

## Cure Fraction Models on Survival Data and Covariates with a Bayesian Parametric Estimation Methods

Umar Yusuf Madaki<sup>1</sup>, Babangida Ibrahim Babura<sup>2,3\*</sup>, Muhammad Sani<sup>3</sup>, Ibrahim Abdullahi<sup>4</sup>

<sup>1</sup>Department of Mathematics and Statistics, Faculty of Science, Yobe State University, Damaturu, Nigeria.

<sup>2</sup>Department of Mathematics, Faculty of Science, Federal University Dutse, Jigawa State, Nigeria.

<sup>3</sup>Department of Mathematical Sciences, Federal university Dutsin-Ma Katsina State, Nigeria.

<sup>4</sup>Department of Mathematics and Statistics, School of Mathematics and Computing, Kampala International University, Uganda.

\*Corresponding author: bibabura@gmail.com

Received: 10 November 2022

Accepted: 10 January 2023

### ABSTRACT

*Cure fraction models are usually meant for survival data that contains a proportion of non-subject individuals for the event under study. In order to get an accurate estimate of the cure fraction model, researchers often used one of two models: the mixture model or the non-mixture model. This study presents both mixture and non-mixed cure fraction models, together with a survival data format that is based on the beta-Weibull distribution. In this body of work, an alternative extension to the Weibull distribution was devised for the purpose of analyzing lifetime data. The beta-Weibull distribution is a four-parameter distribution established in this study as an alternate extension to the Weibull distribution in lifetime data analysis. The suggested addition allows for the inclusion of covariate analysis in the model, with parameter estimation performed using a Bayesian approach and Gibbs sampling methods. In addition, a simulation study was carried out to emphasize the benefits of the new development.*

**Keywords:** Bayesian analysis, Beta-Weibull distribution, Cure fraction models, Survival analysis, MCMC algorithm.

## 1 INTRODUCTION

A suitable distribution is often of interest in the analysis of survival data proposed by [1], as it provides insight into characteristics of failure times and hazard functions such as Weibull, Beta and Gamma distributions respectively given by the probability density function of the 2-parameter Weibull distribution is:

$$f_0(t) = \gamma \lambda t^{\gamma-1} e^{-\lambda t^\gamma}, t \geq \gamma, \lambda > 0 \quad (1)$$

where  $\gamma$  is the model shape parameter and  $\lambda$  is the model scale parameter [2]. Also, the density function of the general Beta distribution is:

$$f_0(t) = \frac{(t-a)^{p-1}(b-t)^{q-1}}{B(\alpha, \beta)(b-1)^{p+q-1}}, \quad a \leq t \leq b; \alpha, \beta \geq 0 \quad (2)$$

where  $\alpha$  and  $\beta$  are the fitted shape parameters, with the lower and upper bounds, given by,  $a$  and  $b$  respectively, of the distribution. We denote  $B(\alpha, \beta)$  as the corresponding parametric beta function [1].

Similarly, the density function for the generalized Gamma distribution is given by:

$$f_0(t) = \frac{\left(\frac{t-\mu}{\beta}\right)^{\gamma-1} \exp\left(-\frac{t-\mu}{\beta}\right)^{\gamma-1}}{\beta\Gamma(\gamma)}, \quad t \geq \mu; \beta, \gamma \geq 0 \quad (3)$$

where  $\gamma$  is the shape parameter,  $\mu$  is the location parameter,  $\beta$  is the scale parameter, and  $\Gamma$  is the gamma function [3] and  $\alpha, \beta, \lambda$  are positive. Weibull distribution is regarded as a well-known distribution, named after its inventor, Waloddi Weibull [2], in 1951 a Swedish physicist. In 1939, Weibull used his proposed model to conduct an analysis on the breaking strength of various materials [4]. Since its inception, it has seen widespread application for the purpose of lifetime data analysis on account of the relative adaptability of its hazard function and the simplicity with which its parameters can be estimated. It is one of the families that is utilized the most frequently for modeling these kinds of data. However, the traditional Weibull distribution with two parameters can only be used to construct hazard functions that are either monotonically increasing or monotonically decreasing [4]. One of the disadvantages of the Beta-Weibull distribution is that both the survival function and hazard function cannot be written in a closed form, especially when more covariates are involved; hence, a numerical approach, namely integration techniques, is required to estimate the parameters in the model.

### 1.1 Cure rate model

The cure fraction model [5] is an extension to classic survival models that accounts for the number of individuals who will not witness the event of interest. In terms of the type of event specified, cure fraction models are also known as long-term survival models [6]. The mixture and non-mixing types are indeed the two cure models that are used most frequently. The standard cure rate model or simply the mixture cure rate model, assumes that the studied population is a mix of predisposed individuals who experience the event of interest "p" which is the proportion of "long-term survivors" or "cured patients" regarding the event of interest ( $0 < p < 1$ ) and non-susceptible individuals who will never exposed to it "(1 - p)".

$S(t)$  denote the survival function for the studied population and is given by

$$S(t) = p + (1-p)S_0(t), \quad t > 0 \quad (4)$$

where  $S_0(t)$  is the standard survival curve function for the vulnerable individuals. The non-mixture cure rate model establishes an asymptote for the cumulative hazard and, thus, for the cure proportion. [7]. Then the survival function is given as:

$$S(t) = p^{F_c(t)} = \exp(\ln(p) F_c(t)), \quad t > 0 \quad (5)$$

## 1.2 Related Work

Normally in any parametric distribution can be incorporated into larger families of distribution through probability integral transform procedures [8], [9]. Thus, the BW density can be re-expressed as a mixture of Weibull density as proposed by [10] who further drive an expression for their moment generating function. He further investigated the potential application of the BW distribution in censored survival data modelling for breast cancer research. A recent research identify some novel extension of the beta-Weibull distribution [11]. The beta modified Weibull distribution is another generalization of the Weibull distribution [10]. The distribution has an edge due to its flexibility upon accommodation of multiple forms of risk function while handling various problems in survival data modeling [8]. Several literature suggest Bayesian formulation of the cure fraction model [5], [12]. Numerous attempt in the literature on new techniques for estimation of cure rates consider the context for the partially observed or missing covariate [9], [13]–[17].

## 2 MATERIAL AND METHODS

### 2.1 Model and Distributional Assumptions.

We denote  $G_0(t)$  as the cumulative distribution function (cdf) of a random variable T is given by:

$$G_0(t) = I_{G_0(t)}(\alpha, \beta) = \frac{B_{G_0(t)}(\alpha, \beta)}{B(\alpha, \beta)} = \frac{\int_0^{G_0(t)} w^{\alpha-1} (1-w)^{\beta-1} dw}{B(\alpha, \beta)} \quad (6)$$

Where  $\alpha > 0, \beta > 0, B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the beta function, with associated gamma function given by  $\Gamma(\alpha) = \int_0^\infty z^{\alpha-1} e^{-z} dz$  and  $B_{G_0(t)}(\alpha, \beta)$  is the incomplete beta function. If  $G_0(t)$  in Equation (6) assumes a cdf of a gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , we then have an expression for beta-normal distribution [18]. A model based on the cdf of the Weibull distribution with shape and scale parameters of  $\gamma$  and  $\lambda$  respectively will assumes:

$$G_0(t) = 1 - \exp\left[-\left(\frac{t}{\lambda}\right)^\gamma\right], \quad t > 0 \quad (7)$$

And from Equation (6), we substitute Equation (7) to have

$$F_0(t) = \frac{1}{B(\alpha, \beta)} \int_0^{1 - \exp\left[-\left(\frac{t}{\lambda}\right)^\gamma\right]} w^{\alpha-1} (1-w)^{\beta-1} dw, \quad t > 0 \quad (8)$$

Now, in a survival analysis modelling framework, the baseline survival function for the susceptible individuals is given by

$$S_0(t) = 1 - F_0(t). \quad (9)$$

We observed that the closed form of the function is undefined with reference to the limitation of the given expression. The baseline pdf of the four parameter BW distribution is of the form

$$f_0(t) = \frac{\gamma}{\lambda^\gamma B(\alpha, \beta)} \exp\left[-\beta \left(\frac{t}{\lambda}\right)^\gamma\right] \left\{1 - \exp\left[-\beta \left(\frac{t}{\lambda}\right)^\gamma\right]\right\}^{\alpha-1}, t > 0 \quad (10)$$

where  $\alpha, \beta, \gamma$  and  $\lambda$  are positive numbers. The corresponding hazard expression is given by

$$h(t) = \frac{f_0(t)}{S_0(t)} = \frac{\gamma t^{\gamma-1} \lambda^\gamma \exp\left[-\beta \left(\frac{t}{\lambda}\right)^\gamma\right] \left\{1 - \exp\left[-\beta \left(\frac{t}{\lambda}\right)^\gamma\right]\right\}^{\alpha-1}}{B(\alpha, \beta) \int_0^{1-\exp\left[-\left(\frac{t}{\lambda}\right)^\gamma\right]} w^{\alpha-1} (1-w)^{\beta-1} dw}, t > 0 \quad (11)$$

In the context of mixture model, the likelihood expression for  $\theta = (\alpha, \beta, \gamma, \lambda, p)$  is given by

$$L_I(\theta) = \prod_{i=1}^n \left[ \frac{(1-p)\gamma}{\lambda^\gamma B(\alpha, \beta)} t_i^{\gamma-1} \exp\left(-\beta \left(\frac{t_i}{\lambda}\right)^\gamma\right) \left(1 - \exp\left[-\beta \left(\frac{t_i}{\lambda}\right)^\gamma\right]\right)^{\alpha-1} \right]^{\delta_i} \\ \times \prod_{i=1}^n [p + (1-p)S_0(t_i)]^{1-\delta_i} . \quad (12)$$

While for the non-mixture model, the likelihood function for  $\theta = (\alpha, \beta, \gamma, \lambda, p)$  is given by

$$L_{II}(\theta) = \prod_{i=1}^n \left[ -\frac{\gamma \ln(p)}{\lambda^\gamma B(\alpha, \beta)} t_i^{\gamma-1} \exp\left(-\beta \left(\frac{t_i}{\lambda}\right)^\gamma\right) \left(1 - \exp\left[-\beta \left(\frac{t_i}{\lambda}\right)^\gamma\right]\right)^{\alpha-1} \right]^{\delta_i} \\ \times \prod_{i=1}^n [p + (1-p)S_0(t_i)]^{1-\delta_i} . \quad (13)$$

## 2.2 Further Incorporation

Implementation of the conventional estimation methods especially maximization or direct methods on the likelihood functions  $L_I(\theta)$  and  $L_{II}(\theta)$  are tedious and usually computationally expensive due to complexity of some distributional expressions. Bayesian Inference based on Markov Chain Monte Carlo (MCMC) estimation methods bring down those complexities without compromise to precision and thus utilized in our implementation in this work which was appropriately justified in . The vector of covariate  $X_i$  which are closely related with proportion  $p$  of cure rate fraction models were incorporated by replacing  $p$  in the likelihood function expressions  $L_I(\theta)$  and  $L_{II}(\theta)$  with

$$p_i(t) = \frac{\exp(x_i' \eta')}{1 - \exp(x_i' \eta')} . \quad (14)$$

where  $x_i' = (1, x_{i1}, \dots, x_{in})$  is  $J$  covariates's vector of observations for the  $i$ th individual and  $\eta' = (\eta_0, \eta_1, \dots, \eta_n)$  is the unknown parameters vector. To study the effect of vector of covariates  $W_i$  on the parameter  $\lambda$ ,  $\lambda$  is replaced in both mixture and non-mixture expression of the likelihood function  $L_I(\theta)$  and  $L_{II}(\theta)$  by

$$\lambda_i(t) = \exp(w_i' \zeta') \quad (15)$$

where  $w_i' = (\mathbf{1}, w_{i1}, \dots, w_{in})$  represent the vector form of K covariates corresponding to the i-th individual and  $\zeta' = (\zeta_1, \zeta_2, \dots, \zeta_n)$  is the vector of unknown parameters.

### 2.3 Bayesian Analysis

We begin with a Bayesian analysis of the long-term survival models without considering covariates [19], we also assume the beta prior for the specified probability of proportion "p" of cure models which is denoted by  $p \sim B(a, b)$  where  $a$  and  $b$  are known hyper parameters [12] [5]. We also assume a gamma prior distribution for the parameters  $\alpha, \beta, \gamma$  and  $\lambda$ . That is  $\alpha \sim \Gamma(c_\alpha, d_\alpha), \beta \sim \Gamma(c_\beta, d_\beta), \gamma \sim \Gamma(c_\gamma, d_\gamma), \lambda \sim \Gamma(c_\lambda, d_\lambda)$  where  $c_\alpha, d_\alpha, c_\beta, d_\beta, c_\gamma, d_\gamma, c_\lambda, d_\lambda$  are known hyperparameters and  $\Gamma(c, d)$  denotes a gamma distribution with mean  $\frac{c}{d}$  and variance  $\frac{c}{d^2}$ . The joint prior distribution is then established in all circumstances by assuming prior independence between the parameters,

$$\begin{aligned} \pi(\theta) &= \pi(\alpha), \pi(\beta), \pi(\gamma), \pi(\lambda) \\ &\propto \alpha^{c_\alpha-1} \beta^{c_\beta-1} \lambda^{c_\lambda-1} \gamma^{c_\gamma-1} \exp\left(-\frac{\alpha}{d_\alpha} - \frac{\beta}{d_\beta} - \frac{\lambda}{d_\lambda} - \frac{\gamma}{d_\gamma}\right) p^{\alpha-1} (1-p)^{b-1} \end{aligned} \quad (16)$$

the following covariates for models incorporating are the assumed prior distribution for the unknown parameters:  $\alpha \sim \Gamma(c_\alpha, d_\alpha), \beta \sim \Gamma(c_\beta, d_\beta), \gamma \sim \Gamma(c_\gamma, d_\gamma), \lambda \sim \Gamma(c_\lambda, d_\lambda), \zeta_j \sim N(c_{\zeta_j}, d_{\zeta_j}^2), j = 0, 1, \dots, J,$  and  $\eta_k \sim N(c_{\zeta_k}, d_{\zeta_k}^2), k = 0, 1, \dots, K.$  where  $N(c, d^2)$  denotes a gaussian distribution with  $c$  and  $d^2$  as the mean and variance respectively, normally referred as the hyper parameters. The situation necessitate to focus solely on the independence among the prior distributions as implemented by [19].

### 2.4 Log Pseudo Maximum Likelihood

Log Psuedo Marginal Likelihood (LPML) and the Pseudo Factor is an efficient tool for comparison of mixture and non-mixture models based on varied distributional assumption. The derivation of LPML  $D, D[i]$  is done through conditional predictive ordinate (CPO) statistics [20]. That is, for the ith observation,  $CPO_i$  is given by

$$f(D_i/y_{|i}t) = \int f(D_i/\theta) f(f(\theta/D_i)) d\theta \quad (17)$$

where  $\theta$  is the incomplete parametric vector,  $D_i$  is each instance of the data  $D$  without the current observation  $i$ , and  $f(\theta/D_i)$  is the posterior density  $D[i]; i = 1, 2, \dots, n$ : An MCMC estimate of  $CPO_i$  is given by

$$\widehat{CPO}_i = \left[ \frac{1}{B} \sum_{b=1}^B \frac{1}{f(D_i/\theta_{(b)})} \right]^{-1}, i = 1, 2, \dots, n \quad (18)$$

such that,  $B$  is the iteration count for the MCMC implementation procedure after burn-in period and  $\Theta_{(b)}$  is vector of the obtained samples at 4th and 5th iterations [13]. Thus, for a given model, the LPML value is given by

$$\widehat{LPML} = \sum_{i=1}^n \log(\widehat{CPO}_i), i = 1, 2, \dots, n \quad (19)$$

For an increase value of LPML, we have a better fit for the model [20]. Alternatively, the derived Pseudo Bayes factor (PMF) for comparing multiple models  $m$  and  $m'$  is

$$PBF_{mm'} = \exp(\widehat{LPML}_m - \widehat{LPML}_{m'}) \quad (20)$$

So also, each parameter of interest require an estimate for the highest probability density (HPD) intervals [20]. Assuming  $100(1 - \omega)\%$  HPD interval associated with a generic parameter  $\theta$  an arbitrary subset of the parametric space  $C\theta$  given by  $C = \{\theta; \pi(\theta/D) \geq k\}$  where  $\pi(\theta/D)$ . Then the posterior distribution for  $\theta$  given the data  $D$  and  $k$  as the largest number such that the expression

$$\int_{\pi(\theta/D) \geq k} \pi(\theta/D) = 1 - \omega \quad (21)$$

### 3 RESULTS AND DISCUSSION

#### 3.1 Simulation Study for the Analysis of Result

A sample of 30,000 was generated for each parameter of interest based on each cases under consideration. Assuming a burn-in sample of 10,000 data size which can minimize the initialization effect on the simulation process. However, a 15,000 sample size, with each of the 200th sample having approximately uncorrelated values was utilized to achieve a posterior summaries of interest.

In the estimation procedure we use the MCMC estimation with Bayesian Approach where the LMPL considers the highest value to be the best models parameter selection in described in Table 1 and Table 2. The simulation process utilizes some R-code implementation using the baseline model that is Beta Weibull distribution and an optimization package 'optim' in R. As starting values we used the estimator obtained from a Weibull model based on the censoring indicator (as a surrogate for the unobserved cure indicator). However, due to the non-concavity of our likelihood function and due to the inconsistency of this vector of starting values, the procedure optim often ends up in a local maximum instead of the global maximum. To avoid this problem, we added the some intermediate step to the estimation procedure on the initial starting values, we then estimate all the parameters from a BW model based on the parametric estimate that maximize the log-likelihood globally. Since this log-likelihood as starting value for our likelihood function of BW model having a close form property.

Table 1 demonstrate results according to the described simulation procedure. For the generated samples, the result show the bias and mean squared error (MSE) of model for several values of  $c$ ,

namely  $c = 2, 3$  and  $4$ . However, the cross-validation(CV) procedure proposed by [19] for MCMC bootstrap estimators indicate that the convergence of the MCMC algorithm was not obtainable for values less than 1 for the selected hyper parameters, even when using very large burn-in -period for the algorithm. Considering the beta-Weibull (BW), exponentiated Weibull (EW), beta-Exponentiated (BE) and Weibull distributions respectively in Table1 and Table2. Estimated parameters were obtained as median estimate of Gibbs samples drawn as a join posterior distribution. Median is preferred here over mean due to the skewed nature of the distribution in the simulation process. The p values from Heidelberger and Welch (HW) convergence diagnostic criteria do not reject the null hypothesis of stationary of the chains, for being larger or equal than 0.10. In the case of Geweke's p value which also suggests convergence. The result further suggests that, among the models in consideration, Weibull distribution has the lowest Log pseudo marginal likelihood (LPML) value unlike BW, BE and EW distributions all having similar LPML value.

Table 1: The posterior summaries of the model parameters excluding a cure fraction while considering GBC study dataset.

Model	Parameter	Posterior median	95% $HDP^a$	$LPML^b$	$HW^c$ p value	Geweke's p value
BW	$\alpha$	3.8116	(1.5334,6.7999)	-864.148	0.314	0.112
	$\beta$	0.0806	(0.0166,0.2681)		0.392	0.251
	$\gamma$	0.9084	(0.5611,1.2922)		0.618	0.307
	$\lambda$	2.3126	(1.1376,4.3364)		0.099	0.223
EW	$\alpha$	7.4618	(4.5975,11.2487)	-835.774	0.232	0.076
	$\gamma$	0.4409	(0.3456,0.5518)		0.324	0.234
	$\lambda$	4.0185	(1.4213,7.6297)		0.122	0.166
BE	$\alpha$	3.3499	(1,9664,5.2017)	-844.909	0.463	0.334
	$\beta$	0.0590	(0.0250,0.1167)		0.818	0.772
	$\lambda$	2.2978	(1.2235,4.1378)		0.699	0.394
Weibull	$\gamma$	1.654	(1.4712,1.8571)	-825.444	0.736	0.757
	$\lambda$	3.850	(2.788,3.293)		0.710	0.842

However, an additional evidence of a better fit is the non-convergence of the MCMC estimation on fitting BW distribution in presence of cure fraction as against standard Weibull distribution [21].

The inferences for the non-mixture and mixture model which are based on the Beta-Weibull distribution with its special cases are clearly demonstrated in Table 2. Based on highest LPML of the models, mixture models get a better fit [22]. Furthermore, the 95% credible interval for  $\eta_2$  based on its non zero value inclusion suggest that the subjects in the AML high risk and low risk groups have

contrasting cure fractions. The Bayesian estimates for the cure fractions for every risk group can be obtain by considering the simulated samples for  $\eta_0, \eta_1$  and  $\eta_2$  and the relation

$$P(AML\ lowrisk) = exp(\eta_0), P(AML\ highrisk) = exp(\eta_0 + \eta_1) \text{ and}$$

$$P(all) = exp(\eta_0 + \eta_2).$$

Therefore, the estimated results obtained for the cure fractions of the patients classified as AML low risk, ALL and AML high risk are shown above respectively. The graphs in Figure 2, show that Kaplan - Meier survival curves for bone marrow transplant patients based on the BW distribution fit the mixture model at all levels of risks. Note that the curves obtained from the model are close to those estimated by Kaplan-Meier method, a great indication of good fit based on the models for the [22].

Table 2: The posterior summaries assuming the mixture model with covariate and considering the data set of 137 bone marrow transplant patients

Model	Parameter	Posterior median	95% HDP <sup>a</sup>	LPML <sup>b</sup>	HW <sup>c</sup> p value	Geweke's p value
BW	$\alpha$	1.0129	(0.3684,2.1799)	-67.403	0.453	0.729
	$\beta$	1.2932	(0.1107,3.2681)		0.224	0.251
	$\gamma$	1.0381	(0.5337,1.6854)		0.338	0.607
	$\zeta_0$	-0.3526	(-1.1376,0.3364)		0.069	0.533
	$\zeta_1$	-0.6706	(-0.0166,0.2681)		0.152	0.651
	$\zeta_2$	-0.9889	(-0.5611,-.2922)		0.148	0.307
	$\eta_0$	-0.1337	(-1.1376,0.3364)		0.799	0.343
	$\eta_1$	-0.4843	(-0.2611,0.2922)		0.618	0.707
	$\eta_2$	-0.9124	(-1.1376, 0.3364)		0.779	0.243
EW	$\alpha$	0.9930	(4.5975,11.2487)	-67.374	0.232	0.076
	$\gamma$	1.0456	(0.5611,1.2922)		0.618	0.307
	$\zeta_0$	-0.5640	(1.1376,0.3364)		0.099	0.223
	$\zeta_1$	-0.6876	(0.0166,0.2681)		0.372	0.351
	$\zeta_2$	-1.0025	(0.5611,1.2922)		0.618	0.307
	$\eta_0$	-0.1466	(-1.1376,1.3364)		0.099	0.256
	$\eta_1$	-0.4688	(-0.0350,1.1427)		0.834	0.435
	$\eta_2$	-0.9530	(0.0150,0.2147)		0.568	0.872
BE	$\alpha$	1.0649	(1,9664,5.2017)	-66.561	0.483	0.007



	$\beta$	1.2006	(0.0166,0.2681)		0.392	0.251
	$\zeta_0$	-0.4126	(-1.1376,0.3364)		0.099	0.223
	$\zeta_1$	-0.6561	(-.0166,0.2681)		0.692	0.281
	$\zeta_2$	-0.9876	(-1.5611,1.2922)		0.618	0.307
	$\eta_0$	-0.1365	(-0.1376,1.9864)		0.099	0.223
	$\eta_1$	-0.4753	(-.5611,1.2922)		0.618	0.337
	$\eta_2$	-0.9146	(-1.1376,0.2324)		0.099	0.223
Weibull	$\gamma$	1.654	(1.4712,1.8571)	-65.557	0.10	0.265
	$\zeta_0$	-0.3126	(-1.1376,4.3364)		0.023	0.123
	$\zeta_1$	-0.0806	(-1.0166,-0.2881)		0.362	0.281
	$\zeta_2$	-1.9084	(-1.5611,-0.4922)		0.618	0.237
	$\eta_0$	-0.3126	(-0.1376,0.3364)		0.099	0.243
	$\eta_1$	-0.9084	(-1.4621,0.2322)		0.228	0.417
	$\eta_2$	-0.3126	(-1.2346,-0.5364)		0.565	0.287

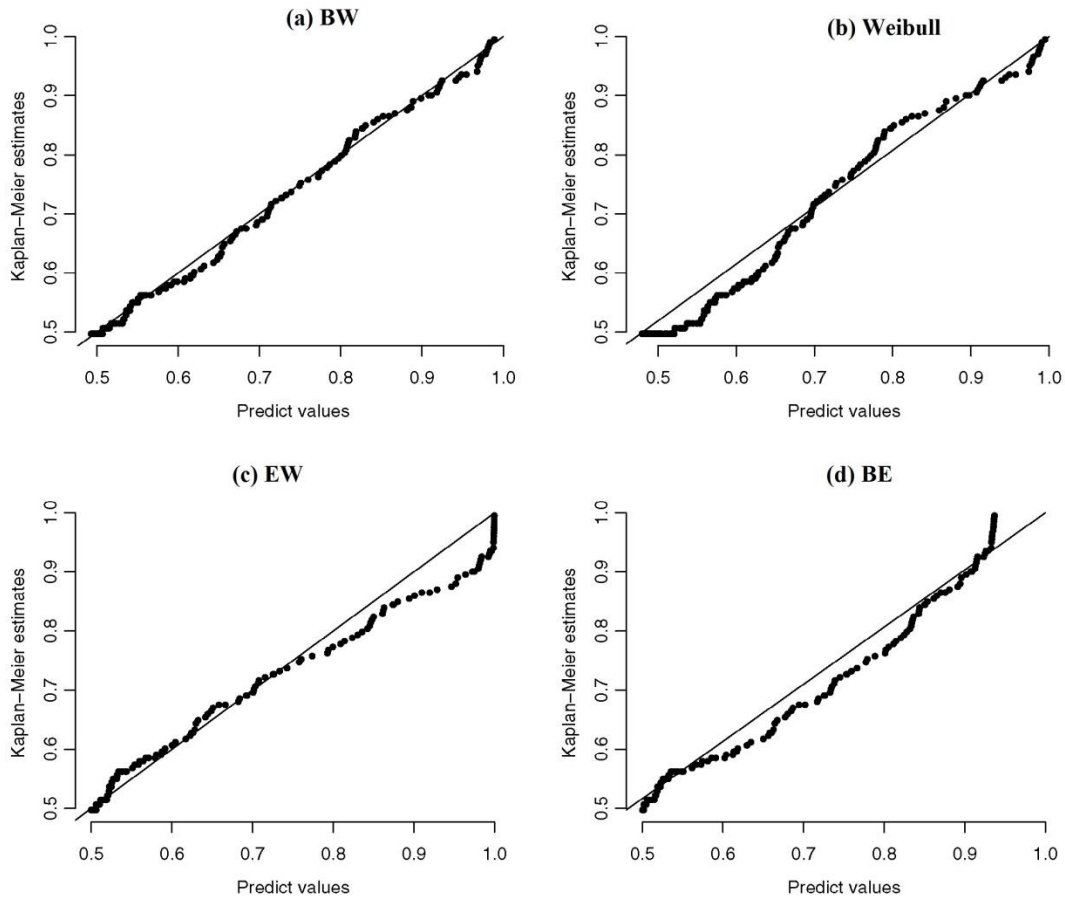


Figure 1: A Kaplan-Meier estimates for survival function plotted against respective values generated from the parametric mixture models for each of the distribution of interest: (a) BW (b) Weibull, (c) EW, (d) BE.

### 3.2 Application to German Breast Cancer and Bone Marrow Transplant Data

We first consider the case where the cure fraction parameter and covariates are not included in the model which can favorably be applied to the popular data set from the German Breast Cancer (GBC) study dataset [23]. The data set comprises 686 patients under 65 years of age where 299 had an event recurrence-free survival and 171 died. We used the time to death as the event of interest. It is estimated that the maximum follow-up time available was 7 years. Figure 2(a) is the plots of the survival functions estimated by Kaplan - Meier method and from the models based on the BW and Weibull distributions. While in Figure 2(b) shows the hazard functions based on the Bone-Marrow Transplant data [23], where (AML) indicate a low risk scenario.

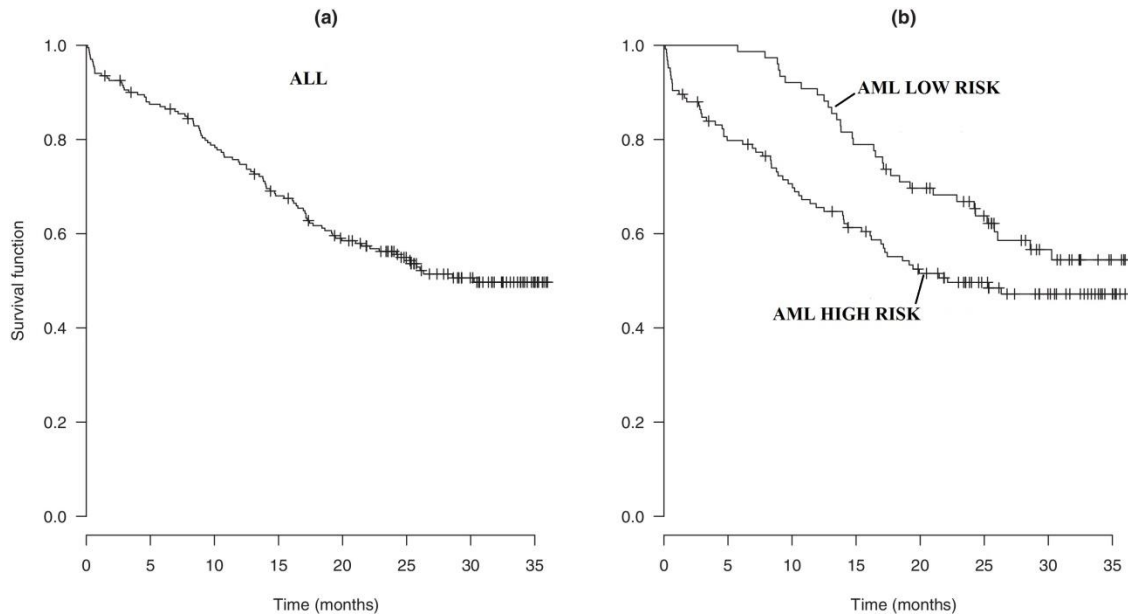


Figure 2: The panel (a), plots of the survival functions estimated by Kaplan - Meier method and from the models based on the BW and Weibull distributions Panel (b), shows the hazard functions based on the Bone-Marrow Transplant data.

#### 4 CONCLUSION

The cure fraction model and covariates analysis are strong features of a lifetime data analysis. Deployment of different parametric formulations for the analysis of such data can be done as mixture or non-mixture models. This paper establish a parametric models approach based on BW distribution with special cases useful in analyzing medical data set [8]. Also, the Bayesian methodology using MCMC methods was demonstrated in this work as a suitable tool to establish certain inferences about parameters of the model. As highlighted by [24], the limitation of the BW distribution is that the survival function has no closed form of expression and thus numerical integration techniques were utilize for parameter estimate of the model. Same limitations were more critical in terms of covariates because the likelihood function became more complex. An advantage of Bayesian approach over other conventional methods is it explicit incorporation of an expert prior opinion for the parameters. In clinical application, the knowledge of a specialist of the expected proportion of patience who are immune to the event of interest can be added into a prior distribution for the cure fraction  $p$  to have a more precise inference.

#### REFERENCES

- [1] M. Pal and M. Tiensuwan, "The Beta Transmuted Weibull Distribution," *Austrian J. Stat.*, vol. 43, no. 2, pp. 133–149, Jun. 2014, doi: 10.17713/AJS.V43I2.37.
- [2] W. Weibull, "A Statistical Distribution Function of Wide Applicability," *J. Appl. Mech.*, vol. 18, no. 3, pp. 293–297, Sep. 1951, doi: 10.1115/1.4010337.

- [3] E. W. Stacy, "A generalization of the gamma distribution," *The Annals of Mathematical Statistics*, vol. 33, no. 3, pp. 1187–1192, 1962. Accessed: Feb. 01, 2022. [Online]. Available: <https://www.jstor.org/stable/2237889>
- [4] G. Biau and E. Scornet, "A random forest guided tour," *TEST*, vol. 25, pp. 197–227, 2016, doi: 10.1007/s11749-016-0487-1.
- [5] J. A. Achcar, E. A. Coelho-Barros, and J. Mazucheli, "Cure fraction models using mixture and non-mixture models," *Tatra Mt. Math. Publ.*, vol. 51, no. 1, pp. 1–9, Apr. 2012, doi: 10.2478/v10127-012-0001-4.
- [6] F. Felizzi, N. Paracha, J. Pöhlmann, and J. Ray, "Mixture Cure Models in Oncology: A Tutorial and Practical Guidance," *PharmacoEconomics - Open*, vol. 5, no. 2, pp. 143–155, Jun. 2021, doi: 10.1007/S41669-021-00260-Z/TABLES/4.
- [7] M. Amico and I. Van Keilegom, "Cure Models in Survival Analysis," *Annu. Rev. Stat. Its Appl.*, vol. 5, pp. 311–342, Mar. 2018, doi: 10.1146/ANNUREV-STATISTICS-031017-100101.
- [8] A. S. Wahed, T. M. Luong, and J. H. Jeong, "A new generalization of Weibull distribution with application to a breast cancer data set," *Stat. Med.*, vol. 28, no. 16, pp. 2077–2094, Jul. 2009, doi: 10.1002/SIM.3598.
- [9] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin, "Bayesian Data Analysis," *Bayesian Data Anal.*, Nov. 2013, doi: 10.1201/B16018.
- [10] G. M. Cordeiro, S. Nadarajah, and E. M. M. Ortega, "General results for the beta Weibull distribution," *J. Stat. Comput. Simul.*, vol. 83, no. 6, pp. 1082–1114, Jun. 2013, doi: 10.1080/00949655.2011.649756.
- [11] G. M. Cordeiro, A. E. Gomes, C. Q. Da-Silva, and E. M. M. Ortega, "The beta exponentiated weibull distribution," *J. Stat. Comput. Simul.*, vol. 83, no. 1, pp. 114–138, 2013, doi: 10.1080/00949655.2011.615838.
- [12] J. A. Achcar, E. A. Coelho-Barros, and J. Mazucheli, "Cure fraction models using mixture and non-mixture models," *Tatra Mt. Math. Publ.*, vol. 51, no. 1, Mar. 2012, doi: 10.2478/TATRA.V51I1.140.
- [13] A. D. Tsodikov, J. G. Ibrahim, and A. Y. Yakovlev, "Estimating Cure Rates from Survival Data: An Alternative to Two-Component Mixture Models," *J. Am. Stat. Assoc.*, vol. 98, no. 464, pp. 1063–1078, Dec. 2003, doi: 10.1198/01622145030000001007.
- [14] M. R. A. Bakar, K. A. Salah, N. A. Ibrahim, and K. Haron, "Bayesian approach for joint longitudinal and time-to-event data with survival fraction," *Bull. Malaysian Math. Sci. Soc.*, vol. 32, no. 1, pp. 75–100, 2009. Accessed: Feb. 01, 2022. [Online]. Available: <http://emis.math.tifr.res.in/journals/BMMSS/pdf/v32n1/v32n1p8.pdf>
- [15] M. Bebbington, C. D. Lai, and R. Zitikis, "A flexible Weibull extension," *Reliab. Eng. Syst. Saf.*, vol. 92, no. 6, pp. 719–726, 2007, doi: 10.1016/j.ress.2006.03.004.

- [16] G. S. Mudholkar, D. K. Srivastava, and G. D. Kollia, "A Generalization of the Weibull Distribution with Application to the Analysis of Survival Data," *J. Am. Stat. Assoc.*, vol. 91, no. 436, pp. 1575–1583, Dec. 1996, doi: 10.1080/01621459.1996.10476725.
- [17] C. Lee, F. Famoye, and O. Olumolade, "Beta-Weibull distribution: Some properties and applications to censored data," *J. Mod. Appl. Stat. Methods*, vol. 6, no. 1, pp. 173–186, 2007, doi: 10.22237/jmasm/1177992960.
- [18] N. Eugene, C. Lee, and F. Famoye, "Beta-normal distribution and its applications," *Commun. Stat. - Theory Methods*, vol. 31, no. 4, pp. 497–512, 2002, doi: 10.1081/STA-120003130.
- [19] E. Z. Martinez, J. A. Achcar, A. A. A. Jácome, and J. S. Santos, "Mixture and non-mixture cure fraction models based on the generalized modified Weibull distribution with an application to gastric cancer data," *Comput. Methods Programs Biomed.*, vol. 112, no. 3, pp. 343–355, Dec. 2013, doi: 10.1016/j.cmpb.2013.07.021.
- [20] A. E. Gelfand, D. K. Dey, and H. Chang, "Model determination using predictive distributions, with implementation via sampling-based methods (with discussion)," *Bayesian Stat. 4*, 1992, Accessed: Feb. 01, 2022. [Online]. Available: <https://apps.dtic.mil/sti/citations/ADA258777>
- [21] W. Sauerbrei, P. Royston, H. Bojar, C. Schmoor, and M. Schumacher, "Modelling the effects of standard prognostic factors in node-positive breast cancer," *Br. J. Cancer*, vol. 79, no. 11–12, pp. 1752–1760, 1999, doi: 10.1038/sj.bjc.6690279.
- [22] M. B. Rao, J. P. Klein, and M. L. Moeschberger, "Survival Analysis Techniques for Censored and Truncated Data," *Technometrics*, vol. 40, no. 2, p. 159, 1998, doi: 10.2307/1270658.
- [23] D. W. Hosmer, S. Lemeshow, and S. May, *Applied Survival Analysis: Regression Modeling of Time to Event Data*, 2nd ed. Wiley Blackwell, 2011. doi: 10.1002/9780470258019.
- [24] G. M. Cordeiro, A. B. Simas, and B. D. Stošić, "Closed form expressions for moments of the beta Weibull distribution," *An. Acad. Bras. Cienc.*, vol. 83, no. 2, pp. 357–373, 2011, doi: 10.1590/S0001-37652011000200002.