

# Imputing Partial Status and Estimating Incidence Rate in an Illness-death Model with Application to a Phase IV Cancer Trial

**Deo Srivastava**

St. Jude Children's Research Hospital

**Jianmin Pan**

University of Louisville

**Chen Qian**

University of Louisville

**Melissa Hudson**

St. Jude Children's Research Hospital

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

University of Louisville

---

## Research Article

**Keywords:** Phase IV clinical trial, Imputation, Cross-sectional survey data, Interval censored data, K-M Method, Missing Value

**Posted Date:** December 2nd, 2020

**DOI:** <https://doi.org/10.21203/rs.3.rs-108160/v1>

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

[Read Full License](#)

---



14 **Abstract**

15 **Background**

16 Phase IV clinical trials are designed to monitor the long-term toxic effects of drugs in cancer  
17 survivors. Evaluations to study the long-term effects of the cancer treatment are often made with  
18 cross-sectional surveys. This leads to interval censored data since the exact time of the onset of  
19 toxicity is not known. In addition to finding prognostic factors for long-term survival outcome,  
20 estimating and comparing the cumulative incidence rates for adverse outcomes of interest for  
21 interval censored data is also desired. However, the analysis of such data is further complicated by  
22 many issues, such as incomplete data, competing risks and selection bias. For example, one such  
23 study was designed by Hudson et al. to study the effect of anthracyclines exposure, received as part  
24 of treatment for childhood cancer, to cardiotoxicity. Rai et al. had utilized a parametric approach for  
25 assessing the effect of anthracycline on the cumulative incidence of cardiotoxicity but excluded the  
26 patients with missing information on the parameters used for assessing cardiotoxicity.

27 **Methods**

28 In this paper our focus is on imputing the missing data and then using the current status regression  
29 methods, previously described in Rai et al. for estimating and comparing cumulative incidence rates  
30 in an illness-death/failure model.

31 **Results**

32 We undertook a comprehensive simulation study to evaluate the performance of our imputation  
33 approach and applied it to a Phase IV clinical trial to evaluate the effect of anthracycline exposure  
34 on long-term cardiotoxicity in childhood cancer survivors, which had missing cardiotoxicity  
35 information.

36 **Conclusions**

37 Our simulations suggest that the results obtained by imputing the missing values using regression  
38 methods are significantly more efficient than those obtained without imputation. The proposed  
39 approach is easy to implement, and we demonstrate its usefulness by applying it to the data reported  
40 in Rai et al. and compare the results reported there to our approach that utilizes imputation.

41 **Keywords:** Phase IV clinical trial; Imputation; Cross-sectional survey data; Interval censored data;  
42 K-M Method; Missing Value.

## 43 **Background**

44 As a result of more modern therapies and better supportive care the 5-year survival rate for  
45 childhood cancer has improved significantly and currently exceeds 80%<sup>3-6</sup>. In 2011 there were  
46 390,000 survivors of childhood cancer living in the US and it is expected that by 2020 there will be  
47 more than 500,000 childhood cancer survivors<sup>7</sup>. This improvement in survival rate comes at a price  
48 as these survivors are at an elevated risk of experiencing long-term morbidity and early mortality as  
49 a result of their cancer and its treatment. The purpose of Phase IV clinical trials is to monitor long-  
50 term sequela and develop interventions to mitigate their effect in long-term. A chemotherapy agent  
51 Anthracycline has served as the backbone for many pediatric malignancies because of its  
52 therapeutic effects but it is also well known to be cardiotoxic<sup>8-9</sup>. One such study was undertaken by  
53 Hudson et al.<sup>1</sup> to evaluate the effect of anthracycline exposure on cardiotoxicity using non-invasive  
54 modalities.

55 Cardiotoxicity is the occurrence of heart electrophysiology dysfunction or/and muscle damage. The  
56 heart becomes weaker and is not as efficient in pumping, and therefore, circulating blood. There are  
57 many measures of electrophysiology dysfunction or/and muscle damage, including shortening  
58 fraction, afterload, QTc interval, and ejection fraction<sup>1, 9, 10</sup>. It is not economic and feasible to  
59 evaluate patients very frequently to estimate the onset time of cardiotoxicity, and hence estimate the  
60 incidence rates. Usually patients are followed longitudinally in the clinics, but not all follow a  
61 routine pattern. Therefore, it is convenient to design cross-sectional surveys for estimating the effect  
62 of long-term side effect of treatments and its predictors. We only know the current status of the  
63 patient with onset prior to current status but not the actual onset times of these events. These types  
64 of incomplete data are referred to as interval censored data since the actual onset time of the events  
65 are unknown<sup>2, 8-11</sup> and our interest is in estimating the onset rate or the cumulative incidence rate.

66 Nonparametric procedures for analyzing interval censored failure time data have been extensively  
67 studied and discussed in the literature.<sup>11-15</sup> Another issue in the cross-section survey study is that  
68 results need to be generalized to the specific population. There can be competing toxic effects from  
69 the same drug. Sun<sup>11</sup> provides an extensive survey of non-parametric methods of estimation using  
70 EM algorithm in studies involving interval censored data. In this paper, we have the same interest,  
71 as Rai et al.<sup>2</sup>, in estimating the cumulative incidence rates in a parametric setting but focus on  
72 improving the accuracy by imputing the missing observations using multivariable regression  
73 method.

74 In practice, most investigators exclude observations with missing values and incomplete cases.  
75 While using only complete cases has its simplicity, one may lose the important information in the  
76 incomplete cases and ignore the possible systematic differences between the complete and  
77 incomplete cases. Hence, the resulting inference may not be applicable to the population of all  
78 cases, especially with a smaller number of complete cases. It is well known that imputation is a  
79 widely used method for handling missing data. Little and Rubin<sup>16</sup> and Buuren<sup>17</sup> provide an excellent  
80 overview of the methods for conducting analyses with missing data. For further information on  
81 multiple imputations see Rubin<sup>18-19</sup>, and Rubin, Stern, and Vehovar<sup>20</sup> discuss imputation of missing

82 discrete data. King et al.<sup>21</sup> review many of the practical costs and benefits of multiple imputations.  
83 For routine imputation of missing data, Schafer<sup>22</sup> presents a method based on multivariate normal  
84 distribution. Liu<sup>23</sup> uses the  $t$  distribution, and Van Buuren, Boshuizen, and Knook<sup>24</sup> use  
85 interlocking regressions. Furthermore, Troxel, Ma, and Heitjan<sup>25</sup> present a method to study the  
86 sensitivity of inferences to missing-data assumptions.

87 This paper is organized as follows. In the following subsection, we provide the details of the  
88 motivating example to introduce the problem. In Methods section, we give a brief description of the  
89 procedure introduced in Rai, et al.<sup>2</sup> and construct corresponding likelihood function. The data from  
90 the motivation example is analyzed by imputing the missing values and compared with the results  
91 obtained without imputation in the Results section. An extensive simulation experiment to study the  
92 performance of the imputation approach are summarized as well. The Discussion section is devoted  
93 to miscellaneous remarks.

#### 94 *Motivation Example*

95 A study was undertaken by Hudson et al.<sup>1</sup> to evaluate the effect of anthracycline on cardiotoxicity  
96 using 12-lead ECG and echocardiography, non-invasive technique. For the study, the cancer  
97 survivors were recruited from St. Jude Children's Research Hospital After Completion of Therapy  
98 Clinic. The survivors were classified into two groups; the first group of survivors consisted of  
99 survivors that received cardiotoxic therapy (anthracycline and/or thoracic radiation) and the other  
100 group did not receive cardiotoxic therapy (no anthracycline nor thoracic radiation). The details of  
101 the study can be found in Hudson et al.<sup>1</sup>. The study was approved by the institutional review boards  
102 at St Jude Children's Research Hospital and Stanford University (Stanford, CA). All study  
103 participants or their parents provided informed consent. Survivors with cardiotoxic therapy were  
104 designated as At-Risk (AR) and those without cardiotoxic therapy as Not At-Risk (NR).

105 Onetime clinical assessment was made by the primary oncologists to identify the survivors with  
106 signs of heart failure using New York Heart Association classification. Along with the clinical  
107 evaluation non-invasive testing based on 12-lead ECG and echocardiography within 24 hours of the  
108 clinical assessment.

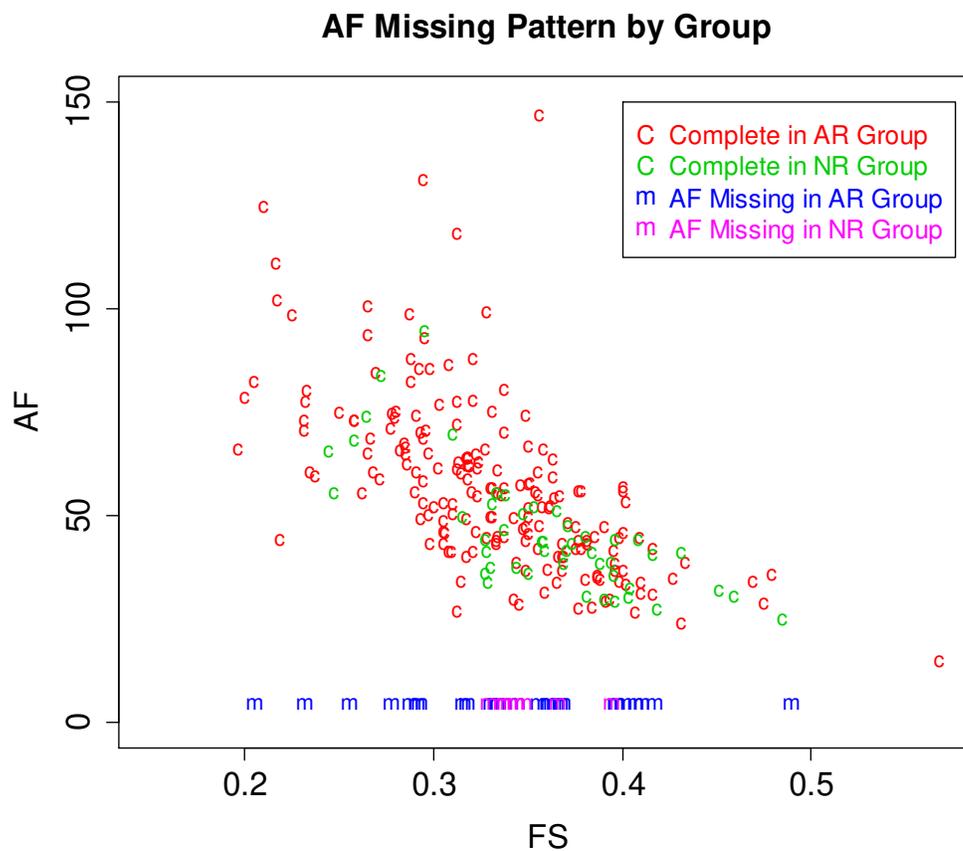
109 When the cardiac measures, discussed above, are out of the normal range, patients are declared to  
110 have cardiotoxicity. Following Hudson et al.<sup>1</sup>, we consider two outcome measures, fractional  
111 shortening ( $FS$ ) and afterload ( $AF$ ). The main measure is defined as  $FS = (LVEdD - LVEsD) / LVEdD$ ,  
112 where  $LVEdD$  is the left ventricular end-diastolic diameter,  $LVEsD$  is the left ventricular end-  
113 systolic diameter. The other measure  $AF$  can be described as the pressure that the chamber of the  
114 heart has to generate in order to eject blood out of the chamber. The study was planned to enroll  
115 almost equal number of patients from each group to detect a medium effect size<sup>26</sup> increase in mean  
116  $AF$  (or decrease in mean  $FS$ ) at  $\alpha = 0.05$  and  $\beta = 0.20$ , without adjusting for multiple outcomes or  
117 multiple comparisons.

118 A short summary about the cohort is summarized in Table 1. Further description is available in  
119 Hudson et al.<sup>1</sup> On closer examination of the data it was seen that there were missing data in  
120 outcome measures  $AF$  and  $FS$ . A total of 40 survivors had missing  $AF$  values and 6 of those were  
121 also missing  $FS$  values. To keep the discussion simple and straight forward these 6 observations

122 were deleted from our analysis and we focused our attention on imputing 34 missing *AF*  
 123 observations. Although, the imputation approach discussed could easily be applied to data missing  
 124 in several variables in a recursive manner. The scatter plot of *AF* and *FS* displayed in Figure 1, does  
 125 not show any clear missing pattern; the missing values of *AF* are in the entire range of values of *FS*.  
 126 Also note that the missing proportion of *AF* values in the NR ( $7/54 = 13\%$ ) and AR ( $27/218 = 12\%$ )  
 127 were almost similar.

<b>Table 1. Characteristics of 278 Patients Enrolled Onto the Noninvasive Cardiac Study</b>						
	<b>At-Risk Group (n=223)</b>		<b>Not At-Risk Group (n=55)</b>		<b>Total (n=278)</b>	
	<b>N</b>	<b>%</b>	<b>N</b>	<b>%</b>	<b>N</b>	<b>%</b>
<b>Demographics</b>						
<b>Sex</b>						
Male	108	51.6	31	56.4	139	50
Female	115	48.4	24	43.6	139	50
<b>Treatment group</b>						
Anthracycline	157	70.4	0	0	157	56.5
Anthracycline + Radiation	60	26.9	0	0	60	21.6
Radiation	6	2.7	0	0	6	2.1
None	0	0	55	100	55	19.8
<b>Race/ethnicity</b>						
White	183	82.1	44	80.0	227	81.7
Black	30	13.5	11	20.0	41	14.7
Other	10	4.4	0	0	10	3.6
<b>Diagnosis</b>						
Leukemia	67	30.0	10	18.2	77	27.7
Sarcomas	60	26.9	14	25.4	74	26.6
Lymphoma	54	34.2	2	3.6	56	20.1
Embryonal tumors	42	18.8	29	52.7	71	25.6
<b>Age at Cancer Diagnosis, Years</b>						

N	223	55	278
Mean	7.37	5.77	7.05
Median	5.46	3.11	4.68
Range	0.01-23.56	0.29-20.06	0.01-23.56
<b>Afterload</b>			
N	191	47	238
Mean	57.50	45.73	55.18
Median	55.43	42.18	51.88
Range	15.38-147.32	25.66-95.02	15.38-147.32
<b>Fractional Shortening</b>			
N	218	54	272
Mean	0.33	0.36	0.34
Median	0.33	0.36	0.34
Range	0.20-0.57	0.24-0.49	0.20-0.57



128

129 **Figure 1:** Scatter plot of *AF* and *FS* within AR and NR group of patients. Complete data (labeled c)  
 130 displays strong correlation between *AF* and *FS*. Missing values of *AF* are in almost entire range of  
 131 *FS* values (labeled m).

132 Using actual measures of these dependent variables *FS* and *AF*, the threshold values were used to  
 133 classify patients as abnormal if ( $FS < 0.28$ ) or ( $AF > 74 \text{ g/cm}^2$ ). Threshold values for *FS* and *AF* were  
 134 determined based on published normative data<sup>27-28</sup>; these are well accepted norms. Let *AFS* and  
 135 *AAF* denote the indicators of these abnormalities. The 278 patients participated in the study; 223  
 136 were designated AR and 55 were designated NR based on treatment. Data on each individual also  
 137 included demographics, date of cancer diagnosis, time since treatment completion, disease related  
 138 variables (such as type, histology, and stage of cancer), treatment related variables (such as  
 139 chemotherapy drugs, doses and irradiation). In the AR group, noninvasive assessment identified  
 140 subclinical dysfunction with *FS* in 37 (13.6%) of 272 and *AF* in 33 (13.9%) of 238; prolonged QTc  
 141 interval in 11 (4.0%) of 273. These are the estimates of prevalence of cardiac abnormalities. Among  
 142 others, one main objective of the study is to estimate cumulative incidence rates of *AFS* and *AAF*.

143 In this study echocardiography was performed as a research measure and not in response to clinical  
 144 symptoms. Individuals with previously established cardiac disease were excluded from  
 145 participation. As formally assessed by New York Heart Association classification, none of the study  
 146 participants reported clinical symptoms of cardiac dysfunction at enrollment. The imaging quality

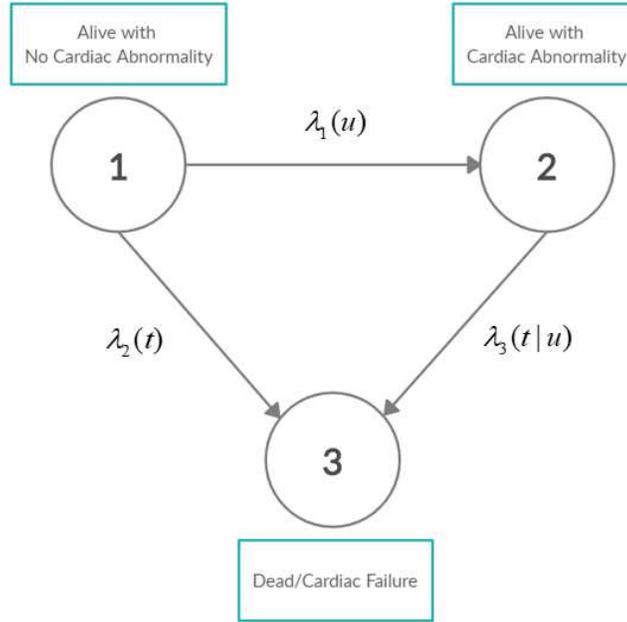
147 in echocardiography is dependent on obtaining a clear acoustic window from which ventricular  
148 volumes are estimated based on geometric assumptions. Operator experience and variations in  
149 thoracic structures can contribute to difficulties in obtaining technically satisfactory data in a given  
150 study. These factors randomly contributed to missing data among study participants. *AF* is not a  
151 standardly used assessment in clinical practice, thus, despite training of ultrasonographers for this  
152 study, this factor may have contributed to a higher prevalence of missing *AF* measurements  
153 compared to *FS*. In other words, *AF* is either under detected or over detected, but operators might  
154 miss it, and therefore, causing missing values. Thus, we feel that there is no selection bias related to  
155 those with and without abnormal *FS* and *AF* identified as part of the study. The missing values of  
156 *AF* are displayed in Figure 1, also do not show any pattern.

157 A crude approach to estimating the incidence rates and obtaining confidence intervals is to apply  
158 the Kaplan-Meier estimator with the assumption of the evaluation time as the onset time. Then,  
159 incidence rate of each type of toxicity is estimated. Some of these toxicity measures could be  
160 missing. In this paper, we impute the missing measurements using regression method first and then  
161 use a parametric approach to estimate incidence rates of specific toxicity. For the cardiotoxicity data  
162 the 34 missing *AF* values were imputed using a multivariable regression with *FS* and other  
163 covariates, such as age, diagnosis, risk status, BMI as predictors. The estimates of cumulative  
164 incidence rates were derived based on the data after imputation and then compared with those  
165 derived in Rai, et al.<sup>2</sup> without imputation.

## 166 **Methods**

### 167 **Cardio-Measures Abnormality Model**

168 The descriptive statistics of all participants can be found in Table 1. In Hudson et al.,<sup>1</sup> the subjects  
169 who died or had cardiac failure during the treatment or during the follow-up after completion of  
170 therapy were excluded because the number of deaths at the time of the analysis were too few.  
171 However, this information is available from the medical record abstraction and with longer follow-  
172 up the number of deaths would increase. Hence, we present the general theory here for a cross-  
173 sectional data with indicators of cardiac abnormality and death/cardiac failure, and time since the  
174 treatment to the survey or the death/cardiac failure, as depicted in Figure 2. We also assume cardiac  
175 abnormality is the precursor for cardiac failure.



176

177 **Figure 2:** An abnormal cardiac measure-death/cardiac failure model involving three states. State 1  
 178 corresponds to patients who are alive with normal value. Patients who are alive with abnormal value  
 179 are in state 2. State 3 is an absorbing state and corresponds to death or cardiac failure.

180 Let stochastic process  $\{X(t)\}$  identify the state occupied by a patient at time  $t$ . For simplicity, we  
 181 suppose that  $n$  patients in state 1 at time  $t = 0$  are those who are identified with different disease  
 182 groups and are planned for treatment, where we have assumed that no patient has cardiac abnormality  
 183 at time  $t = 0$ . Let the random variable  $T$  denote the observation time (survey, death/cardiac failure)  
 184 from the study evaluation and  $U$  the time of *AFS* or *AAF* from the study evaluation. Thus, at any  
 185 time  $t$ ,  $X(t) = 1$ ,  $X(t) = 2$  and  $X(t) = 3$  indicate the patient alive with normal cardiac measure,  
 186 alive with abnormal cardiac measure and died or had cardiac failure with or without cardiac  
 187 abnormality, respectively. We also assume that the development of abnormality is an irreversible  
 188 event without the treatment for cardiotoxicity, and therefore, transitions from state 2 to state 1 do  
 189 not occur, as illustrated in Figure 2. According to practice in this study, the patients are chosen for  
 190 survey independent of their health status, which ensures that the survey results can be regarded as  
 191 independent of the times of the events of interest. The intensities  $\lambda_1(u)$ ,  $\lambda_2(t)$  and  $\lambda_3(t|u)$ , shown  
 192 in Figure 2, are corresponding transitions rates, where  $t$  is the observation time and  $u$  is the time of  
 193 *AAF* or *AFS*.

194 The survival function and the cardiac abnormality prevalence function are derived in Rai et al.<sup>2</sup> as  
 195 follows

196 
$$S(t) = Q(t) + \int_0^t \lambda_1(u)Q(u)Q_3(t|u)du \quad \text{and} \quad \pi(t) = \frac{\int_0^t \lambda_1(u)Q(u)Q_3(t|u)du}{S(t)},$$

197 where  $Q(t) = Q_1(t)Q_2(t)$  is the probability that the time to the first event— alive with abnormal  
 198 value or death with normal value — exceeds  $t$ , and

199

$$Q_i(t) = \exp\left\{-\int_0^t \lambda_i(v)dv\right\}$$

200 for  $i = 1$  and  $2$ , and

201

$$Q_3(t|u) = \exp\left\{-\int_u^t \lambda_3(v|u)dv\right\}.$$

202 are pseudo-survival functions corresponding to the intensities  $\lambda_1(u)$ ,  $\lambda_2(t)$  and  $\lambda_3(t|u)$ .

203 We are interested in estimating  $\Lambda_1(t) = \int_0^t \lambda_1(u)du$ , the cumulative incidence function (CIF), but  
 204 focus on the comparison of CIFs between AR and NR groups based on the original data and  
 205 imputation data.

206

207 Table 2 identifies the various types of observations which occur in this illness-death/failure model  
 208 and the corresponding contribution to the likelihood, denoted as  $L_1(t)$  to  $L_4(t)$ , which are functionals  
 209 of intensities and pseudo-survival functions. Rai et al.<sup>2</sup> derive the explicit form of  $L_1(t)$  to  $L_4(t)$  for  
 210 both constant and piecewise exponential model and the likelihood functions, which are summarized  
 211 in appendix.

212

Table 2: Likelihood Contributions for Anthracycline Cardiac Toxicity Study		
Observation Type	Outcome	Likelihood Contribution
Death with No Cardiac Abnormality	$T = t, X(t^-) = 1$	$L_1(t) = \lambda_2(t)Q(t)$
Alive with No Cardiac Abnormality	$T > t, X(t) = 1$	$L_2(t) = Q(t)$
Death/Cardiac Failure with Cardiac Abnormality	$T = t, X(t^-) = 2$	$L_3(t) = \int_0^t \lambda_1(u)Q(u)\lambda_3(t u)Q_3(t u)du$
Alive with Cardiac Abnormality	$T > t, X(t) = 2$	$L_4(t) = \int_0^t \lambda_1(u)Q(u)Q_3(t u)du$

213

### 214 Imputation Model

215 Let the cardiac measure, such as  $AF$ , be denoted by  $Y$ . Assume that  $Y = (Y_1, Y_2)^T$  be a  $n \times 1$   
 216 response vector with  $Y_1$  ( $n_1 \times 1$ ) observed and  $Y_2$  ( $n_2 \times 1$ ) missed, and  $X = (X_1, X_2)^T$  be  
 217 corresponding  $n \times p$  matrix comprised of covariates including other cardiac measures (other  
 218 response variables).

219 There are several methods for imputation which can be broadly classified as single imputation or  
 220 multiple imputation (MI). In MI approach several copies of the complete data set are created and  
 221 then the appropriate statistical method is applied to each data set and the results from these analyses  
 222 are then combined to provide the final results. Usually, MI approaches are preferred over single

223 imputation as they incorporate variability due to imputation<sup>17, 29-30</sup>. There are many MI approaches  
224 discussed in literature, but two most commonly used approaches based on joint multivariate  
225 modeling or fully conditional specification perform quite well in the regression setting, as seen in  
226 Huque et al<sup>31</sup>. It may be noted that PROC MI can perform imputations for data that have monotone  
227 or arbitrary missing patterns. PROC MI with FCS option, a standard feature in SAS version 9.4<sup>32</sup>,  
228 utilizes the conditional distribution and can incorporate both continuous and categorical variables  
229 appropriately, see Liu and De<sup>33</sup>. In our setting we had missing values only in  $AF$  and we wanted to  
230 take advantage of the relationship between  $AF$  and other covariates of interest that included  
231 categorical variables. Therefore, we preferred to perform the imputations using PROC MI in SAS  
232 with FCS option. The method can be briefly described as follows:

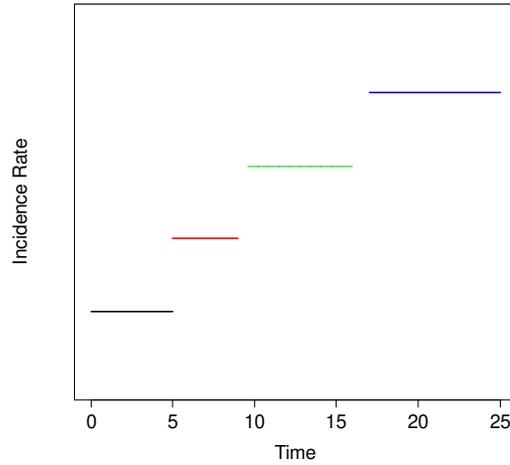
233 A multivariable regression model  $y = x\beta + \varepsilon$  is fitted based on the complete data  $Y_1$  and  $X_1$ , and  
234 the least squared estimator  $\hat{\beta}$  of  $\beta$  ( $p \times 1$ ) and associated variance-covariance matrix is obtained.  
235 Then, missing values in  $Y_2$  are imputed using the posterior predictive distributions, see PROC MI in  
236 SAS<sup>32</sup> for details. It is natural to use the imputation data  $(Y_1, \hat{Y}_2)^T$  instead of only  $Y_1$  and is  
237 anticipated that the imputed information in  $\hat{Y}_2$  will improve the related results in statistical analysis.

238 To each complete data set likelihood ratio test was applied to compare the two risk groups (AR and  
239 NR). In the regression framework one could use PROC MIANALYZE in SAS to combine the  
240 results from multiple imputations to conduct inference that incorporates inherent variability  
241 introduced due to imputations<sup>17,29-30</sup>. However, in our setting p-values associated with each  
242 imputation are obtained based on the likelihood ratio test. Then, the overall conclusion can be based  
243 on some type of summary measure of all the p-values such as mean or median. We prefer to report  
244 the results based on median as that would be much more robust than mean.

## 245 **Results**

### 246 **Application: Cancer Survivor Study**

247 In this section we obtained the imputed data for the cardiotoxicity example and applied the theory  
248 for the exponential model described in appendix to evaluate the effect of anthracyclines on  
249 cardiotoxicity. Furthermore, the results obtained using the imputation approach are then compared  
250 with those obtained without imputation, reported in Rai et al.<sup>2</sup>, under the assumption of no  
251 deaths/cardiac failures. The simplest model is the Parametric-1, which is one parameter Exponential  
252 model. Since there are very few events before 5 years and after 10 years, we also fit two piecewise  
253 Exponential models; Parameter-2, based on two incidence rates, one up-to year five and the second  
254 for year 5 and above, and Parameter-3, based on three incidence rates one up-to year 5, second  
255 between years 5 and 10 and the last one for year 10 and above, (see Figure 3).



256

257

Figure 3: Piecewise Incidence Rate

258 In the cardiotoxicity example, there are 278 subjects, Leukemia (n=77), Sarcoma (n=74),  
 259 Lymphoma (n=56) and Embryonal (n=71), and 34 measurements were missing for *AF* and 6 were  
 260 missing in both measurements *AF* and *FS*, but no covariate information was missing. We exclude  
 261 the 6 with both missing. Hence, we have 272 subjects and employ the multivariable regression  
 262 method to estimate the 34 missing measurements in *AF* based on the values of *FS* and  
 263 corresponding covariates, like age at diagnosis, race, gender, BMI, QTC, diagnosis group and risk  
 264 group (AR/NR)<sup>27-28</sup>. Based on 4 diagnosis groups, we define three dummy variables as follows:

265 
$$\text{Diag1} = \begin{cases} 1 & \text{Leukemia} \\ 0 & \text{Otherwise,} \end{cases} \quad \text{Diag2} = \begin{cases} 1 & \text{Sarcoma} \\ 0 & \text{Otherwise,} \end{cases} \quad \text{Diag3} = \begin{cases} 1 & \text{Lymphoma} \\ 0 & \text{Otherwise.} \end{cases}$$

266 Before conducting the regression analysis, the Shapiro-Wilk test of normality was applied to  
 267 original *AF* and *FS* measurements and a few commonly used transformations for making the  
 268 underlying distributions of *AF* and *FS* more normal.  $\text{Log}(AF)$  was normally distributed ( $p=0.701$ )  
 269 but original *AF* ( $p<0.001$ ) was not. On the other hand,  $\text{Log}(FS)$ ,  $\text{Logit}(FS)$  and *FS* were not  
 270 normally distributed with  $p$  values  $<0.001$ ,  $0.007$  and  $0.026$ , respectively. This suggests fitting the  
 271 regression model using logarithm transformation of *AF* and original *FS*. The significant predictors  
 272 with coefficients, their  $p$ -values and  $R^2$  in the model are presented in Table 3a. That is,

273 
$$\log(AF) = 5.534 - 3.891FS + 0.008Age + 0.081Risk + 0.098Diag2 - 0.010BMI + \varepsilon \quad (4.1)$$

274 where Risk=1 for patient in AR group and 0 in NR group. Based on the regression model, the  
 275 values of *FS* and the covariates, the 34 missing *AF* values are imputed. Thus, after imputing the  
 276 missing values the total sample size is 272.

Table 3a: Coefficients and P-values in Regression Model			
Variable	Estimator	SE	p-value
Intercept	5.534	0.106	<0.001

<b><i>FS</i></b>	-3.891	0.274	<0.001
<b>Age at Diagnosis</b>	0.008	0.003	0.005
<b>Risk Group</b>	0.081	0.038	0.033
<b>Diag2</b>	0.098	0.035	0.006
<b>BMI</b>	-0.010	0.003	<0.001
$R^2$	0.586		

277

278 **Remark:** It may be noted that the  $R^2$  for the above model is 0.586 which represents a reasonable fit  
 279 but may not be the best model fit. However, in general, it would not be unreasonable to expect that  
 280 the imputation process would be more efficient if a better model with higher value of  $R^2$  can be  
 281 obtained.

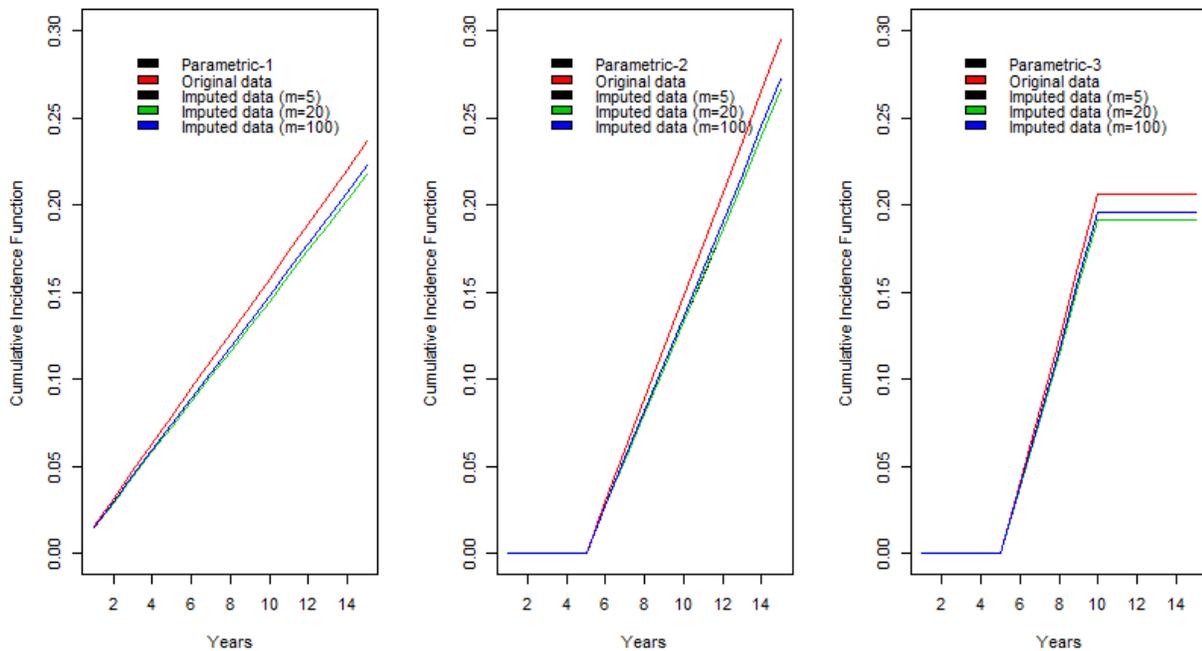
282 Because *AF* is the only variable with missing values in our data (n=272), we consider the data  
 283 including age at diagnosis, BMI, diagnosis group, risk group, *FS*. Based on each imputation data,  
 284 we calculate the Cumulative Incidence Function (CIF) for constant exponential model, two or three  
 285 piecewise exponential model using the method. Then, the methods described in Appendix were  
 286 applied to the imputed data sets for each group, AR and NR, and both groups combined. Then,  
 287 based on the likelihood ratio test, the corresponding p-value for group effect for the variable *AF* is  
 288 0.014 (median of  $m = 20$  imputations). On the other hand, the p-value without imputation is  
 289  $0.020^2$ . The results with imputation seem to be little bit more sensitive for the group effect in *AF*  
 290 compared to that without imputation. Note that not using the illness/death model as proposed and,  
 291 instead, using logistic regression, without imputation, the group effect was only marginally  
 292 significant ( $p=0.065$ )<sup>1-2</sup>. Thus, the approach base on illness/death model led to a better  
 293 understanding of this data and motivated the current development. A summary of the results is  
 294 given in Table 3b. The cumulative incidence function (CIF) was derived for exponential and  
 295 piecewise exponential models for the imputed data using the above regressions model and SAS  
 296 procedure (PROC MI, with  $m=5, 10$  and  $100$  imputations) and were compared to those based on  
 297 the original data (without imputation).

298 In Table 3b, the group effects are reported for both data without imputation and with imputation for  
 299  $m=5, 20$  and  $100$ . For imputed data, we reported the p-values of group effect as mean, minimum,  
 300 maximum, and median. A comparison of cumulative incidence rates can be found in Figure 4.

<b>Table 3b. AF p-values for Group Effect using Likelihood Ratio Test</b>					
Procedure	Without Imputation	With Imputation (MI)			
		mean	median	min	max
Parameter-1					
m=5	0.0199	0.0125	0.0121	0.0098	0.0187
m=20		0.0180	0.0138	0.0098	0.0561
m=100		0.0160	0.0123	0.0051	0.0696

Parameter-2					
m=5	0.0118	0.0082	0.0077	0.0063	0.0126
m=20		0.0120	0.0091	0.0063	0.0397
m=100		0.0105	0.0080	0.0031	0.0487
Parameter-3					
m=5	0.0782	0.0564	0.0519	0.0469	0.0805
m=20		0.0743	0.0621	0.0469	0.1983
m=100		0.0656	0.0565	0.0206	0.2164

301



302

303 **Figure 4:** Cumulative Incidence Comparison for AF Based on Original Data and Imputed Data  
 304 (Using the mean of intensity estimates) Corresponding to  $m = 5, 20$  and 100.

305 Note: The lines for  $m=5$  and 20 are almost overlapped.

### 306 Simulation Study

307 To assess the performance of the imputation approach we conducted simulation studies as described  
 308 below.

309 The primary focus is on assessing the performance of imputing AF in Anthracycline Cardiac  
 310 Toxicity data for comparing the cumulative incidences of cardiac toxicity in the illness-death model  
 311 as discussed above. The detail steps are described as follows.

312 **Step 1:** From Table 3a it is clear that  $\log(AF)$  values are associated with the risk group (AR and  
 313 NR), diagnosis group (Sarcomas vs. others), Age, FS, and BMI and using the equation (4.1) we first  
 314 filled all the missing values of  $AF$  with the mean predicted values and obtained a complete copy of  
 315 the data set.

316 **Step 2:** Then, for simulation studies we first created four subgroups:

317 Group 1: The patients which are in AR group diagnosed with Sarcomas (sample size  $n_1$ )

318 Group 2: The patients which are in AR group diagnosed with other cancers (sample size  $n_2$ )

319 Group 3: The patients which are in NR group diagnosed with Sarcomas (sample size  $n_3$ )

320 Group 4: The patients which are in NR group diagnosed with other cancers (sample size  $n_4$ )

321 Now to keep the covariance structure consistent with the observed data the sample mean and  
 322 variance-covariance matrix were obtained for the variables  $LAF=\log(AF)$ , FS, Age, BMI and  
 323  $LTime=\log(Time)$ , where Time is the length of follow-up from diagnosis to the time of survey and  
 324 the individual can be in one of the three states as shown in Figure 2.

325 **Step 3: (Generate multi-normal data by group):** Generate a sequence of random vectors,  
 326  $(LAF, FS, Age, BMI, LTime)_j, j = 1, \dots, n_i$ , from multi-normal distribution for Group  $i$  ( $i =$   
 327  $1, 2, 3, 4$ ). That is, we assumed that,

$$328 \quad (LAF, FS, Age, BMI, LTime)_j \sim i. i. d. \quad MVN(v_i, \Sigma_i) \text{ for } j = 1, 2, \dots, n_i, i = 1, 2, \dots, 4$$

329 where  $n_1 = n_2 = n/3$  and  $n_3 = n_4 = n/6$  for  $n = 180, 240, \text{ or } 300$  (the sample sizes for three  
 330 simulation studies) and  $v_i$  and  $\Sigma_i$  are the sample mean vector and covariance matrix. In this  
 331 simulation study we used unequal samples size to reflect higher proportion of At-risk survivors in  
 332 our cohort.

333 **Note:** We noticed the means of FS and LAF were significantly low (high) for Group 1 which would  
 334 have led to highly significant p-values for comparing the two groups (AR vs. NR). Therefore, we  
 335 adjusted the mean values for these two variables in the simulation studies as follows:

Mean of FS and LAF in Anthracycline Cardiac Toxicity data	
	Group

		1	2	3	4
Original	FS	0.307	0.342	0.359	0.358
	LAF	4.209	3.932	3.865	3.786
Adjusted	FS	<b>0.330</b>	0.342	0.359	0.358
	LAF	<b>4.000</b>	3.932	3.865	3.786

336

337 **Step 4: (Incomplete Data in AF):** From the sample size generated in Step 2, we randomly deleted  
338 R% (R=20 or 30) of AF values, and got incomplete data with sample sizes (100-R)%n.

339 **Step 5: (Imputed data):** Using SAS procedure PROC MI with FCS option we imputed AF values  
340 and obtained a complete copy of the data set.

341 **Step 6 (Calculate p-value for group Effect):** For the one parameter exponential distribution, the p-  
342 values for group effect (comparing AR with NR) were obtained using likelihood ratio test for  
343 complete (originally generated), incomplete and imputed data set.

344 **Step 7:** The imputation process (Step 5) was repeated 20 times to obtain 20 copies of complete data  
345 sets, which resulted in 20 p-values. A description of the p-values in terms of mean, median,  
346 minimum, and maximum is summarized in Table 4.

347 **Step 8:** Steps 2 – 7 were repeated 10 times to assess the performance of the imputation approach on  
348 10 independently generated data sets. The results of the simulation study are summarized in Table  
349 4.

350 From Table 4, it is seen that, in general, the median of the p-values is much closer to the p-value  
351 obtained from the complete data compared to those obtained from the incomplete data. However,  
352 there are some extreme situations where the results from the complete data and those obtained from  
353 imputations and incomplete data are not in agreement and this could be due to chance that more  
354 observations were deleted from a particular group and the regression is not able to completely  
355 exploit the underlying correlation structure. For example, for the 4<sup>th</sup> simulation, when the sample  
356 size is 180 and 30% observations are imputed the p-value for the complete data set is 0.007 but  
357 those corresponding to incomplete data and imputed data are 0.059 and 0.469, respectively. This  
358 clearly suggests that the manner in which data are generated and the observations are randomly  
359 deleted might have changed the underlying structure particularly those who are in the AR and NR

360 groups. However, for larger sample size (n=300) we see that in all simulations the median p-value  
 361 based on imputation is much closer to the p-value obtained from the complete data. Thus, it is clear  
 362 that with imputation approach we are able to exploit the underlying correlation structure and obtain  
 363 nearly unbiased conclusions.

364

<b>Table 4. p-values for Group Effects using Likelihood Ratio Test</b>									
			Complete Data	Incomplete Data	Imputed Data (p-values)				
n	R%	Simulation #	p-value	p-value	Median	Min	Max	Mean	SD
180	20	1	<.001	<.001	<.001	<.001	0.011	0.001	0.003
		2	0.285	0.322	0.256	0.127	0.628	0.300	0.143
		3	0.004	0.019	0.006	0.002	0.014	0.006	0.003
		4	0.007	0.057	0.025	0.003	0.202	0.045	0.050
		5	0.003	0.003	0.004	<.001	0.025	0.007	0.008
		6	<.001	0.001	0.005	<.001	0.073	0.011	0.016
		7	0.013	0.031	0.020	0.002	0.088	0.027	0.027
		8	0.761	0.937	0.770	0.539	0.995	0.801	0.135
		9	0.007	0.001	<.001	<.001	0.100	0.008	0.024
		10	0.011	0.099	0.026	0.004	0.191	0.045	0.053
180	30	1	<.001	0.130	0.004	<.001	0.091	0.015	0.027
		2	0.285	0.002	0.035	0.005	0.130	0.046	0.037
		3	0.004	0.003	<.001	<.001	0.003	<.001	0.001
		4	0.007	0.059	0.469	0.115	0.889	0.493	0.230
		5	0.003	0.033	0.314	0.022	0.935	0.330	0.253
		6	<.001	0.040	0.001	<.001	0.060	0.007	0.015
		7	0.013	0.262	0.008	<.001	0.090	0.020	0.026
		8	0.761	0.003	0.035	0.003	0.150	0.050	0.044
		9	0.007	0.005	<.001	<.001	0.032	0.005	0.009
		10	0.011	0.250	0.140	0.011	0.372	0.140	0.100
240	20	1	<.001	<.001	<.001	<.001	0.003	<.001	0.001
		2	0.087	0.075	0.087	0.011	0.248	0.094	0.072
		3	0.020	0.041	0.040	0.003	0.117	0.043	0.026
		4	0.011	0.005	0.007	<.001	0.086	0.017	0.022
		5	0.002	0.016	0.009	0.001	0.74	0.025	0.060
		6	<.001	<.001	<.001	<.001	0.001	<.001	<.001
		7	0.001	0.001	0.001	<.001	0.008	0.001	0.002
		8	0.221	0.390	0.347	0.109	0.933	0.388	0.197
		9	0.001	<.001	<.001	<.001	0.003	0.001	0.001
		10	<.001	0.001	<.001	<.001	0.033	0.002	0.007
		1	<.001	<.001	<.001	<.001	<.001	<.001	<.001

240	30	2	0.087	0.327	0.105	0.102	0.499	0.154	0.141
		3	0.020	0.060	0.023	0.004	0.135	0.037	0.038
		4	0.011	0.028	0.011	0.001	0.140	0.027	0.033
		5	0.002	0.014	0.008	<.001	0.042	0.011	0.013
		6	<.001	<.001	<.001	<.001	0.004	0.001	0.001
		7	0.001	0.002	0.001	<.001	0.081	0.011	0.022
		8	0.221	0.121	0.041	0.004	0.201	0.051	0.047
		9	0.001	0.019	0.012	0.001	0.079	0.022	0.026
		10	<.001	0.005	0.015	<.001	0.159	0.041	0.054
300	20	1	<.001	<.001	<.001	<.001	0.001	<.001	<.001
		2	0.268	0.107	0.155	0.023	0.809	0.212	0.184
		3	0.026	0.011	0.014	0.001	0.077	0.022	0.022
		4	<.001	<.001	<.001	<.001	0.003	<.001	0.001
		5	<.001	0.001	0.001	<.001	0.005	0.001	0.002
		6	<.001	0.028	0.007	0.001	0.056	0.014	0.015
		7	0.001	0.003	0.002	<.001	0.016	0.004	0.005
		8	0.107	0.193	0.118	0.031	0.370	0.148	0.098
		9	<.001	0.001	0.002	<.001	0.017	0.004	0.005
		10	0.001	0.013	0.007	0.001	0.336	0.047	0.090
300	30	1	<.001	<.001	<.001	<.001	<.001	<.001	<.001
		2	0.268	0.759	0.383	0.078	0.889	0.389	0.231
		3	0.026	0.024	0.019	0.003	0.103	0.032	0.026
		4	<.001	0.003	<.001	<.001	0.007	0.001	0.002
		5	<.001	0.009	0.001	<.001	0.056	0.006	0.013
		6	<.001	0.015	0.003	<.001	0.042	0.009	0.014
		7	0.001	0.015	0.003	<.01	0.022	0.002	0.006
		8	0.107	0.169	0.110	0.007	0.620	0.142	0.142
		9	<.001	<.001	<.001	<.001	0.003	<.001	0.001
		10	0.001	0.002	0.003	<.001	0.114	0.015	0.027

## 365 Discussion

366 In this paper, we have employed a well-established methodology of illness-death/Failure model<sup>2</sup>  
367 and imputed the missing observations for a phase IV clinical trial study as an example. Although,  
368 we assumed a very simple parametric model, it is straightforward to expand to other parametric or  
369 semi-parametric models. From a clinician's point of view the simple approaches such as log rank  
370 test and KM survival curves<sup>34</sup> are most commonly used and understood. However, in the setting of  
371 interval censored data the approach proposed here is simple to use and can be implemented easily to  
372 estimate fixed-time cumulative incidence function with or without imputation<sup>2</sup>.

373 When studying the long-term effects of treatment, there can be multiple unwanted events identified  
374 at the time of observation. Some of these events can be competing and others are not correlated.  
375 This leads to multivariate time-to-event data. One simple approach is to study the incidence of first  
376 event and then incidence of specific event. In our example, cardiotoxicity measures included  
377 abnormal *AF* and *FS* but there are some other measures to evaluate cardiotoxicity. For some reason,  
378 not all patients had both measures and the models based on bi-variate time-to-event outcomes  
379 would include only those patients who have data on both outcomes and this would reduce the  
380 sample size and potentially ignore important information. Based on this consideration, the  
381 multivariable regression method was used to impute the missing observations and to apply the  
382 parametric method to the imputed data and compare the results with those obtained without  
383 imputation.

384 As stated before, the problem of evaluating possible toxic effects of cancer therapies in a Phase IV  
385 trial setting is an important problem. Among the many issues in such studies, missing data is a key  
386 aspect that can influence the inferences. It is also important to understand the nature of impact of  
387 missing data on the analysis and the interpretation of the study data. Also note that imputation  
388 increases the sample size, and thus increases statistical power to detect the same effect size. But if  
389 the model assumptions are not correct, the inference may not be valid. Thus, it is recommended to  
390 report the p-values with and without imputation. However, with higher absolute correlations  
391 between two outcome measures (the primary outcome measure, *AF*, with higher missing and the  
392 secondary outcome measure, *FS*, with little or no missing), produced efficient results, a rigorous  
393 simulation study with different amount of correlations between two outcome measures, amount of  
394 missing and model uncertainty is underway to consider this aspect and will be reported elsewhere;  
395 this is along the lines our work for a randomized clinical study<sup>33</sup>.

396 The study involved all the patients visiting the clinic in a pre-specified time frame (such as 1 year of  
397 accrual) and represents a somewhat unbiased survey of patients. Since the outcome measure may  
398 depend on disease type, an almost equal allocation was used to enroll patients. It has been reported  
399 in Hudson et al.<sup>1</sup> that the prevalence depends on disease type; hence it will be another research  
400 direction to adjust the sampling allocation and variability due to sampling and modeling for  
401 generalizing the results for the entire patient population<sup>35</sup>.

402 Another limitation of this study is that this is a cross-sectional survey to estimate the long-term  
403 effect of cardiotoxicity of the primary treatment of cancer. Dodge<sup>36</sup> defined cross-sectional survey  
404 “A method of data collection whereby a battery of questions is asked of participation at one single  
405 point or in a relatively small interval of time. Inferences about a population must be anchored to the  
406 time period in which the sample was taken. Data from cross-sectional surveys are typically unable  
407 to be used to prove the existence of cause-and-effect relationships.” Even though this is based on  
408 enrolling consecutive eligible patients in a very homogeneous environment (St. Jude Children’s  
409 Hospital treats patients without charge to patients), effect of this limitation is minimized but cannot  
410 be reduced to zero. Generalization of the results to a general population should be done with  
411 caution.

412 **Conclusions**

413 Based on simulation output, it is suggested that the results obtained by imputing the missing values  
414 using regression methods are significantly more efficient than those obtained without imputation. It  
415 is recommended to carefully impute missing values in certain dataset in order to retain most of the  
416 information from the incomplete cases.

417 **List of Abbreviations**

418 **AR:** At-Risk

419 **NR:** Not At-Risk

420 **FS:** Fractional shortening

421 **AF:** Afterload

422 **LVEdD:** Left ventricular end-diastolic diameter

423 **LVEsD:** Left ventricular end-systolic diameter

424 **BMI:** Body mass index

425 **CIF:** Cumulative incidence function

426 **Declarations**

427 **Ethics approval and consent to participate**

428 The study was approved by the institutional review boards at St Jude Children's Research Hospital  
429 and Stanford University (Stanford, CA).

430 All study participants or their parents provided informed consent.

431 All methods were carried out in accordance with relevant guidelines and regulations.

432 **Consent to publish**

433 Not applicable

434 **Availability of data and materials**

435 The datasets generated and/or analyzed during the current study are not publicly available due to  
436 privacy restrictions.

437 **Competing interests**

438 The authors declare that they have no competing interests.

439 **Funding**

440 D.K. Srivastava was in part supported by the Grant CA21765 from the National Institutes of Health  
441 (NIH) and by the American Lebanese Syrian Associated Charities (ALSAC).

442 C. Qian was supported by the University of Louisville Fellowship.

443 S. N. Rai was partly supported with Wendell Cherry Chair in Clinical Trial Research Fund, multiple  
444 National Institutes of Health (NIH) grants (5P20GM113226, PI: McClain; 1P42ES023716, PI:  
445 Srivastava; 5P30GM127607-02, PI: Jones; 1P20GM125504-01, PI: Lamont; 2U54HL120163, PI:  
446 Bhatnagar/Robertson; 1P20GM135004, PI: Yan; 1R35ES0238373-01, PI: Cave; 1R01ES029846,  
447 PI: Bhatnagar; 1R01ES027778-01A1, PI: States;), and Kentucky Council on Postsecondary  
448 Education grant (PON2 415 1900002934, PI: Chesney).

449 **Author's contributions**

450 Conception: All Authors

451 Data analysis: JP

452 Original Draft: All Authors

453 Review & Editing: All Authors

454 Approval of the Final Version: All Authors

455 **Acknowledgements**

456 We thank Dr. Donald Miller for his support for this research work.

457 **References**

- 458 1. Hudson MM, Rai SN, Nunez C, Merchant TE, Marina NM, Zalamea N, Cox C, Phipps S,  
459 Pompeu R and Rosenthal D. Noninvasive evaluation of late anthracycline cardiac toxicity in  
460 childhood cancer survivors. *Journal of Clinical Oncology* 2007; **25**: 3635-3643.
- 461 2. Rai SN, Pan J, Sun J, Hudson MM and Srivastava Dk. Estimating incidence rate on current  
462 status data with application to a Phase IV cancer trial. *Communications in Statistics: Theory*  
463 *and Methods* 2013; **42**:2417-2433.

- 464 3. Armstrong GT, Chen Y, Yasui Y, et al. Reduction in late mortality among 5-year survivors of childhood cancer. *N Engl J Med* 2016;**374**(9):833–42.  
465
- 466 4. Howlader N, Noone AM, Krapcho M, et al. SEER Cancer Statistics Review (CSR), 1975-  
467 2012. Bethesda, MD: National Cancer Institute (internet). 2015.  
468 [http://seer.cancer.gov/csr/1975\\_2012/](http://seer.cancer.gov/csr/1975_2012/) (accessed November 18, 2020).
- 469 5. Phillips SM, Padgett LS, Leisenring WM, et al. Survivors of childhood cancer in the United  
470 States: prevalence and burden of morbidity. *Cancer Epidemiol Biomarkers Prev*  
471 2015;**24**(4):653–63.
- 472 6. Fidler MM, Reulen RC, Winter DL, et al. Long term cause specific mortality among 34 489  
473 five year survivors of childhood cancer in Great Britain: population based cohort study.  
474 *BMJ* 2016;354: i4351.
- 475 7. Robison LL, Hudson MM. Survivors of childhood and adolescent cancer: life-long risks and  
476 responsibilities. *Nat Rev Cancer* 2014;**14**(1):61–70.
- 477 8. Krischer JP, Epstein S, Cuthbertson DD, et al. Clinical cardiotoxicity following an-  
478 thracycline treatment for childhood cancer: The Pediatric Oncology Group experience.  
479 *Journal of Clinical Oncology* 1997; **15**: 1544-1552.
- 480 9. Sorensen K, Levitt GA, Bull C, et al. Late anthracycline cardiotoxicity after childhood  
481 cancer: A perspective longitudinal study. *Cancer* 2003; **97**: 1991-1998.
- 482 10. Pein F, Sakiroglu O, Dahan M, et al. Cardiac abnormalities 15 years and more after  
483 adriamycin therapy in 229 childhood survivors of a solid tumor at the Institut Gustave  
484 Roussy. *British Journal of Cancer* 2004; **91**: 37-44.
- 485 11. Sun J. *The Statistical Analysis of Interval-Censored Failure Time Data*. Springer: New  
486 York, 2006.
- 487 12. Ayer M, Brunk HD, Ewing, GM, Reid WT, and Silverman E. An empirical distribution  
488 function for sampling with incomplete information. *Ann. Math. Statist.* 1955; **26**: 641-647.
- 489 13. Van Eeden C. Maximum likelihood estimation of ordered probabilities. *Indagationes*  
490 *Mathematicae* 1956; **18**: 444-455.
- 491 14. Peto R. Experimental survival curves for interval-censored data. *Appl. Statist.* 1973; **22**: 86-  
492 91.
- 493 15. Turnbull BW. The empirical distribution function with arbitrarily grouped, censored and  
494 truncated data. *J. R. Statist. Soc. B* 1976; **38**: 290-295.
- 495 16. Little RJ, Rubin DB. *Statistical Analysis with Missing Data*, 2<sup>nd</sup> edition, New York: John  
496 Wiley, 2002.
- 497 17. Burren Sv. *Flexible Imputation of Missing Data*, 2<sup>nd</sup> edition, Chapman & Hall/CRC, Florida,  
498 2012.

- 499 18. Rubin DB. *Multiple Imputation for Nonresponse in Surveys*. J. Wiley & Sons, New York,  
500 1987.
- 501 19. Rubin DB. Multiple imputation after 18+ years (with discussion). *Journal of the American*  
502 *Statistical Association*, **91**: 473-489
- 503 20. Rubin DB, Stern H, and Vehovar V. (1995), Handling “don’t know” survey responses: The  
504 case of the Slovenian plebiscite. *Journal of the American Statistical Association* 1995; **90**:  
505 822-828.
- 506 21. King G, Honaker J, Joseph A, and Scheve K. Analyzing incomplete political science data:  
507 an alternative algorithm for multiple imputation. *The American Political Science Review*  
508 2001; **95**: 49–69.
- 509 22. Schafer JL. *Analysis of Incomplete Multivariate Data*. Chapman & Hall, London, 1997.
- 510 23. Liu C. Missing data imputation using the multivariate t distribution. *Journal of Multivariate*  
511 *Analysis* 1995; **53**: 139-158.
- 512 24. Buuren Sv, Boshuizen HC and Knook DL. Multiple imputation of missing blood pressure  
513 covariates in survival analysis. *Statistics in Medicine* 1999; **18**: 681-694.
- 514 25. Troxel A, Guoguang M and Heitjan DF. An index of local sensitivity to nonignorability.  
515 *Statistica Sinica* 2004; **14**: 1221-1237.
- 516 26. Cohen, J. *Statistical power analysis for the behavioral sciences*, 2<sup>nd</sup> edition. Hillsdale, NJ:  
517 Lawrence Erlbaum Associates, 1988.
- 518 27. Henry WL, Gardin JM and Ware JH. Echocardiographic measurements in normal subjects  
519 from infancy to old age. *Circulation* 1980; 62:1054-1061.
- 520 28. Colan SD, Parness IA, Spevak PJ, et al. Developmental modulation of myocardial  
521 mechanics: Age- and growth-related alterations in afterload and contractility. *Journal of the*  
522 *American College of Cardiology*, 1992; 19:619-629.
- 523 29. Molenberghs G and Fitzmoaurice G. Incomplete data: Introduction and overview. In  
524 *Handbooks of Modern Statistical Methods; Longitudinal Data Analysis*, (Chapter 17, pages  
525 395-408), Chapman & Hall/CRC, Florida, 2009.
- 526 30. Kenward MG and Carpenter JR. Multiple Imputation. In *Handbooks of Modern*  
527 *Statistical Methods; Longitudinal Data Analysis*, (Chapter 21, pages 477-500), Chapman &  
528 Hall/CRC, Florida, 2009.
- 529 31. Huque, M.H., Carlin, J.B., Simpson, J.A. *et al*. A comparison of multiple imputation  
530 methods for missing data in longitudinal studies. *BMC Med Res Methodol* **18**, 168 (2018).
- 531 32. SAS/STAT®14.1 User’s Guide: The MI Procedure, SAS Institute Inc., Cary, NC, USA,  
532 2015.

- 533 33. Liu Y, and De A. Multiple imputations by fully conditional specification for dealing with  
534 missing data in a large epidemiologic study. *International Journal of Statistics in Medical*  
535 *Research*, 4(3): 287-295, 2015.
- 536 34. Kaplan EL and Meier P. Nonparametric estimation from incomplete observation. *Journal of*  
537 *the American Statistical Association* 1958; **53**: 457-481.
- 538 35. Kovacevic MS and Rai SN. Log-Linear Modeling of Change Using Longitudinal Survey  
539 Data. *Communications in Statistics—Theory and Methods* 2002; **31**: 1815-1835.
- 540 36. Dodge Y. *The Oxford Dictionary of Statistical Terms*. The International Statistical Institute,  
541 Oxford Press, 2003

542 **Appendix**

543 **A1: Details of the Likelihood Approach**

544 Due to the actual time of onset of abnormality,  $U$ , is not known, the observed quantity for each  
 545 patient includes the observation time,  $T$ , and two indicators of status,  $\delta$  and  $\gamma$ , at the time of  
 546 survey, where  $\delta$  is an indicator of patient alive with no cardiac failure or dead/cardiac failure, and  
 547  $\gamma$  is an indicator of patient with a normal or abnormal value. Let  $t_i$  be the observation time (death,  
 548 cardiac failure or survey) for the  $i^{th}$  subject. That is,

550 
$$\delta_i = \begin{cases} 1, & \text{if unit } i \text{ dead or cardiac failure at time } t_i \\ 0, & \text{if unit } i \text{ alive and no cardiac failure at time } t_i, \end{cases}$$

551 and

552 
$$\gamma_i = \begin{cases} 1, & \text{if unit } i \text{ with abnormal value at time } t_i \\ 0, & \text{if unit } i \text{ with normal value at time } t_i. \end{cases}$$

553 The simplified forms of intensities  $\lambda_i(t) = \lambda_i$  for  $i = 1$  or  $2$  and  $\lambda_3(t|u) = \lambda_3$  lead to  $Q_i(t) =$   
 554  $e^{-\lambda_i t}$  for  $i = 1$  or  $2$ ,  $Q_3(t|u) = e^{-\lambda_3(t-u)}$  and  $Q(t) = e^{-(\lambda_1+\lambda_2)t}$ . We derive the corresponding  
 555 likelihood contributions from  $L_1(t)$  to  $L_4(t)$  for the four observation types in Table 2 and then the  
 556 log-likelihood function as follows

557 
$$l(\lambda_1, \lambda_2, \lambda_3) = \sum_{i=1}^n a_i [\log \lambda_2 - (\lambda_1 + \lambda_2)t_i] - \sum_{i=1}^n b_i (\lambda_1 + \lambda_2)t_i]$$
  
 558 
$$+ \sum_{i=1}^n c_i [\log \lambda_1 + \log \lambda_3 - \log(\lambda_1 + \lambda_2 - \lambda_3) + \log(e^{-\lambda_3 t_i} - e^{-(\lambda_1+\lambda_2)t_i})]$$
  
 559 
$$+ \sum_{i=1}^n d_i [\log \lambda_1 - \log(\lambda_1 + \lambda_2 - \lambda_3) + \log(e^{-\lambda_3 t_i} - e^{-(\lambda_1+\lambda_2)t_i})],$$

560 where  $a_i = \delta_i(1 - \gamma_i)$ ,  $b_i = (1 - \delta_i)(1 - \gamma_i)$ ,  $c_i = \delta_i \gamma_i$  and  $d_i = (1 - \delta_i)\gamma_i$  are the indicators  
 561 corresponding to observation type 1 to type 4 in Table 2. Then the maximum likelihood estimators  
 562  $\hat{\lambda}_1, \hat{\lambda}_2$  and  $\hat{\lambda}_3$  of  $\lambda_1, \lambda_2$  and  $\lambda_3$  are derived from the following equations<sup>2</sup>.

563

564 
$$\frac{D_3 + D_4}{\lambda_1} - \frac{D_3 + D_4}{\lambda_1 + \lambda_2 - \lambda_3} - T_2 + \sum_{i=1}^n \frac{(c_i + d_i)t_i}{e^{(\lambda_1+\lambda_2-\lambda_3)t_i} - 1} = 0$$

565 
$$\lambda_2 = \frac{D_1}{D_3 + D_4} \lambda_1$$

566 
$$\lambda_3 = \frac{D_3}{T_1 \lambda_1 - D_3 - D_4} \lambda_1$$

567 where

568 
$$D_1 = \sum_{i=1}^n a_i, D_2 = \sum_{i=1}^n b_i, D_3 = \sum_{i=1}^n c_i, D_4 = \sum_{i=1}^n d_i$$

569 
$$T_1 = \sum_{i=1}^n (a_i + b_i + c_i + d_i) t_i \text{ and } T_2 = \sum_{i=1}^n (a_i + b_i) t_i.$$

570 It is further extended to the model to allow the intensity  $\lambda_1$  with piecewise constant<sup>2</sup>. Assume two  
 571 intervals: less than  $t_c$  years and above  $t_c$  (including  $t_c$ ) years (say,  $t_c = 5$ ) and let these two rates  
 572 be  $\lambda_{11}$  and  $\lambda_{12}$ . Then the log-likelihood function is derived as

$$\begin{aligned}
 573 \quad l(\lambda_{11}, \lambda_{12}) &= \sum_{i=1}^n [b_i \log L_2(t_i) + d_i \log L_4(t_i)] \\
 574 \quad &= \sum_{i \in S_1} [-b_i t_i \lambda_{11} + d_i \log(1 - e^{-t_i \lambda_{11}})] + \sum_{i \in S_2} \{-b_i [t_c \lambda_{11} + (t_i - t_c) \lambda_{12}] + d_i \log(1 \\
 575 \quad &\quad - e^{-t_c \lambda_{11} - (t_i - t_c) \lambda_{12}}\},
 \end{aligned}$$

576 for a special case with no deaths/cardiac failures, that is,  $a_i = c_i = 0$ ,  $\lambda_2 = \lambda_3 = 0$ , where  $S_1 =$   
 577  $\{i: t_i < t_c\}$  and  $S_2 = \{i: t_i \geq t_c\}$ . Hence, the estimates of  $\lambda_{11}$  and  $\lambda_{12}$  can be derived easily from  
 578 following score equations

$$\begin{aligned}
 579 \quad \sum_{i \in S_1} b_i t_i - \sum_{i \in S_1} \frac{d_i t_i}{e^{\lambda_{11} t_i} - 1} + \sum_{i \in S_2} b_i t_c - \sum_{i \in S_1} \frac{d_i t_c}{e^{\lambda_{11} t_c + (t_i - t_c) \lambda_{12}} - 1} &= 0 \\
 580 \quad \sum_{i \in S_2} b_i (t_i - t_c) - \sum_{i \in S_2} \frac{d_i (t_i - t_c)}{e^{\lambda_{11} t_c + (t_i - t_c) \lambda_{12}} - 1} &= 0.
 \end{aligned}$$

581 For a general model with piecewise constants in parameter  $\lambda_1$ , the log-likelihood function is

$$583 \quad l(\lambda_{11}, \lambda_{12}, \lambda_2, \lambda_3) = \sum_{i=1}^n [a_i \log L_1(t_i) + b_i \log L_2(t_i) + c_i \log L_3(t_i) + d_i \log L_4(t_i)]$$

582 where  $L_1(t) = \lambda_2 Q(t)$ ,  $L_2(t) = Q(t)$ ,  $L_4(t) = \dot{L}_3(t)/\lambda_3$  and if  $t < t_c$ ,

$$585 \quad L_3(t) = \frac{\lambda_{11} \lambda_3}{\lambda_{11} + \lambda_2 - \lambda_3} e^{-\lambda_3 t} (1 - e^{-(\lambda_{11} + \lambda_2 - \lambda_3) t})$$

586 if  $t \geq t_c$ ,

$$\begin{aligned}
 587 \quad L_3(t) &= \frac{\lambda_{11} \lambda_3}{\lambda_{11} + \lambda_2 - \lambda_3} e^{-\lambda_3 t_c} (1 - e^{-(\lambda_{11} + \lambda_2 - \lambda_3) t_c}) \\
 588 \quad &+ \frac{\lambda_{12} \lambda_3}{\lambda_{12} + \lambda_2 - \lambda_3} e^{-(\lambda_{11} - \lambda_{12}) t_c - \lambda_3 t} (e^{-(\lambda_{12} + \lambda_2 - \lambda_3) t_c} - e^{-(\lambda_{12} + \lambda_2 - \lambda_3) t}),
 \end{aligned}$$

589 in which  $Q_1(t) = e^{-\lambda_{11} t}$  if  $t < t_c$ , and  $Q_1(t) = e^{-(\lambda_{11} - \lambda_{12}) t_c - \lambda_{12} t}$  if  $t \geq t_c$ ;  $Q_2(t) = e^{-\lambda_2 t}$ ,  
 591  $Q_3(t|u) = e^{-\lambda_3(t-u)}$  and  $Q(t) = Q_1(t) Q_2(t)$ . Based on the log-likelihood equation, the  
 592 maximum likelihood estimates of  $\lambda_{11}$ ,  $\lambda_{12}$ ,  $\lambda_2$  and  $\lambda_3$  can be computed from the score equations  
 593 using numerical method. It is similar to derive the likelihood function for exponential model if  $\lambda_1$   
 594 has three or more pieces.

# Figures

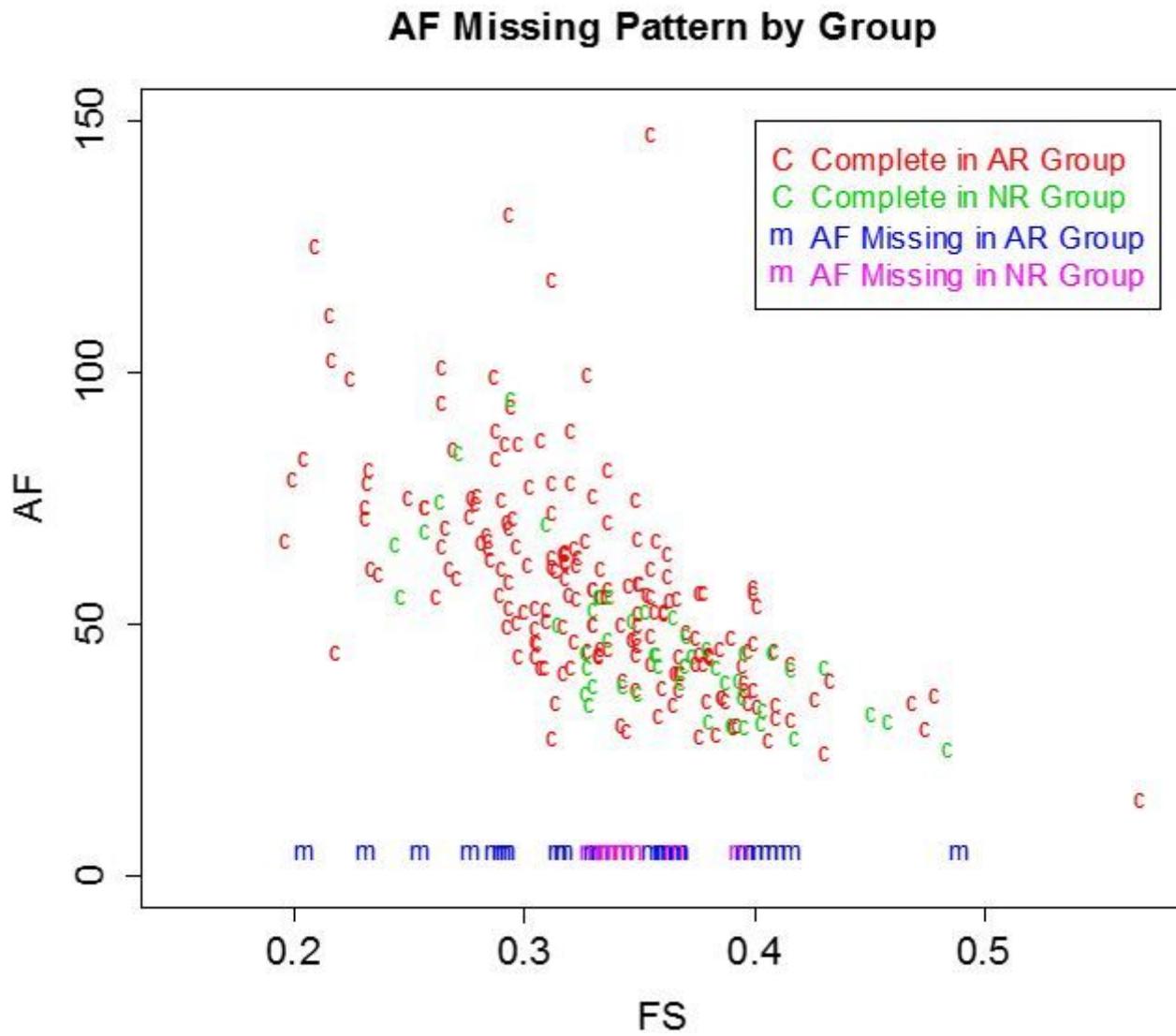
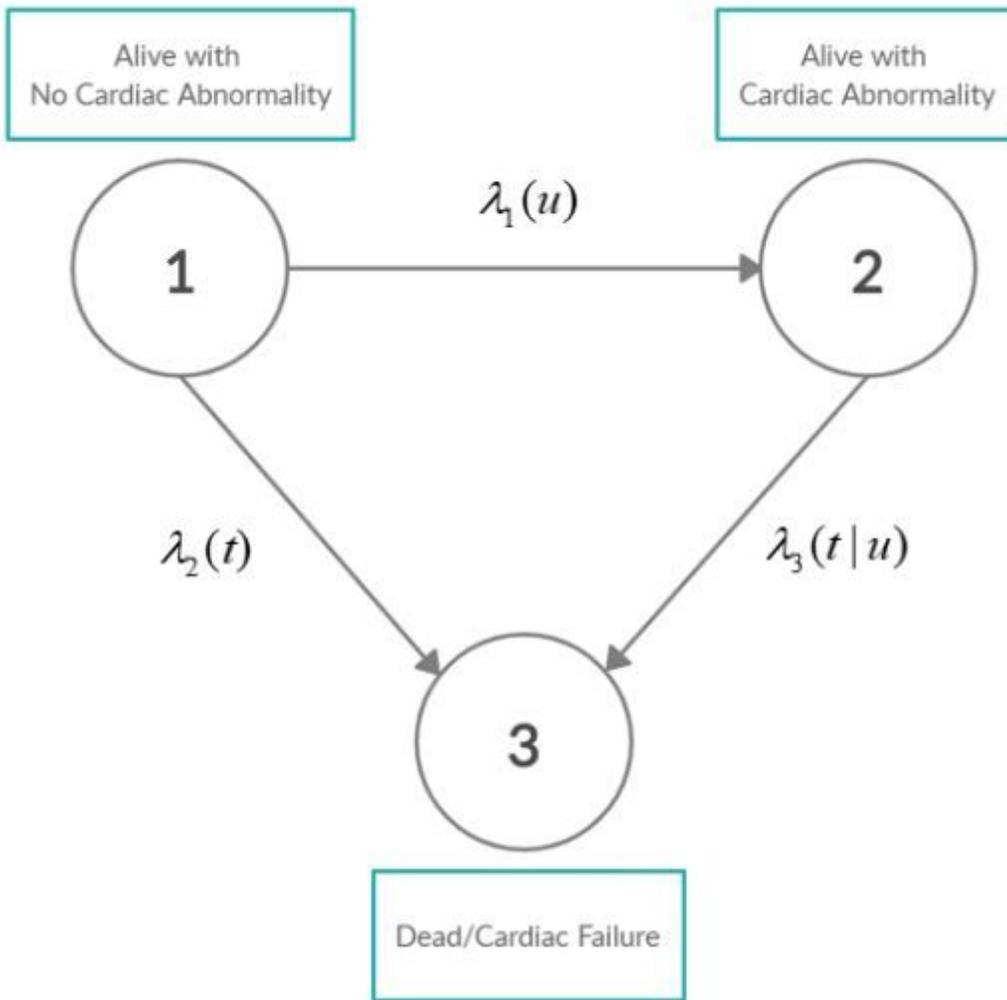


Figure 1

Scatter plot of AF and FS within AR and NR group of patients. Complete data (labeled c) displays strong correlation between AF and FS. Missing values of AF are in almost entire range of FS values (labeled m).



**Figure 2**

An abnormal cardiac measure-death/cardiac failure model involving three states. State 1 corresponds to patients who are alive with normal value. Patients who are alive with abnormal value are in state 2. State 3 is an absorbing state and corresponds to death or cardiac failure.

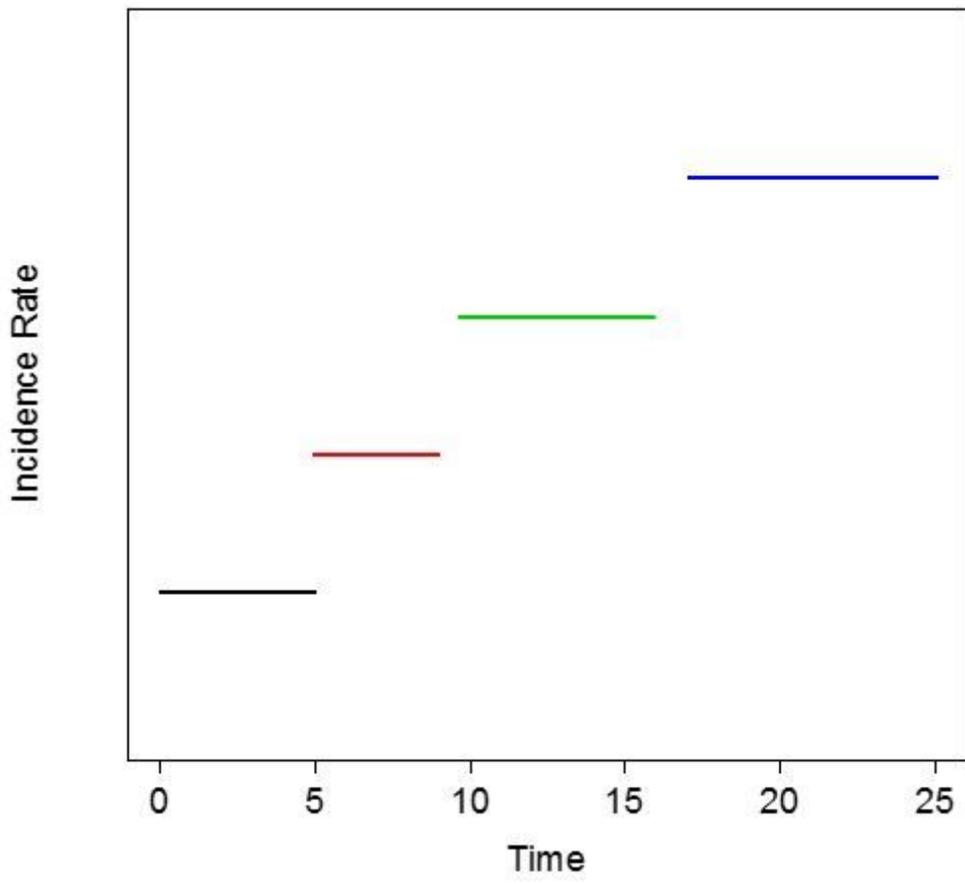
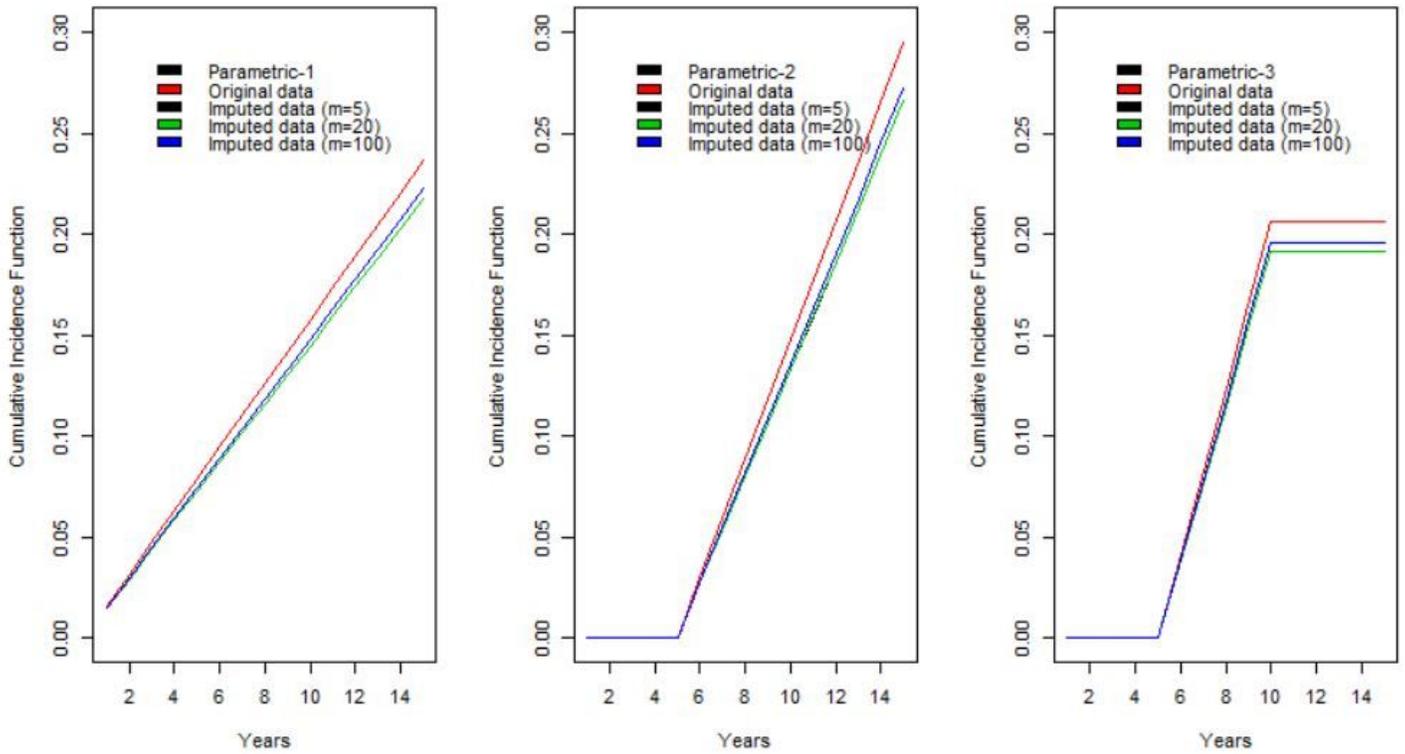


Figure 3

Piecewise Incidence Rate



**Figure 4**

Cumulative Incidence Comparison for AF Based on Original Data and Imputed Data (Using the mean of intensity estimates) Corresponding to  $m = 5, 20$  and  $100$ . Note: The lines for  $m=5$  and  $20$  are almost overlapped.