

The log-beta Weibull regression model with application to predict recurrence of prostate cancer

**Edwin M. M. Ortega, Gauss M. Cordeiro
& Michael W. Kattan**

Statistical Papers

ISSN 0932-5026

Stat Papers

DOI 10.1007/s00362-011-0414-1



Your article is protected by copyright and all rights are held exclusively by Springer-Verlag. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your work, please use the accepted author's version for posting to your own website or your institution's repository. You may further deposit the accepted author's version on a funder's repository at a funder's request, provided it is not made publicly available until 12 months after publication.

The log-beta Weibull regression model with application to predict recurrence of prostate cancer

Edwin M. M. Ortega · Gauss M. Cordeiro · Michael W. Kattan

Received: 8 August 2010 / Revised: 7 November 2011
© Springer-Verlag 2011

Abstract We study the properties of the called log-beta Weibull distribution defined by the logarithm of the beta Weibull random variable (Famoye et al. in *J Stat Theory Appl* 4:121–136, 2005; Lee et al. in *J Mod Appl Stat Methods* 6:173–186, 2007). An advantage of the new distribution is that it includes as special sub-models classical distributions reported in the lifetime literature. We obtain formal expressions for the moments, moment generating function, quantile function and mean deviations. We construct a regression model based on the new distribution to predict recurrence of prostate cancer for patients with clinically localized prostate cancer treated by open radical prostatectomy. It can be applied to censored data since it represents a parametric family of models that includes as special sub-models several widely-known regression models. The regression model was fitted to a data set of 1,324 eligible prostate cancer patients. We can predict recurrence free probability after the radical prostatectomy in terms of highly significant clinical and pathological explanatory variables associated with the recurrence of the disease. The predicted probabilities of remaining free of cancer progression are calculated under two nested models.

E. M. M. Ortega (✉)
Departamento de Ciências Exatas, Universidade de São Paulo,
Piracicaba, SP 13418-900, Brazil
e-mail: edwin@esalq.usp.br

G. M. Cordeiro
Departamento de Estatística, Universidade Federal de Pernambuco,
Recife, PE 50740–540, Brazil
e-mail: gausscordeiro@uol.com.br

M. W. Kattan
Department of Quantitative Health Sciences, Cleveland Clinic,
Desk JN3-01, 9500 Euclid Avenue, Cleveland, OH 44195, USA
e-mail: kattanm@ccf.org

Keywords Beta Weibull distribution · Censored data · Log-beta Weibull distribution · Log-Weibull regression model · Survival function

1 Introduction

Standard lifetime distributions usually present very strong restrictions to produce bathtub curves, and thus appear to be inappropriate for interpreting data with this characteristic. Some distributions were introduced to model this kind of data, as the generalized gamma distribution (Stacy 1962), the exponential power family (Smith and Bain 1975), the beta integrated model (Hjorth 1980), and the generalized log-gamma distribution (Lawless 2003), among others. A good review of these models is described, for instance, in Rajarshi and Rajarshi (1988). In the last decade, new classes of distributions for modeling this type of data based on extensions of the Weibull distribution were developed. See, for example, the exponentiated Weibull (EW) (Mudholkar et al. 1995), the additive Weibull (Xie and Lai 1995), the modified Weibull (Lai et al. 2003), the beta Weibull (BW) (Famoye et al. 2005; Lee et al. 2007) and the generalized modified Weibull (Carrasco et al. 2008) distributions. Further, Cordeiro et al. (2011) investigated several mathematical properties of the BW geometric distribution, which is a highly flexible lifetime model to cope with different degrees of kurtosis and asymmetry. The BW distribution, due to its flexibility in accommodating the four types of the risk function (i.e. increasing, decreasing, unimodal and bathtub) depending on its parameters, can be used in a variety of problems in modeling survival data. The main motivation for the use of the BW model is that it contains as special sub-models several distributions such as the EW, exponentiated exponential (EE) (Gupta and Kundu 1999) and generalized Rayleigh (GR) (Kundu and Raqab 2005) distributions, among others.

Prostate cancer is the second most common cancer in American men and also the second leading cause of cancer death. The American Cancer Society estimates (in 2010) 217,730 new cases, 32,050 deaths per year and a ten year relative survival rate of 91% for all stages combined. A man with a localized prostate cancer may have a high probability of full recovery if he receives a radical prostatectomy (surgical removal of the prostate gland). Radical prostatectomy provides excellent control of prostate cancer confined to the prostate gland. However, when the cancer breaches the capsule, the cancer recurrence after this surgery is quite higher.

Accurate models to predict cancer recurrence after radical prostatectomy for clinically localized prostate patients are important for the rational application of adjuvant therapy and patient counseling. Previous studies by Kattan et al. (1999) and Stephenson et al. (2005) indicate that some individual patient characteristics such as the PSA value before surgery, biopsy Gleason sum, extracapsular extension, surgical margins, seminal vesicle invasion, lymph node involvement, neo-adjuvant hormone, experience of the surgeon, year of the surgery, among others variables, are very important to predict the risk of prostate cancer recurrence after open radical prostatectomy. Patient follow-up was conducted according to accepted clinical practice, and prostate cancer recurrence is defined as a PSA level > 0.4 ng/ml.

For the first time, we propose a log-beta Weibull (LBW) regression model to predict the t months biochemical recurrence free probability after radical prostatectomy in

terms of highly significant clinical and pathologic variables associated with disease recurrence after surgery. The study cohort comprises 1,324 patients with clinically localized prostate cancer treated by open radical prostatectomy between 1987 and 2003. Patient data were obtained from the Cleveland Clinic from a single surgeon. Patients with clinical stage T1a or T1b disease, who received neoadjuvant therapy, adjuvant therapy or who had missing data for prostate specific antigen were excluded. All information was obtained with appropriate Institutional Review Board waivers.

In this article, we propose a location-scale regression model based on the LBW distribution, referred to as the LBW regression model, which is a feasible alternative for modeling the four existing types of failure rate functions. Some inferential issues were carried out using the asymptotic distribution of the maximum likelihood estimators (MLEs). The sections are organized as follows. In Sect. 2, we define the LBW distribution. Mathematical properties of this distribution are investigate in Sect. 3. In Sect. 4, we obtain the order statistics. We propose a LBW regression model for censored data and discuss inferential issues in Sect. 5. In Sect. 6, a prostate cancer data set is analyzed to show the flexibility, practical relevance and applicability of our regression model. Section 7 ends with some concluding remarks.

2 The log-beta Weibull distribution

Most generalized Weibull distributions have been proposed in reliability literature to provide better fitting of certain data sets than the traditional two and three parameter Weibull models. The BW density function (Famoye et al. 2005) with four parameters $a > 0$, $b > 0$, $c > 0$ and $\lambda > 0$ is given by (for $t > 0$)

$$f(t) = \frac{c}{\lambda^c B(a, b)} t^{c-1} \exp \left\{ -b \left(\frac{t}{\lambda} \right)^c \right\} \left[1 - \exp \left\{ - \left(\frac{t}{\lambda} \right)^c \right\} \right]^{a-1}, \quad (1)$$

where $B(a, b) = [\Gamma(a)\Gamma(b)]/\Gamma(a+b)$ is the beta function and $\Gamma(\cdot)$ is the gamma function. Here, a and b are two additional shape parameters to the Weibull distribution to model the skewness and kurtosis of the data.

The important characteristic of the BW distribution is that it contains, as special sub-models, the EE (Gupta and Kundu 1999), EW (Mudholkar et al. 1995) and GR (Kundu and Raqab 2005) distributions, and some other distributions (see, for example, Cordeiro et al. 2011). The survival and hazard rate functions corresponding to (1) are

$$S(t) = 1 - \frac{1}{B(a, b)} \int_0^{1-\exp\{-(t/\lambda)^c\}} w^{a-1} (1-w)^{b-1} dw = 1 - I_{1-\exp\{-(t/\lambda)^c\}}(a, b)$$

and

$$h(t) = \frac{c(1/\lambda)^c t^{c-1} \exp\{-b(t/\lambda)^c\} [1 - \exp\{-(t/\lambda)^c\}]^{a-1}}{B(a, b) [1 - I_{1-\exp\{-(t/\lambda)^c\}}(a, b)]},$$

respectively, where $I_y(a, b) = B(a, b)^{-1} \int_0^y w^{a-1} (1-w)^{b-1} dw$ is the incomplete beta function ratio.

Let T be a random variable having the BW density function (1). We study the mathematical properties of the LBW distribution defined by the random variable $Y = \log(T)$. The density function of Y , parameterized in terms of $\sigma = c^{-1}$ and $\mu = \log(\lambda)$, can be expressed as

$$f(y; a, b, \sigma, \mu) = \frac{1}{\sigma B(a, b)} \exp \left\{ \left(\frac{y - \mu}{\sigma} \right) - b \exp \left(\frac{y - \mu}{\sigma} \right) \right\} \left\{ 1 - \exp \left[- \exp \left(\frac{y - \mu}{\sigma} \right) \right] \right\}^{a-1}, \quad (2)$$

where $-\infty < y < \infty$, $\sigma > 0$ and $-\infty < \mu < \infty$. We refer to the new model (2) as the LBW distribution, say $Y \sim \text{LBW}(\mu, \sigma, a, b)$, where μ is a location parameter, σ is a dispersion parameter and a and b are shape parameters. The following results holds

$$\text{if } T \sim \text{BW}(\lambda, a, b, c) \text{ then } Y = \log(T) \sim \text{LBW}(\mu, \sigma, a, b).$$

We emphasize that the LBW distribution could also be called the beta extreme value (BEV) distribution, since they are identical. The survival function corresponding to (2) is

$$\begin{aligned} S(y) &= 1 - \frac{1}{B(a, b)} \int_0^{1 - \exp \left[- \exp \left(\frac{y - \mu}{\sigma} \right) \right]} w^{a-1} (1-w)^{b-1} dw \\ &= 1 - I_{1 - \exp \left[- \exp \left(\frac{y - \mu}{\sigma} \right) \right]}(a, b). \end{aligned} \quad (3)$$

3 Properties of the LBW distribution

Here, we study some properties of the standardized LBW random variable defined by $Z = (Y - \mu)/\sigma$. The density function of Z reduces to

$$\pi(z; a, b) = \frac{1}{B(a, b)} \exp[z - b \exp(z)] \{1 - \exp[-\exp(z)]\}^{a-1}, \quad -\infty < z < \infty. \quad (4)$$

The associated cumulative distribution function (cdf) is $F_Z(z) = I_{1 - \exp[-\exp(z)]}(a, b)$. The basic exemplar $a = b = 1$ corresponds to the standard extreme-value distribution.

3.1 Linear combination

By expanding the binomial term in (4), we can write

$$\pi(z; a, b) = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} (-1)^j \binom{a-1}{j} \exp[z - (b+j) \exp(z)]. \quad (5)$$

The density function $h_b = b \exp[z - b \exp(z)]$ (for $b > 0$) gives the Kumaraswamy extreme value (KumEV) distribution (Cordeiro and Castro 2011) with parameters one and b . Its associated cumulative function is $H_a(x) = 1 - [1 - \exp(-e^x)]^a$. Thus,

$$\pi(z; a, b) = \sum_{j=0}^{\infty} w_j h_{b+j}(z),$$

where the coefficients are

$$w_j = \frac{(-1)^j \binom{a-1}{j}}{(b+j)B(a, b)}.$$

So, the LBW density function can be expressed as a linear combination of KumEV densities. For $a = 1$, the LBW distribution reduces to the KumEV distribution with parameters one and b . For $b = 1$, it becomes the log exponentiated Weibull, which is a new model defined here. The LBW random variable Z can be generated directly from the beta variate V with parameters $a > 0$ and $b > 0$ by $Z = \log[-\log(1 - V)]$.

3.2 Moments

The s th ordinary moment of the LBW distribution (4) is

$$\mu'_s = E(Z^s) = \frac{1}{B(a, b)} \int_{-\infty}^{\infty} z^s \exp[z - b \exp(z)] \{1 - \exp[-\exp(z)]\}^{a-1} dz.$$

By expanding the binomial term and setting $w = e^z$, we obtain

$$\mu'_s = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} \binom{a-1}{j} (-1)^j \int_0^{\infty} \log(w)^s \exp[-(b+j)w] dw.$$

The above integral can be calculated from Prudnikov et al. (1986, Vol. 1, Eq. 2.6.21.1) as

$$I(s, j) = \left(\frac{\partial}{\partial p} \right)^s [(b+j)^{-p} \Gamma(p)] \Big|_{p=1}$$

and then

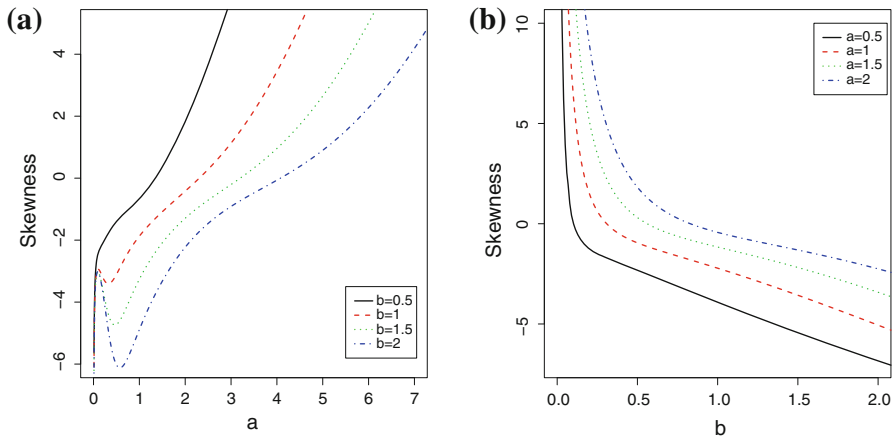


Fig. 1 Skewness of the LBW distribution. **a** Function of a for some values of b . **b** Function of b for some values of a

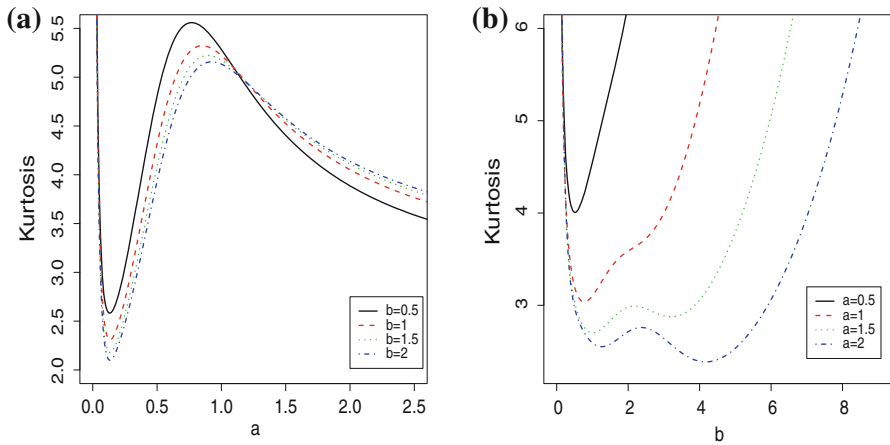


Fig. 2 Kurtosis of the LBW distribution. **a** Function of a for some values of b . **b** Function of b for some values of a

$$\mu'_s = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} (-1)^j \binom{a-1}{j} I(s, j). \quad (6)$$

Equation 6 gives the moments of the LBW distribution. The skewness and kurtosis measures can be calculated from the ordinary moments using well-known relationships. These measures are controlled mainly by the parameters a and b . Plots of the skewness and kurtosis for selected values of b as function of a , and for selected values of a as function of b , for $\mu = 0$ and $\sigma = 1$, are shown in Figs. 1 and 2, respectively. These plots reveal that the skewness for fixed b (a), as function of a (b) decreases and then increases (decreases), whereas the kurtosis for fixed b (a) as function of a (b) decreases, increases and then decreases (decreases and then increases).

3.3 Moment generating function

The moment generating function (mgf) of Z , say $M(t) = E(e^{tZ})$, follows from (4) as

$$M(t) = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} (-1)^j (a-1)^j \int_0^{\infty} w^t \exp[-(b+j)w] dw$$

and then

$$M(t) = \frac{\Gamma(t+1)}{B(a, b)} \sum_{j=0}^{\infty} (-1)^j \binom{a-1}{j} (b+j)^{-(t+1)}. \quad (7)$$

Clearly, the moments (6) can be obtained from (7) by simple differentiation.

3.4 Quantile function

We now give an expansion for the quantile function $q = F^{-1}(p)$ (given p) of the LBW distribution. First, we have $p = F(q) = I_s(a, b)$, where $s = 1 - \exp[-\exp(q)]$. It is possible to obtain s as function of p from some expansions for the inverse of the beta incomplete function $s = I_p^{-1}(a, b)$. One of them can be found in Wolfram website¹ as

$$s = I_p^{-1}(a, b) = w + \frac{b-1}{a+1} w^2 + \frac{(b-1)(a^2+3ba-a+5b-4)}{2(a+1)^2(a+2)} w^3 \\ + \frac{(b-1)[a^4+(6b-1)a^3+(b+2)(8b-5)a^2+(33b^2-30b+4)a+b(31b-47)+18]}{3(a+1)^3(a+2)(a+3)} w^4 + O(p^{5/a}),$$

where $w = [a p B(a, b)]^{1/a}$ for $a > 0$. Hence, $q = \log[-\log(1-s)]$ and the above expansion defines the LBW quantile function.

3.5 Mean deviations

The amount of scatter in Z is evidently measured to some extent by the totality of deviations from the mean μ'_1 and median m . These are known as the mean deviations about the mean and the median—defined by

$$\delta_1(Z) = \int_{-\infty}^{\infty} |x - \mu| \pi(z; a, b) dz \quad \text{and} \quad \delta_2(Z) = \int_{-\infty}^{\infty} |x - m| \pi(z; a, b) dz,$$

respectively. From (6) with $s = 1$, we obtain

¹ <http://functions.wolfram.com/06.23.06.0004.01>.

$$\mu'_1 = E(Z) = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} \frac{(-1)^{j+1} \binom{a-1}{j}}{(b+j)} [\gamma + \log(b+j)],$$

where γ is Euler's constant. The median m is calculated from the nonlinear equation $I_{1-\exp[-\exp(m)]}(a, b) = 1/2$. The measures $\delta_1(Z)$ and $\delta_2(Z)$ can be expressed as

$$\delta_1(Z) = 2\mu'_1[F_Z(\mu'_1) - 1] + 2T(\mu'_1) \quad \text{and} \quad \delta_2(Z) = 2T(m) - \mu'_1,$$

where $T(q) = \int_q^{\infty} z \pi(z; a, b) dz$. We obtain $T(q)$ as

$$\begin{aligned} T(q) &= \frac{1}{B(a, b)} \int_q^{\infty} z \exp[z - b \exp(z)] \{1 - \exp[-\exp(z)]\}^{a-1} \\ &= \frac{1}{B(a, b)} \sum_{j=0}^{\infty} (-1)^j \binom{a-1}{j} \int_{e^q}^{\infty} \log(w) \exp[-(b+j)w] dw. \end{aligned}$$

For $b > 0$ and $p > 0$, using a result in [Prudnikov et al. \(1986, Vol. 1, Eq. 1.6.10.3\)](#), namely

$$K(p, a) = \int_p^{\infty} \log(x) e^{-bx} dx = b^{-1} [e^{-bp} \log(p) - E_i(-bp)],$$

where $E_i(x) = \int_{-\infty}^x t^{-1} e^t dt$ is the exponential integral, we obtain

$$T(q) = \frac{1}{B(a, b)} \sum_{j=0}^{\infty} \frac{(-1)^j \binom{a-1}{j}}{(b+j)} [q e^{-(b+j)e^q} - E_i(-(b+j)e^q)].$$

This equation for $T(q)$ can be used to determine Bonferroni and Lorenz curves that have applications in economics to study income and poverty, reliability, demography, insurance and medicine and other fields. They are defined by

$$B(p) = \frac{\mu'_1 - T(q)}{p\mu'_1} \quad \text{and} \quad L(p) = \frac{\mu'_1 - T(q)}{\mu'_1},$$

respectively, where $q = F^{-1}(p)$ can be calculated for given p from the quantile function.

4 Order statistics

Order statistics make their appearance in many areas of statistical theory and practice. The density $f_{i:n}(x)$ of the i th order statistic ($Z_{i:n}$) for $i = 1, \dots, n$ from i.i.d. LBW random variables Z_1, \dots, Z_n is simply given by

$$f_{i:n}(z) = \frac{\pi(z; a, b)}{B(i, n-i+1)} \sum_{j=0}^{n-i} (-1)^j \binom{n-i}{j} I_{1-\exp[-\exp(z)]}(a, b)^{i+j-1}. \quad (8)$$

We now obtain an expansion for the density function of the LBW order statistics. First, we use the incomplete beta function expansion for $b > 0$ real non-integer

$$I_{1-\exp[-\exp(z)]}(a, b) = \frac{1}{B(a, b)} \sum_{m=0}^{\infty} \frac{(1-b)_m \{1 - \exp[-\exp(z)]\}^{a+m}}{(a+m) m!},$$

where $(f)_k = \Gamma(f+k)/\Gamma(f)$ is the ascending factorial. We have

$$I_{1-\exp[-\exp(z)]}(a, b) = \sum_{k=0}^{\infty} d_k \exp[-k \exp(z)], \quad (9)$$

where the coefficients d_k (for $k = 0, 1, \dots$) are

$$d_k = \frac{(-1)^k}{B(a, b)} \sum_{m=0}^{\infty} \frac{(1-b)_m \binom{a+m}{k}}{(a+m) m!}.$$

Using the identity $(\sum_{k=0}^{\infty} a_k x^k)^n = \sum_{j=0}^{\infty} c_{n,k} x^k$ for n positive integer (see [Gradshteyn and Ryzhik 2000](#)) in $I_{1-\exp[-\exp(z)]}(a, b)^{i+j-1}$, we readily obtain

$$I_{1-\exp[-\exp(z)]}(a, b)^{i+j-1} = \sum_{k=0}^{\infty} c_{i+j-1,k} \exp[-k \exp(z)], \quad (10)$$

where $c_{i+j-1,0} = d_0^{i+j-1}$ and, for $k = 1, 2, \dots$,

$$c_{i+j-1,k} = (k d_0)^{-1} \sum_{r=1}^k [(i+j)r - k] d_i c_{i+j-1,k-r}. \quad (11)$$

Substituting (5) and (10) in Eq. 8, we have

$$f_{i:n}(z) = \sum_{m,k=0}^{\infty} (-1)^m \binom{a-1}{m} v_k \exp[z - (b+m+k) \exp(z)], \quad (12)$$

where

$$v_k = \frac{\sum_{j=0}^{n-i} (-1)^j \binom{n-i}{j} c_{i+j-1,k}}{B(i, n-i+1) B(a, b)}.$$

The moments, mgf, mean deviations of the LBW order statistics are easily obtained from (12) using the same calculations for those quantities of the LBW distribution. For example, the s th ordinary moment of $Z_{i:n}$ is expressed as

$$E(X_{i:n}^s) = \sum_{m,k=0}^{\infty} (-1)^m \binom{a-1}{m} v_k I(s, m+k),$$

where $I(s, m+k)$ is defined just before (6).

5 The log-beta Weibull regression model

In many practical applications, the lifetimes are affected by explanatory variables such as the cholesterol level, blood pressure, weight and many others. Parametric models to estimate univariate survival functions and for censored data regression problems are widely used. A parametric model that provides a good fit to lifetime data tends to yield more precise estimates of the quantities of interest. Based on the LBW density function, we propose a linear location-scale regression model linking the response variable y_i and the explanatory variable vector $\mathbf{x}_i^T = (x_{i1}, \dots, x_{ip})$ as follows

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \sigma z_i, \quad i = 1, \dots, n, \quad (13)$$

where the random error z_i has density function (4), $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$, $\sigma > 0$, $a > 0$ and $b > 0$ are unknown parameters. The parameter $\mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ is the location of y_i . The location parameter vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^T$ is represented by a linear model $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}$, where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ is a known model matrix. The LBW model (13) opens new possibilities for fitting many different types of data. It contains as special sub-models the following well-known regression models. For $a = b = 1$, we obtain the classical Weibull regression model (see, Lawless 2003). If $\sigma = 1$ and $\sigma = 0.5$, in addition to $a = b = 1$, it coincides with the exponential and Rayleigh regression models, respectively. For $b = 1$, it reduces to the log-exponentiated Weibull regression model (Cancho et al. 1999, 2009; Ortega et al. 2006; Hashimoto et al. 2010). If $\sigma = 1$, in addition to $b = 1$, the LBW model yields the log-exponentiated exponential regression. If $\sigma = 0.5$, in addition to $b = 1$, it becomes the log-generalized Rayleigh regression model. For $\sigma = 1$, we have a new model called the log-beta exponential regression model.

Consider a sample $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$ of n independent observations, where each random response is defined by $y_i = \min\{\log(t_i), \log(c_i)\}$. We assume non-informative censoring such that the observed lifetimes and censoring times are independent. Let F and C be the sets of individuals for which y_i is the log-lifetime and

log-censoring, respectively. The log-likelihood function for the vector of parameters $\theta = (a, b, \sigma, \beta^T)^T$ from model (13) has the form $l(\theta) = \sum_{i \in F} \log[f(y_i)] + \sum_{i \in C} \log[S(y_i)]$, where $f(y_i)$ is the density function (2) and $S(y_i)$ is the survival function (3) of Y_i . The log-likelihood function for θ reduces to

$$l(\theta) = -r \log \{\log(\sigma) + \log[B(a, b)]\} + \sum_{i \in F} z_i - b \sum_{i \in F} \exp(z_i) + (a-1) \sum_{i \in F} \log \{1 - \exp[-\exp(z_i)]\} + \sum_{i \in C} \log \{1 - I_{1-\exp[-\exp(z_i)]}(a, b)\}, \quad (14)$$

where r is the number of uncensored observations (failures) and $z_i = (y_i - \mathbf{x}_i^T \beta) / \sigma$.

The MLE $\hat{\theta}$ of the vector θ of unknown parameters can be calculated by maximizing the log-likelihood (14). We use the subroutine NLMixed in SAS to calculate $\hat{\theta}$. Initial values for σ and β can be taken from the fit of the log-Weibull (LW) regression model with $a = b = 1$. The fitted LBW model gives the estimated survival function of Y for any individual with explanatory vector \mathbf{x}

$$S(y; \hat{a}, \hat{b}, \hat{\sigma}, \hat{\beta}^T) = 1 - I_{1-\exp\left[-\exp\left(\frac{y-\mathbf{x}^T \hat{\beta}}{\hat{\sigma}}\right)\right]}(\hat{a}, \hat{b}). \quad (15)$$

The invariance property of the MLEs yields the survival function for $T = \exp(Y)$

$$S(t; \hat{a}, \hat{b}, \hat{c}, \hat{\lambda}) = 1 - I_{1-\exp\{-(t/\hat{\lambda})^{\hat{c}}\}}(\hat{a}, \hat{b}), \quad (16)$$

where $\hat{c} = 1/\hat{\sigma}$ and $\hat{\lambda} = \exp(\mathbf{x}^T \hat{\beta})$.

Under conditions that are fulfilled for the parameter vector θ in the interior of the parameter space but not on the boundary, the asymptotic distribution of $\sqrt{n}(\hat{\theta} - \theta)$ is multivariate normal $N_{p+3}(0, K(\theta)^{-1})$, where $K(\theta)$ is the information matrix. The asymptotic covariance matrix $K(\theta)^{-1}$ of $\hat{\theta}$ can be approximated by the inverse of the $(p+3) \times (p+3)$ observed information matrix $-\ddot{L}(\hat{\theta}) = \{\mathbf{L}_{r,s}\}$, whose elements $\mathbf{L}_{r,s}$ are given in Appendix A.

The approximate multivariate normal distribution $N_{p+3}(0, -\ddot{L}(\hat{\theta})^{-1})$ for $\hat{\theta}$ can be used in the classical way to construct approximate confidence regions for some parameters in θ . We can use the likelihood ratio (LR) statistic for comparing some special sub-models with the LBW model. We consider the partition $\theta = (\theta_1^T, \theta_2^T)^T$, where θ_1 is a subset of parameters of interest and θ_2 is a subset of remaining parameters. The LR statistic for testing the null hypothesis $H_0 : \theta_1 = \theta_1^{(0)}$ versus the alternative hypothesis $H_1 : \theta_1 \neq \theta_1^{(0)}$ is given by $w = 2\{\ell(\hat{\theta}) - \ell(\tilde{\theta})\}$, where $\tilde{\theta}$ and $\hat{\theta}$ 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 θ_1 of parameters of interest.

6 Application: prostate cancer recurrence data

In this section, we develop an application of the LBW regression model to a prostate cancer data. The study cohort comprises 1,324 patients with clinically localized prostate cancer treated by open radical prostatectomy between 1987 and 2003. Patient data were obtained from the Cleveland Clinic from a single surgeon. The data consist of the random response variable given by the number of months (y_i) without detectable disease after prostatectomy. Uncensored observations correspond to patients having cancer recurrent time computed. Censored observations correspond to patients who were not observed to have cancer recurrence at the time the data were collected. The numbers of censored and uncensored observations are 1,096 and 228, respectively, of the total of 1,324 patients. The following explanatory variables were associated with each patient (for $i = 1, \dots, 1,324$):

- δ_i : is the event indicator where 1 represents the event and 0 is censored;
- $neoad_i$: is whether the patient received neo-adjuvant hormones, i.e., treated with hormone therapy prior to radical prostatectomy (yes = 1 and no = 0);
- psa_i : is the PSA value (in ng/mL) from the laboratory report before receiving prostatectomy;
- $ece\ t_i$: is the extracapsular extension on path report (yes = 1, no = 0);
- $svi\ t_i$: is the seminal vesicle invasion on path report (yes = 1, no = 0);
- pgx : is the pathology report Gleason sum 4–7, 7, 8–10. We construct two dummy random variables: ($pgx\ t1$: [4,7) versus 7 and $pgx\ t2$: [4,7) versus [8,10]);
- $lni\ t_i$: is the lymph node involvement on path report (neg = 1, pos = 0);
- $sm\ t_i$: is surgical margin status (yes = 1, no = 0).

Now, we present results by fitting the model

$$y_i = \beta_0 + \beta_1 neoad_i + \beta_2 psa_i + \beta_3 ece\ t_i + \beta_4 svi\ t_i \\ + \beta_5 lni\ t_i + \beta_6 pgx\ t_{1i} + \beta_7 pgx\ t_{2i} + \beta_8 sm\ t_i + \sigma z_i,$$

where the dependent variable y_i follows the LBW density function (2) for $i = 1, \dots, 1,324$. The MLEs of the model parameters are calculated using the procedure NLMixed in SAS. Iterative maximization of the logarithm of the likelihood function (14) starts with initial values for β and σ taken from the fit of the LW regression model with $a = b = 1$.

Table 1 lists the MLEs of the parameters for the LBW and LW regression models fitted to the current data. The LR statistic for testing the hypotheses $H_0: a = b = 1$ versus $H_1: H_0$ is not true, i.e., to compare the LW and LBW regression models, is $w = 2\{-716.45 - (-730.80)\} = 28.70$ (p -value < 0.0001), which gives favorable indications toward to the LBW model. The LBW model involves two extra parameters which gives it more flexibility to fit the data. The fitted LBW regression model indicates that all explanatory variables are significant at 5%.

Cox (1972) proposed a very useful regression model for analyzing censoring failure times, where the random variable of interest represents failure time and the failures times are assumed identically distributed in some specified form. He noted that if the

Table 1 MLEs of the parameters for the LBW and LW regression models fitted to the recurrence prostate cancer data

θ	LBW regression model				LW regression model			
	Estimate	S.E	<i>p</i> -value	95%CI	Estimate	S.E	<i>p</i> -value	95%CI
<i>a</i>	267.08	0.11	—	(266.85; 267.30)	1	—	—	—
<i>b</i>	21.63	0.12	—	(21.40; 21.86)	1	—	—	—
σ	24.12	1.21	—	(21.74; 26.50)	1.24	0.07	—	(1.11; 1.37)
β_0	−16.00	1.04	<0.0001	(−18.05; −13.96)	7.40	0.42	<0.0001	(6.56; 8.23)
β_1	−0.59	0.23	0.0085	(−1.04; −0.15)	−0.72	0.21	0.0006	(−1.13; −0.31)
β_2	−0.02	0.007	0.0017	(−0.04; −0.01)	−0.01	0.004	0.0040	(−0.02; −0.003)
β_3	−0.84	0.20	<0.0001	(−1.23; −0.45)	−0.93	0.21	<0.0001	(−1.35; −0.51)
β_4	−1.01	0.27	0.0002	(−1.54; −0.48)	−0.76	0.23	0.0013	(−1.22; −0.30)
β_5	0.67	0.25	0.0075	(0.18; 1.16)	0.67	0.29	0.0227	(0.09; 1.25)
β_6	−0.90	0.19	<0.0001	(−1.27; −0.52)	−1.01	0.23	<0.0001	(−1.46; −0.56)
β_7	−2.09	0.30	<0.0001	(−2.68; −1.51)	−2.00	0.30	<0.0001	(−2.59; −1.42)
β_8	−1.09	0.18	<0.0001	(−1.46; −0.74)	−0.88	0.19	<0.0001	(−1.25; −0.51)

proportional hazards assumption holds (or, is assumed to hold) then it is possible to estimate the effect parameter(s) without any consideration of the hazard function (non-parametric approach). This approach to survival data is called proportional hazards model. The Cox model may be specialized if a reason exists to assume that the baseline hazard follows a parametric form. In this case, the baseline hazard can be replaced by a parametric density. Typically, we can then maximize the full likelihood which greatly simplifies model-fitting and provides interpretability at the cost of flexibility.

Let $R(t_i)$ be the set of individuals at risk at time t_i . Conditionally on the risk sets, the required likelihood $L(\beta)$ can be expressed as

$$L(\beta) = \prod_{i=1}^n \left[\frac{\exp(\mathbf{x}_i^T \beta)}{\sum_{j \in R(t_i)} \exp(\mathbf{x}_j^T \beta)} \right]^{\delta_i}, \quad (17)$$

where δ_i is the censoring indicator.

The MLE $\hat{\beta}$ of β can be calculated by maximizing the likelihood function (17) using the matrix programming language SAS. Table 2 provides the estimates, corresponding standard errors and *p*-values for the fitted Cox regression model. All explanatory

Table 2 Estimates for the Cox regression model fitted to the recurrence prostate cancer data

Parameter	Estimate	SE	<i>p</i> -value	95% CI
β_1	0.558	0.168	0.0009	(0.228, 0.887)
β_2	0.008	0.003	0.0122	(0.002, 0.014)
β_3	0.755	0.167	<.0001	(0.428, 1.082)
β_4	0.618	0.186	0.0009	(0.253, 0.982)
β_5	−0.539	0.239	0.0240	(−1.007, −0.071)
β_6	0.797	0.183	<.0001	(0.439, 1.155)
β_7	1.598	0.237	<.0001	(1.134, 2.062)
β_8	0.703	0.147	<.0001	(0.419, 0.986)

Table 3 AIC, BIC and GD statistics for comparing the LBW and LW models

Model	AIC	BIC	GD
LBW	1456.9	1519.1	1432.9
LW	1481.6	1533.5	1461.6
Cox proportional hazards	2742.4	2742.4	2726.4

variables are marginally significant at the 5% significance level. For a prostate cancer patient with explanatory vector \mathbf{x} , the recurrence free probability, say $P(T \geq t; \boldsymbol{\beta}, \mathbf{x}) = S(t; \boldsymbol{\beta}, \mathbf{x})$, can be predicted from Cox regression model by

$$S(t; \hat{\boldsymbol{\beta}}, \mathbf{x}) = [\hat{S}_0(t)]^{\exp(\mathbf{x}^T \hat{\boldsymbol{\beta}})}, \quad (18)$$

where $\hat{S}_0(t) = \exp[-\hat{\Lambda}_0(T)]$, $\hat{\Lambda}_0(T) = \sum_{j:t_j < T} \left[\frac{d_j}{\sum_{l \in R_j} \exp(\mathbf{x}_l^T \hat{\boldsymbol{\beta}})} \right]$ and d_j is the number of failures in t_j .

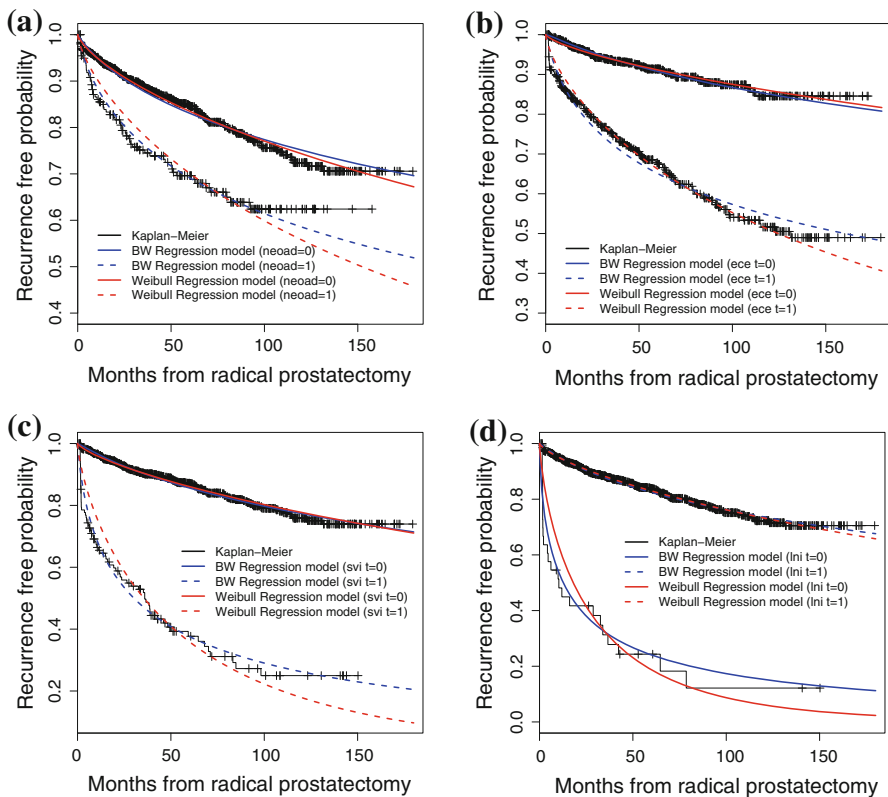


Fig. 3 Kaplan–Meier curves stratified by explanatory variable and estimated survival functions to the recurrence prostate cancer data: **a** *neoad* explanatory variable. **b** *ece* explanatory variable. **c** *svi* explanatory variable. **d** *lni* explanatory variable

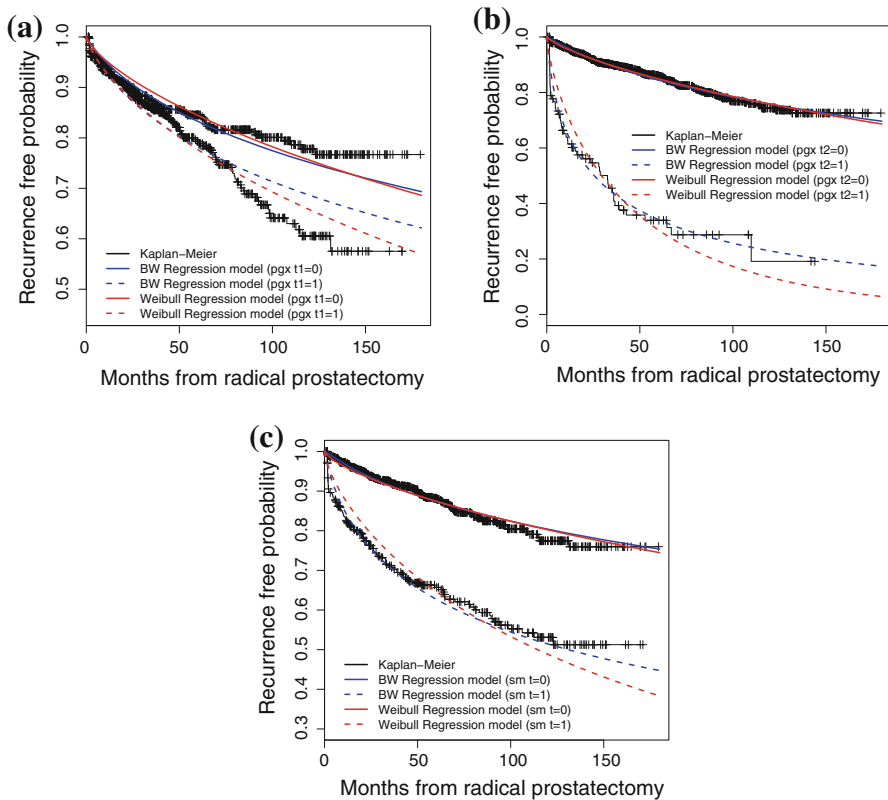


Fig. 4 Kaplan–Meier curves stratified by explanatory variable and estimated survival functions to the recurrence prostate cancer data: **a** $pgxt1$ explanatory variable. **b** $pgxt2$ explanatory variable. **c** smt explanatory variable

Further, Table 3 lists the Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC) and the global deviance (GD) given by $-2 \log\{L\hat{\beta}\}$, to compare the LBW, LW and Cox proportional hazard regression models. The LBW regression model outperforms the other models irrespective of the criteria and it can be used effectively in the analysis of these data. So, the proposed model is a great alternative to model survival data.

In order to assess if the model is appropriate, we fit the LBW and LW regression models for each explanatory variable. In Figs. 3a, b, c, d and 4a, b, c, we plot the empirical survival function and the estimated survival function (16) for each explanatory variable. We conclude that the LBW regression model provides a good fit to these data.

6.1 Prediction

For a prostate cancer patient treated by open radical prostatectomy with explanatory vector \mathbf{x} , we can estimate the recurrence free probability, say $P(T \geq t; a, b, \sigma, \beta, \mathbf{x}) =$

Table 4 Recurrence free probability under the BW regression model

Patient	<i>neoad</i>	<i>psa</i>	<i>ecet</i>	<i>svi t</i>	<i>lni t</i>	<i>pgx t1</i>	<i>pgx t2</i>	<i>sm t</i>
A	0	5	1	0	1	1	0	1
B	0	25	1	0	1	1	0	1
C	1	30	0	1	0	0	1	0
D	1	60	0	1	0	0	1	0

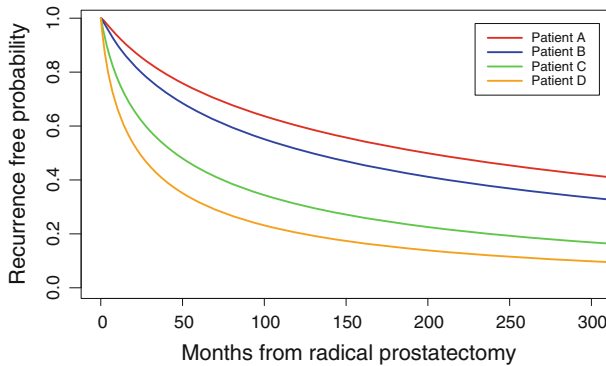


Fig. 5 Estimated recurrence free probability curves for patients A, B, C and D

$S(t; a, b, \sigma, \beta, \mathbf{x})$, by using (16). Evidently, the recurrence free probability converges to zero when the linear predictor $\mu_i = \mathbf{x}_i^T \beta$ tends to $-\infty$ and converges to one when the linear predictor goes to $+\infty$. In other words, the recurrence for patients with clinically localized prostate cancer treated by open radical prostatectomy for a fixing time t after the surgery, approaches one (zero) when the linear predictor μ increases to a very large negative (positive) number.

We can use (16) to predict the recurrence free probability $S(t; \mathbf{x}) = S(t; \hat{a}, \hat{b}, \hat{\sigma}, \hat{\beta}, \mathbf{x})$ of prostate cancer at t months. As an illustration, we consider four hypothetical patients A, B, C and D who underwent radical prostatectomy having fixed values for the explanatory variables given in Table 4. In Fig. 5, we provide the plots of the estimated recurrence free probabilities for these four patients.

7 Concluding remarks

We introduce the called log-beta Weibull (LBW) distribution whose hazard rate function accommodates four types of shape forms, namely increasing, decreasing, bathtub and unimodal. We derive expressions for its moments, moment generating function, quantile function, mean deviations and order statistics. Based on this new distribution, we propose a LBW regression model very suitable for modeling censored and uncensored lifetime data. We provide an application to predict cure of prostate cancer. The new regression model allows to perform goodness of fit tests for some known regression models as special cases. Hence, the proposed regression model serves as a good alternative for lifetime data analysis. Further, the new regression model is much

more flexible than the exponentiated Weibull, Weibull and generalized Rayleigh sub-models. In one application to real prostate cancer data, we show that the LBW model can produce better fit than its sub-models. We compare three fitted models using the AIC, BIC and global deviance criterions to give evidence that the LBW regression model outperforms the other two models.

Acknowledgements This work was supported by CNPq and CAPES.

Appendix A: Matrix of second derivatives $-\ddot{\mathbf{L}}(\theta)$

Here, we give the formulas to obtain the second-order partial derivatives of the log-likelihood function. After some algebraic manipulations, we obtain

$$\mathbf{L}_{aa} = \sum_{i \in F} \left[\psi'(a+b) - \psi'(a) \right] - \sum_{i \in C} \left\{ v_i^{-2} ([\dot{I}_{G(z_i)}(a, b)]_a)^2 + v_i^{-1} \left[\frac{[\psi(a) - \psi(a+b)]^2}{B(a, b)} - \frac{\psi'(a) - \psi'(a+b)}{B(a, b)} + M(a) \right] \right\},$$

$$\mathbf{L}_{ab} = \sum_{i \in F} \psi'(a+b) - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_a [\dot{I}_{G(z_i)}(a, b)]_b + v_i^{-1} \left[\frac{[\psi(a) - \psi(a+b)][\psi(b) - \psi(a+b)]}{B(a, b)} + \frac{\psi'(a+b)}{B(a, b)} + M(ab) \right] \right\},$$

$$\mathbf{L}_{a\sigma} = -\sigma^{-1} \sum_{i \in F} z_i o_i - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_a [\dot{I}_{G(z_i)}(a, b)]_\sigma - v_i^{-1} z_i q_i \log[G(z_i)] \right\},$$

$$\mathbf{L}_{a\beta_j} = -\sigma^{-1} \sum_{i \in F} x_{ij} o_i - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_a [\dot{I}_{G(z_i)}(a, b)]_{\beta_j} - v_i^{-1} x_{ij} q_i \log[G(z_i)] \right\},$$

$$\mathbf{L}_{bb} = \sum_{i \in F} \left[\psi'(a+b) - \psi'(b) \right] - \sum_{i \in C} \left\{ v_i^{-2} ([\dot{I}_{G(z_i)}(a, b)]_b)^2 + v_i^{-1} \left[\frac{[\psi(b) - \psi(a+b)]^2}{B(a, b)} - \frac{\psi'(b) - \psi'(a+b)}{B(a, b)} + M(b) \right] \right\},$$

$$\mathbf{L}_{b\sigma} = \sigma^{-1} \sum_{i \in F} z_i \exp(z_i) - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_b [\dot{I}_{G(z_i)}(a, b)]_\sigma - v_i^{-1} z_i q_i \exp(z_i) \right\},$$

$$\mathbf{L}_{b\beta_j} = \sigma^{-1} \sum_{i \in F} x_{ij} \exp(z_i) - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_b [\dot{I}_{G(z_i)}(a, b)]_{\beta_j} - v_i^{-1} x_{ij} q_i \exp(z_i) \right\},$$

$$\begin{aligned} \mathbf{L}_{\sigma\sigma} = & \sum_{i \in F} \left\{ \sigma^{-2}(1 + 2z_i) - b\sigma^{-2}z_i \exp(z_i) + z_i u_i [2 + z_i(1 - \exp(z_i)) - o_i] \right\} \\ & - \sum_{i \in C} \left\{ v_i^{-2} ([\dot{I}_{G(z_i)}(a, b)]_{\sigma})^2 + v_i^{-1} z_i d_i [z_i(b \exp(z_i) - 1) - \sigma^2 z_i u_i - 2] \right\}, \end{aligned}$$

$$\begin{aligned} \mathbf{L}_{\sigma\beta_j} = & \sum_{i \in F} \left\{ \sigma^{-2} x_{ij} - b\sigma^{-2} x_{ij} \exp(z_i)(1 + z_i) + x_{ij} u_i [1 + z_i(1 - \exp(z_i)) - z_i o_i] \right\} \\ & - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_{\sigma} [\dot{I}_{G(z_i)}(a, b)]_{\beta_j} + v_i^{-1} x_{ij} d_i \right. \\ & \times \left. [z_i(b \exp(z_i) - 1) - \sigma^2 z_i u_i - 1] \right\} \end{aligned}$$

and

$$\begin{aligned} \mathbf{L}_{\beta_j\beta_s} = & - \sum_{i \in F} \left\{ b\sigma^{-2} x_{ij} x_{is} \exp(z_i) - x_{ij} x_{is} u_i [1 - \exp(z_i) - o_i] \right\} \\ & - \sum_{i \in C} \left\{ v_i^{-2} [\dot{I}_{G(z_i)}(a, b)]_{\beta_j} [\dot{I}_{G(z_i)}(a, b)]_{\beta_s} \right. \\ & \left. + v_i^{-1} x_{ij} x_{is} d_i [b \exp(z_i) - 1 - \sigma^2 u_i] \right\}, \end{aligned}$$

where

$$z_i = (y_i - \mathbf{x}_i^T \boldsymbol{\beta})/\sigma, \quad G(z_i) = 1 - \exp[-\exp(z)],$$

$$v_i = 1 - I_{G(z_i)}(a, b), \quad o_i = [G(z_i)]^{-1} \exp[z_i - \exp(z_i)],$$

$$\begin{aligned} M(a) &= \int_0^{G(z_i)} w^{a-1} (1-w)^{b-1} [\log(w)]^2 dw, \\ M(b) &= \int_0^{G(z_i)} w^{a-1} (1-w)^{b-1} [\log(1-w)]^2 dw, \end{aligned}$$

$$\begin{aligned} M(ab) &= \int_0^{G(z_i)} w^{a-1} (1-w)^{b-1} \log(w) \log(1-w) dw, \\ q_i &= \sigma^{-1} [G(z_i)]^{a-1} \exp[z_i - b \exp(z_i)], \end{aligned}$$

$$\begin{aligned}u_i &= \sigma^{-2}[G(z_i)]^{-1}(a-1)\exp[z_i - \exp(z_i)], \\d_i &= \sigma^{-2}[G(z_i)]^{a-1}\exp[z_i - b\exp(z_i)],\end{aligned}$$

$$[\dot{I}_{G(z_i)}(a, b)]_a = [\psi(a+b) - \psi(a)]/B(a, b) + \int_0^{G(z_i)} w^{a-1}(1-w)^{b-1} \log(w)dw,$$

$$[\dot{I}_{G(z_i)}(a, b)]_b = [\psi(a+b) - \psi(b)]/B(a, b) + \int_0^{G(z_i)} w^{a-1}(1-w)^{b-1} \log(1-w)dw,$$

$$[\dot{I}_{G(z_i)}(a, b)]_\sigma = -\sigma^{-1}z_i[G(z_i)]^{a-1}\exp[z_i - \exp(z_i)]$$

and

$$[\dot{I}_{G(z_i)}(a, b)]_{\beta_j} = -\sigma^{-1}x_{ij}[G(z_i)]^{a-1}\exp[z_i - \exp(z_i)].$$

References

- Cancho VG, Bolfarine H, Achcar JA (1999) A Bayesian analysis for the exponentiated-Weibull distribution. *J Appl Stat* 8:227–242
- Cancho VG, Ortega EMM, Bolfarine H (2009) The log-exponentiated-Weibull regression models with cure rate: local influence and residual analysis. *J Data Sci* 7:433–458
- Carrasco JMF, Ortega EMM, Cordeiro MG (2008) A generalized modified Weibull distribution for lifetime modeling. *Comput Stat Data Anal* 53:450–462
- Cordeiro GM, de Castro M (2011) A new family of generalized distributions. *J Stat Comput Simul* 81: 883–898
- Cordeiro GM, Silva GO, Ortega EMM (2011) The beta-Weibull geometric distribution. *Statistics*. doi:[10.1080/02331888.2011.577897](https://doi.org/10.1080/02331888.2011.577897)
- Cox DR (1972) Regression models and life tables (with discussion). *J R Stat Soc Ser B (Stat Methodol)* 34:187–220
- Famoye F, Lee C, Olumolade O (2005) The beta-Weibull distribution. *J Stat Theory Appl* 4:121–136
- Gradshteyn IS, Ryzhik IM (2000) In: Jeffrey A, Zwillinger D (eds) *Table of integrals, series, and products*, 6th edn. Academic Press, New York
- Gupta RD, Kundu D (1999) Generalized exponential distributions. *Austral N Z J Stat* 41:173–188
- Hashimoto EM, Ortega EMM, Cancho VG, Cordeiro GM (2010) The log-exponentiated Weibull regression model for interval-censored data. *Comput Stat Data Anal* 54:1017–1035
- Hjorth U (1980) A reliability distributions with increasing, decreasing, constant and bathtub failure rates. *Technometrics* 22:99–107
- Kattan MW, Wheeler TM, Scardino PT (1999) Postoperative nomogram for disease recurrence after radical prostatectomy for prostate cancer. *J Clin Oncol* 17:1499–1507
- Kundu D, Raqab MZ (2005) Generalized Rayleigh distribution: different methods of estimation. *Comput Stat Data Anal* 49:187–200
- Lai CD, Xie M, Murthy DNP (2003) A modified Weibull distribution. *Trans Reliab* 52:33–37
- Lawless JF (2003) *Statistical models and methods for lifetime data*. Wiley, New York
- Lee C, Famoye F, Olumolade O (2007) Beta-Weibull distribution: some properties and applications to censored data. *J Mod Appl Stat Methods* 6:173–186

- Mudholkar GS, Srivastava DK, Friemer M (1995) The exponentiated Weibull family: a reanalysis of the bus-motor-failure data. *Technometrics* 37:436–445
- Ortega EMM, Cancho VG, Bolfarine H (2006) Influence diagnostics in exponentiated-Weibull regression models with censored data. *Stat Oper Res Trans* 30:172–192
- Prudnikov AP, Brychkov YA, Marichev OI (1986) Integrals and series, vol 1. Gordon and Breach Science Publishers, Amsterdam
- Rajarshi S, Rajarshi MB (1988) Bathtub distributions. *Commun Stat Theory Methods* 17:2521–2597
- Smith RM, Bain LJ (1975) An exponential power life testing distributions. *Commun Stat Theory Methods* 4:469–481
- Stacy EW (1962) A generalization of the gamma distribution. *Ann Math Stat* 33:1187–1192
- Stephenson AJ, Scardino PT, Eastham JA, Bianco FJ Jr, Dotan ZA, DiBlasio CJ, Reuther A, Klein EA, Kattan MW (2005) Postoperative nomogram predicting the 10-year probability of prostate cancer recurrence after radical prostatectomy. *J Clin Oncol* 23:7005–7012
- Xie M, Lai CD (1995) Reliability analysis using an additive Weibull model with bathtub-shaped failure rate function. *Reliab Eng Syst Saf* 52:87–93