Available online at www.elixirpublishers.com (Elixir International Journal)

Bio Technology

Elixir Bio. Tech. 37 (2011) 3726-3731

Comparison of bootstrap survival function with the Kaplan-Meier survival function under the influence of different percentage of censoring

Qamruz Zaman

Department of Medical Statistics, Informatics and Health Economics, Innsbruck Medical University, Department of Statistics, University of Peshawar, Pakistan.

ARTICL	Е	INFO
Article	hi	story:

17 July 2011;

Keywords

Received: 4 June 2011;

Received in revised form:

Accepted: 27 July 2011;

If survival data follows a specific survival distribution (weibull, exponential), survival probabilities can be estimated by using the specific survival function. Alternate easy and assumptions free methods for computing the survival probabilities are Kaplan-Meier and Bootstrap survival functions. The aim of present study is to compare the two non-parametric methods under the influence of different percentage of censoring presented in four different data sets "Leukaemia data set, Stanford Heart Transplant data set, Thalassaemia data set and Lung cancer data set". Results of analysis show that the performance of Kaplan-Meier survival function is better than the Bootstrap survival function. This is due to the easy concept, less time, less time, easily availability as well as in terms of survival probabilities.

© 2011 Elixir All rights reserved.

Kaplan-Meier Survival Function, Bootstrap Survival Function, Censoring, Events.

Introduction

Survival analysis is considered as the backbone of medical research and is based on the observed time of event.

The event of interest may for instance be death or recovery from treatment. Survival analysis is based on three techniques i) parametric ii) Semi-parametric and iii) Non-parametric.

Non-parametric technique is more commonly used technique (Fleming & Harrington, 1984), due to relax conditions about the assumptions of distributions.

Kaplan-Meier survival function (1958) is the most famous and commonly used non-parametric technique of survival analysis. To define the function we proceed as:

Let the survival times $X_1, X_2, ..., X_n$ be independently identically distributed according to the distribution function F (x).

Similarly, let $Y_1, Y_2, ..., Y_n$ be independently identically distributed censoring times according to G(y). In addition, the survival times X_i and censoring times Y_i are assumed to be independent (Miller and Rupert, 1983).

The observable random variables are $T_i = \min \{X_i, Y_i\}$ and $\delta_i = I (X_i \leq Y_i)$ indicates whether the survival time is uncensored or censored. Let $T_1 < \ldots < T_n$ denote the ordered observed survival times, and let $\delta_1,\ ...,\ \delta_n$ be their corresponding (unordered) indicator values.

Let the number of individuals who are alive just before time t_i, including those who are about to die at this time, be denoted by r_i and e_i denotes the number who die at this time.

Using these notations, the Kaplan-Meier estimator is defined as

$$\hat{S}_{KM}(t) = \prod_{i=1}^{n} \left(\frac{r_i - e_i}{r_i} \right)$$

Another non-parametric similar approach for calculating the survival function is the Bootstrap survival function.

A brief introduction of the function is defined below:

Bootstrap Method

The Bootstrap (Effron, 1979) is a method for calculating the approximated biases, standard deviations, confidence intervals etc. Except these applications, Bootstrap also doing a reasonable job under a variety of situations. e.g. in survival analysis, one can use the Bootstrap to estimate the survival probabilities (as by Kaplan-Meier, 1958).

Sampling scheme for Bootstrapping censored data was introduced by Efron for 'Bootstraping individual' k, sample a survival time and censoring indicator pair (t_i, δ_i) with replacement from the data set { (t_i, δ_i) , i=1, ..., n}. If n Bootstrap individuals are generated, a Bootstrap replicate can be found by applying the Kaplan-Meier estimator to the Bootstrap data set. We take M of these replicates and order them as

$$\hat{S}_{[1]}^{B}(t^{*}) \leq \ldots \leq \hat{S}_{[M]}^{B}(t^{*})$$

The procedure is describe as

Draw Bootstrap sample of the pair (ti, δ_i), and note $m_i = \#$ times (t_i, δ_i) appears in the Bootstrap sample, so $m = (m_1, m_2, ..., m_n)$ is an n-category multinomial, n draws, probability 1/n for each category: $m \sim mult(n, 1/n)$ Define

$$M_{j} = \sum_{i=j}^{n} m_{i}$$
, j= 1, 2,..., n

Where

 $M_1 = n$, $M_2 = n - m_1$ and so forth. The survival function based on Bootstrap data is

$$\hat{S}_{m}(t) = \prod_{j=1}^{n} \left(1 - \frac{m_{j}}{M_{j}}\right)^{\delta}$$



ABSTRACT

^{© 2011} Elixir All rights reserved

Aim

"Bootstrap is a way to pull oneself up (from an unfavourable situation) by one's Bootstrap, to provide trustworthy answers despite of unfavourable Circumstances" (Efron, 1979).

The aim of the study is to judge the performance (in terms of bias) of the quote by comparing the curves as well as the survival probabilities of two nonparametric survival functions for four data sets. For small sample Kaplan-Meier gives downward bias results (Whittemore & Keller, 1986). The performance of the two methods will be evaluated through the probabilities and also through graphical presentation on four different data sets containing different percentage of censoring.

For the analysis R-package (2004) is used. Brief Introduction of Data Sets

Leukaemia data set:

The famous leukaemia data set conducted by Freireich (Freireich et. al. 1963), which was reviewed by Gehan (1965) and also used by Borkowf (2005). The data consisted of 21 patients, including 9 events and 12(57%) censored observations. The data set consists of weeks in maintenance of remission. We ignore the placebo controls (containing no censored observation) here.

The weeks in remission are: 6, 6, 6, 6+, 7, 9+, 10, 10+, 11+, 13, 16, 17+, 19+, 20+, 22, 23, 25+, 32+, 32+, 34+ and 35+. Where + denotes a censored observation.

Stanford Heart Transplant Data Analysis

The second data set is the Stanford heart transplant data (Kalbflesch and Printice, 1980) .The data set contains 103 patients 75of were events and 28(27.2%) censored. The data set is given below

1,	2,	2,	2,	3,	3,	3,	5,	5,	6	, 6,
8,	9,	11+,	12,	16,	16,		16,	17,	18,	21,
21,	28,	30,	31	+,	32,	35,	30	5,	37,	39,
39+,	40,	40,		43,	45,	50,	5	1,	53,	58,
61,	66,	68,	68,	69,	72,		72,	77,	78,	80,
81,	85,	90,	96,	100,	102,	,	109+,	11	10,	131+,
149,	153,	165	,	180+,	186,		188,	20	07,	219,
263,	265+	-, 28	5,	285,	308,		334,	34	0,	340+,
342,	370	+, 3	897+,	427	7+,	445	+,	482+	,	515+,
545+,	583	, 596	+,	630+,	670+	-,	675,	73	3,	841+,
852,	915+	, 941	+,	979,	995,	10	32,	1141-	⊦,	1321+,
1386,	1400-	+, 140)7+,	1571+	, 158	86+,	179	99+		
Where	e + den	otes a c	ensor	ed obs	ervatio	n.				

An Application to a Lung Cancer Data

A data set from the Veterans Administration lung cancer trial in which chemotherapy was given to males with advanced inoperable lung cancer (presented by Prentice (1973)). This data set was also used by Gupta (1999). The data set consists of 97 patients out of which 91 were events and 6(6.2%) were censored. The survival times are given in days:

				0	•	
72	228	10	110	314	100*	42
	144	30	384	4	13	
123*	97*	59	117	151	22	18
	139	20	31	52	18	51
	122	27	54	7	63	392
92	35	117	132	162	3	95
	162	216	553	278	260	156
	182*	143	105	103		
112	87*	242	111	587	389	33
25	357	467	1	30		

283	25	21	13	87	7	24
99	8	99	61	25	95	80
29	24	83*	31	51	52	73
8	36	48	7	140	186	19
	45	80	52	53	15	133
	111	378	49			

Where + denotes a censored observation.

An Application to a Breast Cancer Data

A large data set of 1207 breast cancer patient is obtained from SPSS, version 11.5 (SPSS, 2004) containing 1135 censored (94.3 %) and 72 events only.

Method of Bootstrap Sample

Software for the program is prepared in R-package. A Bootstrap of sample "10* size of data set" is drawn with replacement. Sample is selected with the point to give each time enough chances to be included in the sample.

Results

The results of leukaemia data set are shown in table 1 and in Figure 1. For heart transplant data, we prepared table 2 and Figure 2. Table 3 and Figure 3 represent the facts of lung cancer data set. For a very large breast cancer data set, we show the results through Figure 4 and not prepared table of probabilities, due to the same conclusion obtained from the previous data sets.

Figure 1: Survival curves of Kaplan-Meier and Modified Kaplan-Meier for leukaemia data set



Figure 2: Survival curves of Kaplan-Meier and Modified Kaplan-Meier for heart transplant data

Set.



Leukaemia data set containing 21 patients having 57% censored data. Stanford heart transplant data containing 27.2% censored data. Lung cancer data consists of 6.2% censored cases, while the breast cancer data 1207 patients having very high censoring (94.3%). The range of data sets is 21 to 1207 patients and the range of percentage censoring is 6.2 to 94.3. On the basis of these facts and also on the basis of 4 Figures and 3 Tables, we reach to the followings.

· Bootstrap function is a form of Kaplan-Meier survival function.

• Kaplan-Meier gives the underestimate (bias) results for small sample. By comparing the survival probabilities of given data sets and curves, we can say that Bootstrap function gives more bias results than Kaplan-Meier survival function.

• For large data set, Kaplan-Meier gives unbiased results (Maller and Zhou, 1996). If we compare the curves of breast cancer data, the Bootstrap still gives some bias.

• If the last observation is censored, one can not obtain the mean of survival function and the same result we obtained from the Bootstrap survival function.

• Both the survival functions have the same range i.e. 0 to 1.

Figure 3: Survival curves of Kaplan-Meier and Modified Kaplan-Meier for lung cancer data set.



Figure 4: Survival curves of Kaplan-Meier and Modified Kaplan-Meier for breast cancer data set (SPSS data directory).



Discussion and Conclusion

Parametric approach produces better results, if we are able to find a specific parametric survival distribution (Exponential, Weibull etc.) that fits the data. As most of the survival data are skewed and some times, it is very difficult to find the appropriate distribution. The easiest way to avoid the possible error (which may occur due to the application of in appropriate distribution) is the non-parametric approach. The two nonparametric approaches for estimating the survival function are the Kaplan-Meier and Bootstrap approach.

In this article we compared the two approaches by applying them on different data sets and on the basis of analysis we reached to the conclusion that, Kaplan-Meier approach as compared to the Bootstrap approach is easy to understand and to apply. The Kaplan-Meier approach can be easily apply to different data sets by the help of packages e.g. SPSS, R, S, SPLUS, while for applying the Bootstrap technique, there is no option available in commonly used packages. In some situations Kaplan-Meier gives bias results, which are smaller than the results obtained by applying the Bootstrap survival function.

Bootstrap survival function in all the four data sets gives the smaller probabilities of survival (for every censoring percentage).

If comparatively large data set, having moderate censoring is available, then for estimating the probabilities and for drawing the survival curve, Kaplan-Meier is considered as the conventional, easy and less time consuming method. **References**

1. Borkowf, C. B. (2005). A simple hybrid variance estimator fort he Kaplan-Meier survival function. *Statistics in Medicine* 24: 827-851.

2. Efron, B. (1979)."Bootstrap Methods: Another Look at the jacknife," *Annals of Statistics* 7: 1-26.

3. Fleming, T.R., & Harrington, D. P. (1984). Nonparametric estimation of the survival distribution in censored data. *Communication in Statistics – Simulation and Computation* 13: 1-26.

4. Freireich EJ, Gehan E, Frei III E, Schroeder LR, Wolman IJ, Anbari R, Burgert EO, Mills Sd, Pinkel D, Selwry OS, Moon JH, Gendel BR, Spurr CL, Storrs R, Haurani F, Hoogstraten B, Lee S.(1963). The effect of 6-mercaptopurine on the duration of steroid-induced remissions in acute leukaemia: a model for evaluation of other potentially useful therapy. *Blood* 21(6): 699-716.

5. Gehan, E. A. (1965). A generalized Wilcoxon test for comparing arbitrarily singly-censored samples. *Biometrika* 52(1, 2): 203-223.

6. Gupta, R.C. (1999). A Study of Log-Logistic Model in Survival Analysis. *Biometrical Journal* 41: 431-443.

7. Kalbfleisch J.D. and Prentice R.L. (1980). *The Statistical Analysis of Failure Time Data*. Wiley, New York.

8. Kaplan, E.L. & Meier, P. (1958). Non-parametric estimation from incomplete observations. *Journal of the American Statistical Association* 53: 457-481.

9. Maller, R. & Zhou, X. (1996). Survival Analysis with Longterm survivors. John Wiley & Sons Chichester.

10. Miller, Jr.& Rupert, G. (1983). What Price Kaplan-Meier. *Biometrics* 39: 1077-1081.

11. Prentice, R. L. (1973). Exponential survivals with censoring and explanatory variables. *Biometrika* 60: 279-288.

12. R Development Core Team. *R: a language and environment for statistical computing*. Vienna, Austria, R Foundation for Statistical Computing, 2004.

13. SPSS Corporation. SPSS 11.5 for windows. Chicago, 2004 14. Whittemore, A. S. & Keller, J. B. (1986). Survival Estimation Using Splines. *Biometrics* 42: 495-506.

	1101 41 1101	Implan - Micici	Dootstrup
		Survival Function	Survival Function
3	21	0.8571	0.8382
1	17	0.8067	0.7578
0	16	0.8067	0.7578
1	15	0.7529	0.7179
0	13	0.7529	0.7179
1	12	0.6902	0.6257
1	11	0.6275	0.5772
0	10	0.6275	0.5772
0	9	0.6275	0.5772
0	8	0.6275	0.5772
1	7	0.5378	0.4861
1	6	0.4482	0.3797
0	5	0.4482	0.3797
0	4	0.4482	0.3797
0	2	0.4482	0.3797
0	1	0.4482	0.3797
	$ \begin{array}{c} 3\\1\\0\\1\\0\\1\\1\\0\\0\\0\\1\\1\\0\\0\\0\\0\\0\\0\\0\\0\\$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{tabular}{ c c c c c } \hline Survival Function \\ \hline Survival Function \\ \hline \\ \hline \\ 3 & 21 & 0.8571 \\ 1 & 17 & 0.8067 \\ 0 & 16 & 0.8067 \\ 1 & 15 & 0.7529 \\ 0 & 13 & 0.7529 \\ 1 & 15 & 0.7529 \\ 1 & 12 & 0.6902 \\ 1 & 11 & 0.6275 \\ 0 & 10 & 0.6275 \\ 0 & 9 & 0.6275 \\ 0 & 9 & 0.6275 \\ 0 & 9 & 0.6275 \\ 1 & 7 & 0.5378 \\ 1 & 6 & 0.4482 \\ 0 & 5 & 0.4482 \\ 0 & 4 & 0.4482 \\ 0 & 2 & 0.4482 \\ 0 & 1 & 0.4482 \\ \hline \end{tabular}$

Table 1. Estimated survival functions of Kaplan-Meier and Bootstrap for leukaemia data setTimeEventNo. at riskKaplan-MeierBootstrap

rvival	function	ns of Kapla	an-Meier and I	Bootstrap for heart
Time	Event	No. at risk	Kaplan-Meier	Bootstrap
1	1	102	Survival Function	Survival Function
2	3	103	0.9612	0.9500
3	3	99	0.9320	0.9247
5	2	96	0.9126	0.9062
8	1	94 92	0.8835	0.8697
9	1	91	0.8738	0.8616
11	0	90	0.8738	0.8616
12	1	89	0.8640	0.8515
17	1	85	0.8247	0.8066
18	1	84	0.8149	0.7968
21	2	83 81	0.7952	0.7738 0.7633
30	1	80	0.7756	0.7510
31	0	79	0.7756	0.7510
32 35	1	78 77	0.7657	0.7448
36	1	76	0.7458	0.7155
37	1	75	0.7358	0.7084
39 40	1	74	0.7259	0.6933
43	1	70	0.6956	0.6616
45	1	69	0.6856	0.6519
50	1	68 67	0.6755	0.6449
53	1	66	0.6553	0.6256
58	1	65	0.6452	0.6159
61	1	64	0.6352	0.6045
68	2	63 62	0.6251	0.5940
69	1	60	0.5948	0.5662
72	2	59	0.5747	0.5474
78	1	57	0.5646	0.5423
80	1	55	0.5444	0.5219
81	1	54	0.5343	0.5159
85	1	53 52	0.5243	0.5040 0.4955
96	1	51	0.5041	0.4819
100	1	50	0.4940	0.4700
102	1	49	0.4839	0.4615
110	1	47	0.4736	0.4513
131	0	46	0.4736	0.4513
149	1	45 44	0.4631	0.4417
165	1	43	0.4426	0.4181
180	0	42	0.4421	0.4181
186	1	41 40	0.4313	0.4073
207	1	39	0.4097	0.3866
219	1	38	0.3989	0.3748
263 265	1	37	0.3882	0.3577
285	2	35	0.3660	0.3303
308	1	33	0.3549	0.3187
334 340	1	32 31	0.3438 0.3327	0.3088 0.3008
342	1	29	0.3212	0.2909
370	0	28	0.3212	0.2909
397 427	0	27	0.3212 0.3212	0.2909
445	0	25	0.3212	0.2909
482	0	24	0.3212	0.2909
515 545	0	23	0.3212	0.2909
583	1	21	0.3059	0.2780
596	0	20	0.3059	0.2780
630 670	0	19	0.3059	0.2780
675	1	17	0.2879	0.2652
733	1	16	0.2699	0.2410
841 852	0	15	0.2699	0.2410
915	0	13	0.2507	0.2247
941	0	12	0.2507	0.2247
9/9 995	1	11	0.2279	0.2010 0.1757
1032	1	9	0.1823	0.1622
1141	0	8	0.1823	0.1622
1321 1386	0	6	0.1823	0.1622 0.1403
1400	0	5	0.1519	0.1403
1407	0	4	0.1519	0.1403
1571	0	3 2	0.1519	0.1403
1799	ő	1	0.1519	0.1403

Table 2. Estimated sur transplant data set

d	surviv	al_func	tions of Ka	plan-Meier and	Bootstrap for lung	g
	Time	Event	No. at risk	Kaplan-Meier	Bootstrap Survival Function	
	1	1	97	0.9897	0.9897	
	3	1	96	0.9794	0.9784	
	4	1	95	0.9691	0.9650	
	7	2	94	0.9485	0.9343	
	8 10	2	91 89	0.9276	0.9161	
	13	2	88	0.8963	0.8771	
	15	1	86	0.8859	0.8652	
	18	0	85	0.8859	0.8652	
	19	1	83	0.8753	0.8512	
	20	1	81 81	0.8539	0.8402	
	22	1	80	0.8432	0.8162	
	24	2	79	0.8219	0.7923	
	25	3	77	0.7899	0.7602	
	27	1	74	0.7792	0.7515	
	30	2	73	0.7472	0.7124	
	31	2	70	0.7258	0.6863	
	33	1	68	0.7151	0.6734	
	35	1	67	0.7045	0.6595	
	36	1	66 65	0.6938	0.6484	
	42 45	1	64	0.6724	0.6280	
	48	1	63	0.6618	0.6169	
	49	1	62	0.6511	0.5984	
	51	2	61	0.6298	0.5710	
	52 53	3 1	59 56	0.5977	0.5391	
	53 54	1	55	0.5764	0.5095	
	59	1	54	0.5657	0.5008	
	61	0	53	0.5657	0.5008	
	63	1	52	0.5548	0.4946	
	72	1	51	0.5440	0.4847	
	80	1	49	0.5222	0.4613	
	83	1	47	0.5111	0.4506	
	87	2	46	0.4889	0.4363	
	92	1	44	0.4778	0.4240	
	95 07	2	43	0.4555	0.4048	
	99	2	40	0.4222	0.3681	
	100	1	38	0.4111	0.3623	
	103	1	37	0.4000	0.3448	
	105	1	36	0.3889	0.3331	
	111	2	34	0.3555	0.3232	
	112	1	32	0.3444	0.2890	
	117	2	31	0.3222	0.2640	
	122	1	29	0.3111	0.2540	
	123	1	28	0.3000	0.2487	
	132	1	26	0.2778	0.2265	
	139	1	25	0.2667	0.2181	
	140	1	24	0.2555	0.2120	
	143	1	23	0.2444	0.2012	
	144	0	22	0.2333	0.1921	
	156	1	20	0.2217	0.1807	
	162	2	19	0.1983	0.1601	
	182	1	17	0.1867	0.1525	
	186	1	16 15	0.1750	0.1394	
	210	1	13	0.1517	0.1323	
	242	1	13	0.1400	0.1088	
	260	1	12	0.1283	0.0965	
	278	1	11	0.1167	0.0896	
	283 314	1	10	0.1050	0.0820	
	314	1	9 8	0.0955	0.0598	
	378	1	7	0.0700	0.0513	
	384	1	6	0.0583	0.0398	
	389	1	5	0.0467	0.0307	
	392 167	1	4	0.0350	0.0238	
	407 553	1 1	5 2	0.0233	0.0158	
	587	1	1	0.0000	0.0000	

Table 2. Estimated g cancer data set