

7-2-2018

# Robust Estimation and Inference on Current Status Data with Applications to Phase IV Cancer Trial

Deo Kumar Srivastava

*St. Jude Children's Research Hospital*, [kumar.srivastava@stjude.org](mailto:kumar.srivastava@stjude.org)

Liang Zhu

*University of Texas Health Science Center at Houston*, [liang.zhu@uth.tmc.edu](mailto:liang.zhu@uth.tmc.edu)

Melissa M. Hudson

*St. Jude Children's Research Hospital*, [melissa.hudson@stjude.org](mailto:melissa.hudson@stjude.org)


Jianmin Pan

*University of Louisville*, [jianmin.pan@louisville.edu](mailto:jianmin.pan@louisville.edu)

Shesh N. Rai

*University of Louisville*, [Shesh.Rai@Louisville.Edu](mailto:Shesh.Rai@Louisville.Edu)

Follow this and additional works at: <https://digitalcommons.wayne.edu/jmasm>

 Part of the [Applied Statistics Commons](#), [Social and Behavioral Sciences Commons](#), and the [Statistical Theory Commons](#)

## Recommended Citation

Srivastava, Deo Kumar; Zhu, Liang; Hudson, Melissa M.; Pan, Jianmin; and Rai, Shesh N. (2018) "Robust Estimation and Inference on Current Status Data with Applications to Phase IV Cancer Trial," *Journal of Modern Applied Statistical Methods*: Vol. 17 : Iss. 1 , Article 18.

DOI: [10.22237/jmasm/1530544863](https://doi.org/10.22237/jmasm/1530544863)

Available at: <https://digitalcommons.wayne.edu/jmasm/vol17/iss1/18>

This Regular Article is brought to you for free and open access by the Open Access Journals at DigitalCommons@WayneState. It has been accepted for inclusion in *Journal of Modern Applied Statistical Methods* by an authorized editor of DigitalCommons@WayneState.

---

# Robust Estimation and Inference on Current Status Data with Applications to Phase IV Cancer Trial

## **Cover Page Footnote**

Acknowledgements: The research work of Deokumar Srivastava and Melissa Hudson was in part supported by the Cancer Center Support (CORE) grant CA 21765 and by the American Lebanese Syrian Associated Charities (ALSAC). Shesh Rai's research work was supported by Wendell Cherry Chair in Clinical Trial Research. The authors are thankful to the referees and the editor for their very helpful comments which led to a substantially improved manuscript.

# Robust Estimation and Inference on Current Status Data with Applications to Phase IV Cancer Trial

**Deo Kumar Srivastava**  
St. Jude Children's Research  
Hospital  
Memphis, TN

**Liang Zhu**  
University of Texas Health  
Science Center at Houston  
Houston, TX

**Melissa M. Hudson**  
St. Jude Children's Research  
Hospital  
Memphis, TN

**Jianmin Pan**  
University of Louisville  
Louisville, KY

**Shesh N. Rai**  
University of Louisville  
Louisville, KY

---

The use of piecewise exponential distributions was proposed by Rai et al. (2013) for analyzing cardiotoxicity data. Some parametric models are proposed, but the focus is on the Weibull distribution, which overcomes the limitation of piecewise exponential.

*Keywords:* Current status data, change point, constant hazard rate, Weibull distribution, phase IV clinical trial

---

## Introduction

With significant advancements in cancer treatment an increasing number of cancer survivors are living many years following a successful treatment. About a decade ago it was estimated nearly 13 million Americans were cancer survivors and over 379,000 were survivors of childhood and adolescent cancers (Mariotto et al., 2009). With this encouraging success comes the realization that survivors are at an increased risk of late adverse effects and of late mortality many years following cancer treatment. As described in Chow and Liu (2004), phase IV clinical trials are often used to document such long-term safety, toxicity, and mortality in cancer survivors.

Armstrong et al. (2009) showed cardiovascular events are the leading non-malignant cause of death among childhood cancer survivors with a 7-fold higher

---

risk of cardiovascular mortality compared to age-matched controls. Both chemotherapy and radiation therapy are known to be toxic to cardiomyocytes and contribute to early mortality. One such chemotherapy agent anthracycline is known to be cardiotoxic (Hudson et al., 2007), but because of its therapeutic benefits it remains one of the key components of the treatment plan. Often, in long-term follow-up studies the interest is not only on estimating the mortality but also on estimating the cumulative incidence of certain types of events, e.g. the onset of cardiotoxicity. With long-term follow-up studies it is common to see that the data from the survivor were not collected continuously in real time but at regular intervals (e.g., every six months or every year). The onset times for the events of interest are unknown, but the current status of each participant is known, such as whether the event of interest occurred in between the observation periods or not. This was characterized as case I interval censored data (Sun, 2006; Rai, 2008). Based on these data the prevalence of toxicity can be estimated, although the incidence rate is not straightforward.

The use of nonparametric methods for interval-censored data was discussed, for example, by Sun (2006), but the development of a parametric approach has lagged behind. Within the framework of parametric modeling and using likelihood theory one can easily use any one of the parametric models such as Exponential, Weibull, Log-normal, Gamma, Generalized Gamma, Log-logistic, and Generalized F, proposed in Kalbfleisch and Prentice (2002). Several more parametric distributions, such as the hypertextastic distribution proposed by Tabatabai, Bursac, Williams, and Singh (2007) or the generalized log-logistic proposed by Singh and Bartolucci (1997), were proposed for modeling survival data. However, the intensity functions for many of these distributions are not available in a nice closed form and may involve incomplete gamma or normal integrals.

J. K. Lindsey (1998) studied the parametric regression models to estimate the location and dispersion parameters, and compared the performance of nine different distributions. Most of the parametric models provided reasonably robust estimates, but the recommendation was against using exponential (unreasonable assumption of constant hazard over time), log-Student, and log-Cauchy distributions for their thick tails. J. C. Lindsey and Ryan (1998) proposed general approaches for interval censored data and concluded that parametric approaches can have satisfactory performance, especially if the Weibull or log-normal family was chosen, that allows for a reasonably wide range of distributional shape. They suggested a piecewise exponential distribution could be used to provide more flexibility in modeling as long as the number of intervals does not become too large. However, in utilizing a piecewise exponential distribution it remains necessary to either visualize the cut

## ROBUST INFERENCE FOR CURRENT STATUS DATA

points or adopt more rigorous approaches to first estimate the number of cut points and identify specific locations of the cut points. As an alternative, adopting the Weibull distribution provides reasonable flexibility in modeling monotone hazard shapes and performs quite well in comparison to other parametric models.

Often, as seen with the cardiotoxicity data to be discussed below, with the long-term follow-up data the cumulative incidence of the event of interest (cardiotoxicity) occurs at a low rate and remains stable over a few years following the end of the treatment, and then increases with longer follow-up. The identification of the change point is based on human input and visual inspection of the data, rather than a rigorous statistical approach. Rai et al. (2013) used a piecewise exponential distribution to model cardiotoxicity data. However, this approach had three key limitations: (1) it required knowledge of the time point when the incidence rate changes, (2) it required knowledge of how many change points may be needed, and (3) the incidence rate was assumed to be constant within each piece. Therefore, in order to avoid those limitations, the purpose of this study is to propose the use of parametric models for modeling such data and, in particular, focus on comparing the performance of Weibull model to the approach based on piece wise exponential distribution discussed in Rai et al.

### **Motivating Example**

A study to investigate cardiotoxic effects of anthracycline exposure during cancer treatment was described in Hudson et al. (2007). Specific diagnostic groups potentially at risk of cardiotoxicity were identified and recruited from a long-term follow-up clinic. The diagnostic group included survivors of childhood leukemia, lymphoma, sarcoma, and embryonal tumors who were all treated with anthracycline chemotherapy and/or radiation involving the heart, denoted by AR (at risk group). The control group comprised of survivors of acute lymphoblastic leukemia, Wilms tumor, and germ cell tumors who did not receive any cardiotoxic treatment, denoted by NR (not at risk group). The cardiotoxicity can be measured by several cardiac measures such as fractional shortening (FS), afterload (AF), QTc interval, and ejection fraction; see Hudson et al. (2007) and Krischer et al. (1997). Clinically,  $AF > 74 \text{ g/cm}^2$  can be used to identify patients with abnormal AF (AAF).

Of 278 patients who agreed to participate on the study 223 were designated as AR and 55 were designated as NR based on the treatment exposure. At the time of survey, data on each individual included demographics, the date of cancer diagnosis, time since treatment completion, disease-related variables (type, histology, and stage of cancer), treatment-related variables (chemotherapy drugs

and their doses, irradiation), and outcome-related cardiac measure (AF) and quality of life measures (general health, vitality, and physical health; see Cox et al., 2008). None of the patients had clinically defined cardiac dysfunction at the time of study evaluation. The AAF prevalence was 13.9%. Further details regarding the study can be found in Hudson et al. (2007).

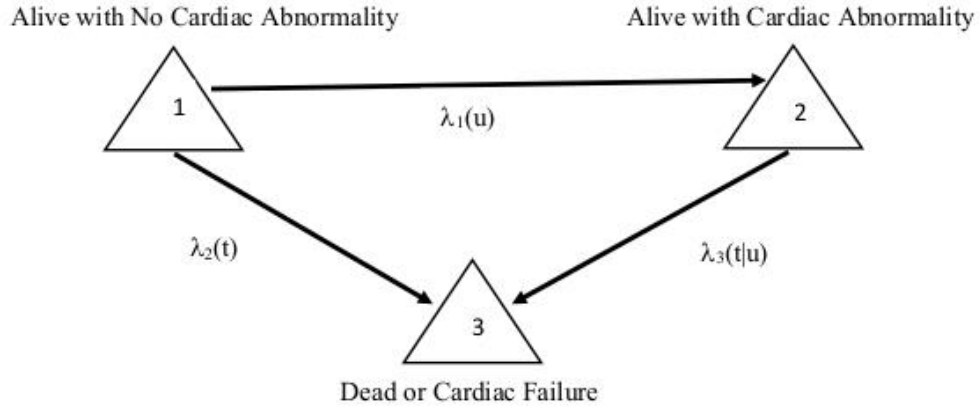
A common practice in estimation of the incidence rates and their confidence intervals is to assume the time of follow-up as the onset time and use the familiar and widely-used Kaplan-Meier approach (Pui, 2003; Kaplan & Meier, 1958). This approach is crude and provides biased estimates of cumulative incidence rates, as noted in Odell, Anderson, and D'Agostino (1992) and J. K. Lindsey (1998). In addition, there are several issues such as missingness among correlated measures of cardiotoxicity and accounting for competing risks (due to death or other toxicities) that make the analysis of such data more interesting and challenging. Furthermore, some patients who were potentially eligible to be enrolled on the study but died were not included in this study. Generalizing the results from the selected group of survivors leads to biased results, but alternative approaches such as those based on sampling weights could be adopted to obtain relatively unbiased estimates. Srivastava, Hudson, Robison, Wu, and Rai (2015) provided a relatively detailed account of the statistical issues involved with the design and analysis of cohort studies. Thus, the focus here will be on estimating the incidence rates of specific toxicities using a robust parametric approach.

## **Likelihood-Estimation for General Model**

Hudson et al. (2007) did not include patients who died or had cardiac failure during the treatment or during the follow-up. However, this information was available from the medical records. Therefore, the general theory for cross-sectional data, with indicators of cardiac abnormality and death, and time since treatment to survey or death is presented here.

Following Rai et al (2013), let  $T$  denote the observation time (death, cardiac failure, or survey) from the date of diagnosis and let  $U$  denote the time of cardiac abnormality (such as AAF) from the date of diagnosis, which is unknown. Note that  $T$  and  $U$  are measured from the date of cancer diagnosis and are not the current age of the patient.

## ROBUST INFERENCE FOR CURRENT STATUS DATA



**Figure 1.** An abnormal cardiac measure-death/cardiac failure model involving three states: State 1 – patients who are alive with no cardiac abnormality; State 2 – patients who are alive with abnormal cardiac measure; and State 3 – (an absorbing state) death or cardiac failure patients; the intensity  $\lambda_1(u)$ ,  $\lambda_2(t)$ , and  $\lambda_3(t|u)$  are the transition rates, where  $t$  is the observation time and  $u$  is the time of cardiac abnormality

The survival and prevalence functions are defined as

$$S(t) = Q(t) + \int_0^t \lambda_1(u) Q(u) Q_3(t|u) du \quad (1)$$

$$\pi(t) = [S(t) - Q(t)] / S(t) \quad (2)$$

The terms in equations (1) and (2) are functions of pseudo-survival functions defined as

$$Q_i(t) = \exp\left\{-\int_0^t \lambda_i(v) dv\right\} \text{ for } i = 1, 2, \quad Q_3(t|u) = \exp\left\{-\int_u^t \lambda_3(v|u) dv\right\},$$

$$Q(t) = Q_1(t) Q_2(t)$$

The parameter of interest is

$$\Lambda_1(t) = \int_0^t \lambda_1(u) du$$

the cumulative hazard function, but we observe  $\pi(t)$ .

Let  $\theta$ , the transition intensities, represent the full parametric vector. To construct the likelihood function, we summarize observations into 4 groups: 1) Death with No Cardiac Abnormality, 2) Alive with No Cardiac Abnormality, 3) Death/Cardiac Failure with Cardiac Abnormality, and 4) Alive with Cardiac Abnormality. Furthermore, let the corresponding contributions to the likelihood be  $L_1(t)$ - $L_4(t)$  and each individual contributes to only one term in the likelihood, where

$$L_1(t) = \lambda_2(t)Q(t); \quad L_2(t) = Q(t); \quad L_3(t) = \int_0^t \lambda_1(u)Q(u)\lambda_3(t|u)Q_3(t|u)du;$$

$$L_4(t) = \int_0^t \lambda_1(u)Q(u)Q_3(t|u)du$$

Note that the likelihood function depends on  $\theta$ , in addition to the observation time and status, but  $\theta$  is suppressed for notational convenience.

For parametric modeling, a variety of distributions can be used and the intensity functions for a selection of commonly-used parametric distributions are provided in Table 1. For the parametric modeling, any of the various parametric forms of the intensities, e.g. those listed in Table 1, can be considered, and the inference based on likelihood approach can be carried out. However, in light of the flexibility provided by the Weibull distribution in modeling monotone intensity

**Table 1.** Some parametric distributions and their corresponding intensity functions

Distribution	Density $f(t)$	Intensity function $\lambda(t)$
Exponential	$\eta e^{-\eta t}$	$\eta$
Weibull	$\eta \alpha (t)^{\alpha-1} \exp(-\eta t^\alpha)$	$\eta \alpha t^{\alpha-1}$
Log-Normal	$(2\pi)^{-1/2} \eta t^{-1} \exp(-\eta^2 (\log \eta t)^2 / 2)$	$f(t) / [1 - (\log t)]$
Gamma	$(t)^{k-1} e^{-t} / \Gamma(k)$	$\eta (\eta t)^{k-1} \exp(-\eta t) / \Gamma(k) [1 - I_k(\eta t)]$
Log-Logistic	$\eta \gamma (\eta t)^{\gamma-1} [1 + (\eta t)^\gamma]^{-2}$	$\eta \gamma (\eta t)^{\gamma-1} / [1 + (\eta t)^\gamma]$

Note:  $I_k(s)$  is the incomplete gamma integral



function and the ease of having a relatively simplified form of the likelihood function, details of the inference procedure under the Weibull distribution are provided.

### Inference under the Weibull Model

For reporting the results of survival analysis to the practitioners, the focus is often on estimating the cumulative incidence at fixed time points, such as 5 or 10 years, along with their corresponding 95% confidence intervals. The data observed for each patient, at a particular time, consists of observation time  $T$  and two status indicators:  $\delta$  (patient's status) and  $\gamma$  (patient's event status). Let  $(t_j, \delta_j, \gamma_j)$  be the triplet observation for the  $j^{\text{th}}$  subject;  $\delta_j = 1$  represents if the  $j^{\text{th}}$  subject had the event at time  $t_j$ , while  $\gamma_j = 1$  represents the  $j^{\text{th}}$  subject's abnormal value (cardiac abnormality) at time  $t_j$ .

The intensity functions defined by  $\lambda_i(t) = \eta_i \alpha_i t^{\alpha_i - 1}$  for  $i = 1, 2$ , leads to  $Q_i(t) = \exp(-\eta_i t^{\alpha_i})$ , and  $Q(t) = \exp(-\eta_1 t^{\alpha_1}) \exp(-\eta_2 t^{\alpha_2})$ . The intensity function  $\lambda_3(t|u)$  can be specified by  $\lambda_3^{\text{SM}}(t|u) = \eta_3 \alpha_3 t^{\alpha_3 - 1}$  under the assumption of a semi-Markov process which leads to  $Q_3^{\text{SM}}(t|u) = \exp(-\eta_3(t^{\alpha_3} - u^{\alpha_3}))$ , or as  $\lambda_3^{\text{M}}(t|u) = \eta_3 \alpha_3 (t-u)^{\alpha_3 - 1}$  under the assumption of a Markov process which leads to  $Q_3^{\text{M}}(t|u) = \exp(-\eta_3(t-u)^{\alpha_3})$ ; see Kalbfleisch and Lawless (1985) and Harezlak, Gao, and Hui (2003). The components of the likelihood contribution are:

$$L_1(t) = \lambda_2(t)Q(t) = \eta_2 \alpha_2 t^{\alpha_2 - 1} \exp(-\eta_1 t^{\alpha_1} - \eta_2 t^{\alpha_2})$$

$$L_2(t) = Q(t) = \exp(-\eta_1 t^{\alpha_1} - \eta_2 t^{\alpha_2})$$

$$\begin{aligned} L_3(t) &= \int_0^t Q(u) \lambda_1(u) Q_3^{\text{SM}}(t|u) \lambda_3^{\text{SM}}(t|u) du \text{ under semi-Markov assumption;} \\ &= \int_0^t Q(u) \lambda_1(u) Q_3^{\text{M}}(t|u) \lambda_3^{\text{M}}(t|u) du \text{ under Markov assumption} \end{aligned}$$

and

$$\begin{aligned}
 L_4(t) &= \int_0^t Q(u)\lambda_1(u)Q_3^{SM}(t|u)du \text{ under semi-Markov assumption;} \\
 &= \int_0^t Q(u)\lambda_1(u)Q_3^{SM}(t|u)du \text{ under Markov assumption}
 \end{aligned}$$

Then the log-likelihood function is presented as

$$\begin{aligned}
 l(\alpha_1, \eta_1, \alpha_2, \eta_2, \alpha_3, \eta_3) \\
 = \sum_{i=1}^n [a_i \log L_1(t_i) + b_i \log L_2(t_i) + c_i \log L_3(t_i) + d_i \log L_4(t_i)] \quad (3)
 \end{aligned}$$

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 corresponding to observations of type 1 to 4 discussed above. The likelihood contributions for rest of the distributions, except exponential, listed in Table 1 are cumbersome. This may be another reason why the Weibull distribution has been extensively used in practice.

The maximum likelihood estimates can be obtained by solving the non-linear equations obtained by taking the first order derivatives of the likelihood function in (3) with respect to the parameters and equating them to 0. The estimates are often obtained using the statistical software packages, such as R, that maximize the likelihood function in (3) directly or solve the non-linear set of equations, but it may be noted that the justification of the asymptotic normality of the estimates is for the solutions of the likelihood equations; see Rao (1973).

## Simulation Study

A simulation study was undertaken to compare the results obtained using piecewise exponential and Weibull distributions. The simulations were conducted using continuous time scale; see Rai (2008). In the continuous scale model with maximum follow-up of 10 years the events can occur at any time between 0 and 10 years. Two sample sizes ( $n = 100$  and  $n = 400$ ) were considered. Three different design settings were considered for the simulation study. For each scenario the estimates of the cumulative incidence (CI) and their standard errors (SEs) were obtained based on 5000 replications.

## ROBUST INFERENCE FOR CURRENT STATUS DATA

**Setting I.** In this setting, the data were generated from piecewise exponential distributions and the estimates of CI (at various time points) were obtained using piecewise exponential and Weibull distributions. Specifically, the data were generated using a piecewise exponential distribution with two pieces characterized by  $\lambda_{11} = \lambda_{11}(t)$  and  $\lambda_{12} = \lambda_{12}(t)$  of the parameter  $\lambda_1 = \lambda_1(t)$  described in Figure 1. The set of parameter values of  $(\lambda_{11}, \lambda_{12})$  chosen was  $(0.15, 0.4)$  to closely resemble the simulation parameters considered by Rai et al. (2013). To mimic the situation of longer follow-up for the majority of patients, 20% of the patients were expected to have shorter follow-up with intensity function  $\lambda_{11}$  and 80% were expected to have longer follow-up with intensity function  $\lambda_{12}$ . For the follow-up time of 10 years the change points considered were 2 and 5 years. Once the data was generated from piecewise exponential with  $\lambda_{11} = 0.15$  and  $\lambda_{12} = 0.4$ , then the estimates of the CI at various time points were estimated using 2-piece exponential and Weibull distributions. The results of the comparison are summarized in Table 2.

**Table 2.** Performance of Weibull model for the data generated from piecewise exponential distributions with  $\lambda_{11} = 0.15$  and  $\lambda_{12} = 0.4$

$n$	Follow-up time/Change point	Time	True CI	Piecewise Exponential		Weibull	
				CI	(SE)	CI	(SE)
100	10 years/ 2 year	1	0.14	0.14	(0.06)	0.14	(0.06)
		2	0.26	0.25	(0.11)	0.32	(0.08)
		3	0.50	0.51	(0.06)	0.50	(0.07)
		4	0.67	0.68	(0.05)	0.65	(0.06)
		5	0.78	0.78	(0.05)	0.76	(0.05)
		6	0.85	0.85	(0.04)	0.85	(0.04)
		7	0.90	0.90	(0.04)	0.90	(0.04)
		8	0.93	0.93	(0.03)	0.94	(0.03)
		9	0.95	0.95	(0.02)	0.96	(0.03)
		10	0.97	0.97	(0.02)	0.98	(0.02)
	10 years/ 5 year	1	0.14	0.14	(0.03)	0.10	(0.05)
		2	0.26	0.26	(0.05)	0.24	(0.07)
		3	0.36	0.36	(0.07)	0.37	(0.07)
		4	0.45	0.45	(0.08)	0.50	(0.06)
		5	0.53	0.52	(0.08)	0.61	(0.06)
		6	0.68	0.69	(0.06)	0.70	(0.05)
		7	0.79	0.79	(0.06)	0.78	(0.05)
		8	0.86	0.86	(0.06)	0.84	(0.05)
		9	0.90	0.90	(0.05)	0.88	(0.05)
10	0.94	0.93	(0.04)	0.91	(0.05)		

**Table 2 (continued).**

$n$	Follow-up time/Change point	Time	True CI	Piecewise Exponential		Weibull	
				CI	(SE)	CI	(SE)
400	10 years/ 2 year	1	0.14	0.14	(0.03)	0.14	(0.03)
		2	0.26	0.26	(0.05)	0.32	(0.04)
		3	0.50	0.51	(0.03)	0.50	(0.04)
		4	0.67	0.67	(0.03)	0.65	(0.03)
		5	0.78	0.78	(0.03)	0.76	(0.03)
		6	0.85	0.85	(0.02)	0.84	(0.02)
		7	0.90	0.90	(0.02)	0.90	(0.02)
		8	0.93	0.93	(0.02)	0.94	(0.02)
		9	0.95	0.95	(0.01)	0.96	(0.01)
		10	0.97	0.97	(0.01)	0.98	(0.01)
	10 years/ 5 year	1	0.14	0.14	(0.02)	0.10	(0.03)
		2	0.26	0.26	(0.03)	0.24	(0.03)
		3	0.36	0.36	(0.03)	0.38	(0.03)
		4	0.45	0.45	(0.04)	0.50	(0.03)
		5	0.53	0.53	(0.04)	0.61	(0.03)
		6	0.68	0.69	(0.03)	0.70	(0.03)
		7	0.79	0.79	(0.03)	0.78	(0.03)
		8	0.86	0.86	(0.03)	0.83	(0.03)
		9	0.90	0.90	(0.03)	0.88	(0.03)
		10	0.94	0.93	(0.02)	0.91	(0.02)

**Setting II.** In this setting, the data were generated from two different Weibull distributions corresponding to two different scenarios: (A)  $h(2) = 0.15$ ,  $h(10) = 0.4$  (slow increasing hazard) and (B)  $h(2) = 0.15$  and  $h(5) = 0.4$  (rapidly increasing hazard), where  $h(t)$  represents the value taken by the Weibull hazard function at time  $t$ . For each generated data set the cumulative incidence estimates were obtained using Weibull and piecewise exponential distributions with change point assumed to be at (a)  $t - 1$ , (b)  $t$ , and (c)  $t + 1$  years, where  $t$  is the time point at which  $h(t) = 0.15$  with the follow-up period of 10 years. The results of the comparison are summarized in Table 3.

**Setting III.** In this phase, the data was generated from piecewise exponential distributions as in Setting 1. However, since the change point is usually unknown, we assumed it to be at (a)  $t - 1$ , (b)  $t$ , and (c)  $t + 1$  years for the data generated with true change points at  $t = 2$  and  $t = 5$  years, respectively. The results of the comparison are summarized in Table 4.

## ROBUST INFERENCE FOR CURRENT STATUS DATA

**Table 3.** Performance of piecewise exponential (PE) with three assumed change points (a), (b), and (c) for the data generated from Weibull distributions corresponding to cases A and B with 10 year follow-up

<i>n</i>	Cases*	Time	True CI	Weibull		PE (a)		PE (b)		PE (c)	
				CI	(SE)	CI	(SE)	CI	(SE)	CI	(SE)
100	A	1	0.06	0.06	(0.03)	0.03	(0.06)	0.06	(0.04)	0.09	(0.04)
		2	0.17	0.17	(0.06)	0.22	(0.04)	0.12	(0.08)	0.17	(0.06)
		3	0.30	0.30	(0.07)	0.37	(0.05)	0.32	(0.05)	0.24	(0.09)
		4	0.43	0.43	(0.06)	0.50	(0.05)	0.48	(0.05)	0.44	(0.06)
		5	0.56	0.56	(0.06)	0.59	(0.05)	0.59	(0.06)	0.58	(0.06)
		6	0.66	0.67	(0.06)	0.67	(0.05)	0.69	(0.06)	0.69	(0.06)
		7	0.75	0.76	(0.06)	0.74	(0.05)	0.76	(0.06)	0.76	(0.06)
		8	0.82	0.83	(0.06)	0.79	(0.05)	0.81	(0.05)	0.82	(0.05)
		9	0.88	0.88	(0.05)	0.83	(0.05)	0.85	(0.05)	0.86	(0.05)
		10	0.92	0.92	(0.05)	0.86	(0.04)	0.88	(0.04)	0.90	(0.04)
	B	1	0.03	0.04	(0.02)	0.10	(0.03)	0.04	(0.04)	0.14	(0.03)
		2	0.13	0.13	(0.05)	0.18	(0.06)	0.07	(0.07)	0.26	(0.05)
		3	0.28	0.28	(0.07)	0.26	(0.08)	0.35	(0.05)	0.36	(0.06)
		4	0.46	0.45	(0.07)	0.33	(0.09)	0.55	(0.05)	0.45	(0.07)
		5	0.62	0.62	(0.06)	0.62	(0.06)	0.68	(0.06)	0.52	(0.08)
		6	0.76	0.76	(0.06)	0.78	(0.06)	0.78	(0.05)	0.59	(0.08)
		7	0.86	0.86	(0.05)	0.87	(0.05)	0.84	(0.05)	0.84	(0.06)
		8	0.92	0.93	(0.04)	0.92	(0.04)	0.89	(0.04)	0.93	(0.04)
		9	0.96	0.96	(0.02)	0.95	(0.03)	0.92	(0.03)	0.96	(0.03)
		10	0.98	0.98	(0.02)	0.98	(0.02)	0.94	(0.03)	0.98	(0.02)
400	A	1	0.14	0.06	(0.02)	0.03	(0.03)	0.06	(0.02)	0.09	(0.02)
		2	0.26	0.17	(0.03)	0.22	(0.02)	0.12	(0.04)	0.17	(0.03)
		3	0.50	0.30	(0.03)	0.37	(0.02)	0.32	(0.03)	0.24	(0.05)
		4	0.67	0.43	(0.03)	0.49	(0.03)	0.47	(0.03)	0.43	(0.03)
		5	0.78	0.56	(0.03)	0.59	(0.03)	0.59	(0.03)	0.58	(0.03)
		6	0.85	0.67	(0.03)	0.67	(0.03)	0.68	(0.03)	0.68	(0.03)
		7	0.90	0.75	(0.03)	0.74	(0.03)	0.75	(0.03)	0.76	(0.03)
		8	0.93	0.82	(0.03)	0.79	(0.03)	0.81	(0.03)	0.82	(0.03)
		9	0.95	0.88	(0.03)	0.83	(0.02)	0.85	(0.02)	0.87	(0.03)
		10	0.97	0.92	(0.02)	0.86	(0.02)	0.88	(0.02)	0.90	(0.02)
	B	1	0.14	0.03	(0.01)	0.10	(0.02)	0.04	(0.02)	0.14	(0.01)
		2	0.26	0.14	(0.03)	0.18	(0.03)	0.07	(0.03)	0.26	(0.02)
		3	0.36	0.28	(0.04)	0.26	(0.04)	0.35	(0.02)	0.36	(0.03)
		4	0.45	0.46	(0.04)	0.33	(0.05)	0.54	(0.03)	0.45	(0.04)
		5	0.53	0.62	(0.03)	0.61	(0.03)	0.68	(0.03)	0.52	(0.04)
		6	0.68	0.76	(0.03)	0.78	(0.03)	0.77	(0.03)	0.59	(0.04)
		7	0.79	0.86	(0.02)	0.87	(0.02)	0.84	(0.02)	0.83	(0.03)
		8	0.86	0.92	(0.02)	0.92	(0.02)	0.89	(0.02)	0.93	(0.02)
		9	0.90	0.96	(0.01)	0.95	(0.01)	0.92	(0.02)	0.97	(0.01)
		10	0.94	0.98	(0.01)	0.97	(0.01)	0.94	(0.01)	0.99	(0.01)

Note: \*For a description of Weibull distributions and other parameters see Setting II in Simulations Study section

**Table 4.** Performance of Weibull model and piecewise exponential PE) distribution for varying change points when the data is generated from piecewise exponential distributions with  $\lambda_{11} = 0.15$  and  $\lambda_{12} = 0.4$

<i>n</i>	Follow-up time/ Change point	Time	Weibull			PE (a)		PE (b)		PE (b)	
			True CI	CI	(SE)	CI	(SE)	CI	(SE)	CI	(SE)
100	10 years/ 2 year	1	0.06	0.14	(0.06)	0.09	(0.09)	0.14	(0.06)	0.17	(0.05)
		2	0.17	0.32	(0.08)	0.36	(0.06)	0.25	(0.11)	0.31	(0.08)
		3	0.30	0.50	(0.07)	0.55	(0.05)	0.52	(0.06)	0.43	(0.10)
		4	0.43	0.65	(0.06)	0.68	(0.05)	0.68	(0.05)	0.65	(0.06)
		5	0.56	0.76	(0.05)	0.78	(0.05)	0.78	(0.05)	0.78	(0.05)
		6	0.66	0.85	(0.04)	0.84	(0.04)	0.85	(0.04)	0.86	(0.05)
		7	0.75	0.90	(0.04)	0.89	(0.04)	0.90	(0.04)	0.91	(0.04)
		8	0.82	0.94	(0.03)	0.92	(0.03)	0.93	(0.03)	0.94	(0.03)
		9	0.88	0.96	(0.03)	0.94	(0.03)	0.95	(0.02)	0.96	(0.03)
		10	0.92	0.98	(0.02)	0.96	(0.02)	0.97	(0.02)	0.97	(0.02)
	10 years/ 5 year	1	0.03	0.10	(0.05)	0.13	(0.03)	0.14	(0.03)	0.15	(0.03)
		2	0.13	0.24	(0.07)	0.25	(0.06)	0.26	(0.05)	0.27	(0.05)
		3	0.28	0.37	(0.07)	0.34	(0.08)	0.36	(0.07)	0.38	(0.06)
		4	0.46	0.50	(0.06)	0.43	(0.09)	0.45	(0.08)	0.47	(0.07)
		5	0.62	0.61	(0.06)	0.60	(0.06)	0.52	(0.08)	0.55	(0.07)
		6	0.76	0.70	(0.05)	0.71	(0.06)	0.69	(0.06)	0.61	(0.08)
		7	0.86	0.78	(0.05)	0.79	(0.06)	0.79	(0.06)	0.77	(0.06)
		8	0.92	0.84	(0.05)	0.85	(0.05)	0.86	(0.06)	0.86	(0.06)
		9	0.96	0.88	(0.05)	0.89	(0.05)	0.90	(0.05)	0.91	(0.05)
		10	0.98	0.91	(0.05)	0.92	(0.04)	0.93	(0.04)	0.94	(0.05)
400	10 years/ 2 year	1	0.14	0.14	(0.03)	0.10	(0.04)	0.14	(0.03)	0.17	(0.03)
		2	0.26	0.32	(0.04)	0.36	(0.03)	0.26	(0.05)	0.32	(0.04)
		3	0.50	0.50	(0.04)	0.55	(0.03)	0.51	(0.03)	0.44	(0.05)
		4	0.67	0.65	(0.03)	0.68	(0.03)	0.67	(0.03)	0.64	(0.03)
		5	0.78	0.76	(0.03)	0.77	(0.02)	0.78	(0.03)	0.77	(0.03)
		6	0.85	0.84	(0.02)	0.84	(0.02)	0.85	(0.02)	0.85	(0.02)
		7	0.90	0.90	(0.02)	0.89	(0.02)	0.90	(0.02)	0.90	(0.02)
		8	0.93	0.94	(0.02)	0.92	(0.02)	0.93	(0.02)	0.94	(0.02)
		9	0.95	0.96	(0.01)	0.94	(0.01)	0.95	(0.01)	0.96	(0.01)
		10	0.97	0.98	(0.01)	0.96	(0.01)	0.97	(0.01)	0.97	(0.01)
	10 years/ 5 year	1	0.14	0.10	(0.03)	0.13	(0.02)	0.14	(0.02)	0.15	(0.01)
		2	0.26	0.24	(0.03)	0.25	(0.03)	0.26	(0.03)	0.27	(0.02)
		3	0.36	0.38	(0.03)	0.35	(0.04)	0.36	(0.03)	0.38	(0.03)
		4	0.45	0.50	(0.03)	0.43	(0.04)	0.45	(0.04)	0.47	(0.03)
		5	0.53	0.61	(0.03)	0.59	(0.03)	0.53	(0.04)	0.55	(0.04)
		6	0.68	0.70	(0.03)	0.71	(0.03)	0.69	(0.03)	0.62	(0.04)
		7	0.79	0.78	(0.03)	0.79	(0.03)	0.79	(0.03)	0.77	(0.03)
		8	0.86	0.83	(0.03)	0.85	(0.03)	0.86	(0.03)	0.86	(0.03)
		9	0.90	0.88	(0.03)	0.89	(0.02)	0.90	(0.03)	0.91	(0.03)
		10	0.94	0.91	(0.02)	0.92	(0.02)	0.94	(0.02)	0.96	(0.02)

## ROBUST INFERENCE FOR CURRENT STATUS DATA

From Table 2, for Setting I, it is clear that when the data are generated from piecewise exponential distributions the estimates obtained using piecewise exponential are almost unbiased, whereas those based on a Weibull distribution are slightly biased with the most bias occurring at the change point. For example, when  $n = 100$ , the true value of the CI at 2 years is 0.26, whereas the estimates of CI corresponding to the piecewise exponential and Weibull distributions are 0.25 and 0.32, respectively. Similarly, for  $n = 400$ , the true value of the CI at 8 years is 0.93 and the estimates corresponding to piecewise exponential and Weibull are 0.93 and 0.94, respectively.

Table 3 summarizes the findings from the simulation experiment for Setting II corresponding to sample sizes 100 and 400. This is the situation where the data are generated from Weibull distributions and both Weibull and piecewise exponential distributions are used to obtain the estimates. For Case B, the rapidly rising hazard situation, when the sample size is 100, the true values of the CI corresponding to 2, 5, and 10 years are 0.13, 0.62, and 0.98, respectively. The corresponding estimates based on Weibull and piecewise exponential with cut point at 5 years are 0.13, 0.62, 0.98, and 0.07, 0.68, 0.94, respectively, suggesting that the estimates obtained using the exponential distribution are biased. Furthermore, if the cut point is not guessed correctly, then these estimates corresponding to cut point at 4 or 6 years are 0.18, 0.62, 0.98 and 0.26, 0.52, 0.98, respectively, suggesting that the bias could be more pronounced if the cut points are not appropriately chosen. The findings are similar for sample size 400.

The simulation results corresponding to Setting III are summarized in Table 4. This is the setting where the data are generated from piecewise exponential distributions with cut point at  $t = 2$  or  $t = 5$  years and the estimates are obtained using piecewise exponential distributions with cut points assumed to be at  $t - 1$ ,  $t$ , and  $t + 1$ . It is seen that if the cut point is not appropriately chosen then the estimates can be significantly biased. For example, when  $n = 100$  and change point is at  $t = 5$ , then the estimates of the CI at 2, 5, and 10 years from the piecewise exponential distribution are 0.25, 0.60, and 0.92 at 4 years, 0.26, 0.52, and 0.93 at 5 years, and 0.27, 0.55, and 0.94 at 6 years, whereas the true values are 0.26, 0.53, and 0.94, respectively.

### **Application to a Phase IV Cancer Trial**

The approach based on a semi-Markov assumption proposed above was applied to the motivating example and then compared with the non-parametric approach as in Sun (2006) and the piecewise exponential proposed by Rai et al. (2013). The study

cohort reported in Hudson et al. (2007) did not include the patients who died or had already experienced cardiac failure. Rai et al. adopted a piecewise exponential distribution of model Cardiac data reported in Hudson et al. and in a subjective manner chose the change point to be at 5 years, i.e.  $t_c = 5$ , and assumed the hazard rates to be  $\lambda_{11}$  for  $t < t_c$  and  $\lambda_{12}$  for  $t \geq t_c$ . In view of the data collected they obtained likelihood solutions for the special case  $\lambda_2 = \lambda_3$ . They also examined the usefulness of having more than one change point and modeled the data with three piecewise exponential distributions. However, the results from the data analysis suggested that the group effects were better detected with two piecewise exponential distributions compared to three piecewise exponential distributions and was used to interpret the findings. Here, the Weibull distribution was used for modeling current status data and applied it to Cardiac data. The maximum likelihood estimates for the special case of  $\lambda_2 = \lambda_3 = 0$ ,  $a_i = c_i = 0$  were obtained using the statistical software package R.

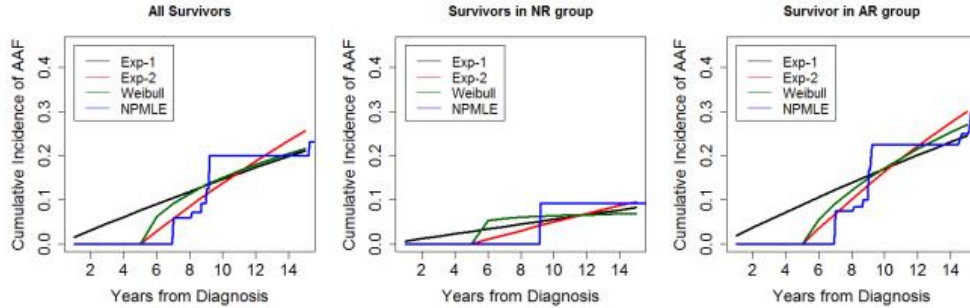
Because Rai et al. (2013) concluded that the piecewise exponential distribution provided better results than the interval censored approach, as implemented in SAS procedure LIFEREG, the current results were not compared with the interval censored approach. Thus, four approaches to compute the incidence rates were used and compared. One is the nonparametric approach based on Sun (2006). Because there are very few events prior to 5 years, two types of

**Table 5.** Cumulative incidence functions for AAF

Year	Nonparametric CI	Weibull		Exp-1		Exp-2	
		CI	SE	CI	SE	CI	SE
1	0.0000	0.0000	0.0000	0.0160	0.0030	0.0000	0.0000
2	0.0000	0.0000	0.0000	0.0320	0.0050	0.0000	0.0000
3	0.0000	0.0000	0.0000	0.0470	0.0080	0.0000	0.0000
4	0.0000	0.0000	0.0000	0.0630	0.0100	0.0000	0.0000
5	0.0000	0.0000	0.0000	0.0790	0.0130	0.0000	0.0000
6	0.0000	0.0620	0.0250	0.0950	0.0160	0.0300	0.0050
7	0.0590	0.0920	0.0280	0.1100	0.0180	0.0590	0.0100
8	0.0590	0.1140	0.0320	0.1260	0.0210	0.0890	0.0150
9	0.1250	0.1330	0.0370	0.1420	0.0230	0.1190	0.0200
10	0.2000	0.1500	0.0430	0.1580	0.0260	0.1490	0.0250
11	0.2000	0.1650	0.0480	0.1740	0.0290	0.1780	0.0300
12	0.2000	0.1790	0.0540	0.1890	0.0310	0.2080	0.0350
13	0.2000	0.1920	0.0600	0.2050	0.0340	0.2380	0.0400
14	0.2000	0.2040	0.0660	0.2210	0.0370	0.2670	0.0450
15	0.2000	0.2150	0.0720	0.2370	0.0390	0.2970	0.0500
20	0.2500	0.2640	0.0980	0.3160	0.0520	0.4460	0.0750



## ROBUST INFERENCE FOR CURRENT STATUS DATA



**Figure 2.** Cumulative incidence function for the four models

exponential models were considered: the first with a constant incidence rate (denoted by Exp-1), the second with two incidence rates (one up to 5 years as zero and the second beyond 5 years as non-zero positive constant, i.e. piecewise exponential with two pieces denoted by Exp-2). The fourth approach is to obtain the incidence rates using the proposed Weibull distribution, denoted by Weibull. The cumulative incidence functions and their standard errors based on these approaches are presented in Table 5 and Figure 2 for AAF.

The cumulative incidence estimates, along with their 95% confidence intervals at specific time points, are provided in Table 6. To avoid the possibility of the lower limits of the confidence intervals going below 0, we obtain the confidence limits on the log-scale using log-transformation in conjunction with the delta method, and then transformed the limits back to obtain the confidence limits in original scale, which are also reported in Table 6.

**Table 6.** Cumulative incidence (CI) functions and 95% Confidence Intervals (CIs) at fixed time points for Weibull fit to AAF data

Year	CI	Confidence intervals based on MLE	Confidence intervals using log-transformation approach
1	0.000	(0.000, 0.000)	(0.000, 0.000)
3	0.000	(0.000, 0.000)	(0.000, 0.000)
5	0.000	(0.000, 0.000)	(0.000, 0.000)
10	0.150	(0.066, 0.234)	(0.086, 0.263)
15	0.215	(0.074, 0.356)	(0.112, 0.414)
20	0.264	(0.072, 0.456)	(0.128, 0.547)

The effect of anthracycline exposure on the cumulative incidence of cardiac abnormality was evaluated using likelihood ratio test by fitting the proposed models for the two groups (AR and NR) independently and then after combing them together. The  $p$ -values corresponding to piecewise exponential (Exp-2), Weibull, and logistic regression are 0.012, 0.044, and 0.065, respectively. The proposed approach supports the finding obtained by the piecewise exponential distribution but not by logistic regression that those at risk (AR group) have a higher CI of developing abnormal afterload (AAF).

## Acknowledgements

Acknowledgements: The research work of Deokumar Srivastava and Melissa Hudson was in part supported by the Cancer Center Support (CORE) grant CA 21765 and by the American Lebanese Syrian Associated Charities (ALSAC). Shesh Rai's research work was supported by Wendell Cherry Chair in Clinical Trial Research. The authors are thankful to the referees and the editor for their very helpful comments which led to a substantially improved manuscript.

## Conclusions

The approach based on the Weibull model is better because it is able to capture the group effect appropriately without having to make any assumptions about the cut point and relaxing the restrictive and unrealistic assumption of constant hazard rates for the two pieces. In addition, it is seen that the CI estimates based on the Weibull model closely match the ones obtained using the nonparametric approach. Thus, the estimates and test based on the Weibull distribution may be capturing the underlying hazard pattern appropriately without making restrictive assumptions about change point or the hazards being constant within particular time period.

There will always be situations for which the use of a piecewise exponential distribution might be more appropriate (e.g., assumptions underlying the use of piecewise exponential distribution would be valid and more appropriate), and in those situations more efficient estimates of the parameters and cumulative incidence rate can be obtained. However, even in such situations, based on the simulation results it is seen that the results obtained from the Weibull distribution are reasonable, with most discrepancies observed near the location of the change points. Thus, this approach is robust in terms of detecting the group effects, providing reasonable fixed term adverse-effect-incidence rates along with

## ROBUST INFERENCE FOR CURRENT STATUS DATA

confidence intervals. This can be implemented as a routine analysis for similar studies at ours and other cancer centers.

Often, when designing similar studies, questions are asked about required sample size for comparing adverse event rates at fixed time-points. The fixed time-point estimates and confidence intervals obtained using Weibull model along with design parameters can be used to justify the required sample size. Another extension of this work is to study other types of long term adverse effects such as kidney stones in long-term survivors of childhood acute lymphoblastic leukemia, Kaste et al, (2009). This can be modelled using a competing risk model or a bivariate model that we plan to study further.

### References

- Armstrong, G. T., Liu, Q., Yasui, Y., Neglia, J. P., Leisenring, W. L., Robison, L. L., & Mertens, A. C. (2009). Late mortality among 5-year survivors of childhood cancer: A summary from the Childhood Cancer Survivor Study. *Journal of Clinical Oncology*, 27(14), 2328-2338. doi: 10.1200/jco.2008.21.1425
- Chow, S. C., & Liu, J. P. (2004). *Design and analysis of clinical trials: Concepts and methodologies* (2nd ed.). Hoboken, NJ: John Wiley & Sons, Inc.
- Harezlak, J., Gao, S., & Hui, S. L. (2003). An illness-death stochastic model in the analysis of longitudinal dementia data. *Statistics in Medicine*, 22(9), 1465-1475. doi: 10.1002/sim.1506
- Hudson, M. M., Rai, S. N., Nunez, C., Merchant, T. E., Marina, N. M., Zalamer, N.,...Rosenthal, D. (2007). Noninvasive evaluation of late anthracycline cardiac toxicity in childhood cancer survivors. *Journal of Clinical Oncology*, 25(24), 3635 - 3643. doi: 10.1200/jco.2006.09.7451
- Kalbfleisch, J. D., & Lawless, J. F. (1985). The analysis of panel data under a Markov assumption. *Journal of the American Statistical Association*, 80(392), 863-871. doi: 10.1080/01621459.1985.10478195
- Kalbfleisch, J. D., & Prentice, R. L. (2002). *The statistical analysis of failure time data*. Hoboken, NJ: John Wiley & Son, Inc. doi: 10.1002/9781118032985
- Kaplan, E. L., & Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282), 457-481. doi: 10.1080/01621459.1958.10501452

- Krischer, J. P., Epstein, S., Cuthbertson, D. D., Goorin, A. M., Epstein, M. L., & Lipshultz, S. E. (1997). Clinical cardiotoxicity following anthracycline treatment for childhood cancer: The Pediatric Oncology Group experience. *Journal of Clinical Oncology*, *15*(4), 1544-1552. doi: 10.1200/jco.1997.15.4.1544
- Lindsey, J. C., & Ryan, L. M. (1998). Tutorial in biostatistics methods for interval-censored data. *Statistics in Medicine*, *17*(2), 219-238. doi: 10.1002/(SICI)1097-0258(19980130)17:2<219::AID-SIM735>3.0.CO;2-O
- Lindsey, J. K. (1998). A study of interval censoring in parametric regression models. *Lifetime Data Analysis*, *4*(4), 329-354. doi: 10.1023/a:1009681919084
- Mariotto, A. B., Rowland, J. H., Yabroff, K. R., Scoppa, S., Hachey, M., Ries, L., & Feuer, E. J. (2009). Long-term survivors of childhood cancers in the United States. *Cancer Epidemiology, Biomarker and Prevention*, *18*(4), 1033-1040. doi: 10.1158/1055-9965.epi-08-0988
- Odell, P. M., Anderson, K. M., & D'Agostino, R. B. (1992). Maximum likelihood estimation for interval-censored data using a Weibull-based accelerated failure time model. *Biometrics*, *48*(3), 951-959. doi: 10.2307/2532360
- Pui, C.-H., Cheng, C., Leung, W., Rai, S. N., Rivera, G. K., Sandlund, J. T., ... Hudson, M. M. (2003). Extended follow-up of long-term survivors of childhood acute lymphoblastic leukemia. *The New England Journal of Medicine*, *349*(7), 640-649. doi: 10.1056/nejmoa035091
- Rai, S. N. (2008). Analysis of occult tumor studies. In W.-Y. Tan & L. Hanin (Eds.), *Handbook of cancer models with applications*. Hackensack, NJ: World Scientific Press. doi: 10.1142/9789812779489\_0018
- Rai, S. N., Pan, J., Yuan, X., Sun, J., Hudson, M. M., & Srivastava, D. K. (2013). Estimating incidence rate on current status data with application to a phase IV cancer trial. *Communications in Statistics – Theory and Methods*, *42*(17), 3117-3135. doi: 10.1080/03610926.2011.620208
- Rao, C. R. (1973). *Linear statistical inference and its applications*. New York: John Wiley & Sons.
- Singh, K. P., & Bartolucci, W. A. A. (1997). A generalized log-logistic model for analysis of environmental pollutant data. In *Proceedings of the International Congress on Modeling and Simulation*. Hobart, Tasmania.
- Srivastava, D. K., Hudson, M. M., Robison, L. L., Wu, X., & Rai, S. N. (2015). Design and analysis of cohort studies: Issues and practices. *Biometrics & Biostatistics International Journal*, *2*(5), 00044. doi: 10.15406/bbij.2015.02.00044

## ROBUST INFERENCE FOR CURRENT STATUS DATA

Sun, J. (2006). *The statistical analysis of interval-censored failure time data*. New York: Springer. doi: 10.1007/0-387-37119-2

Tabatabai, M. A., Bursac, Z., Williams, D. K., & Singh, K. P. (2007). Hypertabastic survival model. *Theoretical Biology and Medical Modeling*, 4(40). doi: 10.1186/1742-4682-4-40