# On Piece-Wise Modelling of Survival Data with Time Changing Covariate Function

[1]**P. E. Omaku** & [2]**J. S Ibinayin**
[1&2]*Department of Mathematics & Statistics,*
*Federal Polytechnic Nasarawa – Nigeria*

## Abstract

Survival analysis involve the set of statistical techniques or procedures used to study time until an event occurs, these techniques are not without some conditions. One of the basic assumptions is that, to enable a straight forward interpretation of hazard rates of subject's covariate(s) on some reference categories or in situations where variables are continuous in nature, the hazard rates must be constant through time "also known as the proportional hazard assumption" for cox regression. This assumption is often violated in medical practice where subject's vital statistics or measures are often time varying, as their medical situations changes with time. This paper under study a modification of Piece wise survival model, where three levels of Weibull distribution were assumed for baseline hazards, the sensitivity of the baselines were assessed under four (4) censoring percentages (0%, 25%, 50%, & 75%) and sample sizes (n=100, n=500 & n=1000), for when models were Single parametric (SPM) and when partitioned – Piece wise Parametric Model (PPM). A Piece-wise Bayesian hazard model with structured additive predictors in which the functional form of time varying covariate was incorporated in a non-proportional hazards framework was developed, capable of incorporating complex situations in a more flexible framework. Analysis was done utilizing MCMC simulation technique. Results revealed on comparison that the PPM outperformed the SPM with smaller DIC values and larger predictive powers with the LPML criterion and consistently so throughout all simulations.

**Keywords:** *Time varying covariates, Proportional hazard, Violation, Piece-wise survival model, Piece-wise parametric model, Single parametric model*

*Corresponding Author:* P. E. Omaku

**Background to the Study**

Survival analysis is a statistical procedure for data analysis for which the outcome variable of interest is time until an event occurs. By time, we mean years, months, weeks, or days from the beginning of follow-up of an individual until an event occurs and by event, we mean death, disease incidence, re-lapse from remission, recovery (e.g., return to work) or any designated experience of interest that may happen to an individual (David and Kleinbaum 2005).

Analysis of survival times data has gained a considerable attention, particularly in the field of medicine, where the conventional denotation 'Survival analysis' arises from (Hennerfeind, 2006). In several other bio-statistical applications on censored follow-up time data, the interest lies mainly on the prognostic role of clinical/biological covariates. To such end, non-parametric and semi-parametric methods have been preferred over parametric ones. The most widely adopted tool is the Cox model, which avoids any assumption of the functional form of the hazard function on time. However, such feature is not useful if the interest lies on investigating the shape of the hazard or in predictive modeling (Kooperberg et al. 1995) when the cox-model is extended to time-varying covariates and time-dependent effects, which combine to give the most general version of the hazard. Again, further progress would require specifying the form of this function of time. In such situation where time is observed to be truly continuous a flexible or semi-parametric strategy is required, where mild assumptions are made about the baseline hazard $\lambda_0(t)$. Specifically, we may subdivide time into reasonably small intervals and assume that the baseline hazard is constant in each interval, leading to a piece-wise survival model.

According to Fabio et al. (2010) the Piecewise Model (PM) arises as a quite attractive alternative to parametric models for the analysis of time to event data. Although parametric in a strict sense, the PM can be thought of as a nonparametric model as far as it does not have a closed form for the hazard function. This nice characteristic of the PM allows us to use this model to approximate satisfactorily hazard functions of several shapes. For this reason, the PM has been widely used to model time to event data in different contexts, such as;
1. Clinical situations including kidney infection (Sahu et al, 1997), heart transplant data (Aitkin et al, 1983).
2. Hospital mortality data Clark and Ryan (2002), and cancer studies including leukemia (Breslow, 1974), gastric cancer (Gamerman, 1991), breast cancer (Ibrahim et al. 2001b) (see also Sinha et al., 1999) for an application to interval-censored data), Melanoma (Kim et al. 2006) and nasopharynx cancer (McKeague and Tighiouart, 2000), among others.
3. The PM has also been used in reliability engineering (Kim and Proschan,1991), (Gamerman, 1994), and economics problems (Gamerman, 1991) and (Bastos et al, 2006).
4. Time-Varying Effect of Tumor Size and Soft Tissue Sarcoma Data by (Marano, et al 2016)

In this paper we shall modify a Piecewise Weibull hazard baseline function of survival model which can cope better with changes in baseline rate over time, leading to a better fit. This paper

investigate, employing three levels of Weibull distributions as baseline; the effects of ignoring time varying effects and regularized estimation of non-linear functions applied often in prognostic factors.

## Materials and Method
The risk data used for this paper was simulated from a Weibull baseline hazard distribution which was used to generate survival times for sample sizes of 100, 500 & 1000 respectively. Various censoring levels or percentages of: no censoring "0%", low "about 25%", moderate "about 50%" and high "about 75%" were used.

## Model Specification
The cox hazard model

$$\lambda_i(t, X) = \lambda_0(t) \exp\left(\sum_{j=1}^{p} \beta_j X_j\right). \hspace{2cm} 1$$

The baseline hazard rate is unspecified, and assumes that covariates $x = (x_1, \ldots, x_p)$ act multiplicatively on the hazard rate through the exponential link function (Abiodun, 2007). An additive representation of model 1

$$\eta(t; w, z, x, s) = f_0(t) + \sum_{j=1}^{p} f_j(t) z_j + x' \gamma \hspace{2cm} 2$$

This is a re-parameterization of the cox model

Where $f_0(t) = log\lambda_0(t)$ which implies, $\exp(f_0(t))$, is the baseline function, other aspects of the models include the functions $f_1(t)z_1 \ldots f_p(t)z_p$ are possibly functional form of time varying covariate $z_1, \ldots, z_p$ and $\gamma$ is the usual linear part of the predictor for some categorical covariates (Abiodun, 2009) and (Hennerfeind *et al.*, 2006)

$$\lambda_{PE}(t; v, x, s) = \{I(t \epsilon T_h(f_h(t)) + \sum_{j=1}^{p} f_j(t) z_j + f_{spat}(s_{ih}). \hspace{2cm} 3$$

With its various terms defined as
The function $f_h = log\lambda_h$ is the baseline effect for the kth interval of PEM
The functions $f_1(z_{1h}), \ldots, f_p(z_{ph})$ are functional forms of time varying covariates $z_{1h}, \ldots, z_{ph}$ in the $h^{th}$ interval and $f_{spat}(s_{ih})$ is a structured spatial effect, where s, s =1, . . . ,S is either a spatial index, with s = $s_i$ if subject i in the $h^{th}$ bit (interval) is from area s or it is an exact spatial coordinate s = $(x_i, y_s)$, e.g. for centriods of regions or if exact locations of individuals are known.

## Model Likelihood Function

$$L_{PE}\left(\underline{\lambda}, \underline{\beta}; D, \Delta, X, s\right) = \prod_{i=1}^{n} \prod_{h=1}^{H_i} (\lambda_h \exp(X_i^T \underline{\beta} + s_{ih})^{d_{ih}}. \exp(\lambda_h \exp\left(X_i^T \underline{\beta} + s_{ih}\right) \Delta_{ih}). \hspace{1cm} 4$$

where for each subject i there is a product of $h_i$ terms, $H_i$ being the number of intervals in which the subject is followed. In the expression above, $d_{ih}$ is the status of the $i^{th}$ subject within the interval $T_h$ (0 = alive or censored, 1 = failed); $\Delta_{ih}$ is the time spent in $T_h$ by the subject. From

expression (3) it may be seen that L... is proportional to the product of Poisson likelihoods for $D_{ih}$ with mean parameters: $\mu_{ih} = \lambda_h \exp\left(X_i^T \underline{\beta} + s_{ih}\right)\Delta_{ih}$. As a consequence, the expression of the Poisson regression model is:

$$D_{ij} \sim POISSON(\mu_{ih}); \log(\mu_{ih}) = \underline{\alpha_h} + X_0^T \underline{\beta} + s_{ih} + \log(\Delta_{ih}) .$$  5

Where $h(i)$ indicate the interval where $t_i$ falls, i.e. the interval where individual $i$ died or was censored.

where $\underline{\alpha_h} = \log(\lambda_h)$ are log-hazard parameters, and the term $\log(\Delta_{ih})$ is an offset.

The expression of the Piecewise model with regularized effects is the following:

$$\begin{cases} V_{ij} \sim POISSON(\mu_{ih}) \\ \log(\mu_{ih}) = \underline{B_0^T \alpha} + \sum_{j=1}^p Z_{1j,i} \underline{B_0^T} \underline{\gamma_{1j}} + v_{ij} + \log(\Delta_{ih}) \\ \qquad \qquad \square \\ (\underline{\alpha}|\tau^2) \sim RW(\tau^2, P_d); \quad \tau^2 \sim \pi_{\tau^2} \\ \left(\underline{\gamma_{ij}}\Big|\tau_j^2\right) \sim RW\left(\tau_j^2, P_d^{(j)}\right); \quad \tau_j^2 \sim \pi_{\tau^2 j}; j = 1, \dots, p \\ \qquad \qquad \square \\ V_i/\{v_j\}_{j\neq i} \sim N\left(-\Sigma_{\{j;j\neq i\}} P_{ij}v_{ij}/P_{ii}, \tau^2/P_{ii}\right) \end{cases}$$  6

The time-dependent effects for each covariate are: $Z_{1j,i}\underline{B_0^T}\gamma_{1j}$; $j = 1, \dots, p$. Thus, for each $Z_{1j}$, its values multiplied for a piecewise constant function: $B_0^T\gamma_{1j}$; in the parameters.
$\gamma_{1j} = (\gamma_{1,j,1}, \dots, \gamma_{1,j,H})^T$. This enables the effect of each $Z_{1j}$ to vary in each interval $T_h$ of the original partition of the follow-up:

$$\underline{B_0^T\alpha} + Z_{1j,i}\underline{B_0^T} = \alpha_h + z_{1j,i}\gamma_{1j,h} \text{ for } t\in T_{h|}.$$

### Gaussian Random Field (GRF) priors

For georeferenced data, it is commonly assumed that $v_i = v(s_i)$ arises from a Gaussian random field (GRF) $\{v(s), s\epsilon S\}$ such that $v = (v_1, \dots, v_m)$ follows a multivariate Gaussian distribution as $v\sim N_m(0, \tau^2 R)$, where $\tau^2$ measures the amount of spatial variation across locations and the (i,j) element of R is modeled as $R[i,j] = \rho(s_i, s_j)$. Here $\rho(.,.)$ is a correlation function controlling the spatial dependence of v(s). In "survregbayes" package in R, the powered exponential correlation function $\rho(s_i, s_j) = \rho(s_i, s_j, \varphi) = \exp\{-(\varphi\|s - s'\|)^v\}$ is used, where $\varphi > 0$ is a range parameter controlling the spatial decay over distance, $v\epsilon (0,2]$ is a prespecied shape parameter, and $\|s - s'\|$ refers to the distance (e.g., Euclidean, great-circle) between s and s' Therefore, the prior $GRF(\tau^2, \emptyset)$ is defined as

$$V_i/\{v_j\}_{j\neq i} \sim N\left(-\Sigma_{\{j;j\neq i\}} P_{ij}v_{ij}/P_{ii}, \tau^2/P_{ii}\right) \qquad (7)$$

$i = 1, \dots, m$ where $P_{ij}$ is the (I,j) element of $R^{-1}$ (Zhou et al, 2017).

**Test for Non-Proportionality**

To test the hypothesis that the proportional hazard assumption is valid, the following statement of hypothesis is made.

$H_0: \delta_1 = \delta_2 = \cdots = \delta_p$ (Assumption is valid)

$H_1: at\ least\ one\ of\ the\ \delta_i's\ is\ not\ equal\ to\ zero$ (Assumption violated)

Decision rule: Reject $H_0$ if $p - value \leq \alpha$ (level of significance)

Residual measures are used to investigate the departure from the proportional hazard assumption. Schoenfeld residuals are used to test the assumption of proportionality. Schoenfeld residuals are usually calculated at every failure of time under the proportional hazard assumption, and usually not defined for censored observations. The overall significance test is called the global test (sighted in Adeniyi and Akinrefon, 2018)

**Data Analysis**

The simulations apply the functional form of time varying covariate by Bender, Augustin and Blettner (2005) given as

$$f(t) = 0.5\sqrt{t} * y. \qquad y \sim binom(N, 1, 0.5) \qquad (8)$$

For spatial frailty we propose, S= pnorm(v) and $v \sim mvrnorm(1, \Sigma)$; if S=pnorm(v) then $S \sim mvrnorm$, enhanced in simulations via the Mass package in R. Where $\Sigma$ is the covariance matrix for spatial correlation in form frailty model

Co-ordinates for spatial correlations follow the uniform distribution. $s_1 = runif(N, 0, 40)$ and $s_2 = runif(N, 0, 100)$.

(Ulviya, 2011), obtained the shape and scale parameters of the Weibull distribution from the formulas below

$$\eta = \frac{1}{\Gamma(1+\frac{1}{\alpha})} \qquad (9)$$

And

$$\left( \frac{\Gamma\left(1+\frac{2}{\alpha}\right)}{(\Gamma(1+\frac{1}{\alpha}))^2} - 1 \right) = 0.5 \qquad (10)$$

for a convenience choice of mean 1 and variance 0.5. Using the uniroot function in R. parameters were given to be approximately = 1.435523 and = 1.101321. We considered studying the impact of increasing and decreasing the variance of the Weibull distribution while keeping the mean at 1. The result is displayed in table 1 below

**Table 1:** Shape and scale parameters of the Weibull distributions

| E(T) | Var(T) | α | η |
|---|---|---|---|
| 1 | 0.25 | 2.101377 | 1.129063 |
| 1 | 0.5 | 1.435523 | 1.101321 |
| 1 | 0.75 | 1.157975 | 1.052847 |

The simulation study is to investigate:
1. How the baseline hazards behave under functional forms of time varying effect and continuous covariates in the presence of spatial correlations and
2. Investigate the performance of Single hazard models or Single Parametric models (SPM) and the modified Piece-wise model extension or Piece-wise Parametric models (PPM) under various censoring percentages and sample sizes employing their levels of Weibull distributions as baseline.

**Model Specification to advance Simulation**

Model1: $\lambda_{PI}(t; z) = f_0(t) + f(t)z_j + f_{spat}(s_{ih})$.

Model2: $\lambda_{PD}(t; z) = \{I(t \epsilon T_h(f_h(t))\} + f_h(t)z_j + f_{spat}(s_{ih})$.

Where $\lambda_{PI}$ is the hazard function when Partitioning is Ignored (PI) or Single Parametric model (SPM)

Where $\lambda_{PD}$ is the hazard function when Partitioning is done (PD) or Piece wise Parametric model (PPM)

Simulations and analysis were carried out in R using the coda package for spBayesSurv, version 3.6.2. Comparisons were done using Deviance Information Criterion (DIC) (smaller is better) which places emphasis on the relative quality of model fitting and log pseudo marginal likelihood (LPML) (larger is better) focuses on the predictive performance. Both criteria are readily computed from the MCMC output.

**Results and Interpretation of Simulation Study**

## Results and Interpretation of Simulation Study

**Table 1** : DIC and LPML of $\beta(t)$ by three (3) levels of Weibull baseline hazard and level of censoring for all sample sizes and $\beta=0.5\sqrt{t}$ executed for models I & II

**n=100**

| Weibull baseline with low variance of 0.25 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Partitioning is ignored (PI) | | | | PPM (PD) | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_3$ |
|  | 1.083 | 900.3891 | -450.8276 | 878.2437 | -444.161 | 2.1120 | -1.6191 | -2.8592 | -29.978 |
| 25% | 0.9342 | 671.157 | -336.9246 | 664.8024 | -334.399 | 0.8049 | -0.0379 | 0.8172 | -9.0184 |
| 50% | 1.062 | 503.4064 | -252.7583 | 500.7776 | -251.824 | 0.7330 | 0.6288 | 1.1786 | -88.591 |
| 75% | 1.147 | 284.3667 | -143.1534 | 287.7841 | -148.394 | 0.4906 | 1.0942 | 2.4062 | -14.162 |
| Weibull baseline with intermediate variance of 0.5 | | | | | | | | | |
| PI | | | | PD | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_3$ |
|  | 1.388 | 1008.247 | -504.8157 | 986.830 | -494.999 | 2.4477 | -1.3555 | -2.7600 | -3.4787 |
| 25% | 1.305 | 732.8461 | -367.927 | 725.413 | -364.1833 | 1.2463 | -0.2178 | 0.5533 | -91.743 |
| 50% | 1.296 | 540.0572 | -271.5071 | 538.297 | -270.9676 | 0.9155 | 0.5045 | 1.6572 | -107.66 |
| 75% | 1.665 | 296.9573 | -149.734 | 287.729 | -145.1475 | 0.4918 | 0.4660 | 10.515 | 5.8117 |
| Weibull baseline with high variance of 0.75 | | | | | | | | | |
| PI | | | | PD | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
|  | 1.598 | 1078.553 | -540.2496 | 1065.927 | -536.175 | 2.3533 | -1.1581 | -2.5135 | -8.2776 |
| 25% | 1.542 | 772.0216 | -387.4007 | 755.6687 | -381.620 | 1.6048 | -0.4414 | 1.4104 | 1.2518 |
| 50% | 1.432 | 570.3035 | -286.657 | 570.6291 | -289.261 | 1.0935 | -0.1834 | 2.9965 | 2.1824 |
| 75% | 1.586 | 311.8024 | -157.3363 | 210.647 | -105.213 | 1.1233 | 1.13243 | 2.8323 | 2.3014 |

**n=500**

| Weibull baseline with low variance of 0.25 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Partitioning is ignored | | | | Partitioning is done | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
|  | 0.7422 | 4033.102 | -2017.044 | 4020.982 | -2011.72 | 0.8626 | -0.5374 | -0.1370 | -2.0944 |
| 25% | 0.99 | 3278.719 | -1639.99 | 3269.356 | -1636.23 | 1.0717 | -0.0804 | -0.6797 | 0.9899 |
| 50% | 1.113 | 2364.521 | -1182.72 | 2369.647 | -1186.77 | 1.2679 | -0.1969 | 0.30105 | -0.6474 |
| 75% | 1.201 | 1243.873 | -622.8476 | 1248.3 | -625.933 | 0.9743 | 0.2131 | 0.1425 | 1.2370 |
| Weibull baseline with intermediate variance of 0.5 | | | | | | | | | |
| PI | | | | PD | | Parameter Estimates for PEM | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
|  | 1.061 | 4686.231 | -2343.574 | 4664.45 | -2333.295 | 1.2189 | -0.7703 | 0.7253 | -2.9976 |
| 25% | 1.35 | 2519.796 | -1260.203 | 2518.78 | -1261.124 | 1.6592 | -0.5801 | -1.2489 | 0.06127 |
| 50% | 1.297 | 2638.277 | -1319.477 | 2624.79 | -1313.904 | 1.5359 | -0.6076 | -0.8941 | 3.3212 |
| 75% | 1.397 | 1364.274 | -682.8361 | 1349.48 | -676.838 | 1.4953 | -0.3659 | -0.9157 | 15.9160 |
| Weibull baseline with high variance of 0.75 | | | | | | | | | |
| PI | | | | PD | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
|  | 1.314 | 5117.817 | -2559.707 | 5109.984 | -2556.26 | 1.7013 | -0.7955 | -1.1526 | -0.5874 |
| 25% | 1.398 | 3898.273 | -1949.766 | 3869.554 | -1935.94 | 1.5475 | -0.3632 | -0.8892 | 4.3118 |
| 50% | 1.544 | 2651.088 | -1326.259 | 2645.09 | -1323.97 | 1.9363 | -1.0053 | -1.0276 | 0.3476 |
| 75% | 1.739 | 1390.576 | -696.2893 | 1385.082 | -694.934 | 1.8493 | -1.1994 | 1.0844 | 1.8478 |

| **n=1000** | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Weibull baseline with low variance of 0.25 | | | | | | | | | |
| Partitioning is ignored | | | | Partitioning is done | | Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
| | 0.825 | 9506.566 | -4754.146 | 9466.548 | -4734.02 | 1.1574 | -0.6495 | -0.7254 | -1.396 |
| 25% | 0.7595 | 6310.06 | -3155.877 | 6306.177 | -3154.69 | 0.8024 | -0.2136 | -0.1869 | 2.0150 |
| 50% | 0.8459 | 4573.932 | -2287.615 | 4570.375 | -2286.98 | 1.0371 | -0.4995 | -0.6433 | 2.5102 |
| 75% | 0.9339 | 2461.439 | -1231.467 | 2472.358 | -1238.24 | 1.2121 | -0.6227 | -0.2574 | -0.8896 |
| Weibull baseline with intermediate variance of 0.5 | | | | | | | | | |
| PI | | | | PD | | PM Parameter Estimates | | | |
| No censoring | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
| | 1.123 | 10820.7 | -5411.503 | 10776.3 | -5389.516 | 1.4273 | -0.8297 | 0.9314 | -0.6932 |
| 25% | 0.9805 | 7083.88 | -3542.516 | 7086.044 | -3544.98 | 1.0542 | -0.3835 | -0.2667 | 1.65377 |
| 50% | 1.101 | 4991.19 | -2495.927 | 4998.826 | -2501.537 | 1.1879 | -0.2026 | -0.8005 | 0.24082 |
| 75% | 1.26 | 2616.16 | -1308.904 | 2619.639 | -1312.382 | 1.5240 | -0.6691 | -0.3263 | -0.7864 |
| Weibull baseline with high variance of 0.75 | | | | | | | | | |
| PI | | | | PM | | Parameter Estimates | | | |
| 0% | $\beta$ | DIC | LPML | DIC | LPML | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ |
| | 1.359 | 11679.5 | -5841.049 | 11615.29 | -5808.81 | 1.487 | 6.407e-01 | 7.365e-01 | 2.760e+04 |
| 25% | 1.212 | 7489.133 | -3744.906 | 7432.219 | -3789.62 | 1.2854 | -0.3711 | 0.11473 | 0.8471 |
| 50% | 1.317 | 5213.994 | -2607.245 | 5204.697 | -2603.91 | 1.4699 | -0.3473 | -0.6845 | -0.6357 |
| 75% | 1.593 | 2647.961 | -1324.38 | 2632.603 | -1327.79 | 1.8139 | -0.5563 | -0.6096 | 0.3036 |

**Interpretation**

Table 2, present the mean posterior estimates, DIC and LPML across all sample sizes and censoring percentages for single models and for the modified Piece wise models in the presence of the functional form of Time changing covariate, we observed that the values of estimates when models were fitted with data partitioning having observed the graph of beta against time for appropriate cut points are different (not constant), which indicate a change of effect parameters over time. We observed that the PPMs perform better than the single models throughout the simulations, for all censoring percentages & sample sizes.

When variance parameters for the Weibull baseline hazard were examined for low at 0.25, moderate or Intermediate at 0.5 & high at 0.75, estimates become worse with increase in variance and sample sizes, reflective in high DIC values and weak predictive power. In all of these, the Piece wise models out-performed the single ones; we again, noticed that the mean posterior estimates were better with increase in censoring percentages.

## Conclusion

We observed that the mean posterior estimates when the PPM - Model II was fitted, indicates change in effect parameters over time in all four intervals, with DIC and LPML values suggesting that PPM performs better than the Single model, for all censoring percentages, sample sizes & for the three (3) levels Weibull baseline. When the Weibull baseline hazard gain spread estimates were worse. In all of these, the PPM out-performed the SPM.

## Recommendations

The researcher recommends that:
1. other life distributions should be assumed as baseline to study the behavior of the models
2. combinations of baseline distributions to study competing risk problems

## References

Abiodun, A. A. (2007). *Analyzing competing risk survival time data using Cox and parametric proportional hazards models*, JNSA. 19,74-79.

Abiodun, A. A. (2009). *A Bayesian approach to exploring unobserved heterogeneity in Clustered survival and competing risk d*ata, JNSA 20

Adeniyi, O. I. & Akinrefon, A. A. (2018). *First birth interval: Cox regression model with time varying covariates*, CJPS.

Aitkin, M., Laird, N. & Francis, B. (1983). A reanalysis of the Stanford heart transplant data (with discussion), *J Am Stat Assoc 78*, 264–292.

Arjas, E. & Gasbarra, D. (1994). Nonparametric Bayesian inference from right censored survival data, *Stat Sinica 4*: 505–524.

Bastos, L. S. & Gamerman, D. (2006). Dynamic survival models with spatial frailty, *Lifetime Data Anal 12*, 441–460.

Bender R., Augustin, T. & Blettner. M. (2005). Generating survival times to simulate cox proportional hazards models, *Statistics in Medicine 24* (11) 1713-1723

Breslow, N. E. (1974). Covariance analysis of censored survival data, *Biometrics 30*, 89–99.

Clark, D. E. & Ryan, L. M. (2002). Concurrent prediction of hospital mortality and length of stay from risk factors on admission, *Health Services Res 37*, 631–645.

David, G. & Kleinbaum, M. K. (2005). *Survival analysis: A self-learning text*, New York, NY: Springer 2005

Fabio, N. D, Rosangela, H. L, Enrico, A, & Dipak, K. (2010). *Extensions of the piecewise exponential model,* Corpus ID: 53392910

Fan, J. & Gijbels, I. (1996). *Local polynomial modelling and its applications*, Chapman & Hall. London

Friedman, J. & Silverman, B. (1989). Flexible parsimonious smoothing and additive modelling (with discussion), *Technometrics, 31*, 3-39.

Gamerman, D. (1994). *Bayes estimation of the piece-wise exponential distribution*, IEEE Trans Reliab 43: 128–131.

Gamerman, D. (1997). *Efficient sampling from the posterior distributions in generalized linear models*, Statistics and Computing. 7, 57-68.

Hastie, T. & Tibshirani, R. (1990). *Generalized additive models*, Chapman and Hall London.

Hennerfeind, A, Brezger, A. & Fahrmeir, L. (2006). Geoadditive survival models, *Journal of the American Statistical Association, 101* (475), 1065–1075.

Ibrahim, J. G., Chen, M. H., & Sinha, D. (2001). *Bayesian survival analysis*, Springer-Verlag.

Kim, J. S. & Proschan, F. (1991). *Piecewise exponential estimator of the survival function*, IEEE Trans Reliab 40: 134–139.

Kim, S., Chen, M. H., Dey, D. K. & Gamerman, D. (2006). *Bayesian dynamic models for survival data with a cure fraction*, Lifetime Data Anal 13: 17–35.

Kooperberg, C. & Intrator. N. (1995). Trees and splines in survival analysis, *Statistical Methods in Medical Research, 4* 237–261.

Marano, G., Boracchi, P. & Biganzoli, E. M. (2016). Estimation of the piecewise exponential model by Bayesian P-Splines via Gibbs sampling: Robustness and reliability of posterior estimates. *Open Journal of Statistics, 6,* 451-468

McKeague, I. W. & Tighiouart, M. (2000). *Bayesian estimators for conditional hazard functions, Biometrics 56*, 1007–1015.

Sahu, S. K., Dey, D. K., Aslanidu, H. & Sinha, D. (1997). Aweibull regression model with gamma frailties for multivariate survival data, *Lifetime Data Anal 3,* 123–137.

Sinha, D., Chen, M. H. & Gosh, S. K. (1999). Bayesian analysis and model selection for interval-censored survival data, *Biometrics 55*: 585–590.

Stone, C., Hamsen, M., Kooperberg, C. & Truong, Y. (1997). Polynomial splines and their tensor products in extended linear modelling, *The Annals of Statistics, 25*, 1371-1470.

Zhou, H, & Hanson, T. (2017). A unified framework for fitting Bayesian Semiparametric Models to Arbitrarily Censored Survival Data, Including Spatially-Referenced Data, *Journal of the American Statistical Association*.