

# P8160 Project 1 Final Report

## A Simulation Study to Compare Three Survival Models

Jingchen Chai, Yi Huang, Ruihan Zhang

March 3, 2023

### Contents

<b>Abstract</b>	<b>2</b>
<b>Introduction</b>	<b>2</b>
<b>Objective</b>	<b>2</b>
<b>Background</b>	<b>2</b>
Survival function, CDF of survival time, hazard function . . . . .	2
Intro to Proportional hazard model . . . . .	3
Three Proportional-hazards model . . . . .	3
<b>Methodology</b>	<b>3</b>
Data Generation: derive survival time T . . . . .	3
Simulation Design . . . . .	4
<b>Results and Discussion</b>	<b>5</b>
Scenario 1: Exponential baseline hazard function . . . . .	5
Scenario 2: Weibull baseline hazard function . . . . .	7
Scenario 3: Gompertz baseline hazard function . . . . .	8
Additional Cases . . . . .	10
<b>Conclusion</b>	<b>11</b>
<b>References</b>	<b>11</b>
<b>Appendix</b>	<b>11</b>

## Abstract

Survival analysis or time-to-event analysis is commonly used in medical, economics, and engineering research. In survival analysis, each observation unit will be an individual or a patient, the analysis of a variable that corresponds to the time from a defined baseline until the occurrence of an event of interest. The objective of the simulation study is to evaluate the robustness of three survival models against the misspecified baseline hazard functions. To achieve this goal, we assess the accuracy and efficiency of the estimated treatment effects ( $\beta$ ) for each model under various baseline hazard functions. This project will perform simulations to compare the parametric regression model (Exponential and Weibull) to the semi-parametric regression model (Cox) for survival data. We apply the inverse transformation method to generate survival data with censored observations from Exponential, Weibull, Gompertz, and Two-component mixture distributions. Despite the different models, we observed that appropriate baseline hazard results in the best performance.

## Introduction

Survival analysis is a statistical method used to analyze time-to-event data, such as the time until a patient's death. It is commonly used in medical research to estimate the probability of an event occurring over time and to identify factors that may affect the risk of the event. The analysis also considers time-dependent covariates that vary over time and can be used to compare the survival times between different groups of patients. The uncensored data is recorded when the event occurs at the exact observed event time, and it is coded as 1 in the status indicator variable. Conversely, censored data occurs when patients are lost to follow-up or the event happens after the study period, and it is coded as 0 in the status indicator variable. Time-dependent covariates are also considered in the analysis, which can vary over time and can be used to compare the survival times of different participant groups, such as a treatment group and a control group.

## Objective

The objective is to design a simulation study to compare and contrast the efficiency and accuracy of the estimated treatment effects under different baseline hazard functions and assess their robustness against misspecified baseline hazard functions. A practical and effective recommendation will be provided to general users to select a suitable model on the basis of the numerical investigations.

## Background

### Survival function, CDF of survival time, hazard function

Let  $T \in (0, \infty)$  denote the survival time for a patient. Survival function,  $S(t)$  is the probability of observing individual survival time  $T$  beyond a certain time  $t$ . To analyze survival data, we define a **survival function**  $S(t)$  as

$$S(t) = \Pr(T > t) = \int_t^\infty f(s)ds$$

The Cumulative distribution function of the random variable survival time  $T$ ,  $F(t)$  is the probability of observing individual survival time less than a certain time  $t$ , we define **CDF of survival time**  $F(t)$  as

$$F(t) = \Pr(T \leq t) = 1 - S(t) = 1 - \int_t^\infty f(s)ds$$

The Hazard function,  $h(t)$  measures the instantaneous risk of failure at time  $t$  giving that a patient has survived until time  $t$ , defined as the ratio of the probability density function of time variable and Survival function.

$$h(t) = \lim_{\Delta_t \rightarrow 0} \frac{\Pr((T \in (t, t + \Delta_t)] | T > t)}{\Delta_t} = \frac{f(t)}{S(t)}$$

Cumulative Hazard Function  $H(t)$  can be expressed as

$$H(t) = \int_t^{\infty} h(s)ds$$

## Intro to Proportional hazard model

Proportional hazard model is the primary regression model to investigate the effectiveness of treatment  $X$  over survival time  $T$ , where the  $i$ -th patient at a time  $t$  is

$$h_i(t) = h_0(t)e^{x_i\beta}$$

where

- $h_0(t)$  is the baseline hazard function
- $x_i$  is the treatment indicator variable (control = 0, treatment = 1)
- $\beta$  is the parameter of interest, which is the log hazard ratio for the treatment effect.  $\beta$  measures the relative hazard reduction due to treatment in comparison to the control.

The **proportional hazard** can be expressed as ratio of two hazard functions at time  $t$  given individuals in different treatment groups, and does not depend on  $t$ .

$$\frac{h(t|x_0)}{h(t|x_1)} = e^{\beta(x_0-x_1)}$$

There are different ways to formulate the baseline hazard function  $h_0(t)$ , which lead to different models and estimations.

## Three Proportional-hazards model

**An exponential proportional-hazards model** assumes the baseline hazard function is a constant

$$h_0(t) = \lambda$$

**A Weibull proportional-hazards model** assumes the baseline hazard function follows Weibull distribution, where

$$h_0(t) = \lambda\gamma t^{\gamma-1}$$

for  $\gamma > 0$

**A Cox proportional-hazards model** leaves  $h_0(t)$  unspecified.

Note that exponential distribution is a special case of Weibull distribution where  $\lambda = 1$ . Hence, among the three models, the exponential proportional-hazards model is the most restrictive model, while the Cox model is the most general one.

## Methodology

### Data Generation: derive survival time $T$

Given:

$$H(t) = -\log(S(t))$$

$$S(t) = e^{-H(t)} = e^{\int_t^\infty h(s)ds}$$

Utilize Inverse Transformation Method

$$T = F^{-1}(U) = H_0^{-1}\left(\frac{-\log(U)}{e^{X\beta}}\right)$$

where  $U \sim U(0, 1)$

Thus, we can derive the following table

Table 1: Characteristic of Exponential, Weibull and Gompertz Distributions

	Distribution		
	Exponential	Weibull	Gompertz
Scale Parameter	$\lambda > 0$	$\lambda > 0$	$\lambda > 0$
Shape Parameter		$\gamma > 0$	$\alpha \in (-\infty, \infty)$
Baseline Hazard function	$h_0(t) = \lambda$	$h_0(t) = \lambda\gamma t^{\gamma-1}$	$h_0(t) = \lambda \exp(\alpha t)$
Cumulative Baseline Hazard Function	$H_0(t) = \lambda t$	$H_0(t) = \lambda t^\gamma$	$H_0(t) = \frac{\lambda}{\alpha}(e^{\alpha t} - 1)$
Inverse Cumulative Hazard Function	$H^{-1}(t) = \lambda^{-1}t$	$H^{-1}(t) = (\lambda^{-1}t)^{\frac{1}{\gamma}}$	$H^{-1}(T) = \frac{1}{\alpha} \log(1 + \frac{\alpha}{\lambda}t)$
Cumulative Distribution Function	$F(t) = 1 - e^{-\lambda e^{X\beta}}$	$F(t) = 1 - e^{-\lambda t^\gamma e^{X\beta}}$	$F(t) = 1 - e^{-\frac{\lambda}{\alpha}(e^{\alpha} - 1)e^{X\beta}}$
Survival Time $T$	$T = -\frac{\log(U)}{\lambda e^{X\beta}}$	$T = \left(-\frac{\log(U)}{\lambda e^{X\beta}}\right)^{\frac{1}{\gamma}}$	$T = \frac{1}{\alpha} \log[1 - \frac{\alpha \log(U)}{\lambda e^{X\beta}}]$

## Simulation Design

The goal of this simulation study is to evaluate how misspecifying the baseline hazard function influences the estimated treatment effect from 3 models. We generate 4 types of data, Exponential, Weibull, Gompertz, and Mixture distribution to evaluated the model robustness. All simulation data was generated using inverse transformation method and use `simsurv` function in the `simsurv` package to double check the result. The resulting dataset contains time of event (`eventtime`), status indicator (`status`), and treatment group (`trt`) with censored observations.

- Define Random Variable  $X$ , treatment, from a binomial distribution with  $p = 0.5$
- Generate survival time  $T$ , time to event, using  $X$  and  $\beta$
- Randomly generate censoring time  $C$ , from an exponential distribution
- We observe either the survival time  $T$  or else the censoring time  $C$ . Specifically, we observe the random variable  $Y$  takes the minimum value between survival time and censoring time

$$Y = \min(T, C)$$

Create a status indicator variable ( $1 = \text{event}$ ,  $0 = \text{censored}$ ), if the patient dies before censoring time, then status equals one, vice versa.

$$Status = \begin{cases} 1, & T_i \leq C_i \\ 0, & T_i > C_i \end{cases}$$

## Define true treatment effect and parameter

Before running the simulation, also need to define the true treatment effect and parameters

1. True treatment effect  $\beta = 2$
2. 7 different sample size  $N$  ranging from 100 to 400 increasing by 50
3. Simulate exponential distribution with  $\lambda = 0.5$
4. Simulate weibull distribution with  $\lambda = 0.5$  and  $\gamma = 2$
5. Simulate gompertz distribution with  $\lambda = 0.5$  and  $\alpha = 2$

## Simulation Times

A large number of simulations can improve the accuracy and efficiency of estimated treatment effect, thus we simulate 1000 times for N sample. Each datasets will be used to fit in three proportional hazard models to obtain the estimated treatment effects

## Performance Measures

To compare accuracy and efficiency of estimated treatment effects, we use performance measures such as Bias, Variance, Mean square error, and confidence interval to compare 3 models.

$k$ : number of independent datasets

## Estimation

$$\begin{aligned} Bias &= \frac{1}{k} \sum_{i=1}^k (\hat{\beta}^{(k)} - \beta) = \frac{1}{k} \sum_{i=1}^k (\hat{\beta}^{(k)} - 2) \\ Variance &= \frac{1}{k-1} \sum_{i=1}^k (\hat{\beta}^{(k)} - \beta)^2 = \frac{1}{k-1} \sum_{i=1}^k (\hat{\beta}^{(k)} - 2)^2 \\ MSE &= \frac{1}{k} \sum_{i=1}^k (\hat{\beta}^{(k)} - \beta)^2 = \frac{1}{k} \sum_{i=1}^k (\hat{\beta}^{(k)} - 2)^2 \end{aligned}$$

## Confidence Interval

$$Coverage : \frac{1}{k} \sum_{i=1}^k I\{\beta \in CI^{(k)}\}$$

## Results and Discussion

The parameters in each model were held constant with  $\lambda = 0.5$ ,  $\gamma = 2$ ,  $\alpha = 2$ ,  $\beta = 2$ . Each scenario are simulated for 1000 times, sample sizes from 100 to 400. The variance was calculated to examine the efficiency of the models. Additionally, MSE and bias were calculated to demonstrate the performance of survival models as the sample size increases. After running the three models in each scenario, a set of  $\beta$ 's were extracted and used to calculate the 95% confidence interval for comparison.

### Scenario 1: Exponential baseline hazard function

Figure 1 shows the variance from an exponential distribution. We can obviously see from the plot that the exponential model has the lowest versatality, whereas the cox model has the largest.

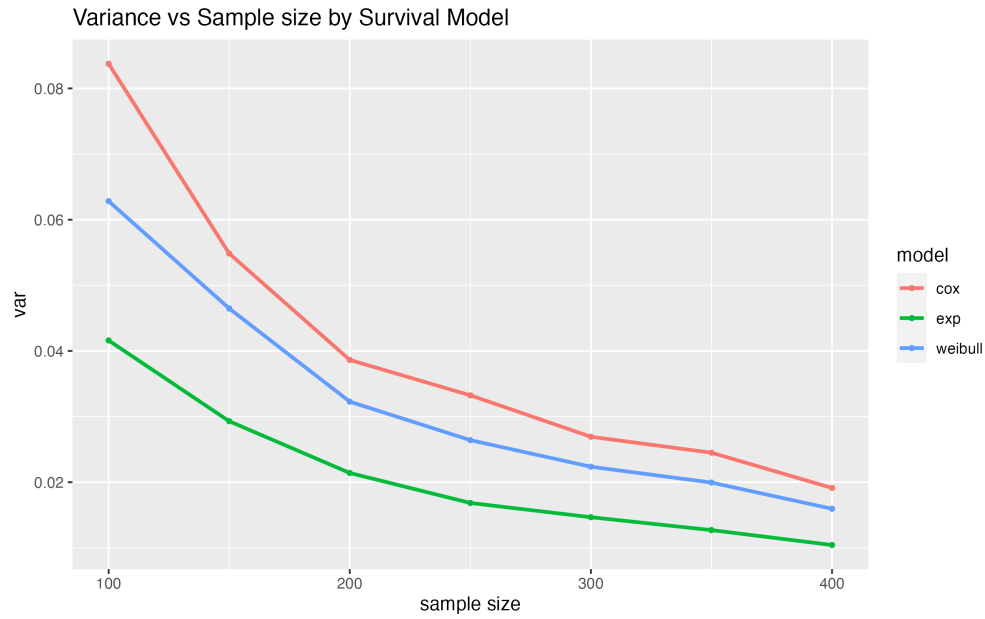


Figure 1

Figure 2-3 shows the bias and MSE from the exponential distribution. As the sample size increased, the MSE of all three models decreased. The Cox model had the MSE, whereas the exponential model had the least. The trend of the bias of the three models switched between increasing and decreasing with the increasing sample size. The Weibull model had the largest bias, and the exponential model had the least bias. Therefore, the exponential model is the most suitable model when fitting the data generated from an exponential distribution.

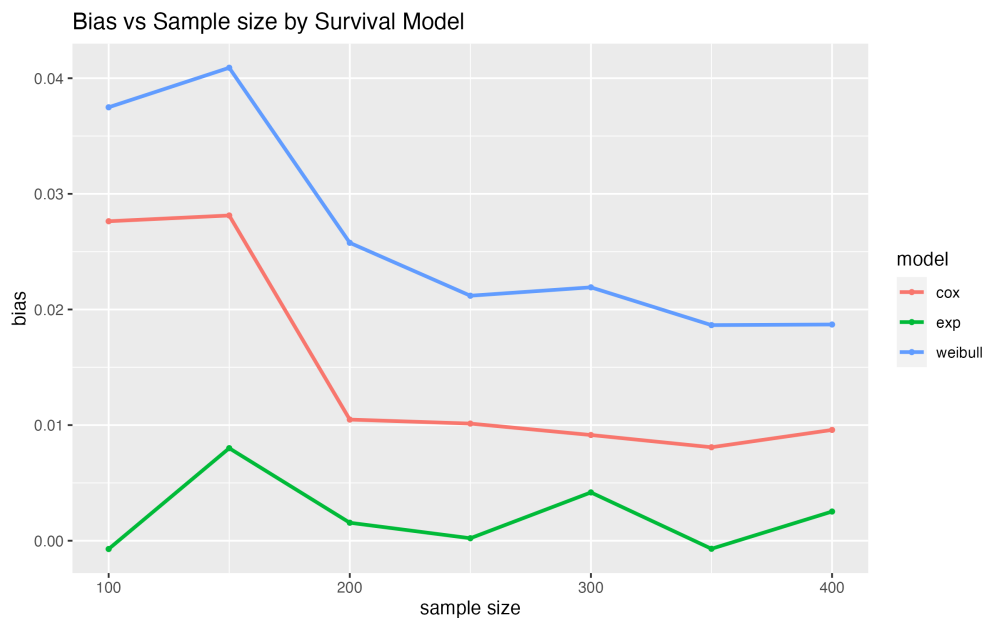


Figure 2

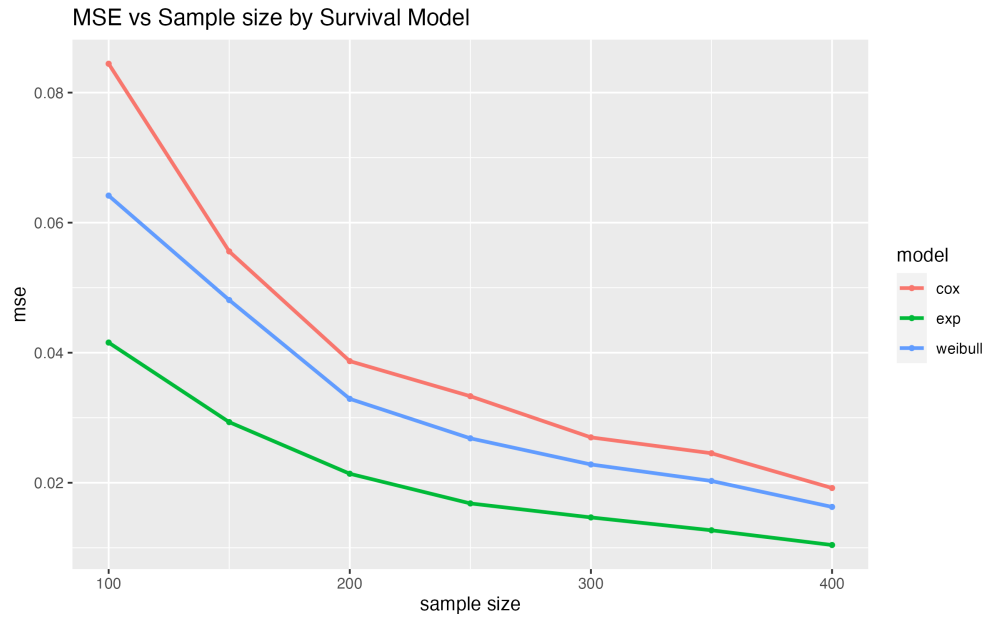


Figure 3

## Scenario 2: Weibull baseline hazard function

Figure 4 shows the variance from a weibull distribution. The exponential model shows to have the lowest variability, while the cox model has the largest.

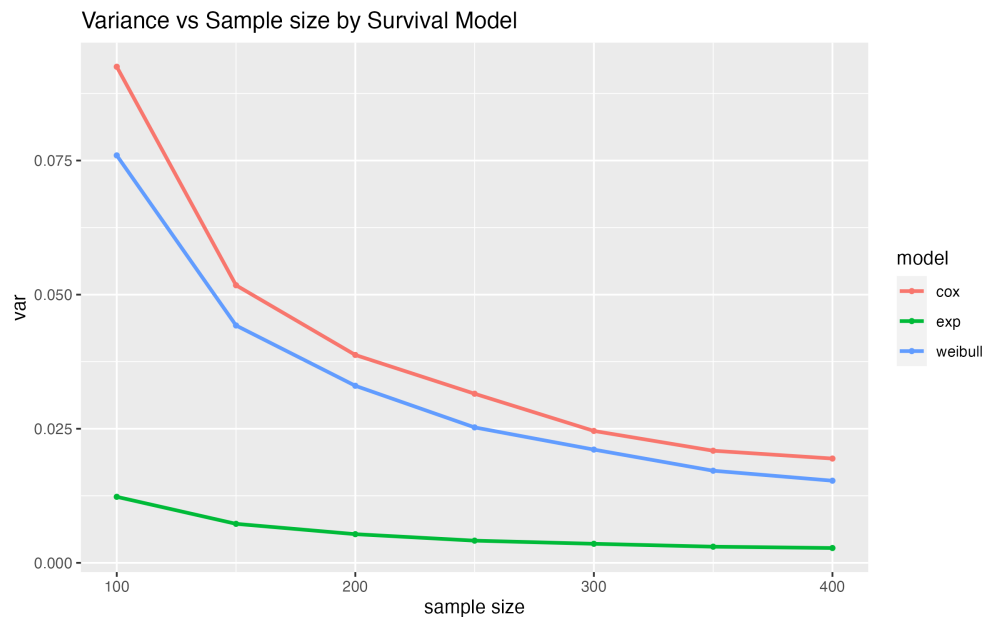


Figure 4

Figure 5-6 shows the bias and MSE from the weibull distribution. When looking at the bias plot, the exponential model has the largest bias, while the cox model and Weibull model both perform well with a low bias. From mse plot, weibull fits the data better than cox model when sample size is less than 300, for more than 300, cox and weibull model has a nearly equivalent prediction accuracy. Thus, weibull and cox both fit data from weibull distribution. Cox model outperforms weibull for large sample sizes.

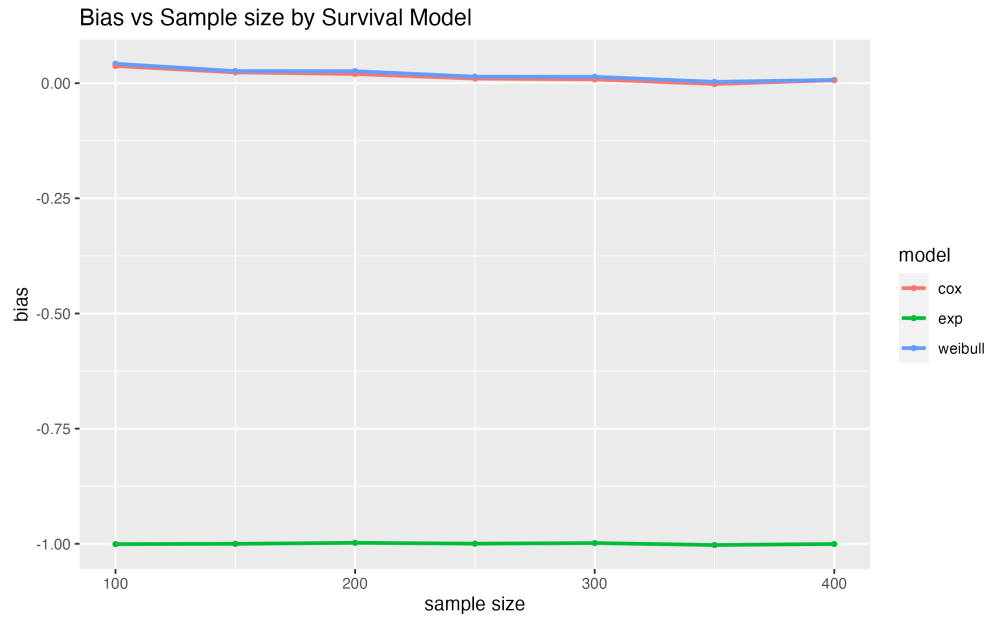


Figure 5

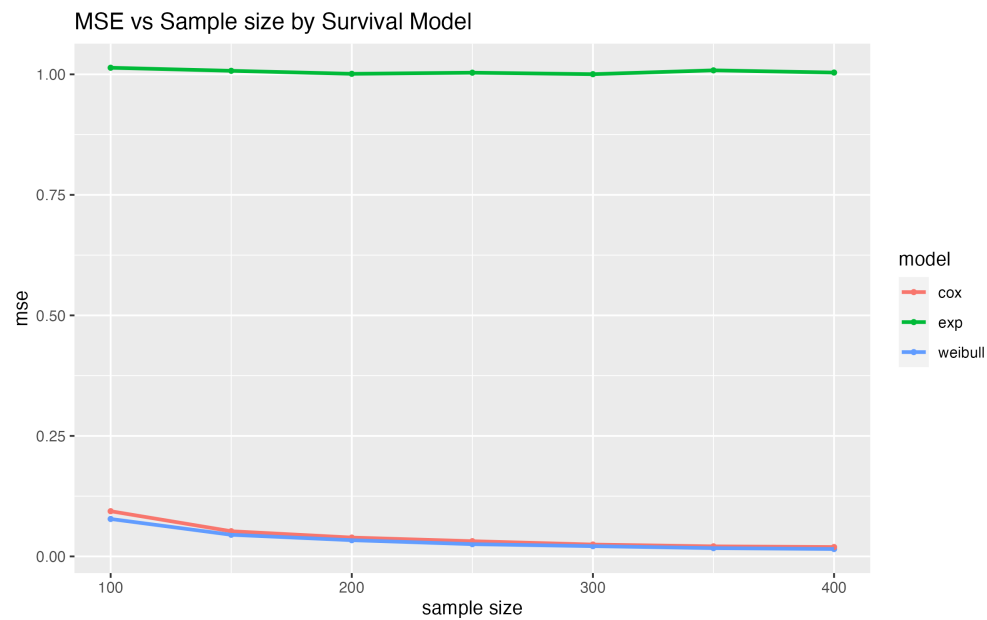


Figure 6

### Scenario 3: Gompertz baseline hazard function

Figure 7 shows the variance from a gomperz distribution. The exponential model shows to have the lowest versatility, while the cox model has the largest.



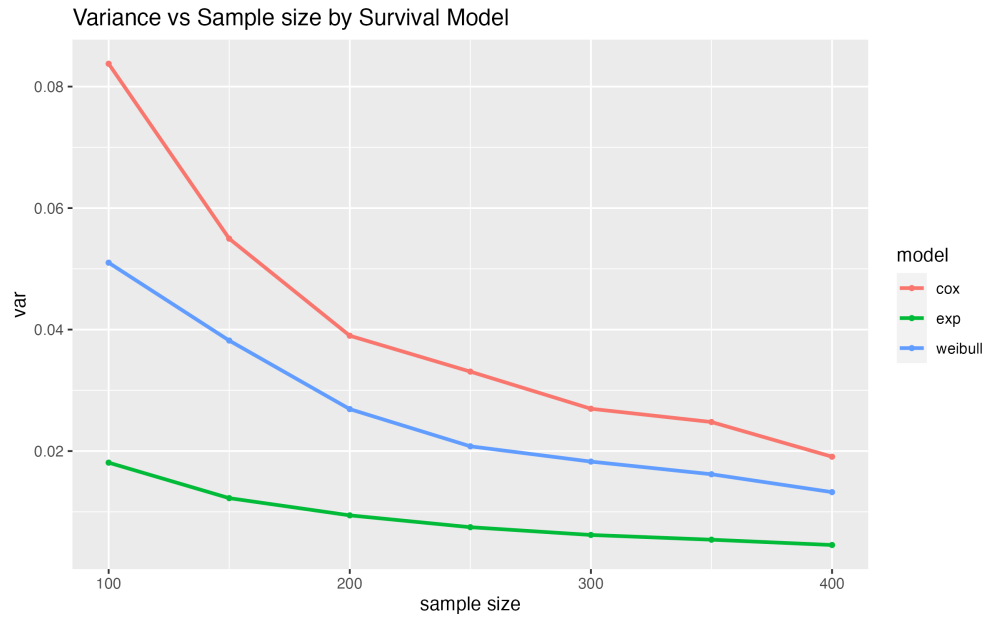


Figure 7

Figure 8-9 shows the bias and MSE from the weibull distribution. The mse and bias has similar trends under the three models, the cox model performs well with the lowest bias and mse, while exponential performs the worst.

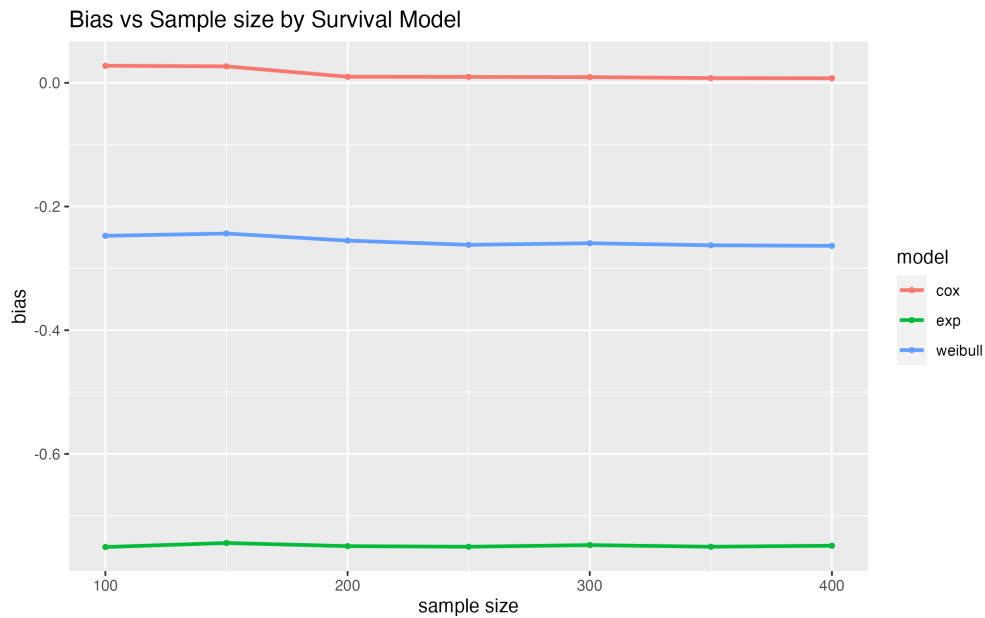


Figure 8

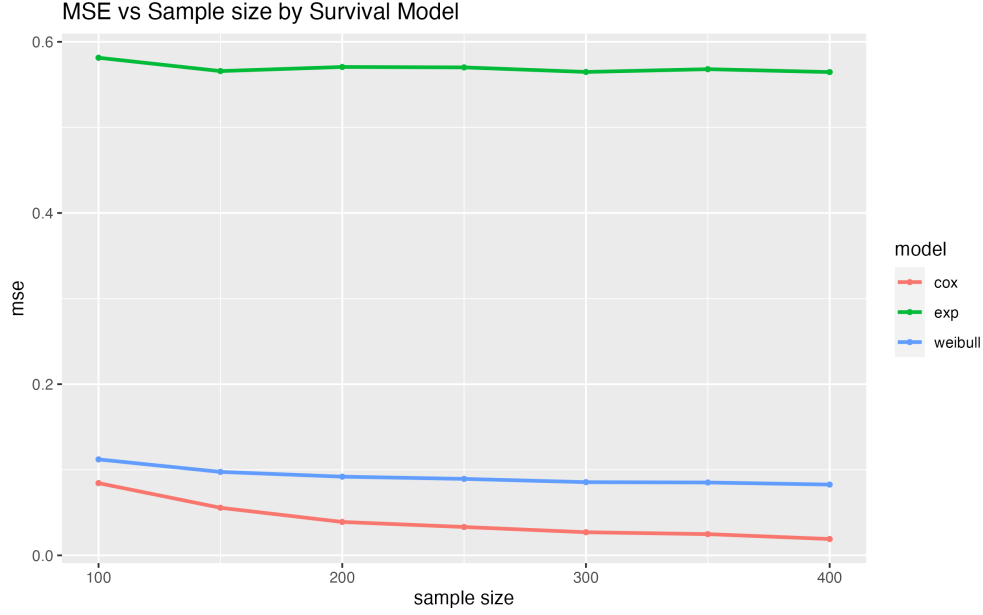


Figure 9

From above scenarios, we found that for an exponential model, the only case it has the highest prediction accuracy is when the true distribution is a purely exponential. The big difference in performances of exponential models in different distributions shows that it is the least robust against a misspecified distribution. On the other hand, the variance plot shows that the exponential model constantly has the lowest versatility. When the distribution is weibull, both the weibull and cox model performs well. However, when the true distribution is gompertz, only the cox model shows to have the best performance.

Table 2 shows the 95% confidence intervals of three models using 1000 simulations with sample size of 400 and  $\beta = 2$ .

Table 2: 95% CI of Estimated Treatment Effect from Three Models

Model/Distribution	Exponential	Weibull	Gompertz
Exponential	(1.80, 2.2)	(0.9, 1.1)	(1.12, 1.39)
Weibull	(1.79, 2.28)	(1.76, 2.29)	(1.53, 1.98)
Cox	(1.76, 2.31)	(1.74, 2.31)	(1.76, 2.31)

With the true baseline hazard as exponential, the exponential proportional model has a narrower confidence interval than the other 2 models, and all contain the true beta. With a weibull baseline hazard, the weibull model provides the most accurate estimation of the treatment effect as expected, with a narrow confidence interval. Cox performs the second best, with a wider confidence interval, but the exponential model significantly deviates from the true treatment and its confidence interval fails to capture the true beta. Under the Gompertz baseline hazard, both exponential and weibull proportional models deviate from the true beta and their 95% confidence intervals fail to include the true beta. The cox model, which makes no assumption about the baseline hazard, fits well.

## Additional Cases

### Mix Distribution

In this case, we tried to mix exponential and weibull distribution together, each distribution constitute half. Similar with the above 3 three scenarios, we simulated the mix distribution for 1000 yimes, sample sizes from 100 to 400, parameters in each model were held constant with  $\lambda = 0.5$ ,  $\gamma = 2$ ,  $\alpha = 2$ ,  $\beta = 2$ .

Figure 10 shows the bias and MSE of the mix distribution. From the bias plot, the exponential model has the

poorest performance, while cox model has a relatively low bias. When looking at the MSE plot, exponential still shows to have a high MSE. Weibull model tend to have a better performace when the sample size is small, when sample size exceed 250, the cox model show tp have better prediction accuracy than weibull model.

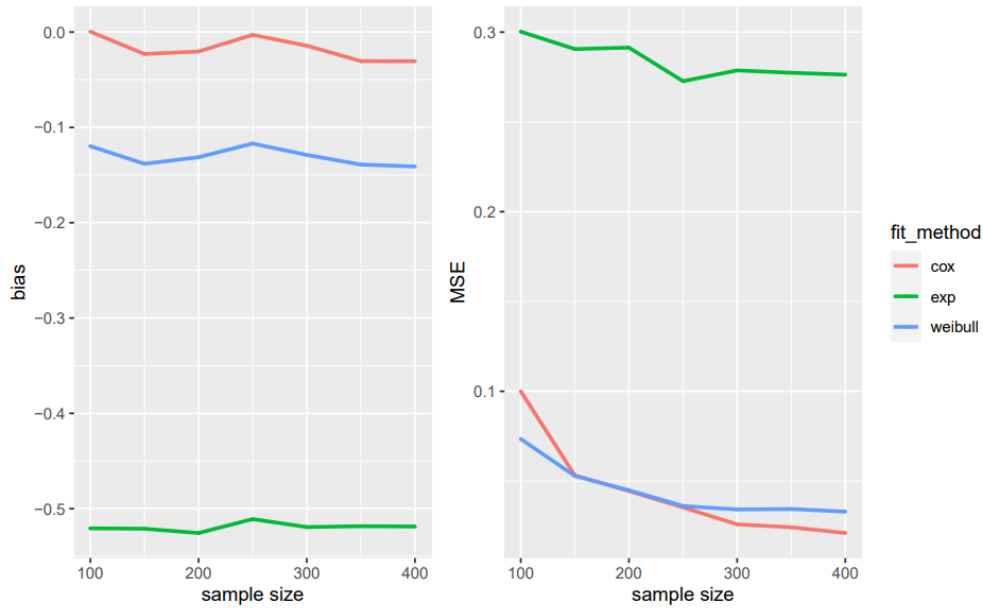


Figure 10

## Conclusion

In conclusion, using the appropriate baseline hazard result in the best performance. The exponential model is the most restrictive one, with only one parameter  $\lambda$ , while the Weibull has an additional  $\gamma$ . The Cox does not specify a certain baseline hazard function, being the most flexible model. Exponential models tend to have high bias and low variance, in other words, they have lower prediction accuracy and higher efficiency when compared to the other two models. In contrast, the cox model is the most robust and flexible one, it can fit to any kinds of distribution, especially when sample size is large. Apparently, the cox model performed well under scenario 2 and 3. However, cox model also has its weakness in the case of an exponential baseline hazard function as we saw in Scenario 1. We noticed that Cox models needed a large sample size (greater than 300) to perform well.

There are still areas that needed for improvements and exploration. Unlike fitted  $\alpha$ ,  $\gamma$ , and  $\beta$ , there is possibilities that the results may vary for different values of the above parameters. Moreover, in our study, we only included one categorical variable, which is the treatment assignment. In real world situations, there also exists other categorical variables, for example, sex that might affect the survival time.

## References

1. Cox, David R. (1972). "Regression Models and Life-Tables". Journal of the Royal Statistical Society, Series B. 34 (2): 187–220.
2. Hosmer, D. W., Lemeshow, S., & May, S. (2008). Applied survival analysis: Regression modeling of time to event data. Wiley.
3. James, Gareth, et al. (2021). An Introduction to Statistical Learning: With Applications in R. Springer.

## Appendix