

# Improved Confidence Intervals for Fixed Term Survival Probabilities in a Small Two-Arm Trial

**Chen Qian**

University of Louisville

**Jianmin Pan**

University of Louisville

**Craig McClain**

University of Louisville

**Shesh Rai** (✉ [shesh.rai@louisville.edu](mailto:shesh.rai@louisville.edu))

University of Louisville

---

## Research Article

**Keywords:** Fixed-term confidence interval, Adjusted covariates, Kaplan-Meier, Cox regression

**Posted Date:** November 24th, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-109278/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 Improved Confidence Intervals for Fixed Term Survival Probabilities in a Small Two-Arm Trial

2 Chen Qian<sup>a,b</sup>, Jianmin Pan<sup>a</sup>, Craig J. McClain<sup>c,d,e,f</sup>, Shesh N. Rai<sup>a,b,e,f\*</sup>

3 <sup>a</sup> *Biostatistics and Bioinformatics Facility, James Graham Brown Cancer Center, University of*  
4 *Louisville, Louisville, Kentucky, 40202, USA*

5 <sup>b</sup> *Department of Biostatistics and Bioinformatics, University of Louisville, Louisville, Kentucky,*  
6 *40202, USA*

7 <sup>c</sup> *Department of Medicine, University of Louisville, Louisville, Kentucky, 40202, USA*

8 <sup>d</sup> *Robley Rex Louisville VAMC, Louisville, Kentucky, 40206, USA*

9 <sup>e</sup> *University of Louisville Alcohol Research Center, University of Louisville, Louisville, Kentucky,*  
10 *40202, USA*

11 <sup>f</sup> *University of Louisville Hepatobiology & Toxicology Center, University of Louisville,*  
12 *Louisville, Kentucky, 40202, USA*

13 <sup>\*</sup> *Correspondence to Dr. Shesh N. Rai, Biostatistics and Bioinformatics Facility, James Graham*  
14 *Brown Cancer Center, University of Louisville, Louisville, Kentucky, 40202, USA*

15 <sup>†</sup> *E-mail: [shesh.rai@louisville.edu](mailto:shesh.rai@louisville.edu)*

16 **Abstract**

17 **Background**

18 The confidence interval for survival probability at a fixed time point provides valuable  
19 information on how the subject performs in terms of survival rate. However, in a two-arm trial  
20 when the sample size in each group is small or when the distribution of events that occurred  
21 within the group is skewed, the confidence interval might become very unstable, and thus may  
22 not provide accurate information for estimating survival rate. In addition, when there are other  
23 covariates available in the dataset, it is important to select those significant variables and include  
24 them in the model. On the other hand, researchers such as physicians who pay more attention to  
25 the final result often analyze the treatment group and control group separately, which may lead to  
26 inaccurate prediction.

27 **Methods**

28 In this study, two treatment groups are combined, and the group indicator variable is considered  
29 as a covariate and is included in the model for computation. Yuan and Rai's adjusted effective  
30 sample size methods are further extended along with Cox proportional hazard model, Weibull  
31 model, and log-logistic model to compute predicted fixed-term overall survival probabilities and  
32 corresponding confidence intervals with other covariates adjusted. Simulations are conducted to  
33 obtain coverage probability.

34 **Results**

35 In a single model, Wilson-Peto provides better confidence intervals than Kaplan-Meier,  
36 especially in the middle and later stages. In addition, AC-Peto produces better coverage  
37 probability at all time points. In a multivariate model, the log-logistic method provides both  
38 better confidence intervals and coverage probability than Cox regression model at all stages.

39 **Conclusions**

40 This paper provides a guideline on how one should correctly analyze survival data with the most  
41 appropriate method. Depending on the dataset, it is important to consider methods other than  
42 traditional Kaplan-Meier and Cox regression models when evaluating survival outcomes.

43 **Keywords:** Fixed-term confidence interval; Adjusted covariates; Kaplan-Meier; Cox regression.

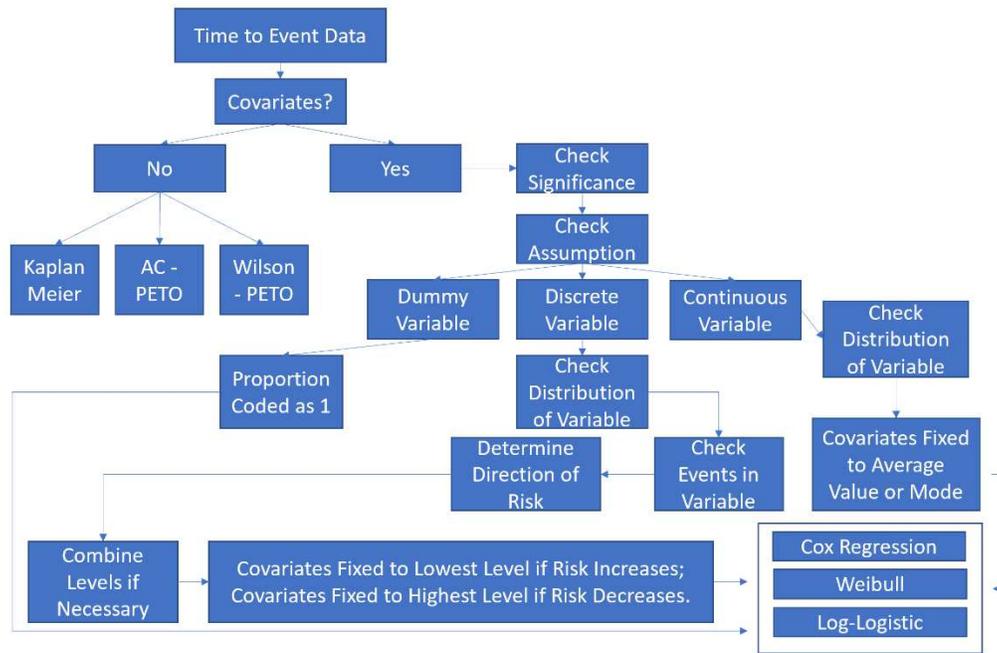
## 44 **Background**

45 When comparing multiple treatments, researchers often want to determine the effectiveness of a  
46 drug by looking at the p-value for the specific statistical test that has been conducted. Their bias  
47 is toward a small p-value so that they could make a successful conclusion on their study.  
48 However, as some articles have pointed out, in clinical research, solely looking at p-value is not  
49 enough to declare success. Alternatively, employing confidence intervals will provide solid  
50 evidence to support the conclusion [2]. The confidence interval can help researchers to decide  
51 whether one treatment is better than the other, which is extremely important for clinical decision  
52 making [3]. The calculation of confidence intervals has been widely adapted in comparing  
53 different treatments because it is more informative than a p-value [4]. The confidence interval is  
54 also highly important in survival analysis. Previous literature suggested the use of the confidence  
55 interval for a binomial proportion to calculate the interval for the survival function [5]. Multiple  
56 articles further examine how to estimate the best confidence interval with appropriate methods in  
57 order to better represent the survival data.

58 When there are several treatments involved in a study, researchers often look at the confidence  
59 intervals among each treatment separately. This allows researchers to make a prediction based on  
60 the observed range of the confidence interval of a certain treatment method. For example, when  
61 there are two treatment groups in a study, one could predict a certain patient's survival  
62 probability given the group to which that patient belongs. However, there will be a problem  
63 when the sample size is small or events are skewed. These conditions happen quite often in  
64 clinical researches due to the cost of trial and difficulty of patient follow up. Often, a trial will  
65 end with a small sample size. Thus, the estimated value which is derived from the data might not  
66 be the true value. Especially when the sample size is small, the true value might be quite  
67 different from the estimated value. The confidence interval provides a range that should contain  
68 the true value. But when the sample size is small, the standard error becomes large, and the  
69 confidence interval will become wider. As a result, the confidence interval will not produce a  
70 precise range to contain the true value. As Henderson and Keiding [6] noted, physicians face  
71 questions from patients for survival prediction after the diagnosis of a terminal disease, but they  
72 often cannot provide accurate survival time prediction since the variation is large among patients.  
73 Henderson further mentioned that the predictive interval could be a way to provide estimation  
74 since the actual survival might fall within the interval. Unfortunately, the predictive interval has  
75 rarely been used because the intervals are too wide. For example, a doctor can tell a patient's  
76 survival probability given a certain condition to be, for example, 50% in treatment A and 60% in  
77 treatment B. However, if the clinical trial information on which this prediction is based is from a  
78 trial with small sample sizes, that estimation is questionable.

79 In this paper, we are going to combine both two treatment groups and label the group indicator  
80 variable as a covariate. We include group covariate in all of our parametric and semi-parametric  
81 models. Our aim is to see if grouping has any impact on the model. We also want to see which  
82 method will provide the best predictive estimation with this improved confidence interval  
83 calculation method. Our overall aim is to provide a guideline on how basic survival data should  
84 be analyzed. Figure 1 below shows a schematic of the process.

85 Figure 1: Flowchart of the Data Analysis.



86 **Data**

87

88 The data used in this paper come from a randomized clinical trial conducted by the Radiation  
 89 Therapy Oncology Group [7]. The dataset is publicly available, and therefore, neither ethical  
 90 approval nor informed consent is needed for our study. The entire trial contains data from 15  
 91 sites with 16 participating institutions; however, in this paper, only the data on three sites with  
 92 the six largest institutions will be used. At the beginning of this study, 193 patients were  
 93 randomly assigned into two treatment groups. Group one (only radiation therapy) has 99 patients  
 94 with 27 censored subjects. Group two (radiation therapy with a chemotherapeutic agent) has 94  
 95 patients with 26 censored subjects. Other variables including *sex*, *age*, *condition*, *T-staging*, and  
 96 *N-staging*. Summary statistics can be found in Table 1.

97 **Methods**

98 Six methods are demonstrated in this paper, including Kaplan-Meier, Agresti-Coull with Peto's  
 99 effective sample size adjusted method, Wilson with Peto's effective sample size adjusted  
 100 method, Cox proportional hazard model, Weibull model, and log-logistic model. In each method,  
 101 fixed-term overall survival probabilities and confidence intervals are calculated, then coverage  
 102 probabilities for each method are obtained through simulation.

103 To have a good understanding of the data, *survdiff* in *R* is used to calculate whether or not there  
 104 is any difference between the two treatment groups. It appears that according to the p-value  
 105 ( $p=0.3$ ), the two treatment groups are not significantly different. In this case, a closer look at the

106 confidence interval becomes necessary. Zhu et. al [8] evaluate several test procedures for  
 107 survival functions comparison when data is interval-censored and the distribution of censoring is  
 108 unequal. Similar approach could be tested for right-censored data in future research.

109 **Kaplan-Meier**

110 Kaplan-Meier curve can be easily produced with the help of *R*. Its confidence interval can also  
 111 be obtained. Note that the default method for *R* in calculating confidence interval is the  
 112 Greenwood (log) method, which can be treated as a Wald confidence interval, and has been  
 113 proven not to be robust regardless the size of the sample [9]. The Kaplan-Meier estimate  $S(t)$  is

$$\hat{S}(t) = \prod_{j=1}^k \left(1 - \frac{d_j}{r_j}\right). \quad (1)$$

114  $d_j$  and  $r_j$  are the number of deaths and the number of patients at risk at time  $t$  respectively. The  
 115 R-code is as follows:

116 `survdif(Surv(Time,Status) ~ Group, data=oropharynx)`  
 117 `km=survfit(Surv(Time,Status)~Group,data= oropharynx)`

118 **Agresti-Coull-Peto**

119 Brown et al. [9] recommend the Agresti-Coull interval when the sample size is greater than 40. It  
 120 is a score interval and appears to be a better way to calculate the confidence interval. Moreover,  
 121 Yuan and Rai further suggest that the combination of Agresti-Coull interval with Peto's adjusted  
 122 effective sample size provides better coverage probability [1]. To construct the AC confidence  
 123 interval, the formula can be written as:

$$\tilde{p} \pm Z_{1-\alpha/2} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{n + Z_{1-\alpha/2}^2}}, \quad (2)$$

124 where  $\tilde{p} = \frac{M + Z_{1-\alpha/2}^2/2}{n + Z_{1-\alpha/2}^2}$  and  $Z_{1-\alpha/2}$  is the critical value at 95% confidence level [10].

125 Here  $M$  is defined as the number of estimated events. Also, the sample size  $n$  needs to be  
 126 adjusted by using Peto's effective sample size.  $n$  will be replaced by  $n_p$  [11].  $n_p$  can be easily  
 127 obtained by

$$n_p = \frac{r_t - d_t}{\hat{S}(t)}. \quad (3)$$

128  $n_p$  is defined as the number of observations that remain at risk at time  $t$  divided by the survival  
 129 probability at  $t$ . Replace  $n$  with  $n_p$  in equation 2 and it will generate the new Peto's adjusted  
 130 confidence interval.

131 The *R* code is as follows:

```

132 kmout=summary(km)           #km from previous output
133 time=kmout$time           #death time
134 risk=kmout$n.risk         #number of risk
135 deathn=kmout$n.event     #number of events
136 surv=kmout$surv          #kaplan-meier estimate of survival function
137 sd=kmout$std.err         #standard error of each S(t) by greenwood formula
138 m=length(surv)
139 #calculate Estimated Effective Sample Size
140 np=(risk-deathn)/surv     #ESS from Peto
141 yp=np*surv
142 a=qnorm(1-0.05/2)
143 padjnp=(yp+0.5*a^2)/(np+a^2)   #adjusted p-hat for AC method using Peto's ESS
144 ci1=numeric(m)
145 ci2=numeric(m)
146 ci1=padjnp-a*sqrt( padjnp*(1-padjnp)/(np+a^2) )
147 ci2=padjnp+a*sqrt( padjnp*(1-padjnp)/(np+a^2) )

```

### 148 **Wilson-Peto**

149 Wilson method takes a similar approach. The variance is estimated differently. It proposes to  
150 provide a shorter confidence interval compare to AC method. It is calculated by using the  
151 following formula:

$$\tilde{p} \pm \frac{Z_{1-\alpha/2}}{1+Z_{1-\alpha/2}^2/n} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{n} + \frac{Z_{1-\alpha/2}^2}{4n^2}}. \quad (4)$$

152 Here, similarly,  $n$  will be replaced by  $n_p$ .

153 For Wilson-Peto's confidence interval, the R code is as follows:

```

154 ci1w<-padjnp-(a/(np+a^2))*sqrt(np*surv*(1-surv)+a^2/4)
155 ci2w<-padjnp+(a/(np+a^2))*sqrt(np*surv*(1-surv)+a^2/4)

```

### 156 **Cox Regression Model**

157 Since there are several covariates in the model, the above models might not best represent the  
158 data, and thus implementing other semi-parametric or parametric methods such as Cox  
159 proportional hazard model is necessary. The first step is to test whether those variables are  
160 significant. Univariate regression tests are used here, and only significant variables will be  
161 included in the multivariate regression. Next, it is important to make sure that the dataset does  
162 not violate the proportional hazard assumption. If the data does not violate the assumption, then  
163 the analysis will turn out to be good. If the data does violate the assumption, then the robust  
164 method should be used. There are three variables in the multivariate regression model including  
165 *T-staging*, *Condition*, and *Group* indicator. The purpose of this paper is to see whether making  
166 *Group* as a covariate will improve the predictive outcome. To look at the impact of *Group* on the  
167 model, the other covariates must be adjusted. Since *T-staging* and *Condition* are both categorical  
168 variables, it is essential to determine which level shall the covariates to be fixed at. Distributions

169 of the variables are checked to see whether they are evenly distributed. Next, the proportion of  
 170 events under each level in each variable will be tested to see whether proportions are similar.  
 171 Besides, based on the information of variable, the direction of risk can be determined. In other  
 172 words, if a variable determines the size of tumor from level 1 to 5 in which 1 means the smallest,  
 173 then one shall know that the direction of risk is increasing from 1 to 5. Note that if the variable is  
 174 not evenly distributed and the proportion of event is also skewed, combining levels become  
 175 necessary. Normally, if a level has the percentage of the number of subjects less than 10%, then  
 176 combining levels will help. After combining levels, the covariate is fixed to the lowest level if  
 177 risk increases, and to the highest level if risk decreases. *R* code as follow:

```
178 res.cox=coxph(Surv(Time,Status)~Group+factor(T)+factor(Condition),data= oropharynx)
179 res.zph=cox.zph(res.cox) #Check assumption
180 print(res.zph)
181 group_df <- with(oropharynx,
182 data.frame(Group = c(1, 2)
183 T = c(3, 3),
184 Condition= c(1,1)
185 )
186 )
187 fit=survfit(res.cox,newdata=group_df)
188 surv_summary(fit)
```

189 Predicted survival probability and confidence interval can be easily obtained from the summary  
 190 output.

### 191 Weibull Model

192 Other than Cox regression, people usually look at parametric methods such as Weibull and log-  
 193 logistic. Those methods are more often recommended when failure time distribution cannot be  
 194 determined by a particular distribution family. Weibull distribution can model multiple types of  
 195 data, regardless of the presence of skew or not. In addition, Weibull distribution can easily model  
 196 hazard function for survival analysis. In *R*, *survreg* can be used and distribution can be specified  
 197 as *weibull* to get scale and intercept values. To start, the cumulative distribution function of  
 198 Weibull is

$$F(x) = P(X \leq x) = 1 - e^{-(1/\alpha)x^\beta}. \quad (5)$$

199 The survival probability at a fixed time point can be carried out as

$$\hat{S}(x) = P(X \geq x) = e^{-\left(\frac{1}{\alpha}\right)x^\beta}, \quad (6)$$

200 where  $\alpha$  is the scale parameter and  $\beta$  is the shape parameter. The confidence interval can also be  
 201 easily calculated. In *R* there is a package called *flexsurv*. Within it, the *flexsurvreg* command can  
 202 easily model Weibull. The same approach that was used with Cox regression to fix variables at  
 203 certain levels can also be used in Weibull model. The *R* code is as follows:

```

204 b=flexsurvreg(Surv(Time,Status)~factor(T)+Group+factor(Condition),data=oropharynx,dist="
205 weibull")
206 group_df <- with(oropharynx,
207     data.frame(Group = c(1, 2),
208               T = c(3, 3),
209               Condition= c(1,1)
210             )
211           )
212 group_df
213 fit=summary(b,newdata=group_df,dist="weibull")
214

```

### 215 Log-Logistic Model

216 Alakus [12] developed a method to calculate confidence intervals for the log-logistic distribution  
217 survival function at any time point. To calculate survival probability, *survreg* can be used and the  
218 distribution can be specified as log-logistic to obtain scale and intercept values. The cumulative  
219 distribution function of Log-logistic is

$$F(x) = P(X \leq x) = \frac{(\lambda x)^\theta}{1 + (\lambda x)^\theta} \quad (6)$$

220 The survival probability at a fixed time point can be carried out as

$$\hat{S}(x) = P(X \geq x) = \frac{1}{1 + (\lambda x)^\theta}, \quad (7)$$

221 where  $\lambda$  is the scale parameter and  $\theta$  is the shape parameter. To calculate confidence intervals,  
222 the interval for score function needs to be formed first, and then it can be extended to the  
223 survival function. The score function can be defined as  $R_i = \log\theta + \theta \log\lambda$ . The estimated  
224 standard error can be defined as  $se(\hat{R}_i) = \{y_i^T Var(\theta) y_i\}^{1/2}$  where  $y_i^T = [1 \ 1]$  is the unit vector  
225 for the model and  $Var(\theta)$  is the variance-covariance matrix that can be obtained from the  
226 program using *vcov* command. The interval can be shown as  $R_i \pm Z_{1-\alpha/2} * se(R_i)$ .

227 Therefore, the confidence interval for the survival function can be written as

$$CI_{low}: \theta[\theta + e^{R_{upp}} x_i^\theta]^{-1}, \quad (8)$$

$$CI_{upp}: \theta[\theta + e^{R_{low}} x_i^\theta]^{-1}. \quad (9)$$

228 Similarly, *flexsurvreg* can be used to obtain survival probability and confidence interval. Same  
229 approach that were used in Cox regression and Weibull to fix variable at certain level can also be  
230 used in log-logistic model. R code as follow:

```

231 c=flexsurvreg(Surv(Time, Status)~ factor(T)+Group+factor(Condition),data=
232 oropharynx,dist="llogis")
233 group_df <- with(oropharynx,

```

```

234     data.frame(Group = c(1, 2),
235               T = c(3, 3),
236               Condition= c(1,1)
237             )
238         )
239     group_df
240     fit1=summary(c,newdata=group_df,dist="llogis")

```

241 For each of the above models, survival probability, as well as the confidence interval, are  
242 obtained at fixed time points such as 3-months, 6-months, 12-months, and 18-months. In  
243 addition, confidence interval length can be calculated by taking the difference between upper and  
244 lower bounds. Bootstrap sampling with 5000 iterations is used to obtain survival probability,  
245 confidence interval, and coverage probability.

## 246 Results

247 Of the 193 patients, approximately 27% are censored (Table 1). The proportion of events within  
248 each group is relatively similar. For all variables in the dataset, none has an association with  
249 treatment group, as p-values are all greater than 0.05. Methods without considering covariates  
250 were tested first. Note that the AC-Peto method and the Wilson-Peto method carry the same  
251 survival probability and are only different in confidence interval. By comparing both methods  
252 with the usual Kaplan-Meier method, it appears that Kaplan-Meier produces higher survival  
253 probabilities at all time points in both treatment groups (Table 2). In group 1, Kaplan-Meier  
254 produces a relatively shorter confidence interval than the other two methods at an early stage  
255 such as 3-months (Table 3). At 6-months, Wilson produces the shortest interval among all three  
256 methods. Toward the later stages, both the AC-Peto and the Wilson-Peto methods produce better  
257 intervals than Kaplan-Meier, and Wilson-Peto is relatively better than AC-Peto. In the later  
258 stage, AC-Peto and Wilson-Peto improve the confidence interval by about 2% from Kaplan-  
259 Meier. Similarly, coverage probabilities also appear to be improved by a small percent for AC-  
260 Peto and Wilson-Peto (Table 4). In this case, AC-Peto has the highest coverage among all  
261 methods. Group 2 seems to carry a similar result to group 1. Kaplan-Meier has the shortest  
262 confidence interval at the early stage, and Wilson-Peto makes the greatest improvement at a later  
263 stage (Table 3). In terms of coverage probability, AC-Peto has better coverage than Kaplan-  
264 Meier but is close to Wilson-Peto in the early stage, whereas, and in the later time, AC-Peto has  
265 the best coverage among all methods (Table 4).

266 To better predict survival outcome, significant covariates must be taken into consideration.  
267 Univariate Cox regression analysis showed that only variable *Condition* ( $p < 0.01$ ) and *T-Staging*  
268 ( $p = 0.01$ ) are significant (Table 5). Therefore, variables that will be included in the multivariate  
269 regression are limited to *Condition*, *T-Staging*, and *Group*. Proportional hazards assumption  
270 analysis showed that the model does not violate the assumption (Global  $p = 0.052$ ) which  
271 indicates that the Cox regression model would be a good fit for analyzing those variables (Table  
272 6). Since *Condition* and *T-Staging* are categorical variables, it is necessary to take a closer look  
273 at their distribution. Under *T-Staging*, the number of subjects at level 1 is only 4.6% of overall  
274 subjects. Typically, combining levels become necessary when a certain group falls below 10% of  
275 the total. Before combining levels, it is extremely important to check the distribution of events

276 among all levels. In this case, the ratios between the number of events and number of subjects  
277 are relatively close among all levels (Table 7). Therefore, combining level 1 and level 2 becomes  
278 reasonable. In the variable *Condition*, the distribution of subjects is skewed wherein level 1 has  
279 around 73% of the total, but level 3 and level 4 have only 3% and 0.5%, respectively (Table 7).  
280 Taking a closer look at the distribution of events, due to the small sample size of level 4, the  
281 proportion will be either 100% or 0% in this case. This situation does not provide much valuable  
282 information, and based on the definition of this variable, people with a higher level of the  
283 condition tend to have a higher risk. It is therefore reasonable to see that the proportion increases  
284 from level 1 to level 3. In this study, level 3 and level 4 are combined for computation. In group  
285 1, log-logistic has the highest survival probabilities at all time points (Table 8). Weibull has  
286 relatively higher survival probabilities than Cox in the later stage but *vice versa* in the early  
287 stage. In terms of confidence intervals, at 3-months, the log-logistic interval is around 15%  
288 shorter than Cox, and about 21% shorter at 6-months (Table 9). In the later stage, Weibull has  
289 roughly 20% shorter intervals compared to Cox intervals. On the other hand, Weibull has the  
290 least coverage probabilities at all tested time points (Table 11). Log-logistic has slightly better  
291 coverage than Cox at 3-months, 6-months, and 12-months, but not at 18-months. In group 2,  
292 Weibull and log-logistic have higher survival probabilities at most time points (Table 8).  
293 Confidence intervals follow the same pattern as group 1 (Table 9). In terms of coverage  
294 probability, log-logistic has the best improvement at all time points (Table 11). Weibull has  
295 better coverage than Cox at earlier stages but becomes worse at later stages.

296 Further comparisons were made between semi/parametric models and AC/Wilson-Peto methods.  
297 As seen in Table 10, in both groups semi-parametric and parametric models produce shorter  
298 confidence intervals than AC/Wilson-Peto methods in earlier stages, but in the long term, basic  
299 models tend to perform better. In terms of coverage, semi-parametric and parametric models  
300 produce better coverage than basic models only at 3-months in group 1 (Table 12). Survival  
301 curves for Kaplan-Meier, AC-Peto, and Cox regression can be found in Figure 2. All methods  
302 compared to Kaplan-Meier can be found in Figure 3 and Figure 4. In most cases, semi-  
303 parametric and parametric models produce shorter confidence intervals in the early stage, but the  
304 pattern does not hold for later stages. Similarly, coverage is higher at early stages for semi-  
305 parametric and parametric methods but becomes worse in the long term.

## 306 **Discussion**

307 This paper illustrates the group effect with other covariates adjusted in survival calculation. This  
308 method can also be expanded to three or more groups. It can further be expanded to determine if  
309 making group as a covariate will benefit large sample size as well. The method provides more  
310 important and more accurate information to researchers as well as clinicians when making a  
311 survival prediction. Note that when distributions of subjects and events are skewed in certain  
312 variables, it is important to determine the best way to combine levels within the variable. There  
313 are many ways to combine levels, such as making it into two blocks. The best way always  
314 depends on the dataset.

315 In this paper, the predictive survival probability is used rather than directly obtaining the result  
316 from data analysis. The reason for doing so is that for a fixed term estimation, it is more

317 accurate. For example, one event occurred at 9-months, and the next event occurred at 13-  
318 months, if we want to find fixed term survival rate at 12-months (1-year), we know that it is the  
319 same between 9-month and 13-month because it is a stepwise function. But the interval, in this  
320 case, is very wide. Suppose there is no event at 12-months, then the survival function is zero, but  
321 the size of the risk set is not. In order to make sense of the data, we need to calculate the  
322 estimated event at 12-months, and then find its corresponding survival probability and  
323 confidence interval.

324 In this paper, only 3-month, 6-month, 12-month, and 18-month survival are tested. The reason  
325 for doing so is that 1-year survival is a normal clinical indicator for many terminal illnesses, thus,  
326 it is the most important term we want to look at. Dong et al. [13] recommend use generalized  
327 inference approach to calculate confidence interval when normality is satisfied, and use bootstrap  
328 percentile approach when assumption being violated as well as when the probability of detecting  
329 early disease stage is large and sample sizes are small. That approach should be considered in  
330 future research.

331 Based on the result above, in a single model, Wilson-Peto provides better confidence interval  
332 than Kaplan-Meier at 12-months and later. In addition, AC-Peto produces better coverage  
333 probability at all time points. In a multivariate model, the log-logistic method provides both  
334 better confidence intervals and coverage probability than Cox regression model at all stages.  
335 Comparing all methods together, long term confidence intervals all become very wide and lose  
336 coverage, and therefore, the confidence interval is better when  $n$  is small.

337 Fay et. al. proposed a non-parametric method called beta product confidence procedure for right-  
338 censored data with independent censoring that should provide better coverage than the basic  
339 Kaplan-Meier method [14]. Further research could compare this method with methods that are  
340 discussed in this paper. Lee et. al also mentioned that the average covariate method might have  
341 some disadvantage, and instead, the corrected group prognostic curve approach was  
342 recommended [15]. This could also be tested in the future.

343 In summary, we have examined six methods for predicting overall survival probabilities and  
344 confidence intervals. Coverage probabilities for each method are obtained through simulation. In  
345 this paper, we combined both two treatment groups and labeled the group indicator variable as a  
346 covariate. We included group covariate in all of our parametric and semi-parametric models. Our  
347 aim was to see if grouping has any impact on the model. We also wanted to see which method  
348 will provide the best predictive estimation with this improved confidence interval calculation  
349 method. Our overall aim is to provide a guideline on how basic survival data should be analyzed.

## 350 **Conclusion**

351 This paper provides detailed guidance on how to deal with survival data. Predicted overall  
352 survival probability and confidence intervals are calculated by making the group a covariate and  
353 looking at group effect with other covariates adjusted. It also compares different methods of  
354 whether or not covariate is present. In a single model, it proves that Wilson-Peto provides better  
355 confidence intervals than Kaplan-Meier, especially in the middle and later stages. Moreover,  
356 AC-Peto produces better coverage probability at all time points. In a multivariate model, the log-

357 logistic method provides both better confidence intervals and coverage probability than Cox  
358 regression model at all stages.

### 359 **List of Abbreviations**

360 **AC:** Agresti-Coull

### 361 **Declarations**

### 362 **Ethics approval and consent to participate**

363 The dataset is publicly available, and therefore, neither ethical approval nor informed consent is  
364 needed for this study.

### 365 **Consent to publish**

366 Not applicable

### 367 **Availability of data and materials**

368 The datasets used and/or analyzed during the current study are available from the corresponding  
369 author on reasonable request.

### 370 **Competing interests**

371 The authors declare that they have no competing interests.

### 372 **Funding**

373 C. Qian was supported by the National Institute of Health grant 5P50 AA024337 to principal  
374 investigator Dr. Craig J. McClain and the University of Louisville Fellowship. S. N. Rai was  
375 partly supported with Wendell Cherry Chair in Clinical Trial Research Fund.

### 376 **Author's contributions**

377 Conception: SNR, CQ, JP

378 Data analysis: CQ

379 Original Draft: SNR, CQ, JP

380 Critical Input: CJM

381 Review & Editing: All Authors

382 Manuscript Revision: All Authors

383 Approval of the Final Version: All Authors

### 384 **Acknowledgements**

385 Not applicable

## 386 References

- 387 [1] Yuan X, Rai SN. Confidence Intervals for Survival Probabilities: A Comparison Study.  
388 *Communications in Statistics - Simulation and Computation*. 2011;40(7):978–991.
- 389 [2] Jirutek MR, Turner JR. In Praise of Confidence Intervals: Much More Informative than P Values  
390 Alone. *The Journal of Clinical Hypertension*. 2016;18(9):955–957.
- 391 [3] Pandis N. Confidence Intervals rather than P Values. *American Journal of Orthodontics & Dentofacial  
392 Orthopedics*. 2013;143(2):293–294.
- 393 [4] du Prel JB, et al. Confidence Interval or P-Value? *Dtsch Arztebl Int*. 2009;106(19):335–339.
- 394 [5] Rothman KJ. Estimation of Confidence Limits for the Cumulative Probability of Survival in Life Table  
395 Analysis. *Journal of Chronic Diseases*. 1978;31:557–560.
- 396 [6] Henderson R, Keidling N. Individual survival time prediction using statistical models. *Journal of Medical  
397 Ethics*. 2005;31:703-706.
- 398 [7] Kalbfleisch JD, Prentice RL. Data Set II: Clinical Trial in the Treatment of Carcinoma of the Oropharynx.  
399 In: *The Statistical Analysis of Failure Time Data Second Edition*; Hoboken, New Jersey. John Wiley &  
400 Sons, Inc.; 2002. p. 379–384; Wiley Series in Probability and Statistics.
- 401 [8] Zhu C, et al. A Nonparametric Test for Interval-Censored Failure Time Data with Unequal  
402 Censoring. *Communications in Statistics: Theory and Methods*. 2008; 37:1895-1904.
- 403 [9] Brown LD, et al. Interval Estimation for a Binomial Proportion (with discussion). *Statistical Science*.  
404 2001;16:101–133.
- 405 [10] Agresti A, Coull BA. Approximate is better than "exact" for interval estimation of binomial  
406 proportions. *American Statistician*. 1998;52:119–126.
- 407 [11] Peto R, et al. Design and Analysis of randomized clinical trials requiring prolonged observation of  
408 each patient. II. Analysis and examples. *British Journal of Cancer*. 1977; 35:1–39.
- 409 [12] Alakus K, Erilli NA. Confidence Intervals Estimation for Survival Function in Log-Logistic Distribution  
410 and Proportional Odds Regression Based on Censored Survival Time Data. *Biometrics & Biostatistics*.  
411 2011;2(3):116.
- 412 [13] Dong T, et al. Parametric and non-parametric confidence intervals of the probability of identifying early  
413 disease stage given sensitivity to full disease and specificity with three ordinal diagnostic groups. *Statistics in  
414 Medicine*. 2011; 30(30):3532-3545.
- 415 [14] Fay MP, et al. Pointwise confidence intervals for a survival distribution with small samples or heavy  
416 censoring. *Biostatistics*. 2013;14(4):723-736.
- 417 [15] Lee J, et al. Covariance adjustments of survival curves based on Cox's proportional hazards regression model.  
418 *Computer Applications in the Biosciences*. 1992;8(1):23-27.

## Appendix

**Table 1.** Summary Statistics for Oropharynx Data.

Variables	Total (N=193)	Treatment 1 (N=99)	2 (N=94)	P Value
Age				0.72
Mean $\pm$ SE	60.19 $\pm$ 0.79	60.47 $\pm$ 1.18	59.89 $\pm$ 1.05	
Sex				0.98
Male (%)	148 (76.7)	76 (76.8)	72 (76.6)	
Female (%)	45 (23.3)	23 (23.2)	22 (23.4)	
Condition				0.09†
1 (%)	143 (74.1)	79 (79.8)	64 (68.1)	
2 (%)	43 (22.3)	19 (19.2)	24 (25.5)	
3 (%)	6 (3.1)	1 (1.0)	5 (5.3)	
4 (%)	1 (0.5)	0 (0.0)	1 (1.1)	
T-Staging				0.72†
1 (%)	9 (4.7)	3 (3.0)	6 (6.4)	
2 (%)	26 (13.5)	14 (14.1)	12 (12.8)	
3 (%)	92 (47.7)	49 (49.5)	43 (45.7)	
4 (%)	66 (34.2)	33 (33.3)	33 (35.1)	
N-Staging				0.92
0 (%)	38 (19.7)	21 (21.2)	17 (18.1)	
1 (%)	28 (14.5)	15 (15.2)	13 (13.8)	
2 (%)	37 (19.2)	19 (19.2)	18 (19.1)	
3 (%)	90 (46.6)	44 (44.4)	46 (48.9)	
Status				0.95
Censored (%)	53 (27.5)	27 (27.3)	26 (27.7)	
Dead (%)	140 (72.5)	72 (72.7)	68 (72.3)	

†Exact test

**Table 2.** Survival Probability and Confidence Interval at Fixed Time Points.

Group 1		3 Months	6 Months	12 Months	18 Months
Kaplan-Meier	S(t)	0.979	0.847	0.632	0.477
	CI Lower	0.955	0.779	0.544	0.388
	CI Upper	0.999	0.921	0.734	0.587
AC – Peto	S(t)	0.964	0.834	0.627	0.478
	CI Lower	0.931	0.763	0.533	0.381
	CI Upper	0.997	0.906	0.720	0.575
Wilson – Peto	S(t)	0.964	0.834	0.627	0.478
	CI Lower	0.935	0.764	0.533	0.381
	CI Upper	0.993	0.905	0.719	0.575
<b>Group 2</b>					
Kaplan-Meier	S(t)	0.947	0.819	0.524	0.357
	CI Lower	0.904	0.745	0.432	0.272
	CI Upper	0.990	0.899	0.636	0.470
AC - Peto	S(t)	0.929	0.807	0.523	0.363
	CI Lower	0.880	0.729	0.423	0.265
	CI Upper	0.979	0.884	0.623	0.462
Wilson - Peto	S(t)	0.929	0.807	0.523	0.363
	CI Lower	0.883	0.730	0.423	0.265
	CI Upper	0.976	0.883	0.623	0.460

**Table 3.** Confidence Interval Length Comparison.

Group 1	3 Months	6 Months	12 Months	18 Months
Kaplan-Meier	0.044	0.141	0.191	0.199
AC - Peto	0.066	0.143	0.187	0.195
Wilson - Peto	0.058	0.141	0.187	0.195
AC vs KM (%)	49.19	1.57	-1.95	-2.02
Wilson vs KM (%)	31.85	-0.15	-2.11	-2.04
Group 2				
Kaplan-Meier	0.086	0.154	0.204	0.199
AC - Peto	0.099	0.155	0.199	0.197
Wilson - Peto	0.093	0.153	0.199	0.197
AC vs KM (%)	15.50	0.58	-1.99	-0.75
Wilson vs KM (%)	8.17	-0.74	-2.02	-0.96

**Table 4.** Coverage Probability Comparison.

Group 1	3 Months	6 Months	12 Months	18 Months
Kaplan-Meier (%)	87.02	93.12	94.66	94.68
AC - Peto (%)	87.02	94.94	94.96	95.34
Wilson - Peto (%)	87.02	93.80	94.96	95.30
AC vs KM (%)	0	1.95	0.32	0.70
Wilson vs KM (%)	0	0.73	0.32	0.65
Group 2				
Kaplan-Meier (%)	87.62	93.10	95.16	95.02
AC - Peto (%)	95.96	95.12	95.48	95.70
Wilson - Peto (%)	95.96	95.12	95.44	95.60
AC vs KM (%)	9.52	2.17	0.33	0.72
Wilson vs KM (%)	9.52	2.17	0.29	0.61

**Table 5.** Univariate Cox Regression Test.

Covariates	Beta	HR (95% CI)	Wald Test	P Value
Sex	0.17	1.2 (0.8-1.8)	0.72	0.39
Age	0.0042	1 (0.99-1)	0.27	0.60
Condition	0.9	2.4 (1.9-3.2)	40	<0.01
Site	-0.04	0.96 (0.84-1.1)	0.34	0.56
T-Staging	0.31	1.4 (1.1-1.7)	7.1	0.01
N-Staging	0.13	1.1 (0.98-1.3)	3.1	0.08

**Table 6.** Test the Proportional Hazards Assumption of a Cox Regression.

Covariates	Rho	Chi-Square	P-Value
Group	-0.0338	0.164	0.69
T-Staging at Level 2	0.1197	2	0.16
T-Staging at Level 3	0.1132	1.82	0.18
T-Staging at Level 4	0.0571	0.467	0.49
Condition at Level 2	-0.2454	8.13	<0.01
Condition at Level 3	-0.0613	0.526	0.47
Condition at Level 4	0.2492	1.42e-07	0.99
Global	NA	1.4	0.052

**Table 7.** Summary Statistics for Variable *T*-Staging and *Condition*.

<i>T</i> -Staging	1	2	3	4
Number of Subjects	9	26	92	66
Number of Events	6	16	64	54
Proportion	67%	62%	70%	82%
Condition	1	2	3	4
Number of Subjects	143	43	6	1
Number of Events	96	38	6	0
Proportion	67%	88%	100%	0%

**Table 8.** Survival Probability and Confidence Interval at Fixed Time Points.

Group 1		3 Months	6 Months	12 Months	18 Months
Cox Regression	S(t)	0.981	0.903	0.715	0.568
	CI Lower	0.965	0.853	0.606	0.436
	CI Upper	0.997	0.957	0.845	0.743
Weibull	S(t)	0.954	0.894	0.773	0.652
	CI Lower	0.922	0.833	0.666	0.512
	CI Upper	0.975	0.937	0.853	0.764
Log-Logistic	S(t)	0.984	0.943	0.821	0.678
	CI Lower	0.966	0.891	0.704	0.524
	CI Upper	0.994	0.973	0.902	0.803
Group 2					
Cox Regression	S(t)	0.980	0.897	0.699	0.546
	CI Lower	0.962	0.844	0.587	0.412
	CI Upper	0.997	0.953	0.832	0.725
Weibull	S(t)	0.952	0.889	0.763	0.637
	CI Lower	0.918	0.825	0.651	0.493
	CI Upper	0.974	0.934	0.847	0.754
Log-Logistic	S(t)	0.980	0.929	0.784	0.624
	CI Lower	0.959	0.869	0.656	0.467
	CI Upper	0.992	0.965	0.876	0.759

**Table 9.** Confidence Interval Length Comparison.

Group 1	3 Months	6 Months	12 Months	18 Months
Cox Regression	0.032	0.104	0.239	0.307
Weibull	0.053	0.104	0.188	0.252
Log-Logistic	0.028	0.082	0.198	0.279
Weibull vs Cox (%)	62.76	-0.19	-21.37	-18.01
Log-logistic vs Cox (%)	-14.83	-21.00	-17.05	-9.00
Group 2				
Cox Regression	0.034	0.109	0.245	0.313
Weibull	0.055	0.109	0.196	0.261
Log-Logistic	0.033	0.096	0.219	0.292
Weibull vs Cox (%)	60.31	0.13	-19.95	-16.54
Log-logistic vs Cox (%)	-5.39	-12.04	-10.44	-6.66

**Table 10.** Confidence Interval Length Comparison with AC-Peto and Wilson-Peto (Change in Percentage).

Group 1	3 Months	6 Months	12 Months	18 Months
Cox vs Wilson	-44.59	-25.92	27.83	57.51
Cox vs AC	-51.03	-27.17	27.64	57.47
Weibull vs Wilson	-9.81	-26.05	0.52	29.13
Weibull vs AC	-20.29	-27.31	0.36	29.10
Log-logistic vs Wilson	-52.81	-41.48	6.05	43.32
Log-logistic vs AC	-58.29	-42.47	5.88	43.29
Group 2				
Cox vs Wilson	-63.10	-28.82	22.67	59.14
Cox vs AC	-65.44	-29.76	22.64	58.81
Weibull vs Wilson	-40.84	-28.73	-1.80	32.81
Weibull vs AC	-44.60	-29.66	-1.83	32.53
Log-logistic vs Wilson	-65.09	-37.40	9.86	48.53
Log-logistic vs AC	-67.31	-38.22	9.83	48.22

**Table 11.** Coverage Probability Comparison.

Group 1	3 Months	6 Months	12 Months	18 Months
Cox (%)	91.14	92.70	93.80	95.16
Weibull (%)	90.28	92.62	92.92	93.30
Log-Logistic (%)	92.06	94.30	94.22	94.08
Weibull vs Cox (%)	-0.94	-0.09	-0.94	-1.95
Log-logistic vs Cox (%)	1.01	1.73	0.45	-1.13
Group 2				
Cox Regression (%)	91.04	93.64	94.58	95.12
Weibull (%)	92.12	94.02	94.14	93.96
Log-Logistic (%)	93.94	96.38	95.88	95.74
Weibull vs Cox (%)	1.19	0.41	-0.47	-1.22
Log-logistic vs Cox (%)	3.19	2.93	1.37	0.65

**Table 12.** Coverage Probability Comparison with AC-Peto and Wilson-Peto (Change in Percentage).

Group 1	3 Months	6 Months	12 Months	18 Months
Cox vs Wilson	4.73	-1.17	-1.22	-0.15
Cox vs AC	4.73	-2.36	-1.22	-0.19
Weibull vs Wilson	3.75	-1.26	-2.15	-2.09
Weibull vs AC	3.75	-2.44	-2.15	-2.14
Log-logistic vs Wilson	5.79	0.53	-0.78	-1.28
Log-logistic vs AC	5.79	-0.67	-0.78	-1.32
Group 2				
Cox vs Wilson	-5.13	-1.56	-0.90	-0.50
Cox vs AC	-5.13	-1.56	-0.94	-0.61
Weibull vs Wilson	-4.00	-1.16	-1.36	-1.72
Weibull vs AC	-4.00	-1.16	-1.40	-1.81
Log-logistic vs Wilson	-2.11	1.32	0.46	0.15
Log-logistic vs AC	-2.11	1.32	0.42	0.04

Figure 2: Survival curve comparison between Kaplan-Meier, AC-Peto, and Cox Regression.

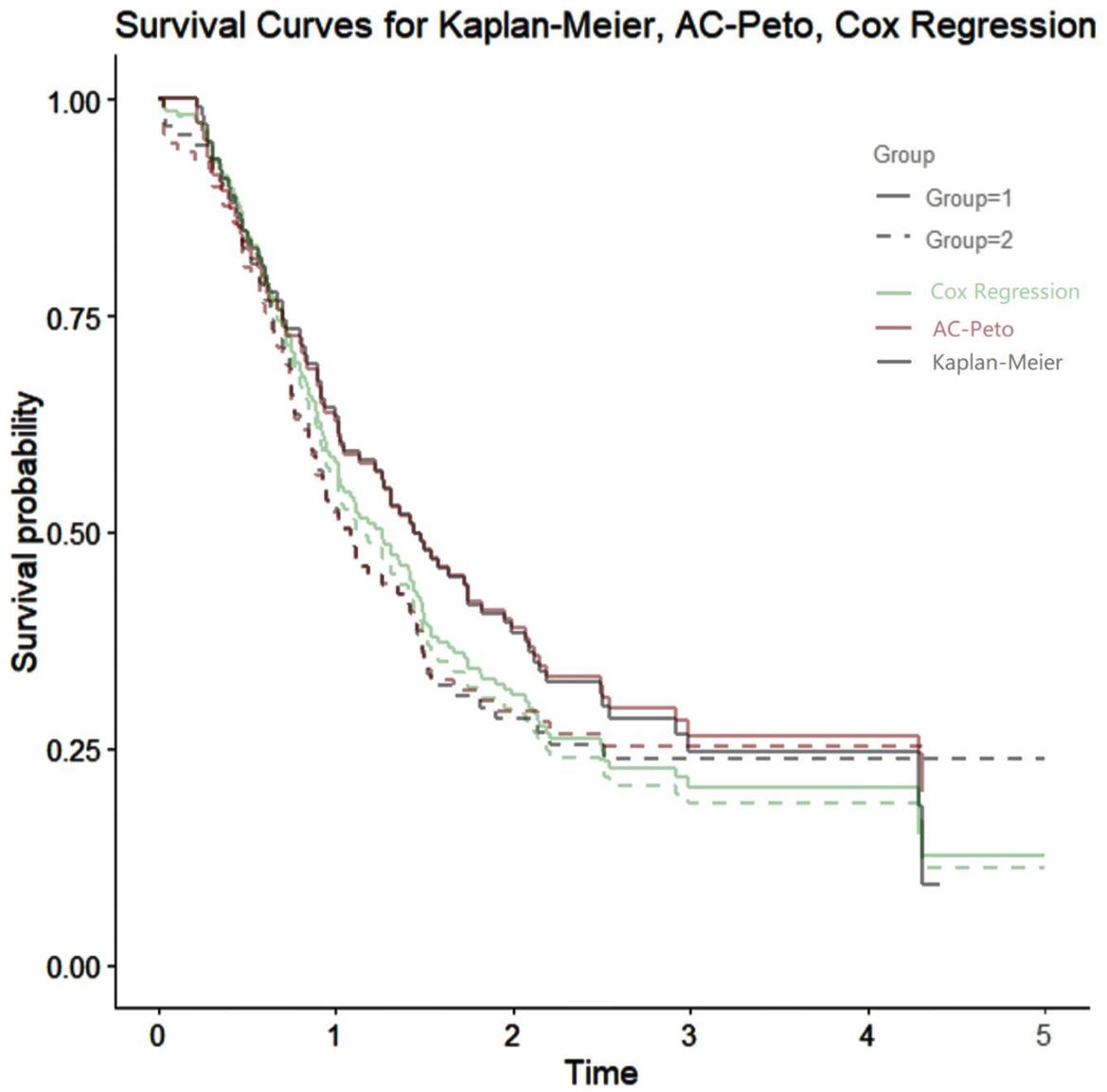


Figure 3: Comparison of confidence interval length changes to the Kaplan-Meier method.

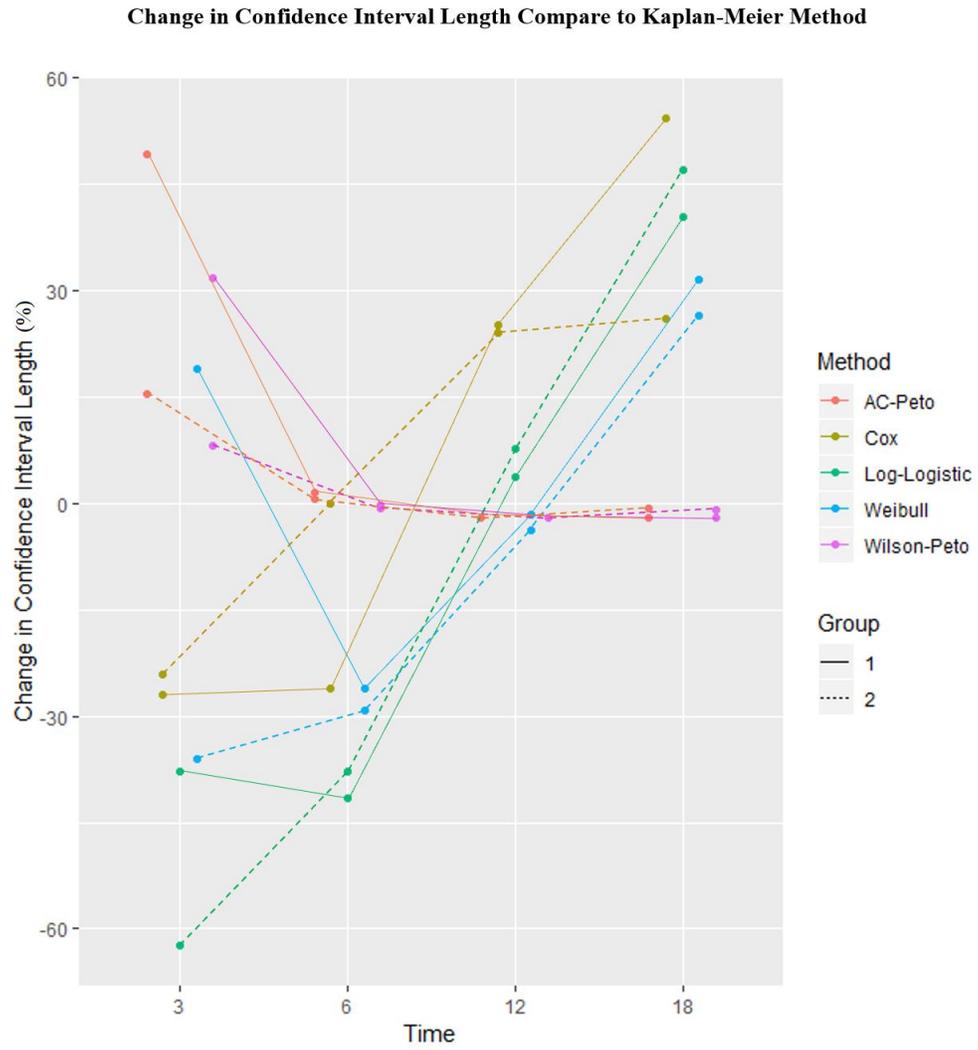
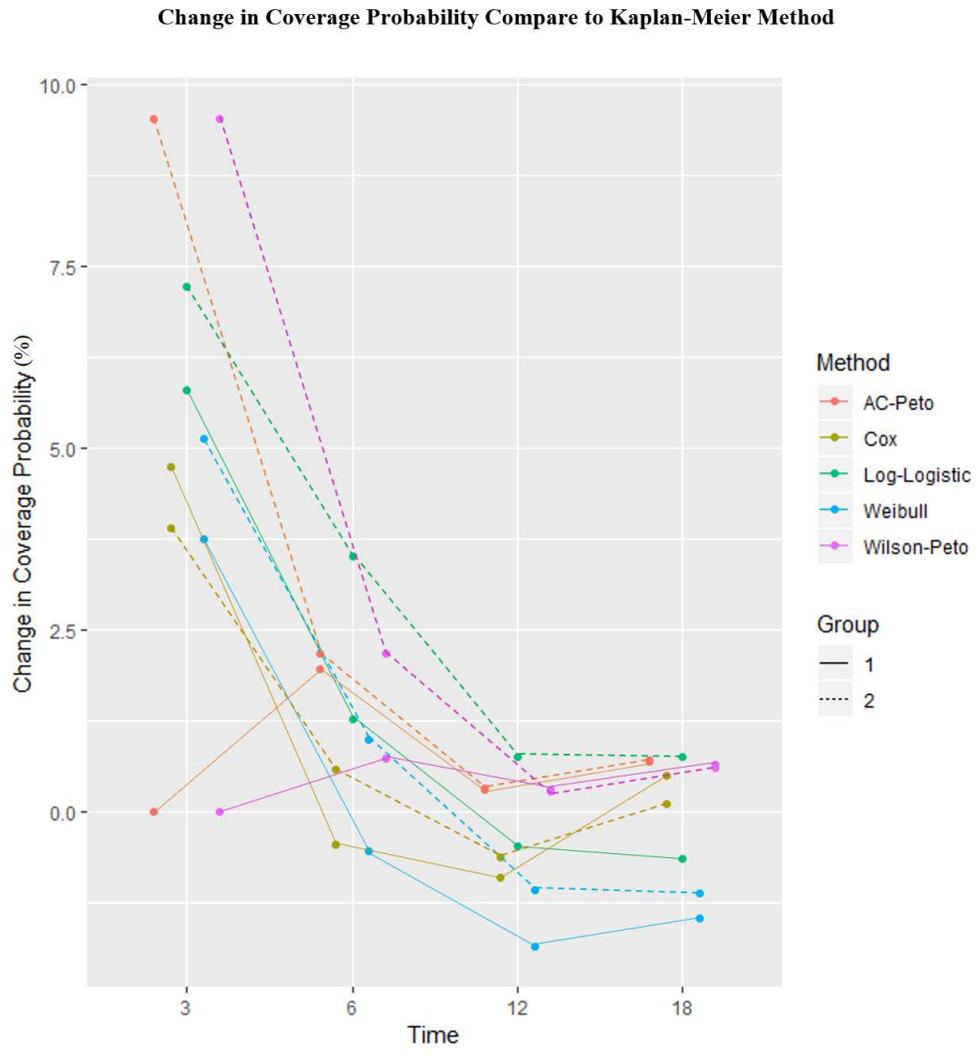


Figure 4: Comparison of coverage probability changes to the Kaplan-Meier method.



# Figures

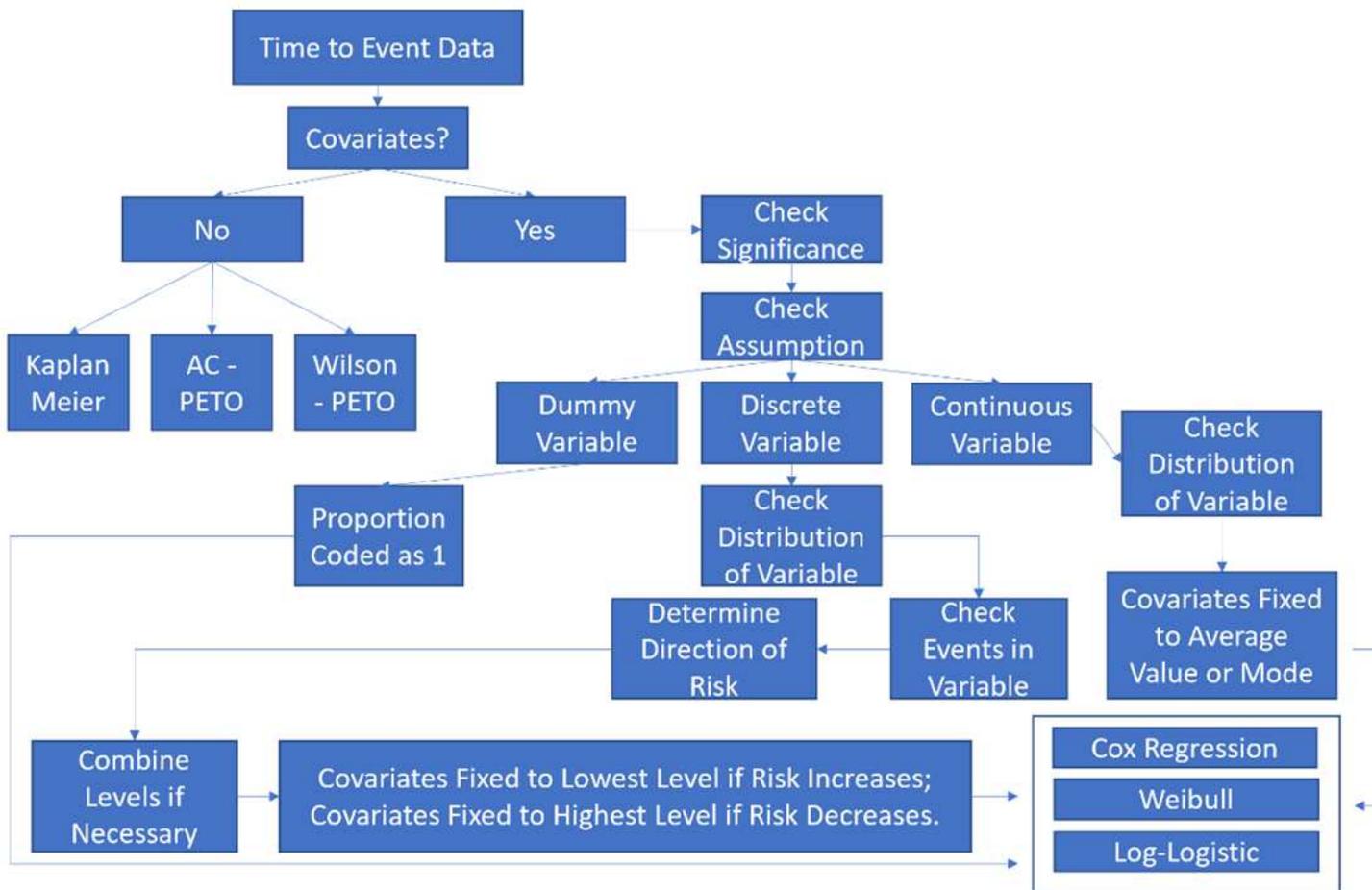


Figure 1

Flowchart of the Data Analysis.

### Survival Curves for Kaplan-Meier, AC-Peto, Cox Regression

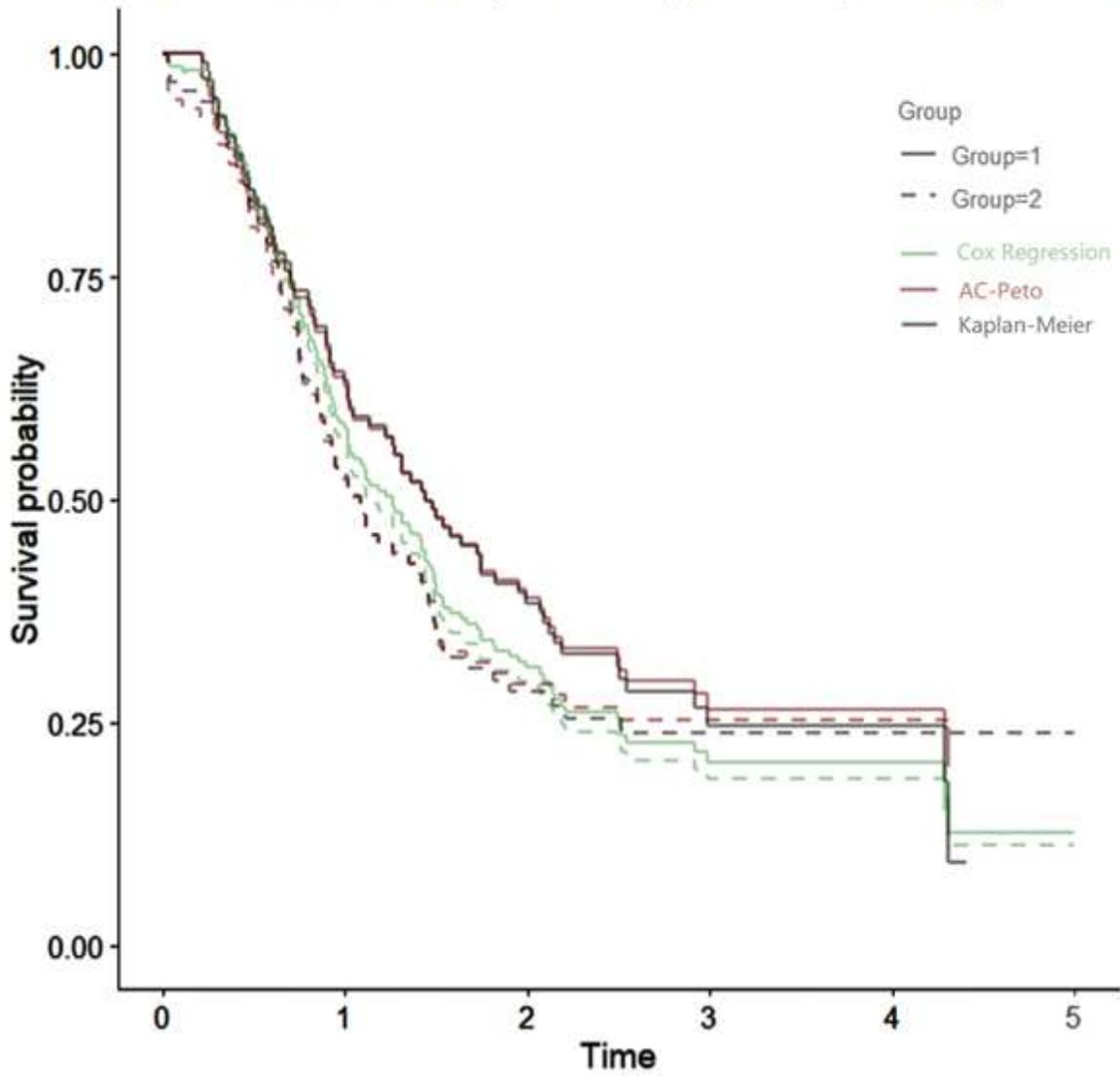


Figure 2

Survival curve comparison between Kaplan-Meier, AC-Peto, and Cox Regression.

### Change in Confidence Interval Length Compare to Kaplan-Meier Method

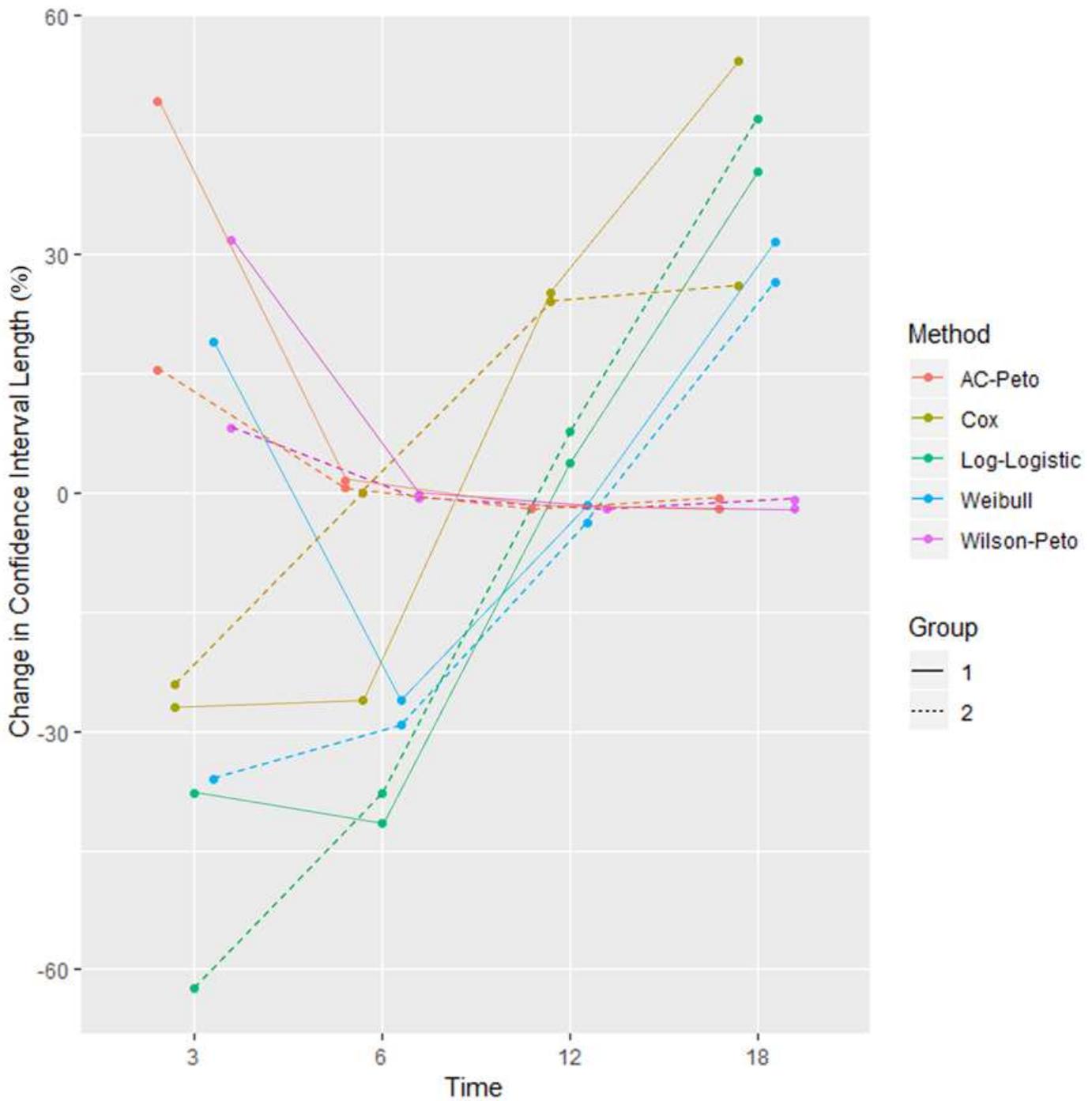


Figure 3

Comparison of confidence interval length changes to the Kaplan-Meier method.

### Change in Coverage Probability Compare to Kaplan-Meier Method

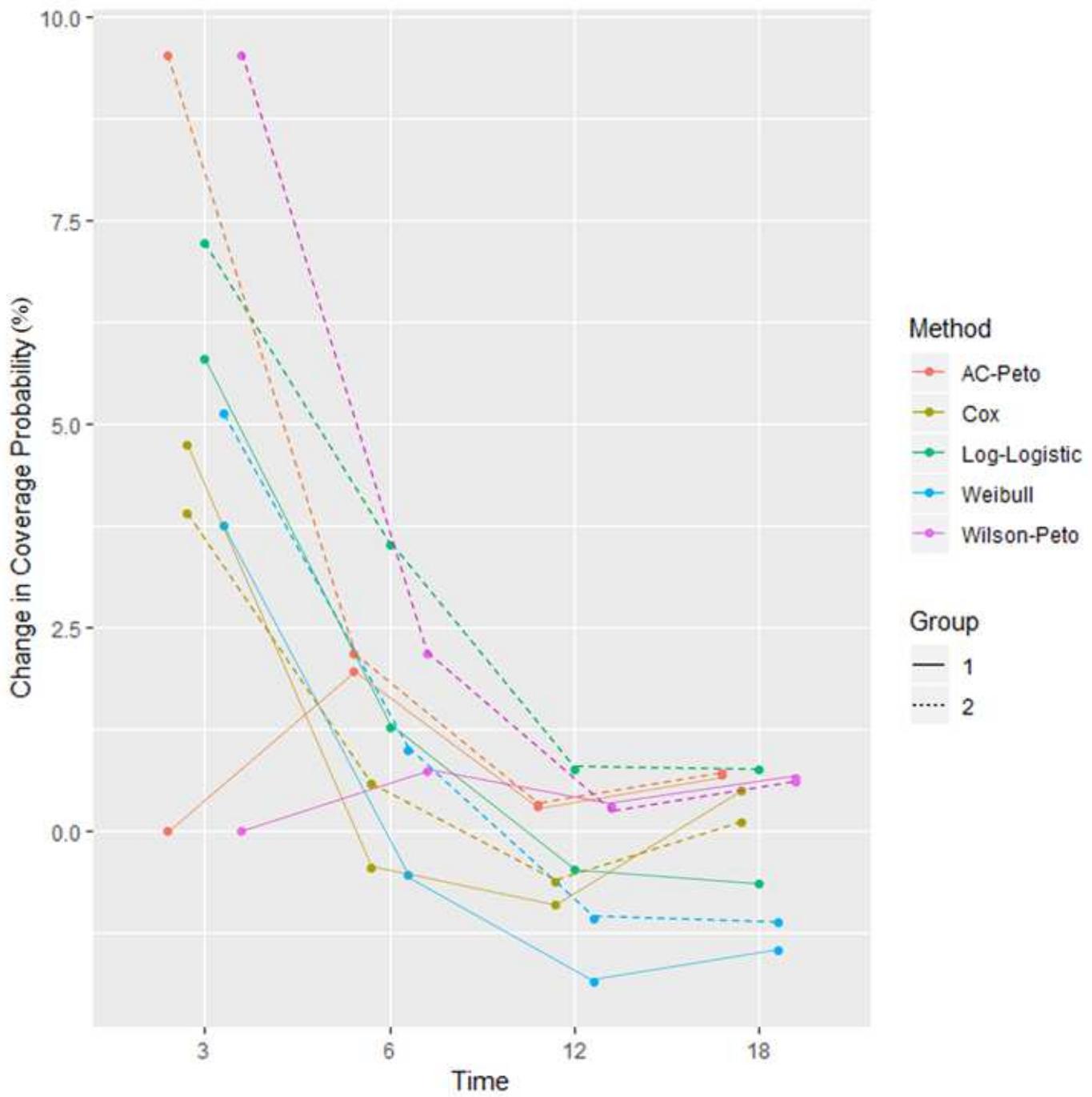


Figure 4

Comparison of coverage probability changes to the Kaplan-Meier method.