

Regression models generated by gamma random variables with long-term survivors

Edwin M. M. Ortega^a, Gauss M. Cordeiro^b, Elizabeth M. Hashimoto^c,
Adriano K. Suzuki^{1,d}

^aDepartamento de Ciências Exatas, USP, Brazil; ^bDepartamento de Estatística, UFPE, Brazil;
^cDepartamento Acadêmico de Matemáticas, UTFPR, Brazil;
^dDepartamento de Matemática Aplicada e Estatística, USP, Brazil

Abstract

We propose a flexible cure rate survival model by assuming that the number of competing causes of the event of interest has the Poisson distribution and the time for the event follows the gamma-G family of distributions. The extended family of gamma-G failure-time models with long-term survivors is flexible enough to include many commonly used failure-time distributions as special cases. We consider a frequentist analysis for parameter estimation and derive appropriate matrices to assess local influence on the parameters. Further, various simulations are performed for different parameter settings, sample sizes and censoring percentages. We illustrate the performance of the proposed regression model by means of a data set from the medical area (gastric cancer).

Keywords: gamma-G family, lifetime data, long-term survivors, Poisson distribution, regression model, sensitivity analysis

1. Introduction

Models for survival data with a surviving fraction (also known as long-term survival models or cure rate models) occupy an outstanding place in reliability and survival analysis. Cure rate models cover situations where there are sampling units insusceptible to the occurrence of the event of interest. The proportion of such units is termed the cured fraction. In clinical studies, the event of interest can be the death of a patient (which can happen due to different competing causes) or a tumor recurrence (which can be attributed to metastasis-component tumor cells left active after an initial treatment). A metastasis-component tumor cell is a tumor cell with a potential to metastasize (Yakovlev and Tsodikov, 1996). The literature on the subject is significant and increasing. Books by Maller and Zhou (1996) and Ibrahim *et al.* (2001), as well as the review paper by Chen *et al.* (1999), Tsodikov *et al.* (2003), and the article by Cooner *et al.* (2007) represent key references. Alternatively, other works dealt with cure rate models. For example, Balakrishnan and Pal (2012) pioneered an EM algorithm-based likelihood estimation for some cure rate models, Balakrishnan and Pal (2013) investigated log-normal lifetimes and likelihood-based inference for flexible cure rate models based on the COM-Poisson family, Balakrishnan and Pal (2016) proposed the EM-based likelihood inference for flexible cure rate models with Weibull lifetimes, Balakrishnan and Pal (2015a) studied likelihood inference for flexible cure rate models with gamma lifetimes, Balakrishnan and Pal (2015b) derived an EM

¹ Corresponding author: Departamento de Matemática Aplicada e Estatística, Universidade de São Paulo, Avenida Trabalhador São-carlense, 400 - Centro, São Carlos-SP 13566-59, Brazil. E-mail: suzuki@icmc.usp.br

algorithm to estimate the parameters of a flexible cure rate model with generalized gamma lifetime and model discrimination using likelihood and information based methods and Balakrishnan *et al.* (2016) proposed piecewise linear approximations for cure rate models and associated inferential issues.

Some proposals have been made recently in the literature by more general classes of distributions to model lifetimes. For example, Zografos and Balakrishnan (2009) and Ristić and Balakrishnan (2012) proposed a family of univariate distributions generated by gamma random variables. For any baseline cumulative distribution function (cdf) $G(t; \boldsymbol{\eta})$ depending on a parameter vector $\boldsymbol{\eta} = (\eta_1, \dots, \eta_q)^\top$, and $t > 0$, they defined the gamma- G family with probability density function (pdf) $f(t)$ and cdf $F(t)$ given by

$$f(t; a, \boldsymbol{\eta}) = \frac{1}{\Gamma(a)} \{-\log[1 - G(t; \boldsymbol{\eta})]\}^{a-1} g(t; \boldsymbol{\eta}) \quad (1.1)$$

and

$$F(t; a, \boldsymbol{\eta}) = \frac{\gamma(a, -\log[1 - G(t; \boldsymbol{\eta})])}{\Gamma(a)} = \frac{1}{\Gamma(a)} \int_0^{-\log[1 - G(t; \boldsymbol{\eta})]} v^{a-1} e^{-v} dv, \quad (1.2)$$

respectively, for $a > 0$, where $g(t) = dG(t)/dt$, $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$ denotes the gamma function, and $\gamma(a, z) = \int_0^z t^{a-1} e^{-t} dt$ denotes the incomplete gamma function. The corresponding hazard rate function (hrf) is given by

$$h(t; a, \boldsymbol{\eta}) = \frac{\{-\log[1 - G(t; \boldsymbol{\eta})]\}^{a-1} g(t; \boldsymbol{\eta})}{\bar{\gamma}(a, -\log[1 - G(t; \boldsymbol{\eta})])},$$

where $\bar{\gamma}(a, z) = \int_z^\infty v^{a-1} e^{-v} dv$ denotes the complementary incomplete gamma function. The gamma- G family has the same parameters of the G distribution plus an additional shape parameter $a > 0$. If T is a random variable with pdf (1.1), we write $T \sim \text{gamma-}G(a)$. Each new gamma- G distribution can be obtained from a specified G distribution. For $a = 1$, the G distribution is a basic exemplar of the gamma- G distribution with a continuous crossover towards cases with different shapes (for example, a particular combination of skewness and kurtosis).

Zografos and Balakrishnan (2009), Ristić and Balakrishnan (2012), and Nadarajah *et al.* (2015) presented several motivations for the gamma- G family of distributions: if $T_{L(1)}, \dots, T_{L(n)}$ are lower record values from a sequence of independent random variables with common pdf $g(\cdot)$, then the pdf of the n^{th} lower record value takes the form (1.1); if Z is a gamma random variable with unit scale parameter and shape parameter $a > 0$, then $X = G^{-1}(1 - \exp(-Z))$ has the pdf (1.1); and, if Z is a log-gamma random variables, then $X = G^{-1}(\exp\{-\exp(Z)\})$ has the pdf (1.1). In order to make the model more flexible, we assume that the time to the event has a gamma- G distribution with density function given by (1.1). In this context, we propose a new family called the Poisson-gamma- G (PG- G) model with cure fraction in competitive-risk structure. Further, we examine statistical inference aspects and formulate the PG- G model with covariates.

After fitting the model, it is important to check the model assumptions and conduct robustness studies to detect possible influential or extreme observations that can cause distortions in the results of the analysis. However, when case-deletion is used, all information from a single subject is deleted at once and it is hard to say if that subject has some influence on a specific aspect of the model. A solution for the earlier problem can be found in the local influence approach, where we reinvestigate how the results of the analysis change under small perturbations in the model or data. Cook (1986)

proposed a general framework to detect the influence of the observations, which indicate how sensitive the analysis is when small perturbations are provoked on the data or in the model. Several authors have applied the local influence methodology in regression analysis with censoring. Ortega *et al.* (2012) considered the problem of assessing local influence in the negative binomial beta Weibull regression model to predict a cure rate for prostate cancer, Hashimoto *et al.* (2013) derived curvature quantities under various perturbation schemes in Neyman type A beta Weibull model for long-term survivors, Fachini *et al.* (2014) used local influence methods to a bivariate regression model with cure fraction and Ortega *et al.* (2015) adapted local influence methods to model a power series beta Weibull regression model to predict breast carcinoma. We propose a similar methodology to detect influential subjects in the PG-G family with cure rate.

The plan of the foregoing sections of the paper is as follows. Section 2 is devoted to model formulation. Some structural properties of the PG-G family for the non-cured population are investigated in Section 3 including moments, generating function and mean deviations. The inference on the model parameters is discussed in Section 4. A simulation study is performed in Section 5. In Section 6, we obtain the normal curvatures for local influence under some usual perturbations. The results of an application to a gastric cancer data set are reported in Section 7. In Section 8, we end with general remarks.

2. The new model

In order to describe the PG-G family with cure fraction, we focus on modeling data, for simplicity, where there are no covariates. The time for the j th competing cause to produce the event of interest (time to event) is denoted by R_j , $j = 1, \dots, M$, where M denotes the unobservable number of competing causes that can produce this event. We consider that, conditional on M , the R_j 's are i.i.d. random variables with cdf $F(t)$ given by (1.2) and survival function $S(t) = 1 - F(t)$. In addition, we assume that R_1, R_2, \dots are independent of M . These latent competing causes M can be assigned to metastasis-component tumor cells left active after an initial treatment. Latent variables cannot be measured directly because they represent a theoretical construction and are not observable. However, they can be represented or measured by other variables. The observable time of the occurrence of the event of interest is defined by $T = \min\{R_1, \dots, R_M\}$ subject to $P(T = \infty | M = 0) = 1$. The i.i.d. assumption on R_1, R_2, \dots is surely a strong one, as remarked by Yakovlev and Tsodikov (1996). This option favors simplicity and analytical tractability at the expense of a more general formulation. Despite this shortcoming, these models have proven to be useful in many real-world applications. Under this setup, Tsodikov *et al.* (2003) and Rodrigues *et al.* (2009) proved that the survival function for the population is given by

$$S_{\text{pop}}(t) = P(T \geq t) = A_p(S(t)), \quad (2.1)$$

where $A_p(\cdot)$ is the probability generating function (pgf) of the number of competing causes (M). The choice of a particular distribution for M implies some consequences. This aspect is investigated in the sequel.

Henceforth, we assume that the number of competing causes has a Poisson distribution with parameter $\tau > 0$ and probability mass function

$$p(m; \tau) = P(M = m | \tau) = \frac{e^{-\tau} \tau^m}{m!}, \quad m = 0, 1, 2, \dots, \quad (2.2)$$

where $\tau > 0$.

The pgf is given by $A_p(s) = \sum_{m=0}^{\infty} p(m; \tau) s^m = \exp[-\tau(1-s)]$, $0 \leq s \leq 1$. Therefore, the improper survival function reduces to

$$S_{\text{pop}}(t; \tau, a, \boldsymbol{\eta}) = \exp \left\{ \frac{-\tau \gamma(a, -\log [1 - G(t; \boldsymbol{\eta})])}{\Gamma(a)} \right\}, \quad (2.3)$$

where $G(t; \boldsymbol{\eta})$ is the baseline cdf. Note that Equation (2.3) is an improper function, since $S_{\text{pop}}(t)$ is not a proper survival function. The cured fraction is $p_0 = \lim_{t \rightarrow \infty} S_{\text{pop}}(t)$, so that, from (2.3), $p_0 = e^{-\tau}$.

The density function corresponding to this model is given by

$$f_{\text{pop}}(t; \tau, a, \boldsymbol{\eta}) = \frac{\tau g(t; \boldsymbol{\eta})}{\Gamma(a)} \{-\log [1 - G(t; \boldsymbol{\eta})]\}^{a-1} \exp \left\{ \frac{-\tau \gamma(a, -\log [1 - G(t; \boldsymbol{\eta})])}{\Gamma(a)} \right\}. \quad (2.4)$$

Equation (2.4) is referred to as the PG-G model with cure fraction in a competitive-risk structure. For $a = 1$, we obtain the Poisson-G (PG) model.

The corresponding population hrf is given by

$$h_{\text{pop}}(t; \tau, a, \boldsymbol{\eta}) = \frac{\tau g(t; \boldsymbol{\eta})}{\Gamma(a)} \{-\log [1 - G(t; \boldsymbol{\eta})]\}^{a-1}.$$

The Weibull, log-logistic (LL) and log-normal distributions are very popular distributions for modelling lifetime data in several areas such as medicine, biology, and engineering. In reliability analysis, the fatigue is a structural damage which occurs when a material is exposed to stress and tension fluctuations. Statistical models allow an analysis of the random variation of the failure time associated to materials exposed to fatigue as a result of different cyclical patterns and strengths. The most popular models used to describe the lifetime process under fatigue are the half-normal (HN), generalized half-normal (GHN) and Birnbaum-Saunders (BS) distributions. Here, we present and study some special cases of this family considering cure fraction and classic Weibull, LL, BS and GHN distributions.

- Poisson-gamma Weibull (PGW) model

The PGW model is defined from (2.3) by taking $G(t; \boldsymbol{\eta})$ and $g(t; \boldsymbol{\eta})$ to be the cdf and pdf of the Weibull distribution. In this case, $G(t; \boldsymbol{\eta}) = 1 - \exp\{-(t/\alpha)^\lambda\}$, with shape parameter $\lambda > 0$, scale parameter $\alpha > 0$ and $\boldsymbol{\eta} = (\alpha, \lambda)^\top$.

- Poisson-gamma log-logistic (PGLL) model

Consider the LL distribution with shape parameter $\alpha > 0$, scale parameter $\lambda > 0$, and cdf given by $G(t; \boldsymbol{\eta}) = 1 - [1 + (t/\alpha)^\lambda]^{-1}$, where $\boldsymbol{\eta} = (\alpha, \lambda)^\top$.

- Poisson-gamma Birnbaum-Saunders (PGBS) model

Let $G(t; \boldsymbol{\eta})$ be the BS distribution with cdf $G(t; \boldsymbol{\eta}) = \Phi[\alpha^{-1}(\sqrt{t/\lambda} - \sqrt{\lambda/t})]$, for $t > 0$, shape parameter $\alpha > 0$, scale parameter $\lambda > 0$, where $\boldsymbol{\eta} = (\alpha, \lambda)^\top$ and $\Phi(\cdot)$ is the standard normal cdf.

- Poisson-gamma generalized half-normal (PGGHN) model

The baseline density $g(t; \boldsymbol{\eta}) = \sqrt{2/\pi}(\alpha/t)(t/\lambda)^\alpha \exp[-1/2(t/\lambda)^{2\alpha}]$ leads to the PGGHN model, where $\alpha > 0$ is the shape parameter, $\lambda > 0$ is the scale parameter and $\boldsymbol{\eta} = (\alpha, \lambda)^\top$.

In Table 1, we list special models of the PG-G family with cure fraction. For $a = 1$, we obtain the following models with cure fraction: Poisson Weibull (PW), Poisson log-logistic (PLL), Poisson Birnbaum-Saunders (PBS) and Poisson generalized-half normal (PGHN) distributions.

Table 1: Some members of the Poisson-gamma-G (PG-G) family with cure fraction, where $\boldsymbol{\eta} = (\alpha, \lambda)^\top$

Model	$f_{\text{pop}}(t; \tau, a, \boldsymbol{\eta})$	$S_{\text{pop}}(t; \tau, a, \boldsymbol{\eta})$
PGW	$\frac{\tau \lambda t^{\lambda-1} \exp\{-(t/\alpha)^\lambda\}}{\alpha^\lambda \Gamma(a)} \left(\frac{t}{\alpha}\right)^{\lambda(a-1)}$ $\exp\left\{\frac{-\tau \gamma [a, (t/\alpha)^\lambda]}{\Gamma(a)}\right\}$	$\exp\left\{\frac{-\tau \gamma [a, (t/\alpha)^\lambda]}{\Gamma(a)}\right\}$
PGLL	$\frac{\tau \lambda t^{\lambda-1} [1+(t/\alpha)^\lambda]^{-2}}{\alpha^\lambda \Gamma(a)} \left\{\log\left[1 + \left(\frac{t}{\alpha}\right)^\lambda\right]\right\}^{a-1}$ $\exp\left\{\frac{-\tau \gamma (a, \log[1+(t/\alpha)^\lambda])}{\Gamma(a)}\right\}$	$\exp\left\{\frac{-\tau \gamma (a, \log[1+(t/\alpha)^\lambda])}{\Gamma(a)}\right\}$
PGBS	$\frac{\tau (t+\lambda) t^{\frac{3}{2}}}{2 \sqrt{2} \lambda \tau \alpha \Gamma(a)} \exp\left\{\frac{-1}{2 \alpha^2} \left[\left(\frac{t}{\lambda}\right) + \left(\frac{\lambda}{t}\right) - 2\right]\right\}$ $\left\{-\log\left\{1 - \Phi\left[\frac{1}{\alpha} \left(\sqrt{\frac{t}{\lambda}} - \sqrt{\frac{\lambda}{t}}\right)\right]\right\}\right\}^{a-1}$ $\exp\left\{\frac{-\tau \gamma (a, -\log[1-\Phi[(1/\alpha)(\sqrt{t/\lambda} - \sqrt{\lambda/t})]]]}{\Gamma(a)}\right\}$	$\exp\left\{\frac{-\tau \gamma (a, -\log[1-\Phi[(1/\alpha)(\sqrt{t/\lambda} - \sqrt{\lambda/t})]]]}{\Gamma(a)}\right\}$
PGGHN	$\frac{\tau \sqrt{2}}{\sqrt{\pi} \Gamma(a)} \left(\frac{\alpha}{t}\right) \left(\frac{t}{\lambda}\right)^\alpha \exp\left[-\frac{1}{2} \left(\frac{t}{\lambda}\right)^{2\alpha}\right]$ $\left\{-\log\left\{2 - 2\Phi\left[\left(\frac{t}{\lambda}\right)^\alpha\right]\right\}\right\}^{a-1}$ $\exp\left\{\frac{-\tau \gamma (a, -\log\{2-2\Phi[(t/\lambda)^\alpha]\})}{\Gamma(a)}\right\}$	$\exp\left\{\frac{-\tau \gamma (a, -\log\{2-2\Phi[(t/\lambda)^\alpha]\})}{\Gamma(a)}\right\}$

PGW = Poisson-gamma Weibull; PGLL = Poisson-gamma log-logistic; PGBS = Poisson-gamma Birnbaum-Saunders; PGGHN = Poisson-gamma generalized half-normal.

2.1. The PG-G model for the non-cured population

Numerous classical distributions have been extensively used for modeling data in several areas. There is a clear need for extended forms of classical distributions in applied areas such as biology, insurance, finance, and lifetime analysis, that is, new distributions which are more flexible to model real data. Hence, several classes of distributions have been introduced in the literature by adding one or more parameters. The PG-G family contains various well-known distributions as special models. Several new distributions can also be easily generated. Some useful distributions in the PG-G family are described below.

The (proper) survival function for the non-cured population (or PG-G survival function), say $S_{\text{PG-G}}$, is given by

$$S_{\text{PG-G}}(t; \tau, a, \boldsymbol{\eta}) = \frac{\exp\left\{\frac{-\tau \gamma (a, -\log[1 - G(t; \boldsymbol{\eta})])}{\Gamma(a)}\right\} - \exp(-\tau)}{1 - \exp(-\tau)}, \quad t > 0.$$

We note that $\lim_{t \rightarrow 0} S_{\text{PG-G}}(t; \tau, a, \boldsymbol{\eta}) = 1$ and $\lim_{t \rightarrow \infty} S_{\text{PG-G}}(t; \tau, a, \boldsymbol{\eta}) = 0$, so that it is a proper survival

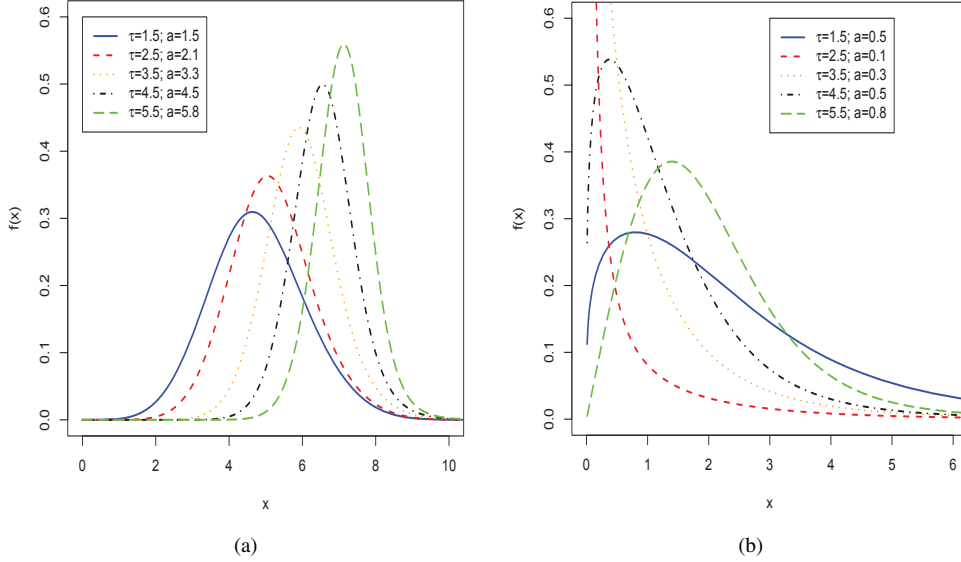


Figure 1: (a) Plots of the Poisson-gamma Weibull (PGW) density function for some values of τ and a with $\alpha = 5.1$ and $\lambda = 3.5$; (b) Plots of the Poisson-gamma log-logistic (PGLL) density function for some values of τ and a with $\alpha = 5.1$ and $\lambda = 2.5$.

function. The pdf for the non-cured population (or the PG-G family) is given by

$$f_{\text{PG-G}}(t; \tau, a, \eta) = \frac{\tau g(t; \eta) \{-\log[1 - G(t; \eta)]\}^{a-1}}{\Gamma(a)[1 - \exp(-\tau)]} \exp\left\{\frac{-\tau \gamma(a, -\log[1 - G(t; \eta)])}{\Gamma(a)}\right\}. \quad (2.5)$$

From Equation (2.5), we note that the parameter a controls its shape. The PG-G family (2.5) can be used to model survival or reliability data. For $a = 1$, the PG-G family reduces to the PG model. The hrf for the non-cured population becomes

$$h_{\text{PG-G}}(t; \tau, a, \eta) = \frac{f_{\text{PG-G}}(t; \tau, a, \eta)}{S_{\text{PG-G}}(t; \tau, a, \eta)}, \quad t > 0.$$

Next, we present some special models of the PG-G family because it extends several widely-known distributions in the literature. The density (2.5) will be most tractable when the cdf $G(t; \eta)$ and pdf $g(t; \eta)$ have simple analytic expressions.

- PGW distribution

If $g(t; \tau) = (\lambda/\alpha^\lambda)t^{\lambda-1} \exp\{-(t/\alpha)^\lambda\}$ is the Weibull pdf, where $\eta = (\alpha, \lambda)^\top$, the PGW density function (for $t > 0$) reduces to

$$f_{\text{PGW}}(t) = \frac{\tau \lambda t^{\lambda-1} \exp\{-(t/\alpha)^\lambda\}}{\alpha^\lambda \Gamma(a)[1 - \exp(-\tau)]} \left(\frac{t}{\alpha}\right)^{\lambda(a-1)} \exp\left\{\frac{-\tau \gamma[a, (t/\alpha)^\lambda]}{\Gamma(a)}\right\}. \quad (2.6)$$

We obtain the PW distribution when $a = 1$. A random variable with density (2.6) is denoted by $T \sim \text{PGW}(a, \tau, \alpha, \lambda)$.

- PGLL distribution

The PGLL distribution is defined from (2.5) by taking $g(t; \tau) = (\lambda/\alpha^\lambda)t^{\lambda-1}[1 + (t/\alpha)^\lambda]^{-2}$ as the LL density, where $\tau = (\alpha, \lambda)^\top$. Its density function (for $t > 0$) is given by

$$f_{\text{PGLL}}(t) = \frac{\tau \lambda t^{\lambda-1} [1 + (t/\alpha)^\lambda]^{-2}}{\alpha^\lambda \Gamma(a) [1 - \exp(-\tau)]} \left\{ \log \left[1 + \left(\frac{t}{\alpha} \right)^\lambda \right] \right\}^{a-1} \times \exp \left\{ \frac{-\tau \gamma(a, \log [1 + (t/\alpha)^\lambda])}{\Gamma(a)} \right\}, \quad (2.7)$$

where α and a are shape parameters and λ is a scale parameter. For $a = 1$, we obtain the PLL distribution. A random variable with density (2.7) is denoted by $T \sim \text{PGLL}(a, \tau, \alpha, \lambda)$. Figure 1 displays the plots of the PGW and PGLL density functions for selected parameter values.

- PGBS distribution

Let $g(t; \tau) = [(t + \lambda)t^{-3/2}/2\alpha \sqrt{2\lambda\pi}] \exp\{-(1/2\alpha^2)(t/\lambda + \lambda/t - 2)\}$ be the BS density. Then, the PGBS density function (for $t > 0$) reduces to

$$\begin{aligned} f_{\text{PGBS}}(t) &= \frac{\tau(t + \lambda)t^{-\frac{3}{2}}}{2\sqrt{2\lambda\pi}\alpha\Gamma(a)[1 - \exp(-\tau)]} \exp \left\{ \frac{-1}{2\alpha^2} \left[\left(\frac{t}{\lambda} \right) + (\lambda t) - 2 \right] \right\} \\ &\times \left\{ -\log \left\{ 1 - \Phi \left[\frac{1}{\alpha} \left(\sqrt{\frac{t}{\lambda}} - \sqrt{\frac{\lambda}{t}} \right) \right] \right\} \right\}^{a-1} \\ &\times \exp \left\{ \frac{-\tau \gamma(a, -\log [1 - \Phi [1/\alpha (\sqrt{t/\lambda} - \sqrt{\lambda/t})]])}{\Gamma(a)} \right\}, \end{aligned} \quad (2.8)$$

where $\lambda > 0$ is a scale parameter and α and a are shape parameters. A random variable with density (2.8) is denoted by $T \sim \text{PGBS}(a, \tau, \alpha, \lambda)$. For $a = 1$, we obtain the PBS distribution.

- PGGHN distribution

Consider the GHN distribution proposed for Cooray and Ananda (2008) with cdf

$$G(t; \tau) = 2\Phi \left[\left(\frac{t}{\lambda} \right)^\alpha \right] - 1 = \text{erf} \left[\frac{(t/\lambda)^\alpha}{\sqrt{2}} \right],$$

where $\Phi(x) = (1/2)[1 + \text{erf}(x/\sqrt{2})]$ and $\text{erf}(x) = (2/\sqrt{\pi}) \int_0^x \exp(-u^2) du$. The PGGHN density function reduces to

$$\begin{aligned} f_{\text{PGGHN}}(t) &= \frac{\tau \sqrt{2} (\alpha/t) (t/\lambda)^\alpha}{\sqrt{\pi} \Gamma(a) [1 - \exp(-\tau)]} \exp \left[-\frac{1}{2} \left(\frac{t}{\lambda} \right)^{2\alpha} \right] \left\{ -\log \left\{ 1 - 2\Phi \left[\left(\frac{t}{\lambda} \right)^\alpha \right] \right\} \right\}^{a-1} \\ &\times \exp \left\{ \frac{-\tau \gamma(a, -\log [1 - 2\Phi [(t/\lambda)^\alpha]])}{\Gamma(a)} \right\}, \end{aligned} \quad (2.9)$$

where $\alpha > 0$ is the shape parameter and $\lambda > 0$ is the scale parameter. For $a = 1$, we obtain the PGHN distribution. A random variable with density (2.9) is denoted by $T \sim \text{PGGHN}(a, \tau, \alpha, \lambda)$. Plots of the PGBS and PGGHN density functions for selected parameter values are displayed in Figure 2.

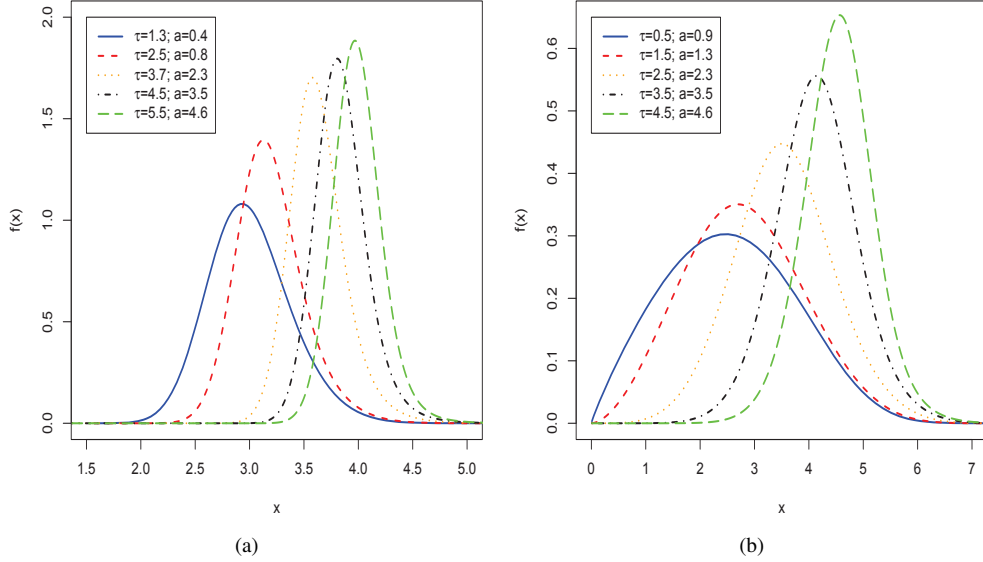


Figure 2: (a) Plots of the Poisson-gamma Birnbaum-Saunders (PGBS) density function for some values of τ and a with $\alpha = 0.1$ and $\lambda = 3.5$; (b) Plots of the Poisson-gamma generalized-half normal (PGGHN) density function for some values of τ and a with $\alpha = 2.1$ and $\lambda = 3.5$.

3. Mathematical properties

The cdf for the non-cured population, say $F_{\text{PG-G}}(t; \tau, a, \eta)$, is given by

$$F_{\text{PG-G}}(t; \tau, a, \eta) = \frac{1 - \exp\{-\tau \gamma(a, -\log[1 - G(t; \eta)])\}}{1 - \exp(-\tau)}, \quad t > 0. \quad (3.1)$$

We omit the dependence on the parameters and write $F_{\text{PG-G}}(t) = F_{\text{PG-G}}(t; \tau, a, \eta)$, $G(t) = G(t; \eta)$, $g(t) = g(t; \eta)$. Based on the power series for the incomplete gamma function ratio, we can write

$$\gamma(a, -\log[1 - G(t)]) = \sum_{m=0}^{\infty} a_m (-\log[1 - G(t)])^{\alpha+m}, \quad (3.2)$$

where $a_m = (-1)^m / [(a+m)m! \Gamma(a)]$. For any real number c and $z \in (0, 1)$, the following formula holds

$$[-\log(1 - z)]^c = z^c + \sum_{i=0}^{\infty} p_i(c) z^{i+c+1}, \quad (3.3)$$

where $p_0(c) = c/2$, $p_1(c) = c(3c+5)/24$, $p_2(c) = c(c^2+5c+6)/48$, $p_3(c) = c(15c^3+150c^2+485c+502)/5760$, etc. Flajolet and Odlyzko (1990, Theorem 3A, p.227) gave the proof of (3.3).

Then, we can write from Equation (3.3)

$$\{-\log[1 - G(t)]\}^{\alpha+m} = G(t)^{\alpha+m} + \sum_{i=0}^{\infty} p_i(\alpha+m) G(t)^{\alpha+i+m+1}$$

and inserting in Equation (3.2)

$$\gamma(a, -\log [1 - G(t)]) = \sum_{m=0}^{\infty} a_m \left[G(t)^{\alpha+m} + \sum_{i=0}^{\infty} p_i(\alpha + m) G(t)^{\alpha+i+m+1} \right].$$

Since $G(t)$ belongs to the interval $(0, 1)$, it is easy to prove

$$G(t)^\beta = \sum_{r=0}^{\infty} s_r(\beta) G(t)^r,$$

where

$$s_r(\beta) = \sum_{j=r}^{\infty} (-1)^{r+j} \binom{\beta}{j} \binom{j}{r}.$$

We can rewrite (3.2) as

$$\gamma(a, -\log [1 - G(t)]) = \sum_{r=0}^{\infty} w_r G(t)^r, \quad (3.4)$$

where

$$w_r = \sum_{m=0}^{\infty} a_m \left[s_r(\alpha + m) + \sum_{i=0}^{\infty} p_i(\alpha + m) s_r(\alpha + i + m + 1) \right].$$

Further, Equation (3.1) reduces to

$$F_{PG-G}(t; \tau, a, \eta) = \frac{1 - \exp\{-\tau \sum_{r=0}^{\infty} w_r G(t)^r\}}{1 - \exp(-\tau)}, \quad t > 0.$$

By expanding the exponential function, we have

$$F_{PG-G}(t; \tau, a, \eta) = \frac{1}{1 - e^{-\tau}} \sum_{j=1}^{\infty} \frac{(-\tau)^j}{j!} \left(\sum_{r=0}^{\infty} w_r G(t)^r \right)^j. \quad (3.5)$$

Next, we adopt an equation of Gradshteyn and Ryzhik (2000) for a power series raised to a positive integer j

$$\left(\sum_{i=0}^{\infty} a_i z^i \right)^j = \sum_{i=0}^{\infty} c_{j,i} z^i, \quad (3.6)$$

where the coefficients $c_{j,i}$ (for $j, i = 1, 2, \dots$) are easily determined from the recurrence equation (for $i = 1, 2, \dots$)

$$c_{j,i} = (i a_0)^{-1} \sum_{l=1}^i [l(j+1) - i] a_l c_{j,i-l}$$

and $c_{j,0} = a_0^j$ (for $j = 0, 1, \dots$). Combining (3.5) and (3.6), we obtain

$$F_{\text{PG-G}}(t) = \sum_{r=0}^{\infty} s_r G(t)^r, \quad (3.7)$$

where (for $r \geq 0$)

$$s_r = \frac{1}{1 - e^{-\tau}} \sum_{j=1}^{\infty} \frac{(-\tau)^j}{j!} e_{j,r},$$

where $e_{j,r} = (r w_0)^{-1} \sum_{m=1}^r [m(j+1) - t] w_m e_{j,r-m}$ (for $j, r = 1, 2, \dots$) and $e_{j,0} = w_0^j$.

A useful representation for the pdf (2.5) can be derived using the concept of exponentiated-G (exp-G) distribution. For an arbitrary baseline cdf $G(x)$, a random variable is said to have the exp-G distribution with parameter $c > 0$, say $Y \sim \text{exp-G}(c)$, if its pdf and cdf are given by $h_c(x) = c G(x)^{c-1} g(x)$ and $H_c(x) = G(x)^c$, respectively.

By differentiating (3.7), we can write

$$f_{\text{PG-G}}(t) = \sum_{r=0}^{\infty} s_{r+1} h_{r+1}(t), \quad (3.8)$$

where $h_{r+1}(t) = h_{r+1}(t; \eta) = (r+1) G(t; \eta)^r g(t; \eta)$ (for $r \geq 0$) is the exp-G density function with power parameter $r+1$. Henceforth, let $Y_{r+1} \sim \text{exp-G}(r+1)$.

3.1. Moments and generating function

A first formula for the n th moment of X can be obtained from (3.8) as

$$\mu'_n = E(X^n) = \sum_{r=0}^{\infty} s_{r+1} E(Y_{r+1}^n). \quad (3.9)$$

Explicit expressions for moments of special PG-G models can be determined from exp-G moments given by Nadarajah and Kotz (2006). A second formula for μ'_n follows from (3.9) in terms of the baseline quantile function (qf) $Q_G(u) = Q_G(u; \eta) = G^{-1}(u; \eta)$. We obtain

$$\mu'_n = \sum_{r=0}^{\infty} (r+1) s_{r+1} \tau_{n,r},$$

where $\tau_{n,r} = \int_{-\infty}^{\infty} x^n G(x)^r g(x) dx = \int_0^1 Q_G(u)^n u^r du$.

For empirical purposes, the shapes of many distributions can be usefully described by the incomplete moments. These types of moments play an important role for measuring inequality, for example, income quantiles and Lorenz and Bonferroni curves. The n th incomplete moment of X is determined as

$$m_n(y) = \int_0^{\infty} t^n f_{\text{PG-G}}(t) dt = \sum_{r=0}^{\infty} (r+1) s_{r+1} \int_0^{G(y)} Q_G(u)^n u^r du. \quad (3.10)$$

The last integral can be computed for most baseline G distributions.

Other kinds of moments such as central and factorial moments and cumulants can also be obtained in closed-form from the ordinary moments using well-known formulae, but we consider only the previous moments for reasons of space.

Next, we obtain the mgf $M(t) = E(e^{tX})$ of X . Let $M_{r+1}(t)$ be the mgf of Y_{r+1} . We can write from (3.8)

$$M(t) = \sum_{r=0}^{\infty} s_{r+1} M_{r+1}(t) = \sum_{r=0}^{\infty} (r+1) s_{r+1} \rho(t, r),$$

where $\rho(t, r) = \int_{-\infty}^{\infty} e^{tq} G(q)^r g(q) dq = \int_0^1 \exp[t Q_G(u)] u^r du$ can be evaluated from $Q_G(u) = G^{-1}(u)$. Hence, $M(t)$ is given by a linear combination of exp- G generating functions.

3.2. Mean deviations

The mean deviations about the mean ($\delta_1 = E(|X - \mu'_1|)$) and about the median ($\delta_2 = E(|X - M|)$) of X can be expressed as $\delta_1 = 2\mu'_1 F(\mu'_1) - 2m_1(\mu'_1)$ and $\delta_2 = \mu'_1 - 2m_1(M)$, respectively, where $\mu'_1 = E(X)$, $M = \text{Median}(X)$ is the median, $F(\mu'_1)$ is evaluated from the cdf (3.1) and $m_1(z) = \int_{-\infty}^z t f_{\text{PG-G}}(t) dt$ can be obtained from (3.10) with $n = 1$. Alternatively, a general equation for $m_1(z)$ can follow from (3.8) as $m_1(z) = \sum_{r=0}^{\infty} s_{r+1} J_{r+1}(z)$, where $J_{r+1}(z) = \int_{-\infty}^z t h_{r+1}(t) dt$.

4. Inference

For many medical problems, the lifetimes are affected by explanatory variables such as the cholesterol level, blood pressure, weight, and others. Parametric models to estimate univariate survival functions and for censored data regression problems are widely used. Different forms of regression models have been proposed in survival analysis. The explanatory variables are commonly used to model the expectation of the number of competing causes. Now, we link the parameter τ in Equation (2.3) to the covariates \mathbf{x}_i by $\tau_i = \exp(\mathbf{x}_i^T \boldsymbol{\beta})$, $i = 1, \dots, n$, where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ denotes the vector of regression parameters.

We consider the situation where the time to the event of interest is not completely observed but subjected to right censoring. Let C_i denote the censoring time. We observe $t_i = \min\{T_i, C_i\}$, for $i = 1, \dots, n$. Let $\boldsymbol{\theta} = (a, \boldsymbol{\eta}^T, \boldsymbol{\beta}^T)^T$ be the parameter vector, from n pairs of times and censoring indicators t_1, \dots, t_n . The full log-likelihood function under non-informative censoring can be expressed as

$$l(\boldsymbol{\theta}) = -r \log[\Gamma(a)] + \sum_{i \in F} \mathbf{x}_i^T \boldsymbol{\beta} + \sum_{i \in F} \log[g(t_i; \boldsymbol{\eta})] + (a-1) \sum_{i \in F} \log\{-\log[1-G(t_i; \boldsymbol{\eta})]\} \\ - \frac{1}{\Gamma(a)} \sum_{i \in F} \exp(\mathbf{x}_i^T \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i; \boldsymbol{\eta})]) - \frac{1}{\Gamma(a)} \sum_{i \in C} \exp(\mathbf{x}_i^T \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i; \boldsymbol{\eta})]), \quad (4.1)$$

where r is the number of failures and F and C denote the uncensored and censored sets of observations, respectively.

The maximum likelihood estimates (MLEs) $\hat{\boldsymbol{\theta}} = (\hat{a}, \hat{\boldsymbol{\eta}}^T, \hat{\boldsymbol{\beta}}^T)^T$ of $\boldsymbol{\theta} = (a, \boldsymbol{\eta}^T, \boldsymbol{\beta}^T)^T$ can be obtained by maximizing (4.1) directly by using the SAS (PROC NLMIXED), R (optim and maxLik functions) and Ox program (sub-routine MaxBFGS). Details for fitting univariate distributions using maximum likelihood in R for censored or non-censored data can be obtained at <http://www.inside-r.org/packages/cran/fitdistrplus/docs/mledist>.

The inference procedures for $\theta = (a, \boldsymbol{\eta}^\top, \boldsymbol{\beta}^\top)^\top$ can be based on the asymptotic normal approximation

$$\left(\hat{a}, \hat{\boldsymbol{\eta}}^\top, \hat{\boldsymbol{\beta}}^\top\right)^\top \sim N_{(p+q+1)}\left\{\left(a, \boldsymbol{\eta}^\top, \boldsymbol{\beta}^\top\right)^\top, -\ddot{\mathbf{L}}^{-1}(\boldsymbol{\theta})\right\}, \quad (4.2)$$

where $-\ddot{\mathbf{L}}(\boldsymbol{\theta}) = \{-\partial^2 l(\boldsymbol{\theta}) / \partial \boldsymbol{\theta} \boldsymbol{\theta}^\top\}$ is the $(p+q+1) \times (p+q+1)$ observed information matrix, which can be evaluated numerically.

The likelihood ratio (LR) statistic can be used to compare some special models with the PG-G model. We consider the partition $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top)^\top$, where $\boldsymbol{\theta}_1$ is the subset of parameters of interest and $\boldsymbol{\theta}_2$ is a subset of remaining parameters. The LR statistic for testing the null hypothesis $H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_1^{(0)}$ versus the alternative hypothesis $H_1 : \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_1^{(0)}$ is given by $w = 2\{\ell(\hat{\boldsymbol{\theta}}) - \ell(\tilde{\boldsymbol{\theta}})\}$, where $\tilde{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ are the estimates under the null and alternative hypotheses, respectively. The statistic w is asymptotically (as $n \rightarrow \infty$) distributed as χ_k^2 , where k is the dimension of the subset of parameters $\boldsymbol{\theta}_1$ of interest. For example, the test of $H_0 : a = 1$ versus $H : a \neq 1$ is equivalent to compare the PG and PG-G models. In this case, the LR statistic is

$$w = 2\left\{l(\hat{a}, \hat{\boldsymbol{\eta}}, \hat{\boldsymbol{\beta}}) - l(1, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\beta}})\right\},$$

where \hat{a} , $\hat{\boldsymbol{\eta}}$ and $\hat{\boldsymbol{\beta}}$ are the MLEs under H and $\tilde{\boldsymbol{\eta}}$ and $\tilde{\boldsymbol{\beta}}$ are the estimates under H_0 .

5. Simulation study

We perform a simulation study in order to evaluate some frequentist properties of the MLE $\hat{\boldsymbol{\theta}}$. The failure time data are simulated from the gamma-G family by taking the Weibull, LL, BS and GHN models for the baseline G distribution. The number of causes of the event of interest for each individual i ($i = 1, \dots, n$), M_i , is generated from a Poisson distribution with parameter τ_i . One explanatory variable is generated from a Bernoulli distribution with parameter 0.5. We link the parameter τ in Equation (4) to the covariate x_i by $\tau_i = \exp(\beta_0 + \beta_1 x_i)$, $i = 1, \dots, n$. The censoring times are sampled from the uniform distribution in the $(0, \nu)$ interval, where ν controls the censoring proportion of the uncured population.

We consider the following values for the parameters: $a = 1.50$, $\alpha = 1.25$, $\lambda = 0.25$, $\beta_0 = -0.50$ and $\beta_1 = 0.70$ for the PGBS model; $a = 0.75$, $\alpha = 2.00$, $\lambda = 0.50$, $\beta_0 = -0.50$ and $\beta_1 = 0.70$ for the PGGHS model; $a = 0.75$, $\alpha = 0.75$, $\lambda = 1.50$, $\beta_0 = -0.50$ and $\beta_1 = 1.25$ for the PGLL model; and $a = 0.50$, $\alpha = 2.00$, $\lambda = 1.50$, $\beta_0 = -0.50$ and $\beta_1 = 1.25$ for the PGW model.

The results are obtained from 3,000 Monte Carlo simulations performed using **R** software with the **Optim** function. In each replication, a random sample of size n is drawn from the four models and parameters are estimated by maximum likelihood. The sample sizes are $n = 100, 200, 300, 400$ and 700 . In this study, the proportion of censored observations is taken approximately to be equal to 55% for the PGBS and PGGHN models and 45% for the PGBS and PGGHN models.

Tables 2–5 list the averages of the MLEs (mean), biases (bias) and mean square errors (MSEs) for the PGBS, PGGHN, PGLL and PGW models, respectively.

We conclude from the figures in Tables 2–5 that the average estimates of the parameters tend to be closer to the true parameters when n increases. We can also note that, even for small sample sizes, the biases and MSEs are small for the estimates of β_0 and β_1 . This fact indicates that asymptotic normal distribution provides an adequate approximation to the finite sample distribution of the MLEs. The normal approximation can be often improved using bias adjustments to the estimators. Approximations to their biases in simple models may be determined analytically. Bias correction typically

Table 2: Summaries of the quantities for the Poisson-gamma Birnbaum-Saunders (PGBS) model

Sample size	Parameter (True value)	Summaries of the parameters		
		Mean	Bias	MSE
100	a (1.50)	1.7020	0.2020	0.5593
	α (1.25)	1.4887	0.2387	0.3342
	λ (0.25)	0.5095	0.2595	1.4779
	β_0 (-0.50)	-0.4009	0.0991	0.1718
	β_1 (0.70)	0.7139	0.0139	0.1020
200	a (1.50)	1.6114	0.1114	0.2923
	α (1.25)	1.4024	0.1524	0.1532
	λ (0.25)	0.3951	0.1451	0.2884
	β_0 (-0.50)	-0.4369	0.0631	0.0820
	β_1 (0.70)	0.7091	0.0091	0.0451
300	a (1.50)	1.5669	0.0669	0.1992
	α (1.25)	1.3526	0.1026	0.0796
	λ (0.25)	0.3433	0.0933	0.1450
	β_0 (-0.50)	-0.4610	0.0390	0.0527
	β_1 (0.70)	0.7058	0.0058	0.0325
400	a (1.50)	1.5769	0.0769	0.1581
	α (1.25)	1.3306	0.0806	0.0524
	λ (0.25)	0.3043	0.0543	0.0801
	β_0 (-0.50)	-0.4769	0.0231	0.0347
	β_1 (0.70)	0.7053	0.0053	0.0241
700	a (1.50)	1.5371	0.0371	0.0917
	α (1.25)	1.3009	0.0509	0.0162
	λ (0.25)	0.2737	0.0237	0.0284
	β_0 (-0.50)	-0.4840	0.0160	0.0149
	β_1 (0.70)	0.6998	-0.0002	0.0136

Mean = averages of the MLEs; Bias = biases; MSE = mean square error.

Table 3: Summaries of the quantities for the Poisson-gamma generalized half normal (PGGHN) model

Sample size	Parameter (True value)	Summaries of the parameters		
		Mean	Bias	MSE
100	a (0.75)	1.2934	0.5434	0.9498
	α (2.00)	1.5624	-0.4376	0.6602
	λ (0.50)	0.4146	-0.0854	0.0316
	β_0 (-0.50)	-0.4979	0.0021	0.0678
	β_1 (0.70)	0.7009	0.0009	0.1034
200	a (0.75)	1.1171	0.3671	0.4919
	α (2.00)	1.6459	-0.3541	0.4718
	λ (0.50)	0.4378	-0.0622	0.0194
	β_0 (-0.50)	-0.5043	-0.0043	0.0350
	β_1 (0.70)	0.7062	0.0062	0.0540
300	a (0.75)	1.0062	0.2562	0.3063
	α (2.00)	1.7352	-0.2648	0.3539
	λ (0.50)	0.4567	-0.0433	0.0131
	β_0 (-0.50)	-0.5011	-0.0011	0.0207
	β_1 (0.70)	0.6990	-0.0010	0.0326
400	a (0.75)	0.9583	0.2083	0.2160
	α (2.00)	1.7673	-0.2327	0.2941
	λ (0.50)	0.4644	-0.0356	0.0100
	β_0 (-0.50)	-0.5022	-0.0022	0.0168
	β_1 (0.70)	0.7011	0.0011	0.0251
700	a (0.75)	0.8612	0.1112	0.0822
	α (2.00)	1.8589	-0.1411	0.1897
	λ (0.50)	0.4823	-0.0177	0.0048
	β_0 (-0.50)	-0.5039	-0.0039	0.0094
	β_1 (0.70)	0.7038	0.0038	0.0144

Mean = averages of the MLEs; Bias = biases; MSE = mean square error.

Table 4: Summaries of the quantities for the Poisson-gamma log-logistic (PGLL) model

Sample size	Parameter (True value)	Summaries of the parameters		
		Mean	Bias	MSE
100	a (0.75)	1.0616	0.3116	1.3849
	α (0.75)	0.6668	-0.0832	0.1360
	λ (1.50)	1.7815	0.2815	0.7338
	β_0 (-0.50)	-0.4841	0.0159	0.1032
	β_1 (1.25)	1.2687	0.0187	0.0915
200	a (0.75)	0.9282	0.1782	0.4558
	α (0.75)	0.6948	-0.0552	0.0962
	λ (1.50)	1.6422	0.1422	0.4122
	β_0 (-0.50)	-0.4659	0.0341	0.0608
	β_1 (1.25)	1.2594	0.0094	0.0447
300	a (0.75)	0.8700	0.1200	0.2288
	α (0.75)	0.7092	-0.0408	0.0725
	λ (1.50)	1.6024	0.1024	0.3001
	β_0 (-0.50)	-0.4719	0.0281	0.0425
	β_1 (1.25)	1.2582	0.0082	0.0302
400	a (0.75)	0.8313	0.0813	0.1441
	α (0.75)	0.7233	-0.0267	0.0547
	λ (1.50)	1.5789	0.0789	0.2131
	β_0 (-0.50)	-0.4772	0.0228	0.0301
	β_1 (1.25)	1.2551	0.0051	0.0215
700	a (0.75)	0.8083	0.0583	0.0797
	α (0.75)	0.7282	-0.0218	0.0345
	λ (1.50)	1.5314	0.0314	0.1233
	β_0 (-0.50)	-0.4844	0.0156	0.0176
	β_1 (1.25)	1.2551	0.0051	0.0131

Mean = averages of the MLEs; Bias = biases; MSE = mean square error.

Table 5: Summaries of the quantities for the Poisson-gamma Weibull (PGW) model

Sample size	Parameter (True value)	Summaries of the parameters		
		Mean	Bias	MSE
100	a (0.50)	0.7289	0.2289	0.3105
	α (2.00)	1.6741	-0.3259	0.7474
	λ (1.50)	1.5217	0.0217	0.5108
	β_0 (-0.50)	-0.5042	-0.0042	0.0773
	β_1 (1.25)	1.2690	0.0190	0.0947
200	a (0.50)	0.6815	0.1815	0.2088
	α (2.00)	1.7685	-0.2315	0.6832
	λ (1.50)	1.4991	-0.0009	0.4328
	β_0 (-0.50)	-0.4800	0.0200	0.0399
	β_1 (1.25)	1.2507	0.0007	0.0431
300	a (0.50)	0.6492	0.1492	0.1582
	α (2.00)	1.8028	-0.1972	0.4584
	λ (1.50)	1.4953	-0.0047	0.3642
	β_0 (-0.50)	-0.4861	0.0139	0.0262
	β_1 (1.25)	1.2561	0.0061	0.0297
400	a (0.50)	0.6294	0.1294	0.1290
	α (2.00)	1.8136	-0.1864	0.3152
	λ (1.50)	1.4952	-0.0048	0.3228
	β_0 (-0.50)	-0.4847	0.0153	0.0202
	β_1 (1.25)	1.2525	0.0025	0.0221
700	a (0.50)	0.5822	0.0822	0.0727
	α (2.00)	1.8880	-0.1120	0.1900
	λ (1.50)	1.5012	0.0012	0.2303
	β_0 (-0.50)	-0.4911	0.0089	0.0122
	β_1 (1.25)	1.2551	0.0055	0.0123

Mean = averages of the MLEs; Bias = biases; MSE = mean square error.

does a very good job to correct MLEs. However, it may also increase the MSEs. Whether bias correction is useful in practice depends basically on the shape of the bias function and on the variance of the MLE.

6. Sensitivity analysis

The approach is one in which the stability of the estimated outputs with respect to the model inputs is studied using various minor model perturbation schemes, such as the local influence approach developed by Cook (1986). In the following section, we describe the background and details of the classical diagnostic methods to detect influential observations.

6.1. Local influence

Another approach is suggested by Cook (1986), where instead of removing observations, weights are given to them. Local influence calculation can be conducted for model (2.3). Let ω_0 be the no perturbation vector. If likelihood displacement $LD(\omega) = 2\{l(\hat{\theta}) - l(\hat{\theta}_\omega)\}$ is used, where $\hat{\theta}_\omega$ denotes MLE under the perturbed model, the normal curvature for θ at the direction \mathbf{d} , $\|\mathbf{d}\| = 1$, is provided by $C_{\mathbf{d}}(\theta) = 2|\mathbf{d}^\top \Delta^\top [\tilde{\mathbf{L}}(\theta)]^{-1} \Delta \mathbf{d}|$, where Δ is a $(p+q+1) \times n$ matrix, which depends on the perturbation scheme. The elements of Δ are given by $\Delta_{vi} = \partial^2 l(\theta|\omega) / \partial \phi_v \partial \omega_i$, for $i = 1, 2, \dots, n$ and $v = 1, 2, \dots, p+q+1$, evaluated at $\hat{\theta}$ and ω_0 . For the PG-G regression model with long-term survivors, the elements of $\tilde{\mathbf{L}}(\theta)$ are calculated numerically. We can also determine normal curvatures $C_{\mathbf{d}}(a)$, $C_{\mathbf{d}}(\eta)$ and $C_{\mathbf{d}}(\beta)$ to perform various index plots such as the index plot of \mathbf{d}_{\max} , the eigenvector corresponding to $C_{\mathbf{d}_{\max}}$, the largest eigenvalue of the matrix $\mathbf{B} = -\Delta^\top [\tilde{\mathbf{L}}(\theta)]^{-1} \Delta$. The index plots of $C_{\mathbf{d}_i}(a)$, $C_{\mathbf{d}_i}(\eta)$ and $C_{\mathbf{d}_i}(\beta)$, called total local influence, where \mathbf{d}_i denotes an $n \times 1$ vector of zeros with one at the i^{th} position. Thus, the curvature at direction \mathbf{d}_i takes the form $C_i = 2|\Delta_i^\top [\tilde{\mathbf{L}}(\theta)]^{-1} \Delta_i|$, where Δ_i^\top denotes the i^{th} row of Δ . It is usual to point out those cases such that $C_i \geq 2\bar{C}$, where $\bar{C} = (1/n) \sum_{i=1}^n C_i$. Another influence measure for the i^{th} observation is $U_i = \sum_{k=1}^{n_1} \lambda_k e_{ki}^2$, where $\{(\lambda_k, \mathbf{e}_k) | k = 1, \dots, n\}$ are the eigenvalue-eigenvector pairs of \mathbf{B} with $\lambda_1 \geq \dots \geq \lambda_{n_1} \geq \lambda_{n_1+1} = \dots = \lambda_n = 0$ and $\{\mathbf{e}_k = (e_{k1}, \dots, e_{kn})^\top\}$ is the associated orthonormal basis. Zhu *et al.* (2007) studied the influence measure u_i systematically under a case weight perturbation. Therefore, this influence measure expresses local sensitivity to the log-likelihood of the perturbations.

6.2. Curvature calculations

Next, we determine for three perturbation schemes, the matrix

$$\Delta = (\Delta_{vi})_{[(p+q+1) \times n]} = \left(\frac{\partial^2 l(\theta|\omega)}{\partial \theta_v \partial \omega_i} \right)_{[(p+q+1) \times n]},$$

where $v = 1, 2, \dots, p+q+1$ and $i = 1, 2, \dots, n$. We should consider the model defined in (2.3), (2.4) and its log-likelihood function given by (4.1).

6.2.1. Case-weight perturbation

First, we consider a case weight perturbation which modifies the weight given to each subject in the log-likelihood. Consider the vector of weights $\omega = (\omega_1, \dots, \omega_n)^\top$.

In this case, the log-likelihood function is given by

$$\begin{aligned}
l(\boldsymbol{\theta}|\boldsymbol{\omega}) &= -\log[\Gamma(a)] \sum_{i \in F} \omega_i + \sum_{i \in F} (\mathbf{x}_i^\top \boldsymbol{\beta}) \omega_i + \sum_{i \in F} \log[g(t_i; \boldsymbol{\eta})] \omega_i \\
&+ (a-1) \sum_{i \in F} \log\{-\log[1-G(t_i; \boldsymbol{\eta})]\} \omega_i \\
&- \frac{1}{\Gamma(a)} \sum_{i \in F} \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i; \boldsymbol{\eta})]) \omega_i \\
&- \frac{1}{\Gamma(a)} \sum_{i \in C} \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i; \boldsymbol{\eta})]) \omega_i,
\end{aligned}$$

where $0 \leq \omega_i \leq 1$, $\boldsymbol{\omega}_0 = (1, \dots, 1)^\top$. The matrix $\boldsymbol{\Delta} = (\boldsymbol{\Delta}_a^\top, \boldsymbol{\Delta}_\eta^\top, \boldsymbol{\Delta}_\beta^\top)^\top$ is evaluated numerically.

6.2.2. Response perturbation

Since t_i values have different variances, we require a scaling of the perturbation vector $\boldsymbol{\omega}$ by an estimator of the standard deviation of t_i . We consider that each t_i is perturbed as $t_{i\omega} = t_i + \omega_i S_t$, where S_t is a scale factor that may be estimated by the standard deviation of t and $\omega_i \in \mathbb{R}$.

Here, the perturbed log-likelihood function is expressed as

$$\begin{aligned}
l(\boldsymbol{\theta}|\boldsymbol{\omega}) &= -r \log[\Gamma(a)] + \sum_{i \in F} \mathbf{x}_i^\top \boldsymbol{\beta} + \sum_{i \in F} \log[g(t_i^*; \boldsymbol{\eta})] \\
&+ (a-1) \sum_{i \in F} \log\{-\log[1-G(t_i^*; \boldsymbol{\eta})]\} \\
&- \frac{1}{\Gamma(a)} \sum_{i \in F} \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i^*; \boldsymbol{\eta})]) \\
&- \frac{1}{\Gamma(a)} \sum_{i \in C} \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) \gamma(a, -\log[1-G(t_i^*; \boldsymbol{\eta})]),
\end{aligned}$$

where $t_i^* = t_i + \omega_i S_t$ and $\boldsymbol{\omega}_0 = (0, \dots, 0)^\top$. The matrix $\boldsymbol{\Delta} = (\boldsymbol{\Delta}_a^\top, \boldsymbol{\Delta}_\eta^\top, \boldsymbol{\Delta}_\beta^\top)^\top$ is evaluated numerically.

7. Application: gastric cancer data

The data set refers to $n = 201$ patients observed with gastric adenocarcinoma. Gastric (stomach) cancer is a disease in which malignant (cancer) cells form in the stomach lining. Almost all gastric cancers are adenocarcinomas (cancers that begin in cells that make and release mucus and other fluids). Other types of gastric cancer are gastrointestinal carcinoid tumors, gastrointestinal stromal tumors, and lymphomas. The response variable is the time t_i in months after surgery until death. Patients who die from other causes and patients that are still alive at the end of the study are censored observations (53%). For more details, see Martinez *et al.* (2013). The only covariate is the type of therapy: x_{i1} (0=adjuvant chemoradiotherapy, $n = 125$; 1 = surgery alone, $n = 76$). We are interested in the effect of the explanatory variable on the cure fraction. Figure 3 displays the estimated Kaplan-Meier survival function with a well-pronounced plateau. According to Yakovlev and Tsodikov (1996) this may be thought of as an indication of the presence of a proportion of patients for whom the gastric adenocarcinoma will never recur; therefore, the patients can be considered cured.

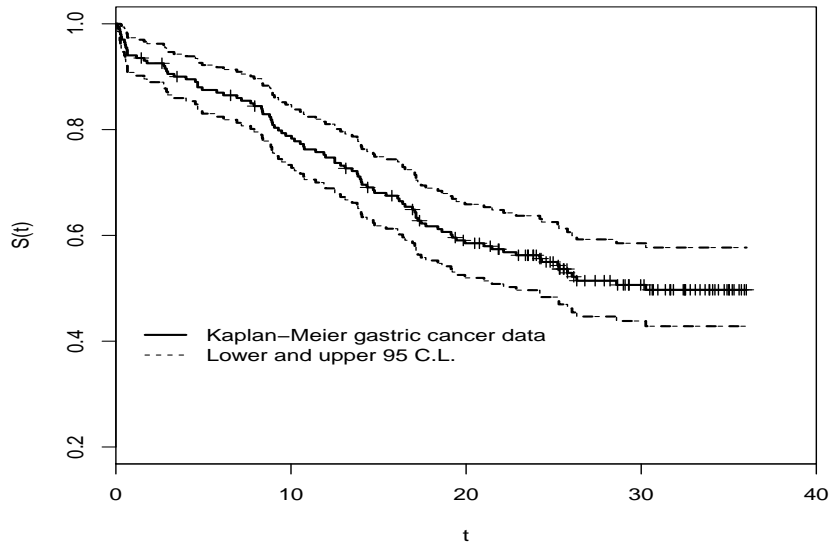


Figure 3: *Kaplan-Meier curves for the gastric cancer data.*

Table 6: Some statistics from the fitted models to the gastric cancer data

Model	Statistics		
	AIC	CAIC	BIC
PGW	900.3	900.6	916.8
PW	898.5	898.7	911.7
PGLL	900.1	900.4	916.7
PLL	898.2	898.4	911.4
PGBS	893.9	894.2	910.4
PBS	970.5	970.7	983.8
PGGHN	892.9	893.2	909.4
PGHN	897.5	897.8	910.8

The explanatory variable is related to the parameter θ according to the following structure. For the PG-G cure rate model, we consider ($i = 1, \dots, 201$):

$$\tau_i = \exp(\beta_0 + \beta_1 x_{i1}).$$

In Table 6, we present the values of the Akaike Information Criterion (AIC), Consistent Akaike Information Criterion (CAIC) and Bayesian Information Criterion (BIC) for all models discussed in Section 2. So, we will have more evidence to be able to discriminate and choose the most suitable model. The lowest values of these information criteria correspond to the PGGHN model, which provides a better fit to the current data than other models.

Table 7 lists the MLEs for the fitted PGGHN regression model. At a 5% significance level, the regression coefficient is significant for the type of therapy (x_1). Table 8 provides LR statistics to compare the PGW and PW, PGLL and PLL, PGBS and PBS, PGGHN and PGHN models to the current data. We reject the null model in favor of the PGBS and PGGHN models.

Next, we analyze the local influence for the stomach cancer data.

1) Influence using case-weight perturbation

Table 7: MLEs for the full Poisson-gamma generalized-half normal (PGGHN) regression model with cure rate fraction fitted to the gastric cancer data

Parameter	Estimate	Standard error	95% C.L.	p -value
a	0.1588	0.0153	(0.1286, 0.1890)	-
α	5.7268	0.1137	(5.5026, 5.9509)	-
λ	28.8416	1.7077	(25.4743, 32.2090)	-
β_0	-0.6451	0.1788	(-0.9976, -0.2926)	0.0004
β_1	0.4485	0.2176	(0.0195, 0.8775)	0.0406

Table 8: Likelihood ratio (LR) tests

Endurance	Hypotheses	Statistic w	p -value
PGW versus PW	$H_0 : a = 1$ vs $H_1 : H_0$ is false	0.2	0.6547
PGLL versus PLL	$H_0 : a = 1$ vs $H_1 : H_0$ is false	0.1	0.7518
PGBS versus PBS	$H_0 : a = 1$ vs $H_1 : H_0$ is false	78.6	<0.0001
PGGHN versus PGHN	$H_0 : a = 1$ vs $H_1 : H_0$ is false	6.6	0.0102

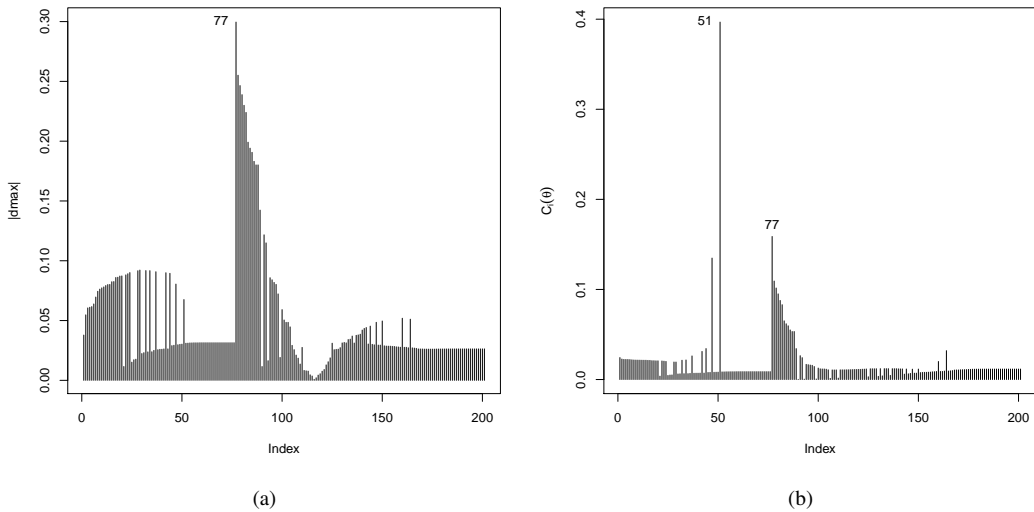


Figure 4: Index plots of the case-weight perturbation for θ on the gastric cancer data ((a) d_{\max} ; (b) C_i).

By applying the local influence methodology developed in Section 6.1, where case-weight perturbation is used, the value $C_{d_{\max}}(\theta) = 1.445$ is found as a maximum curvature. Figure 4 displays the index plots of $d_{\max}(\theta)$ and C_i for all observations. Clearly, the cases #51 and #77 are the most influential observations on $\hat{\theta}$ (Figures 4(a) and (b)).

2) Influence using response variable perturbation

Further, we examine the influence of perturbations on observed survival times. The value for the maximum curvature is $C_{d_{\max}}(\theta) = 599.715$. In Figure 5, we provide plots for $d_{\max}(\theta)$ and C_i for all points. The plots in Figures 5(a) and (b) indicate that the case #82 is the most influential observation on $\hat{\theta}$.

3) Impact of the detected influential observations

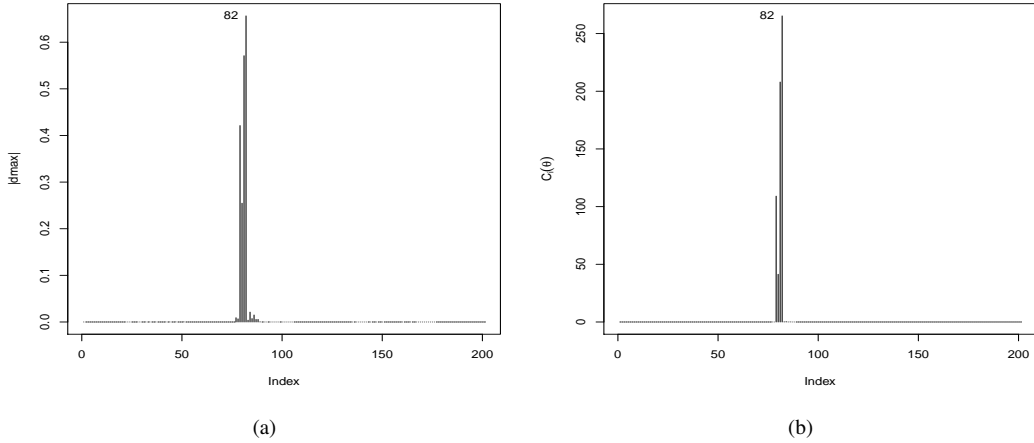


Figure 5: Index plot of the response perturbation scheme for θ on the gastric cancer data ((a) \mathbf{d}_{\max} , (b) C_i).

The diagnostic analysis detected, as potentially influential, the following three cases: #51, #77 and #82. In order to reveal the impact of these observations on the parameter estimates, we refit the model under some situations. First, we individually eliminate each of the three cases. Next, we remove the totality of potentially influential observations from set “A” (original data set).

In Table 9, we provide the relative changes (in percentages) of the parameter estimates given by $\mathbf{RC}_{\theta_j} = [(\hat{\theta}_j - \hat{\theta}_{j(I)})/\hat{\theta}_j] \times 100$, parameter estimates and the corresponding p -values, where $\hat{\theta}_{j(I)}$ denotes the MLE of θ_j after the set “ I ” of observations has been removed. Note that $I_1 = \{\#51\}$, $I_2 = \{\#77\}$, $I_3 = \{\#82\}$, $I_4 = \{\#51, \#77\}$, $I_5 = \{\#51, \#82\}$, $I_6 = \{\#77, \#82\}$ and $I_7 = \{\#51, \#77, \#82\}$.

From Table 9, we can note that the MLEs from the PGGHN regression model with cure rate fraction are not highly sensitive under the deletion of the outstanding observations. The significance of the parameter estimates does not change (at the 5% level) after removing set I . Therefore, we do not have inferential changes after removing the observations provided in the diagnostic plots.

4) Goodness of fit

We adopt a regression structure for the cure probability in long-term survivor models (Section 4). We now estimate the cure rate (p_0). Note that

$$\hat{\tau} = \frac{1}{201} \sum_{i=1}^{201} \hat{\tau}_i = 0.7093,$$

where $\hat{\tau}_i = \exp(-0.6451 + 0.4485x_{i1})$ and then $\hat{p}_0 = e^{-\hat{\tau}} = 0.4920$.

In order to assess if the model is appropriate, Figure 6(a) displays the empirical survival function and the estimated marginal survival functions given by (2.3) from the fitted PGGHN model with long-term survivors.

The estimates of the cure rate for patients stratified by type of therapy (x_1) are:

- For **Chemoradiotherapy** ($x_1 = 0$)

$\hat{\tau}_0 = \exp(-0.6451)$ and the cured fraction is $\hat{p}_{00} = e^{-\hat{\tau}_0} = 0.5918$.

Table 9: Relative changes [RC-in %], estimates and the corresponding p -values in parentheses for the regression coefficients to explain survival times

Drooping	a	α	λ	β_0	β_1
A	[-] 0.1588 (-)	[-] 5.7268 (-)	[-] 28.8416 (-)	[-] -0.6451 (0.0004)	[-] 0.4485 (0.0406)
A- I_1	[3] 0.1547 (-)	[-4] 5.9279 (-)	[5] 27.4986 (-)	[-6] -0.6812 (0.0002)	[-3] 0.4606 (0.0376)
A- I_2	[-2] 0.1626 (-)	[-1] 5.7832 (-)	[0] 28.8386 (-)	[0] -0.6437 (0.0003)	[2] 0.4409 (0.0436)
A- I_3	[4] 0.1625 (-)	[0] 5.7332 (-)	[0] 28.8163 (-)	[0] -0.6438 (0.0003)	[4] 0.4311 (0.0486)
A- I_4	[4] 0.1518 (-)	[-8] 6.1962 (-)	[2] 28.1716 (-)	[-2] -0.6602 (0.0003)	[2] 0.4408 (0.0453)
A- I_5	[-5] 0.1661 (-)	[0] 5.6984 (-)	[5] 27.2627 (-)	[-6] -0.6847 (0.0002)	[0] 0.4499 (0.0419)
A- I_6	[-5] 0.1670 (-)	[-1] 5.7857 (-)	[0] 28.8639 (-)	[-1] -0.6407 (0.0004)	[5] 0.4265 (0.0519)
A- I_7	[13] 0.1388 (-)	[-5] 6.0239 (-)	[4] 27.5831 (-)	[-4] -0.6736 (0.0003)	[2] 0.4385 (0.0491)

- For **Surgery alone** ($x_1 = 1$)

$\hat{\tau}_1 = \exp(-0.6451 + 0.4485)$ and the cured fraction is $\hat{p}_{01} = e^{-\hat{\tau}_1} = 0.4398$.

The estimated survival function and cure fraction stratified by x_1 are also displayed in Figure 6(b), from which a significant fraction of survivors can be observed. Note that the proportion of cured is greater for patients receiving chemoradiotherapy.

The plots of the hrf in Figure 7 corresponding to the type of therapy variable (x_1) under the PGGHN regression model with cure fraction reveal that the hrf is larger for *surgery alone* than *chemoradiotherapy*. There exists a substantial difference between the two hrfs.

More information is provided by a visual comparison of the histogram of the data with the fitted density functions. The plots of the fitted PGGHN and GHN distributions for the non-cured population, see Equation (2.9), are displayed in Figure 8. We also conclude that the PGGHN distribution for the non-cured population provides an adequate fit to the data.

8. Concluding remarks

This study modifies the PG-G regression model to include long-term individuals. Our proposal takes parametric modeling as the basis for survival time. The model estimates the effects of explanatory variables on the acceleration/deceleration of timing in a surviving fraction, that is, the proportion of the population for which the event never occurs. We provide applications of influence diagnostics in the PG-G regression model with long-term survivors. By applying the procedures in a data set from the medical area, we can assess the sensitivity aspects of the MLEs under some perturbation schemes as well as check the goodness-of-fit of the postulated model. The diagnostic plots detect

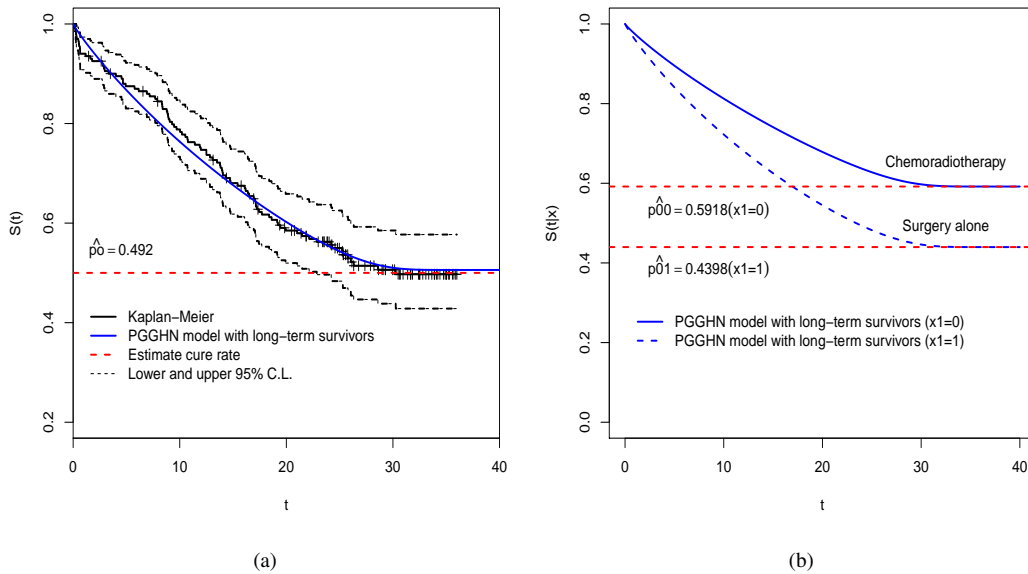


Figure 6: (a) Kaplan-Meier curves (solid lines), the estimated Poisson-gamma generalized-half normal (PGGHN) survival function and the estimated cure fraction for the gastric cancer data; (b) Estimates of the survival function and cure fraction of model stratified by type of therapy for the gastric cancer data.

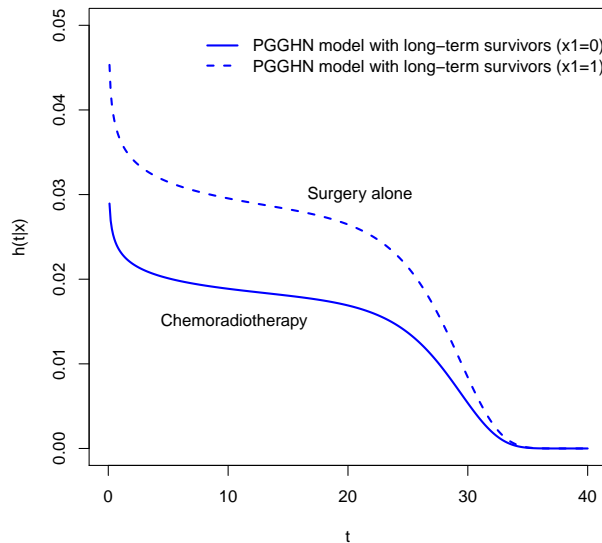


Figure 7: Fitted hrfs using the Poisson-gamma generalized-half normal (PGGHN) regression model with cure fraction for the gastric cancer data.

some possible influential observations; however, their deletion does not cause inferential changes in the results. Future studies can be conducted to compare the Cox regression model with our model. This can be done empirically, considering different simulation scenarios. Therefore, the PG-G model

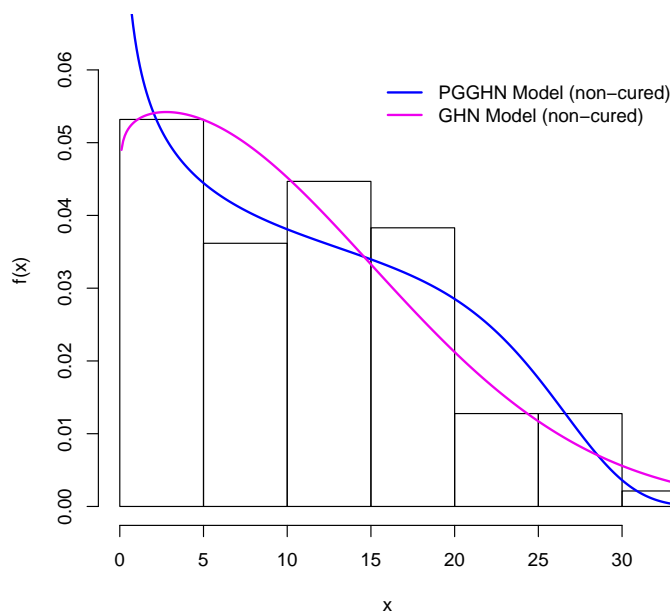


Figure 8: Fitted Poisson-gamma generalized-half normal (PGGHN) and generalized-half normal (GHN) densities (non-cured) for the gastric cancer data.

with long-term survivors represent an interesting option to explain/predict survival times for long-term individuals.

References

- Balakrishnan N and Pal S (2012). EM algorithm-based likelihood estimation for some cure rate models, *Journal of Statistical Theory and Practice*, **6**, 698–724.
- Balakrishnan N and Pal S (2013). Lognormal lifetimes and likelihood-based inference for flexible cure rate models based on COM-Poisson family, *Computational Statistics and Data Analysis*, **67**, 41–67.
- Balakrishnan N and Pal S (2015a). Likelihood inference for flexible cure rate models with gamma lifetimes, *Communications in Statistics-Theory and Methods*, **44**, 4007–4048.
- Balakrishnan N and Pal S (2015b). An EM algorithm for the estimation of flexible cure rate model parameters with generalized gamma lifetime and model discrimination using likelihood- and information-based methods, *Computational Statistics*, **30**, 151–189.
- Balakrishnan N and Pal S (2016). Expectation maximization-based likelihood inference for flexible cure rate models with Weibull lifetimes, *Statistical Methods in Medical Research*, **25**, 1535–1563.
- Balakrishnan N, Koutras MV, Milienos F, and Pal S (2016). Piecewise linear approximations for cure rate models and associated inferential issues, *Methodology and Computing in Applied Probability*, **18**, 937–966.
- Chen MH, Ibrahim JG, and Sinha D (1999). A new Bayesian model for survival data with a surviving fraction, *Journal of the American Statistical Association*, **94**, 909–919.

- Cook RD (1986). Assessment of local influence, *Journal of the Royal Statistical Society Series B (Methodological)*, **48**, 133–169.
- Cooner F, Banerjee S, Carlin BP, and Sinha D (2007). Flexible cure rate modeling under latent activation schemes, *Journal of the American Statistical Association*, **102**, 560–572.
- Cooray K and Ananda MMA (2008). A generalized of the half-normal distribution with applications to lifetime data, *Communications in Statistics - Theory and Methods*, **37**, 1323–1337.
- Fachini JB, Ortega EMM, and Cordeiro GM (2014). A bivariate regression model with cure fraction, *Journal of Statistical Computation and Simulation*, **84**, 1580–1595.
- Flajolet P and Odlyzko A (1990). Singularity analysis of generating functions, *SIAM Journal on Discrete Mathematics*, **3**, 216–240.
- Gradshteyn IS and Ryzhik IM (2000). *Table of Integrals, Series and Products*(6th ed), Academic Press, San Diego, CA.
- Hashimoto EM, Cordeiro GM, Ortega EMM (2013). The new Neyman type A beta Weibull model with long-term survivors, *Computational Statistics*, **28**, 933–954.
- Ibrahim JG, Chen MH, and Sinha D (2001). *Bayesian Survival Analysis*, Springer, New York.
- Maller RA and Zhou X (1996). *Survival Analysis with Long-Term Survivors*, John Wiley & Sons, New York.
- Martinez EZ, Achcar JA, Jácome AAA, and Santos JS (2013). Mixture and non-mixture cure fraction models based on the generalized modified Weibull distribution with an application to gastric cancer data, *Computer Methods and Programs in Biomedicine*, **112**, 343–355.
- Nadarajah S, Cordeiro GM, and Ortega EMM (2015). The Zografos-Balakrishnan-G family of distributions: mathematical properties and applications, *Communications in Statistics - Theory and Methods*, **44**, 186–215.
- Nadarajah S and Kotz S (2006). The beta exponential distribution, *Reliability Engineering and System Safety*, **91**, 689–697.
- Ortega EMM, Cordeiro GM, Campelo AK, Kattan MW, and Cancho VG (2015). A power series beta Weibull regression model for predicting breast carcinoma, *Statistics in Medicine*, **34**, 1366–1388.
- Ortega EMM, Cordeiro GM, and Kattan MW (2012). The negative binomial-beta Weibull regression model to predict the cure of prostate cancer, *Journal of Applied Statistics*, **39**, 1191–1210.
- Ristić MM and Balakrishnan N (2012). The gamma-exponentiated distribution, *Journal of Statistical Computation and Simulation*, **82**, 1191–1206.
- Rodrigues J, Cancho VG, de Castro M, and Louzada-Neto F (2009). On the unification of the long-term survival models, *Statistics and Probability Letters*, **79**, 753–759.
- Tsodikov AD, Ibrahim JG, and Yakovlev AY (2003). Estimating cure rates from survival data: an alternative to two-component mixture models, *Journal of the American Statistical Association*, **98**, 1063–1078.
- Yakovlev AY and Tsodikov AD (1996). *Stochastic Models of Tumor Latency and Their Biostatistical Applications*, World Scientific Publishing, Singapore.
- Zhu H, Ibrahim JG, Lee S, and Zhang H (2007). Perturbation selection and influence measures in local influence analysis, *The Annals of Statistics*, **35**, 2565–2588.
- Zografos K and Balakrishnan N (2009). On families of beta-and generalized gamma-generated distributions and associated inference, *Statistical Methodology*, **6**, 344–362.