# Measuring the performances of covariates using exponential survival analysis with partly-interval censored simulation data

Siti Afiqah Muhamad Jamil ✉; Jessintha Lai; Mohd Asrul Affendi Abdullah

Check for updates

View Online    Export Citation    CrossMark

19 March 2024 04:31:57

# Measuring the Performances of Covariates using Exponential Survival Analysis with Partly-Interval Censored Simulation Data

Siti Afiqah Muhamad Jamil[1, a)], Jessintha Lai[2, b)] and
Mohd Asrul Affendi Abdullah[3, c)]

[1]*School of Mathematical Sciences, College of Computing, Informatics and Media, Universiti Teknologi MARA, 40450, Shah Alam,*
*Selangor, Malaysia*
[2, 3] *Faculty of Applied Sciences and Technology, Universiti Tun Hussein Onn Malaysia, 84600, Pagoh,*
*Johor Malaysia*

a) *Corresponding author: afiqahjamil@uitm.edu.my*
b) *jessinthalai@gmail.com*
c) *afendi@uthm.edu.my*

**Abstract.** In many fields of science, modelling and analyzing survival rates has shown to be a valuable element of statistical study. This paper aims at proposing the partly-interval censored data into the fixed and time-varying covariates and measure the performances of Exponential survival distribution using mean square error (MSE), mean bias error (MBE), mean absolute error (MAE), root mean square error (RMSE), mean absolute percentage error (MAPE) and standard error. As a result, when dealing the data without censored observations, the exponential distribution significantly fit the simulation data since low values of error measurements appeared when the data included the exact and complete types of simulation. Thus, this study proposed that the uncensored data could be applicable towards the Exponential survival distribution compared to other distributions of survival analysis.

## INTRODUCTION

Sequential design of survival analysis has extensively been applied through medical records observation. In conducting time-to-event data, the fundamental approach of survival is focusing on the event of interest where the occurrences of an event will particularly lead to censored observation or missing outcomes [1]. Despite the rapid application of survival analysis in medical studies, partly-interval censored data limitedly being applied. The estimation of survival analysis has been divided into parametric, non-parametric and semi-parametric regression analysis [2].

To start off the analysis of survival, unknown distribution leads to the non-parametric survival which is the Kaplan-Meier regression analysis as one of the most popular application of survival. One of the studies has compared the non-parametric Kaplan Meier, semi parametric cox proportional hazard and parametric methods to estimate the survival and compare the treatment with higher survival rates [3]. Cox-proportional hazard semi parametric survival analysis is applicable because of its simplicity. Nevertheless, the proportionality of hazard function needs to be observed prior to cox-proportional hazards or ultimately, parametric survival analysis could be applied with regards to certain assumptions of data analysis [4].

Consequently, the advanced application of exponential distribution with various extension or modification has broadly been develop not only towards the medical records data but also in other field of study. Besides, several exponential models have been employed towards Bayesian approach of analysis by considering the simple, change point, mixture and survival fraction of exponential model [5]. Besides, comparison of exponential distribution among the parametric method of survival can be seen in many of the previous studies [6]; [7] and [8].

Briefly, partly-interval censored data involves a modification of censored data where the study needs to consider all types of censored to be estimated in the analysis [9]. In order to cope with the medical data, partly-interval censored data has been created to record the observation and improves the precision. Compared to right, left and interval censored data, the partly-interval censored data is rarely applied by the researcher as the function is much more complex and it was found that several papers involving partly-interval censored data had been written by the same author which are [10]; [11], [12], [13] and [14].

Besides, difficulty in identifying which methods of survival is preferable and precise has become one of the problems to analyze the medical records data. Among all parametric survival distributions, Exponential distribution has the constant hazard distribution which is appropriate to start off the analysis through censored and uncensored data. In other words, adjusting the methods of survival analysis towards the partly-interval censored data with the presence of varying effect covariate is to be measured in this study using simulation data. Hence, this study proposes survival time of exponential method of survival with the presence of partly-interval censored data and compared the fixed and time-varying covariate effect of survival using several performance indicators.

## MATERIAL AND METHODS

### Framework of analysis

Based on Figure 1, the analysis starts with generating the simulation data where the data consists of the exact data, observed data and complete data. The exact data is the data without censored observation or fixed data, while observed data represents the fixed and all the censored data (partly-interval censored data) and the complete data is the data which represents the fixed, right censored and the missing values. Besides, 1000 iterations have been applied to test the fixed and time-varying covariate of partly-interval censored data with different assignation of sample sizes. The best fit model towards exponential with either fixed or time varying covariates will be selected based on the smallest value of model performances.
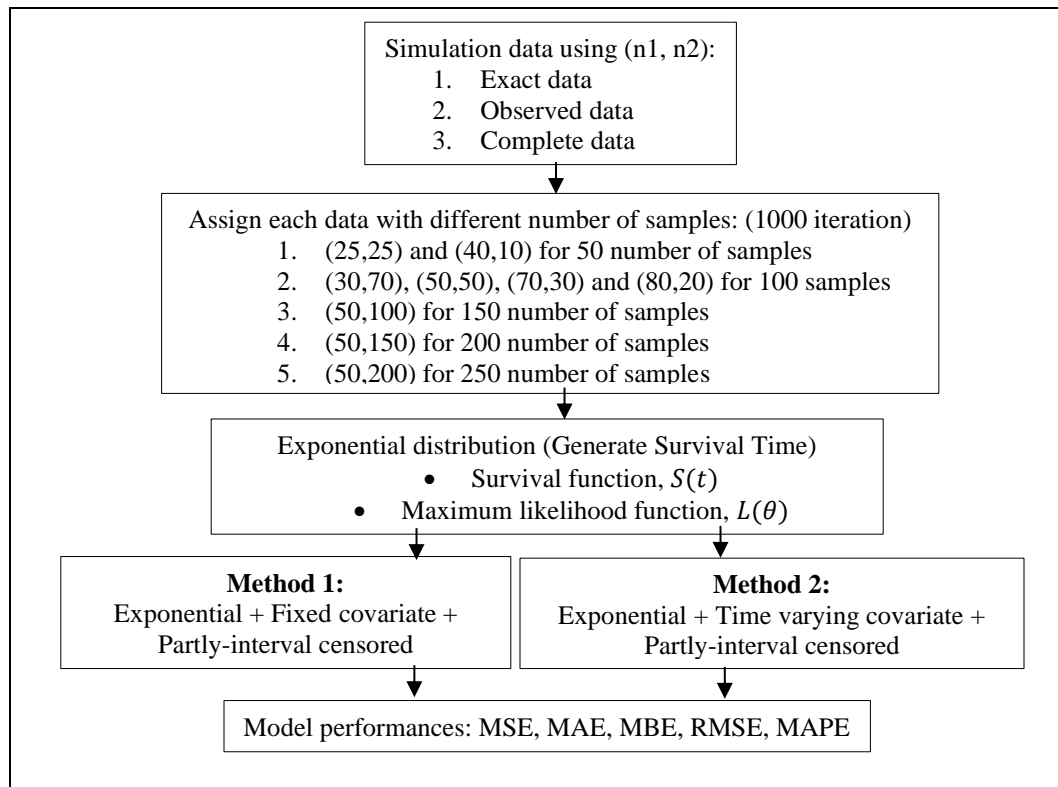


**FIGURE 1.** Study flowchart of simulation data

## Exponential fixed and time varying covariate likelihood estimation

Based on [2], the exponential distribution represents the survival function, $S(t_i)$, hazard function, $h(t_i)$, density function, $f(t_i)$, and likelihood function, $L(\theta)$ which are, $S(t_i) = exp\left(-\frac{t_i}{b}\right)$, $h(t_i) = \frac{1}{b}$, $f(t_i) = \frac{1}{b}exp\left(-\frac{t_i}{b}\right)$, and $L(b) = \prod_{i=1}^{n}\frac{1}{b}exp\left(-\frac{t_i}{b}\right)$ respectively. To be clear, $t$ represents the survival time while $b$ is the scale parameter with $b = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \ldots + \beta_p x_{pi}$. All these functions are modified in the partly-interval censored data of fixed and time-varying covariate effect.

For the time varying covariate, the notation is in binary form where the covariate might either change its value or not changing the value. The code "1" and "0" represent that the covariates changed and did not change respectively as the notation will follow the symbol $c_i$ and $(1 - c_i)$. Thus, on the study by [15], the dependent covariate can be defined by the following notation: $c_i = \begin{cases} 0, \text{covariate is updated} \\ 1, \text{covariate is not updated} \end{cases}$

The likelihood of exponential with fixed covariate and partly-interval censored data can be observed based on equation (1) and (2) below.

$$\ell = \log L(\theta) = \sum_{i=1}^{n_1} \log b^{-1} \exp\left(-\frac{t_i}{b}\right)$$
$$+ \sum_{i=n_1+1}^{n} (1-\delta_i) \log\left[\exp\left(-\frac{R_i}{b}\right)\right] + \sum_{i=n_1+1}^{n} \delta_i \log\left[1 - \exp\left(-\frac{L_i}{b}\right)\right]$$
$$+ \sum_{i=n_1+1}^{n} \rho_i \log\left[\exp\left(-\frac{R_i}{b}\right) - \exp\left(-\frac{L_i}{b}\right)\right] \tag{1}$$

The log likelihood of exponential with partly-interval censored data has been modified as shown in Equation (1) where the exponential only has one parameter, which is $b$ that represents the scale parameter and the first part of likelihood holds the exact data. The censored observation depends on the values of $\delta_i$ and $\rho_i$. When $\delta_i = 1$, the likelihood will follow the left censored observation, and otherwise when $\delta_i = 0$, it will become right censored observation. The interval censored observation will be applied when $\rho_i = 1$. The $R_i$ represent the right censored survival time, $L_i$ represent the left censored survival time, and $t_i$ represent the exact duration of survival time. Meanwhile the equation (2) represents the log-likelihood with partly-interval censored data and time-varying covariate effect.

$$\ell = \log L(\theta) = \sum_{i=1}^{n_1} (1-c_i) \log b_{i_0}^{-1} \exp\left(-\frac{t_i}{b_{i0}}\right)$$
$$+ \sum_{i=n_1+1}^{n} (1-c_i)(1-\delta_i) \log\left[\exp\left(-\frac{R_i}{b_{i0}}\right)\right] + \sum_{i=n_1+1}^{n} (1-c_i)\delta_i \log\left[1 - \exp\left(-\frac{L_i}{b_{i0}}\right)\right]$$
$$+ \sum_{i=n_1+1}^{n} (1-c_i)\rho_i \log\left[\exp\left(-\frac{R_i}{b_{i0}}\right) - \exp\left(-\frac{L_i}{b_{i0}}\right)\right]$$
$$+ \sum_{i=1}^{n_1} (c_i) \log b^{-1} \exp\left(-\frac{t_i}{b_{i1}}\right)$$
$$+ \sum_{i=n_1+1}^{n} (c_i)(1-\delta_i) \log\left[\exp\left(-\frac{R_i}{b_{i1}}\right)\right] + \sum_{i=n_1+1}^{n} (1-c_i)\delta_i \log\left[1 - \exp\left(-\frac{L_i}{b_{i1}}\right)\right]$$
$$+ \sum_{i=n_1+1}^{n} (c_i)\rho_i \log\left[\exp\left(-\frac{R_i}{b_{i1}}\right) - \exp\left(-\frac{L_i}{b_{i1}}\right)\right] \tag{2}$$

## Generating survival time from Exponential distribution

The derivation uses parametric family approach with covariate that is associated with the event time, $t$ that could either be fixed or time varying covariate which is $b = e^{x'\beta}$ or $b = e^{x'\beta + \gamma z(t)}$ respectively. For the simulation

procedure, the time varying covariate for $z(t)$ is in a binary form. The derivation will begin with finding the cumulative hazard function, $H(t)$ by integrating the hazard function of exponential of partly-interval censored from 0 to $t$ of $du$ for the fixed covariate and time varying covariate, $H(t) = \int_0^t h(u)du$ where the hazard function has been mentioned earlier. The hazard function in piecewise form is as follows (Collet, 2015):

$$h(t|a,x,z(t)) = \begin{cases} \left(\dfrac{1}{exp\{x'\beta\}}\right), & t < t_c \\[3mm] \left(\dfrac{1}{exp\{x'\beta + \gamma z(t)\}}\right), & t \geq t_c \end{cases}$$

After the integration of the hazard function of Weibull model, the cumulative hazard function can be defined as follows:

$$H(t|a,x,z(t)) = \begin{cases} \displaystyle\int_0^t \left(\dfrac{1}{exp\{x'\beta\}}\right) du, & t < t_c \\[3mm] \displaystyle\int_0^t \left(\dfrac{1}{exp\{x'\beta + \gamma z(u)\}}\right) du, & t \geq t_c \end{cases}$$

Therefore, in order to get the survival function, the cumulative hazard function will be used by the exponent of the negative cumulative function. The integration has been made in order to find the cumulative hazard with fixed and varying covariate in the model. Besides, the main idea of exponential function was adapted from [1]. Thus, the cumulative hazard function for exponential has been developed in order to find the survival time of fixed and varying effect covariate,

For $t < t_c$,
$$H(t \mid a, x, z(t)) =$$
$$= \int_0^t \left( \frac{1}{\exp\{x'\beta\}} \right) du$$
$$= \left( \frac{1}{\exp\{x'\beta\}} \right) [u]_0^t$$
$$= \left( \frac{t}{\exp\{x'\beta\}} \right)$$

For $t \geq t_c$,
$$H(t \mid a, x, z(t)) =$$
$$= \int_0^t \left( \frac{1}{\exp\{x'\beta + \gamma\}} \right) du$$
$$= \int_0^{t_c} \left( \frac{1}{\exp\{x'\beta\}} \right) du + \int_{t_c}^t \left( \frac{1}{\exp\{x'\beta + \gamma\}} \right) du$$
$$= \left( \frac{t_c}{\exp\{x'\beta\}} \right) + \left( \frac{1}{\exp\{x'\beta + \gamma\}} \right) \int_{t_c}^t u\, du$$
$$= \left( \frac{t_c}{\exp\{x'\beta\}} \right) + \left( \frac{1}{\exp\{x'\beta + \gamma\}} \right) [u]_{t_c}^t$$
$$= \left( \frac{t_c}{\exp\{x'\beta\}} \right) + \left( \frac{1}{\exp\{x'\beta + \gamma\}} \right) [t - t_c]$$
$$= \left( \frac{t_c}{\exp\{x'\beta\}} \right) + \left( \frac{t}{\exp\{x'\beta + \gamma\}} \right) - \left( \frac{t_c}{\exp\{x'\beta + \gamma\}} \right)$$

Next, in order to find the survival time, $t$, the survival function is assumed to follow uniform distribution and extract the survival time from the equation. Thus, the survival time, $T$ could be defined as $S(t|a,x,z(t)) = exp[-H(t|a,x,z(t))]$:

For $t < t_c$,

$$S(t) = \exp\left\{-\wedge(t|x_i)\right\}$$

Let $S(t) \sim U(0,1)$

$$unif(0,1) = \exp\left[-\left(\frac{t}{\exp\{x'\beta\}}\right)\right]$$

$$-\log(u) = \left(\frac{t}{\exp\{x'\beta\}}\right)$$

$$t = (-\log(u)) \cdot \exp\{x'\beta\}$$

For $t \geq t_c$,

$$S(t) = \exp\left\{-\wedge(t|x_i)\right\}$$

Let $S(t) \sim U(0,1)$

$$unif(0,1) = \exp\left[-\left(\frac{t_c}{\exp\{x'\beta\}}\right)+\left(\frac{t}{\exp\{x'\beta+\gamma\}}\right)-\left(\frac{t_c}{\exp\{x'\beta+\gamma\}}\right)\right]$$

$$\left(\frac{t}{\exp\{x'\beta+\gamma\}}\right) = -\log(u) - \left(\frac{t_c}{\exp\{x'\beta\}}\right)+\left(\frac{t_c}{\exp\{x'\beta+\gamma\}}\right)$$

$$t = \exp\{x'\beta+\gamma\} \cdot \left[-\log(u) - \left(\frac{t_c}{\exp\{x'\beta\}}\right)+\left(\frac{t_c}{\exp\{x'\beta+\gamma\}}\right)\right]$$

Therefore, the survival time for exponential of fixed and time varying covariate could be applied in simulation procedure based on the following piecewise function:

$$t = \begin{cases} (-\log(u)) \cdot exp\{x'\beta\}, & t < t_c \\ exp\{x'\beta + \gamma\} \cdot \left[-\log(u) - \left(\frac{t_c}{exp\{x'\beta\}}\right) + \left(\frac{t_c}{exp\{x'\beta+\gamma\}}\right)\right], & t \geq t_c \end{cases} \tag{3}$$

where the equation (3) of survival time, $t$ is to be applied in the R coding so that the analysis of exponential distribution could be done. Through the used of the cumulative hazard function of exponential distribution that has been developed to obtain the survival time with fixed and varying effect covariate, simulation study could be done for this study.

## Simulation Studies

To evaluate the performance of parametric survival analysis of exponential distribution, a simulation strategy has been illustrated by applying the standard value of distributions in modelling the parameters which are, $\lambda = \frac{1}{0.25}$, $x_1 \sim$ Binomial $(1,0.5)$, $x_2 \sim$ Normal $(0,1)$, $x_3 \sim$ time varying covariate. The survival time equation of exponential (3) has been applied in the analysis. Besides, the censoring follows partly interval censored data are the observe data where coding R produce simulation data based on exact, observe and complete data which subsequently means that the simulation data are having an exact duration of survival for exact data, the observe data represent the data while considering all right, left, and interval censored with exact observation while the complete data is data before introducing the censored observation by considering the missing value. Comparison between different sample sizes which are 50, 100, 150, 200 and 250 with 1000 iteration of simulation data has been carried out.

# RESULTS AND DISCUSSION

**TABLE 1.** Simulation results of root mean square error (RMSE) for exponential fixed covariate and time varying covariate

| $(n_1, n_2)$ | Fixed Covariate | | | Time Varying Covariate | | |
|---|---|---|---|---|---|---|
| | **Exact** | **Observed** | **Complete** | **Exact** | **Observed** | **Complete** |
| (25,25) | 1.611 | 3.903 | **1.597** | 1.596 | 14.08 | **1.467** |
| (40,10) | 1.602 | 3.978 | **1.597** | 1.508 | 40.69 | **1.467** |
| (30,70) | 1.614 | 4.043 | **1.588** | 1.564 | 46.42 | **1.354** |
| (50,50) | 1.596 | 3.855 | **1.588** | 1.454 | 10.80 | **1.354** |
| (70,30) | 1.591 | 3.939 | **1.588** | 1.393 | 18.89 | **1.354** |
| (80,20) | 1.590 | 3.948 | **1.588** | 1.372 | 15.21 | **1.354** |
| (50,100) | 1.601 | 3.994 | **1.584** | 1.295 | 12.72 | **1.207** |
| (50,150) | 1.601 | 3.994 | **1.584** | 1.469 | 14.46 | **1.322** |
| (50,200) | 1.596 | 4.058 | **1.583** | 1.466 | 11.10 | **1.313** |

Based on Table 1, the simulation results showed that the lowest values of mean square error for both fixed and time varying covariate were among the complete data of simulation. So, the exponential distribution could likely estimate the data before introducing the censored observation and considering the exact and some missing values.

**TABLE 2.** Simulation results of mean absolute percentage error (MAPE) of exponential fixed covariate and time varying covariate

| $(n_1, n_2)$ | Fixed Covariate | | | Time Varying Covariate | | |
|---|---|---|---|---|---|---|
| | **Exact** | **Observed** | **Complete** | **Exact** | **Observed** | **Complete** |
| (25,25) | **0.995** | 2.395 | 0.998 | 0.815 | 3.762 | **0.770** |
| (40,10) | **0.998** | 2.444 | **0.998** | 0.783 | 8.741 | **0.770** |
| (30,70) | **0.999** | 2.482 | **0.999** | 0.800 | 10.01 | **0.729** |
| (50,50) | **0.997** | 2.377 | 0.999 | 0.767 | 3.270 | **0.729** |
| (70,30) | **0.998** | 2.427 | 0.999 | 0.745 | 6.370 | **0.729** |
| (80,20) | **0.997** | 2.432 | 0.999 | 0.738 | 5.940 | **0.729** |
| (50,100) | 1.001 | 2.457 | **1.000** | 0.701 | 4.947 | **0.667** |
| (50,150) | 1.001 | 2.457 | **1.000** | 0.771 | 5.700 | **0.716** |
| (50,200) | 0.998 | 2.498 | **0.996** | 0.773 | 4.839 | **0.714** |

Based on Table 2, the mean absolute percentage error was likely having lower MAPE among the exact and complete data while the time varying covariate were all having low values of MAPE among the complete simulation data. This would make the results to prefer the exponential distribution to have a better estimation with exact data and complete data without considering the censored observation.

**TABLE 3.** Simulation results of mean absolute error of exponential fixed covariate and time varying covariate

| $(n_1, n_2)$ | Fixed Covariate | | | Time Varying Covariate | | |
|---|---|---|---|---|---|---|
| | **Exact** | **Observed** | **Complete** | **Exact** | **Observed** | **Complete** |
| (25,25) | **1.492** | 3.355 | 1.498 | 1.185 | 5.794 | **1.111** |
| (40,10) | **1.498** | 3.423 | **1.498** | 1.131 | 12.971 | **1.111** |
| (30,70) | 1.501 | 3.470 | **1.499** | 1.160 | 14.85 | **1.050** |
| (50,50) | **1.497** | 3.377 | 1.499 | 1.105 | 5.015 | **1.050** |
| (70,30) | **1.498** | 3.399 | 1.499 | 1.072 | 9.650 | **1.050** |
| (80,20) | **1.497** | 3.406 | 1.499 | 1.061 | 9.038 | **1.050** |
| (50,100) | 1.502 | 3.443 | **1.499** | 1.013 | 7.516 | **0.960** |
| (50,150) | 1.502 | 3.443 | **1.499** | 1.114 | 8.754 | **1.030** |
| (50,200) | 1.497 | 3.495 | **1.496** | 1.114 | 7.395 | **1.025** |

Similarly, the mean absolute error (MAE) was most likely to have lower value towards the exact and complete data. This is because, the Exponential model is considered as the simplest parametric model because of the hazard is constant over time and it does not hold the flexibility of complex functions. Specifically, based on the results, the time-varying covariate of partly-interval censored data would perform better when the data is complete by considering the missing value and before introducing the censored data.

Additionally, the standard error for both fixed and time varying covariate fit much better towards the complete data of simulation. At first, Table 4 showed the result of fixed covariate of exponential distribution which likely preferable with complete data with the lowest value of standard error.

**TABLE 4.** Simulation results of standard error of exponential fixed covariate

| $(n_1, n_2)$ | Standard Error, $b_1$ | | | Standard Error, $b_2$ | | |
|---|---|---|---|---|---|---|
| | Exact | Observed | Complete | Exact | Observed | Complete |
| (25,25) | 0.444 | 0.413 | **0.296** | 0.228 | 0.724 | **0.160** |
| (40,10) | 0.333 | 0.342 | **0.296** | 0.175 | 0.602 | **0.160** |
| (30,70) | 0.402 | 0.319 | **0.206** | 0.204 | 0.514 | **0.102** |
| (50,50) | 0.298 | 0.220 | **0.206** | 0.148 | 0.419 | **0.102** |
| (70,30) | 0.245 | 0.220 | **0.206** | 0.121 | 0.392 | **0.102** |
| (80,20) | 0.229 | 0.199 | **0.206** | 0.114 | 0.367 | **0.102** |
| (50,100) | 0.309 | 0.230 | **0.160** | 0.151 | 0.386 | **0.085** |
| (50,150) | 0.295 | 0.210 | **0.141** | 0.155 | 0.325 | **0.072** |
| (50,200) | 0.285 | 0.196 | **0.126** | 0.148 | 0.266 | **0.064** |

Meanwhile, the results of standard error in Table 5 with considering the time-varying covariate into the model has also showed a more preferable lowest value towards the complete data.

**TABLE 5.** Simulation results of standard error for exponential time-varying covariate

| $(n_1, n_2)$ | Standard Error, $b_1$ | | | Standard Error, $b_2$ | | | Standard Error, $b_3$ | | |
|---|---|---|---|---|---|---|---|---|---|
| | Exact | Observed | Complete | Exact | Observed | Complete | Exact | Observed | Complete |
| (25,25) | 0.453 | 8.070 | **0.314** | 0.256 | 5.794 | **0.162** | 1.463 | 19.61 | **0.865** |
| (40,10) | 0.350 | 19.18 | **0.314** | 0.188 | 24.15 | **0.162** | 1.012 | 57.32 | **0.865** |
| (30,70) | 0.403 | 23.61 | **0.210** | 0.229 | 16.82 | **0.114** | 1.315 | 68.24 | **0.503** |
| (50,50) | 0.298 | 4.873 | **0.210** | 0.171 | 2.849 | **0.114** | 0.788 | 15.06 | **0.503** |
| (70,30) | 0.253 | 7.040 | **0.210** | 0.141 | 3.289 | **0.114** | 0.636 | 22.97 | **0.503** |
| (80,20) | 0.238 | 7.570 | **0.210** | 0.132 | 2.469 | **0.114** | 0.554 | 14.82 | **0.503** |
| (50,100) | 0.304 | 6.998 | **0.174** | 0.169 | 2.121 | **0.094** | 0.633 | 11.27 | **0.331** |
| (50,150) | 0.309 | 8.607 | **0.146** | 0.164 | 1.337 | **0.077** | 0.817 | 12.76 | **0.346** |
| (50,200) | 0.318 | 6.156 | **0.137** | 0.178 | 1.105 | **0.072** | 0.826 | 10.29 | **0.300** |

## CONCLUSION

Based on the results of simulation with different assignations of sample sizes, the comparison of model selection criteria which are the MBE, MSE, RMSE, MAE, MAPE, and standard error had been carried out on the parametric distribution of survival analysis. Based on the exponential distribution, the MBE, MAE, and MAPE clearly showed that this type of distribution is fit when handling the data without censored observation as low values appeared when the data comprised exact and complete types of simulation. In terms of sample sizes, when the number of samples is below 100, the exact data is needed for exponential distribution. Increasing the sample sizes for partly-interval censored data is appropriate with $(n_1 < n_2)$ where $n_1$ represents the exact data, while $n_2$ represents the censored data and will also decrease the values of estimation. Therefore, proposing the partly-interval censored data into fixed and time-varying covariates has been done through Exponential survival functions and it proves that Exponential distribution would fit better towards exact and complete simulation data (without censored data) for both situations when considered the fixed and time-varying covariates. It is recommended that using different values of parameter on the simulation data and different types of distribution of the parametric survival analysis might give further insight on handling the censored observations.

## ACKNOWLEDGMENTS

## REFERENCES

1. D. Collett, Modelling survival data in medical research (CRC press. 2015).
2. J. F. Lawless, *Statistical models and methods for lifetime data* (Vol. **362**). John Wiley & Sons 2003
3. V. Bewick, L. Cheek, & J. Ball, Critical care. **8**(5), 389 (2004).
4. G. Rodrıguez, 2010. Parametric survival models. *Princeton University, Rapport technique, Princeton*
5. Y. Chung, D. K. Dey, M. Kim, & C. Kim, *Communications in Statistics—Theory and Methods*, **34**(12), 2311-2330 (2005).
6. B. Singh, *Lifetime data analysis*, **8**(1), 69-88 (2002).
7. P. C. Austin, *Statistics in medicine.* **31**(29), 3946–3958 (2012).
8. N. Bottomley, L. D. Jones, R. Rout, A. Alvand, I. Rombach, T. Evans, & A. J. Price, *The Bone & Joint Journal*, **98**(10_Supple_B), 22-27 (2016).
9. Y. H. Sparling, N. Younes, J. M. Lachin, & O. M. Bautista, *Biostatistics.* **7**(4), 599–614 (2006).
10. F. A. M. Elfaki, M. Azram, & M. Usman, *International Journal of Applied Physics and Mathematics.* **2**(5), 352 (2012).
11. F. A. M. Elfaki, A. Abobakar, M. Azram, & M. Usman, W*orld Academy of Science, Engineering and Technology.* **7**, 899–902 (2013).
12. N. Zainudin, F. Elfaki, and M. Y. Ali, Confidence Interval for Survival Model based on Partly Interval-Censored data, (2015).
13. A. M. Alharpy, & N. A. Ibrahim, *Journal of Applied Sciences.* **13**(4), 621 (2013).
14. A. M. Alharpy, & N. A. Ibrahim, *Mathematical Problems in Engineering,* (2014).
15. K. Kiani, J. Arasan, & H. Midi, *Sains Malaysiana.* **41**(4), 471–480 (2012).