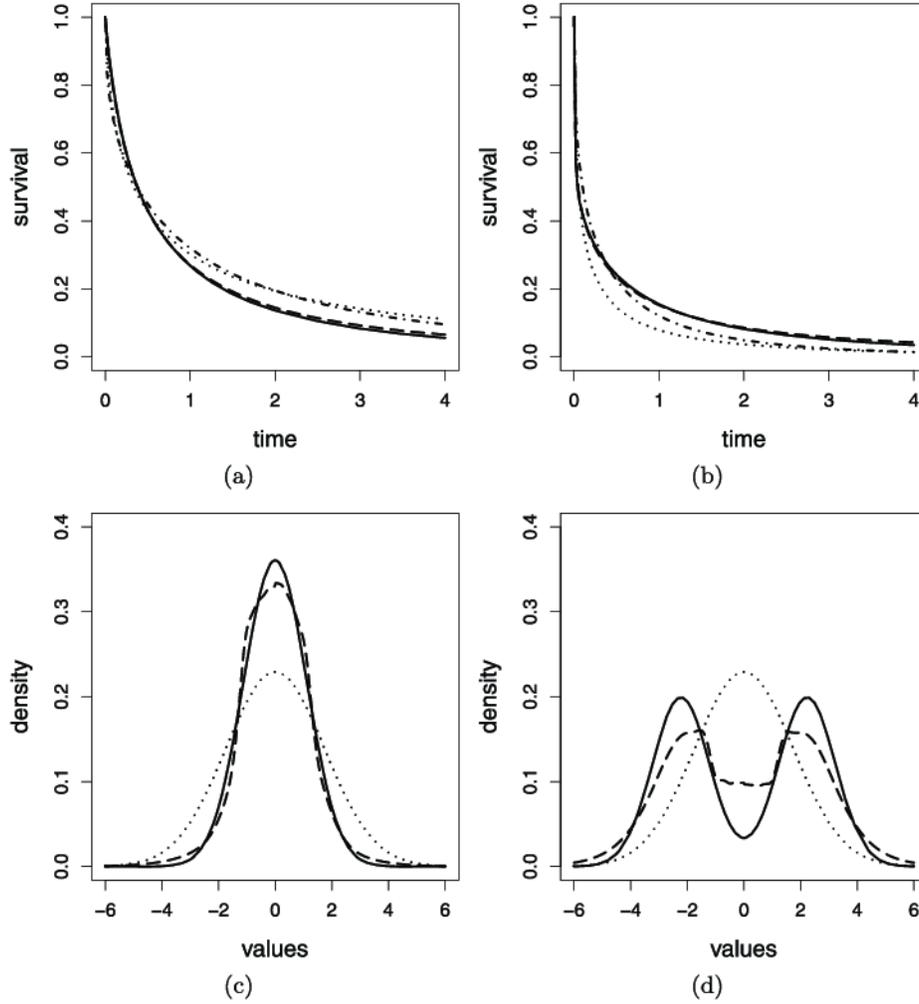


FIG. 4. *Simulated data—Scenario I: Mean, across simulations, of the posterior mean of the survival and frailty density functions under the proposed model. Panels (a) and (b) show the results for the survival functions. Panels (c) and (d) show the results for the frailty densities. Panels (a) and (c) show the results for covariate values (2, 1, -2). Panels (b) and (d) show the results for covariate values (0, 1, 2). The true curves are represented by continuous lines. The results under the proposed model are represented by dashed lines. The results under the exchangeable Gaussian frailty model are represented by dotted lines. In panels (a) and (b) the results obtained under the PSF approach are represented by dot-dashed lines.*

larger sample sizes better estimates of the frailty distributions are obtained and that the approach is not affected by the choice of J in the specification of the LDTFP model. For further comparison, we also fitted the exchange-

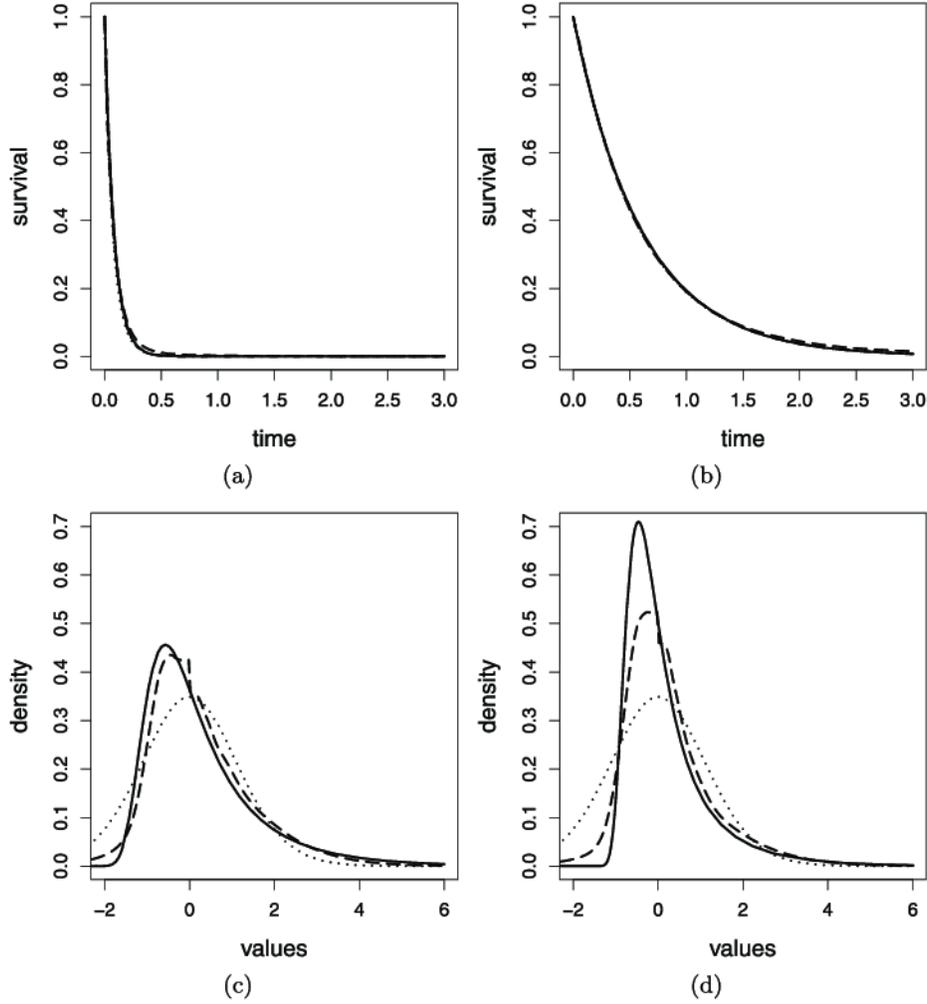


FIG. 5. *Simulated data—Scenario II: Mean, across simulations, of the posterior mean of the survival and frailty density functions under the proposed model. Panels (a) and (b) show the results for the survival functions. Panels (c) and (d) show the results for the frailty densities. Panels (a) and (c) show the results for covariate values (2, 1, 0.5). Panels (b) and (d) show the results for covariate values (0, 1, 1.5). The true curves are represented by continuous lines. The results under the proposed model are represented by dashed lines. The results under the exchangeable Gaussian frailty model are represented by dotted lines. In panels (a) and (b) the results obtained under the PSF approach are represented by dot-dashed lines.*

able mixture of Polya trees (MPT) [Hanson (2006b)] frailty Cox's model using the function `PTglm` available in `DPpackage` [Jara et al. (2011)] under Scenario I, in which the results show that our approach outperforms the

MPT, and considered a third scenario favorable to the GP approach, where the results show that our method pays little price for the extra generality when using the proposed model when normality and exchangeability are valid assumptions. Overall, the proposed approach provides a flexible way to capture the heterogeneity in the frailty distribution, provides superior prediction, and yields an essential improvement for the estimation of population effects, especially when the intra-cluster correlation (or variability in frailties) is relatively large. When the frailty variances are small across clusters, the proposed approach is still recommended due to its flexibility.

5. Concluding remarks. Very limited work has been done on covariate-adjusted frailty survival models for clustered time-to-event data. Liu, Kalbfleisch and Schaubel (2011) proposed a stratified Cox model with positive stable frailties, where the shape parameter of the frailty distribution is allowed to depend on cluster-level covariates. However, they essentially fitted a marginal Cox model, and then utilized the positive stable assumption and some imposed constraints to estimate the parameters in their proposed model. The model proposed in this paper cleanly separates population-level effects from the cluster-level effects, which determine the shape of the frailty distribution. Frailty density shape is modeled using a tractable median-zero LDTFP prior. Other nonparametric density regression approaches could also be considered; however, model identifiability requires a location constraint such as mean-zero or median-zero. The proposed model provides a natural generalization of the conventional PH model with parametric or nonparametric exchangeable frailties, and accommodates frailty distribution “evolution” over cluster-level covariates providing superior prediction, as shown in our simulation studies. When data are truly generated according to an exchangeable Gaussian frailty PH model or the model of Liu, Kalbfleisch and Schaubel (2011), our model does about the same as the underlying true model in terms of fixed effects and/or marginal survival estimations. We illustrate the usefulness of the proposed model with an analysis of a subset of the Iowa SEER breast cancer data, and demonstrate that higher degree of ruralness corresponds to a more bimodal frailty distributional shape with larger variance. In general, the proposed model is more flexible than currently existing frailty PH models, leading to more robust inferences, and thus is recommended. One drawback of the proposed model is that, as currently fit in R, obtaining inference takes longer.

For ease of computation, we have assumed a piecewise constant structure for the baseline hazard function $\lambda_0(t)$ and taken the independent normal prior distributions for $\log(\lambda_k)$'s, so that the baseline hazard heights λ_k and covariate effects ξ can be updated simultaneously. The use of empirically-derived cut-points has permeated much of the Bayesian survival literature for over a decade. Use of Breslow's baseline estimate coupled with the GF

approach led to a greatly increased LPML over the empirical approach. An obvious extension of our current work is to employ a smoothed baseline, for example, using penalized B-splines [Hennerfeind, Brezger and Fahrmeir (2006)], the piecewise exponential with random cut-points [Sahu and Dey (2004)], MPT [Hanson (2006b)], etc. Any of these approaches could improve model fit and prediction, but cannot currently be fitted in the R software. We are currently working on extending the methodology in this paper to other survival models and smoothed baselines.

Acknowledgments. We thank the referees and editors for numerous insightful suggestions which greatly enhanced the readability of this paper.

SUPPLEMENTARY MATERIAL

Supplement to “Modeling county-level breast cancer survival data using a covariate-adjusted frailty proportional hazards model”

(DOI: [10.1214/14-AOAS793SUPP](https://doi.org/10.1214/14-AOAS793SUPP); .pdf). In this online supplemental article we provide (A) technical details on the mixture of linear dependent tailfree processes, (B) a detailed description of the MCMC algorithm, (C) sample R code to analyze the SEER data, (D) additional simulation studies and (E) additional analysis of the SEER data.

REFERENCES

- CHLEBOWSKI, R. T., ANDERSON, G. L., GASS, M., LANE, D. S., ARAGAKI, A. K., KULLER, L. H., MANSON, J. E., STEFANICK, M. L., OCKENE, J., SARTO, G. E. et al. (2010). Estrogen plus progestin and breast cancer incidence and mortality in postmenopausal women. *JAMA* **304** 1684–1692.
- CLAYTON, D. and CUZICK, J. (1985). Multivariate generalizations of the proportional hazards model. *J. Roy. Statist. Soc. Ser. A* **148** 82–117. [MR0806480](#)
- COTTONE, F. (2008). Covariate dependent random effects in survival analysis, Ph.D. dissertation, Univ. degli Studi “Roma Tre.”
- GEISSER, S. and EDDY, W. F. (1979). A predictive approach to model selection. *J. Amer. Statist. Assoc.* **74** 153–160. [MR0529531](#)
- GELFAND, A. E. and DEY, D. K. (1994). Bayesian model choice: Asymptotics and exact calculations. *J. Roy. Statist. Soc. Ser. B* **56** 501–514. [MR1278223](#)
- GOODMAN, L. A. and KRUSKAL, W. H. (1954). Measures of association for cross classifications. *J. Amer. Statist. Assoc.* **49** 732–764.
- GUSTAFSON, P. (1997). Large hierarchical Bayesian analysis of multivariate survival data. *Biometrics* **53** 230–242.
- HANSON, T. E. (2006a). Inference for mixtures of finite Polya tree models. *J. Amer. Statist. Assoc.* **101** 1548–1565. [MR2279479](#)
- HANSON, T. E. (2006b). Modeling censored lifetime data using a mixture of gammas baseline. *Bayesian Anal.* **1** 575–593 (electronic). [MR2221289](#)
- HANSON, T., JOHNSON, W. and LAUD, P. (2009). Semiparametric inference for survival models with step process covariates. *Canad. J. Statist.* **37** 60–79. [MR2509462](#)

- HAUSAUER, A. K., KEEGAN, T. H., CHANG, E. T., GLASER, S. L., HOWE, H. and CLARKE, C. A. (2009). Recent trends in breast cancer incidence in US white women by county-level urban/rural and poverty status. *BMC Medicine* **7** 31.
- HENNERFEIND, A., BREZGER, A. and FAHRMEIR, L. (2006). Geoaddivitive survival models. *J. Amer. Statist. Assoc.* **101** 1065–1075. [MR2324146](#)
- JARA, A. and HANSON, T. E. (2011). A class of mixtures of dependent tail-free processes. *Biometrika* **98** 553–566. [MR2836406](#)
- JARA, A., HANSON, T. E., QUINTANA, F. A., MÜLLER, P. and ROSNER, G. L. (2011). DPpackage: Bayesian semi- and nonparametric modeling in R. *J. Stat. Softw.* **40** 1.
- KALBFLEISCH, J. D. (1978). Non-parametric Bayesian analysis of survival time data. *J. Roy. Statist. Soc. Ser. B* **40** 214–221. [MR0517442](#)
- KRIEGER, N., CHEN, J. T. and WATERMAN, P. D. (2010). Decline in US breast cancer rates after the women’s health initiative: Socioeconomic and racial/ethnic differentials. *American Journal of Public Health* **100** 132–139.
- LAIRD, N. and OLIVIER, D. (1981). Covariance analysis of censored survival data using log-linear analysis techniques. *J. Amer. Statist. Assoc.* **76** 231–240. [MR0624329](#)
- LIU, D., KALBFLEISCH, J. D. and SCHAUBEL, D. E. (2011). A positive stable frailty model for clustered failure time data with covariate-dependent frailty. *Biometrics* **67** 8–17. [MR2898812](#)
- MCCULLOCH, C. E. and NEUHAUS, J. M. (2011). Misspecifying the shape of a random effects distribution: Why getting it wrong may not matter. *Statist. Sci.* **26** 388–402. [MR2917962](#)
- NOH, M., HA, I. D. and LEE, Y. (2006). Dispersion frailty models and HGLMs. *Stat. Med.* **25** 1341–1354. [MR2226790](#)
- QIOU, Z., RAVISHANKER, N. and DEY, D. K. (1999). Multivariate survival analysis with positive stable frailties. *Biometrics* **55** 637–644.
- REICH, B. J., BONDELL, H. D. and WANG, H. J. (2010). Flexible Bayesian quantile regression for independent and clustered data. *Biostatistics* **11** 337–352.
- ROSSOUW, J. E., ANDERSON, G. L., PRENTICE, R. L., LACROIX, A. Z., KOOPERBERG, C., STEFANICK, M. L., JACKSON, R. D., BERESFORD, S. A. A., HOWARD, B. V., JOHNSON, K. C., KOTCHEN, J. M. and OCKENE, J. (2002). Risks and benefits of estrogen plus progestin in healthy postmenopausal women: Principal results from the women’s health initiative randomized controlled trial. *JAMA* **288** 321–333.
- SAHU, S. K. and DEY, D. K. (2004). On a Bayesian multivariate survival model with a skewed frailty. In *Skew-Elliptical Distributions and Their Applications: A Journey Beyond Normality* (M. G. GENTON, ed.) 321–338. Chapman&Hall/CRC, Boca Raton, FL.
- SPIEGELHALTER, D. J., BEST, N. G., CARLIN, B. P. and VAN DER LINDE, A. (2002). Bayesian measures of model complexity and fit. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **64** 583–639. [MR1979380](#)
- SPRAGUE, B. L., TRENTHAM-DIETZ, A., GANGNON, R. E., RAMCHANDANI, R., HAMP-
TON, J. M., ROBERT, S. A., REMINGTON, P. L. and NEWCOMB, P. A. (2011). Socioeconomic status and survival after an invasive breast cancer diagnosis. *Cancer* **117** 1542–1551.
- THERNEAU, T. M. and GRAMBSCH, P. M. (2000). *Modeling Survival Data: Extending the Cox Model*. Springer, New York. [MR1774977](#)
- THERNEAU, T. M., GRAMBSCH, P. M. and PANKRATZ, V. S. (2003). Penalized survival models and frailty. *J. Comput. Graph. Statist.* **12** 156–175. [MR1965213](#)
- TRIPPA, L., MÜLLER, P. and JOHNSON, W. (2011). The multivariate beta process and an extension of the Polya tree model. *Biometrika* **98** 17–34. [MR2804207](#)

- VAUPEL, J. W., MANTON, K. G. and STALLARD, E. (1979). The impact of heterogeneity in individual frailty on the dynamics of mortality. *Demography* **16** 439–454.
- WALKER, S. G. and MALLICK, B. K. (1997). Hierarchical generalized linear models and frailty models with Bayesian nonparametric mixing. *J. Roy. Statist. Soc. Ser. B* **59** 845–860. [MR1483219](#)
- WANG, Z. and LOUIS, T. A. (2004). Marginalized binary mixed-effects models with covariate-dependent random effects and likelihood inference. *Biometrics* **60** 884–891. [MR2133540](#)
- WASSELL, J. T. and MOESCHBERGER, M. L. (1993). A bivariate survival model with modified gamma frailty for assessing the impact of interventions. *Stat. Med.* **12** 241–248.
- WEI, L. J., LIN, D. Y. and WEISSFELD, L. (1989). Regression analysis of multivariate incomplete failure time data by modeling marginal distributions. *J. Amer. Statist. Assoc.* **84** 1065–1073. [MR1134494](#)
- WYSOWSKI, D. K. and GOVERNALE, L. A. (2005). Use of menopausal hormones in the United States, 1992 through June, 2003. *Pharmacoepidemiology and Drug Safety* **14** 171–176.
- YASHIN, A. I. and IACHINE, I. A. (1999). What difference does the dependence between durations make? Insights for population studies of aging. *Lifetime Data Anal.* **5** 5–22. [MR1750339](#)
- ZHAO, L. and HANSON, T. E. (2011). Spatially dependent Polya tree modeling for survival data. *Biometrics* **67** 391–403. [MR2829008](#)
- ZHAO, L., HANSON, T. E. and CARLIN, B. P. (2009). Mixtures of Polya trees for flexible spatial frailty survival modelling. *Biometrika* **96** 263–276. [MR2507142](#)
- ZHOU, H., HANSON, T., JARA, A. and ZHANG, J. (2015). Supplement to “Modeling county level breast cancer survival data using a covariate-adjusted frailty proportional hazards model.” DOI:[10.1214/14-AOAS793SUPP](#).

H. ZHOU
T. HANSON
DEPARTMENT OF STATISTICS
UNIVERSITY OF SOUTH CAROLINA
COLUMBIA, SOUTH CAROLINA 29208
USA
E-MAIL: hansont@stat.sc.edu

A. JARA
DEPARTMENT OF STATISTICS
PONTIFICIA UNIVERSIDAD CATÓLICA DE CHILE
SANTIAGO
CHILE

J. ZHANG
DEPARTMENT OF EPIDEMIOLOGY AND BIostatISTICS
UNIVERSITY OF SOUTH CAROLINA
COLUMBIA, SOUTH CAROLINA 29208
USA