

International Journal of Statistics and Applied Mathematics

ISSN: 2456-1452
Maths 2022; 7(6): 122-128
© 2022 Stats & Maths
www.mathsjournal.com
Received: 16-09-2022
Accepted: 19-10-2022

Pavithra V
Department of Statistics,
Annamalai University, Tamil
Nadu, India

Kannan R
Department of Statistics,
Annamalai University, Tamil
Nadu, India

Corresponding Author:
Pavithra V
Department of Statistics,
Annamalai University, Tamil
Nadu, India

Impact of missing complete at random data in survival analysis

Pavithra V and Kannan R

Abstract

The problem of missing data deserves special attention because it refers to the case where not all data were obtained as intended in the study design. In survival analysis researchers are faced with the problem of identifying subjects for each follow-up visit and in some situations; may not obtain observations about the subject of the study. As a result, there will be a lack of data in some studies and this poses a major challenge for the analysis. In general, missing data analysis deals with replacement and missing data with deletion. In this paper we may try missing complete at random with deletion and Imputation techniques and focuses on the study is to find out the survival probability rate and final outcome of the Breast Cancer patients by the method of Kaplan-Meier, Cox Proportional Hazard Model, Minimum Survival probability, Maximum Survival Probability.

Keywords: Breast cancer, missing data, kaplan-meier test, log rank test, cox proportional hazard model, minimum survival probability, maximum survival probability

Introduction

Breast cancer (BC) is a non-communicable disease that begins in the cells of the breast. BC is one of the leading cancers in Indian women, with over 1.5 lakh new BC patients registered in India in 2018. It accounts for 14% of all cancers in women. The BC is uncommon in men, 1 in 400 men has BC. This is the most common disease among Indian women, 1 in 28 people may develop BC at some point in their lives. Unfortunately, the number of BC cases reported each year is increasing faster than ever. The BC accounts for more than 27% of all new cancer cases. There is an increasing trend in the number of new cancer patients and comparably the risk is higher in urban areas as 1 in 22 women and lower in rural areas as 1 in 60 women. In India, the average age of the high-risk group is between 40-55 years are more prone to BC (American Cancer Society (ACS) 2019-2020). The overall numbers in India are better compared to the number for developed countries like US/UK is less where in 1 in 8 women are diagnosed annually. However, as the awareness level about the disease in developed countries is quite high and there is a lot of government funding which promotes timely detection, most cases are detected and treated at early stages leading to better survival rates.

In India on the other hand, has a very low survival rate due to its large population and low awareness. 1 in 2 women diagnosed with BC will die within the next five year. One of the main reasons for high mortality rates is lack of awareness, late diagnosis and absence of proper BC screening programme. Most of the BCs are diagnosed at advanced stage. Many patients in the urban area are diagnosed at stage-2 and most of the cases from rural areas, these lesions are diagnosed only after they transform to metastatic tumors. The exact cause of BC is still unknown, but years of medical research have identified several risk factors. It is still unclear why some women at very high risk do not develop BC, while some women without risk factors may develop BC. The risk factors for BC include genetics and heritage, late pregnancy, use of oral contraceptives, early onset of menstruation, late menopause, excessive alcohol intake, smoking, adolescent obesity, increases stress and poor eating habits-these factors are due to the increased incidence of BC.

Through cancer, especially BC is a very dangerous disease that is prevalent worldwide (Torre *et al.*, 2015) [29]. Cancer is a group of diseases caused by the uncontrolled growth and spread of abnormal cells throughout the body (Diabate *et al.*, 2018) [6].

BC is expensive that it has received a lot of attention from doctors and statisticians. Mortality with unstable mortality with many different prognostic (Pg) factors (Parkin *et al.*, 2014) [22]. American Joint Committee on Cancer (AJCC) BC staging is associated with survival prognosis (ACS, 2017). This situation is indicated by reduced survival from stage-1 90%, stage-2 65%, stage-3 20% and stage-4 5% (Sinaga *et al.*, 2017) [28].

The majority of BC cases are classified as invasive and non-invasive. Invasive BC has spread throughout the body, but non-invasive did not spread throughout the body (Abay *et al.*, 2018) [1]. Age has a significant effect on whether women get BC. The mortality rate of BC increases with age (Rezaianzadeh *et al.*, 2009) [24]. Drinking alcohol increases the risk of dying from BC in women is about 7% to 12% for every 10g of alcohol consumed per day (Desantis *et al.*, 2013) [5]. A study conducted by Addis Ababa University on the impact of several risk factors on BC and survival showed that stage and type of disease have a significant effect on BC survival (Kantelhardt *et al.*, 2014) [17].

The problem of missing data was dealt with mainly by editing until 1970. Contribution to inference problems in missing data studies was developed by Rubin (1976) [25]. During the 1980's, several methods such as Listwise deletion, pairwise deletion, imputation methods and various models have been widely applied. The missing data problem is very difficult to apply for the common statistical methods (Rubin, 1987) [26]. The studies of Little and Rubin (1987) [26] marked the beginning of the second phase and can be considered a breakthrough in the development of missing data methods. This method is said to be better than the simple imputation method that provides efficient parameter estimates. Pavithra and Kannan (2022) [23] have clearly stated how to handle missing data when using them in survival analysis.

In some cases, missing data can lead to bias and lead to erroneous conclusions about changes in mean responses. Missing data will reduce the efficiency or accuracy of estimates of changes in mean. Unfortunately, the larger the amount of missing data, the greater the loss of model accuracy (Fitzmaurice *et al.*, (2004) [13]; Mohanraj and Srinivasan (2018)) [21]. Several reasons have been attributed to missing data, including equipment failure, audience attrition from treatment design or intrusive questions about a survey or unclear instructions. Fayers and Machin (2007) [10] described a special case of single imputation called hierarchical scale imputation. Hedeker and Gibbons (2006) [14], Fitzmaurice (2003) [12] and Diggle *et al.*, (2002) [9] have made a valuable contribution to the development of the missing data problem. Hesketh and Skrondal (2008) [15] provided a more applicable framework for STATA users to analyze missing data. The analysis of missing data poses a major problem because the estimates of the parameters are mostly biased (Becker and walstad (1990) [8]; Becker and Powers (2001) [7]; Holt (1997) [16]; Rubin (1976) [25]. Studies by several researcher (Anderson *et al.*, (1983); Kim and Curry (1977)) [18] indicate a loss of information and statistical power when conducting analysis of missing data. Scheffer (2002) [27] discussed how the mean and standard deviation are affected by different methods of imputation in dealing with different missingness mechanism.

Types Of Missing Data:

1) **Missing At Random (MAR):** The fact that the missing data is systematically linked to the observed data but not to the unobserved data.

- 2) **Missing Complete At Random (MCAR):** The missing data is independent of observed and unobserved data. In other words, there was no systematic difference between participants with missing data and those with complete.
- 3) **Missing Not At Random (MNAR):** The fact that missing data is systematically linked to unobserved data, i.e., the lack is related to events or factors not measured by the researcher.

In this article, we will calculate the survival analysis and probability by applying the MCAR approach and using the three methods given below.

1.2. Methods Of Handling in Missing Data

In general, the missing data analyze under three different techniques namely,

- a) **Pairwise Deletion:** when the statistical procedure uses instances that contain missing data. A procedure cannot include a particular variable when it has a missing value, but it can still use the case when parsing other variables with a different value.
- b) **Listwise Deletion:** in this method, an entire record is excluded from analysis if any single value is missing and therefore we have the same N (number of records) for all the analysis.
- c) **Single (or) Multiple Imputation:** in this manner, replace the missing value with single or multiple value using a strategy such as: Mean, Median, Most frequent,...,etc.

The goal of this study is to look into the survival and risk of death from the Adyar Cancer Institute in 2013. We study the missing data in the received data from our point of view. Generally, the missing data divided into three types, namely, MACR, MAR and MNAR. Due to the presence of such missing data, the accuracy of the model is greatly reduced during the analysis and the researchers are confused when describing the model. Therefore, researchers have already used three methods namely, pairwise deletion, listwise deletion and imputation method to increase model accuracy when using such missing data. We are going to continue our survival analysis of BC research by using these three methods to overcome the deficiencies caused by missing data. This BC survival analysis implemented using, which includes various models, was employed. The Kaplan-Meier (K-M) with log rank test and Cox Proportion Hazard (PH) models are most commonly utilized models (Lee and Wang (2003)). In addition, we looked at the Minimum survival probability and Maximum survival probability (Felix and Kannan (2007)) [11].

2. Statistical Methods

2.1. Hazard Functions

The hazard function of the hold time X is denoted by $h(x)$ and defined as individual probability fails in the time interval $(x, x + \Delta x)$ that the individual has lived for time x, the hazard function is expressed as:

$$h(x) = \lim_{\Delta x \rightarrow 0} \left[\frac{P(x < X < x + \Delta x | X > x)}{\Delta x} \right] \rightarrow (1)$$

2.2. Cox Proportional Hazard Model

The relationship between the hazard rate and the covariate set can be expressed using the model:

$$\ln[h(t)] = \ln[h_0(t)] + \sum_{i=1}^n x_i \beta_i \rightarrow (2)$$

Where $x_1, x_2, x_3, \dots, x_n$ are covariates. $\beta_1, \beta_2, \beta_3, \dots, \beta_n$ are the regression coefficients to be estimated. t is time and $h_0(t)$ is the baseline hazard rate when all covariates are zero.

2.3. The Survival Function

Individual opportunities to survive for time x are expressed by $S(x) = P(X > x)$. Let X be the continuous random variables, then the survival function is the complement of the Cumulative Distribution function $S(x) = 1 - F(X)$ where $F(X) = P(X \leq x)$. The survival function is the integral of the probability density function $f(x)$:

$$\hat{S}(x) = P(X > x) = \int_x^\infty f(t)dt \rightarrow (3)$$

$$f(x) = -\frac{dS(x)}{dx} \rightarrow (4)$$

Then if X is the discrete random variables, and can be obtained x_j with the probability mass function (p.m.f) $p(x_j) = P(X = x_j)$, $j=1,2,3,\dots$ where x_1, x_2, x_3, \dots then the survival function for the discrete variables X is given by:

$$\hat{S}(x) = P(X > x) = \sum_{x_j > x} p(x_j) \rightarrow (5)$$

2.4. Kaplan-Meier with Log Rank Test

Estimated survival function for K-M Expressed as:

$$\hat{S}(x_{(j)}) = \hat{S}(x_{(j-i)})\hat{P}(X > x_{(j)} | X \geq x_j) \rightarrow (6)$$

In general, log rank is used to compare k-M survival curves formed by the following hypothesis:

H_0 : There is no difference between the survival curves:

H_1 : At least one difference between the survival curves:

$$\text{Log Rank Test} = \frac{(O_i - E_i)^2}{\text{Var}(O_i - E_i)} \rightarrow (7)$$

$$O_i - E_i = \sum_{j=1}^n m_{ij} - e_{ij} \rightarrow (8)$$

m_{ij} denotes the number of individuals who experience the event at time x_j , and e_{ij} is the value of hope. The null hypothesis will be rejected if log rank statistics $\geq \chi^2_\alpha$ with $n-1$ degrees of freedom (df) or p-value $< \alpha$.

2.5. Minimum Survival Probability (MISP)

Survival probabilities are calculated on the assumption that all those that are censored, the result of interest occurred. Then, for any interval i , D_i denotes the number of deaths during i , W_i denotes the number of censored observation during i and N_i denotes the number of subjects at the beginning of i . Then MISP for time interval i is expressed by

$$\text{MISP} = 1 - \frac{(D_i - W_i)}{N_i} \rightarrow (9)$$

2.6. Maximum Survival Probability (MASP):

The survival probabilities are calculated by assuming that all those who are censored at time i are alive till the end of time interval i . Hence the notations of MASP is,

$$\text{MASP} = 1 - \left(\frac{D_i}{N_i}\right) \rightarrow (10)$$

3. Source of Breast Cancer Data:

The data we used in our research were obtained from the Adyar Cancer Institute in Chennai. These data are the newly diagnosed breast cancer for 2013 and where we used the number 257 patients for our study. The data provided by the cancer center for this research: Gender, Age, Medical History, Date of Diagnosis, laterality of the BC, Grade, Stages, Treatments (Surgery, Chemo Therapy, Radiation Therapy, Hormonal Therapy) with dates, follow-up details with dates and Alive Status.

4. Result and Discussion

4.1. Cox Proportional Hazard Model

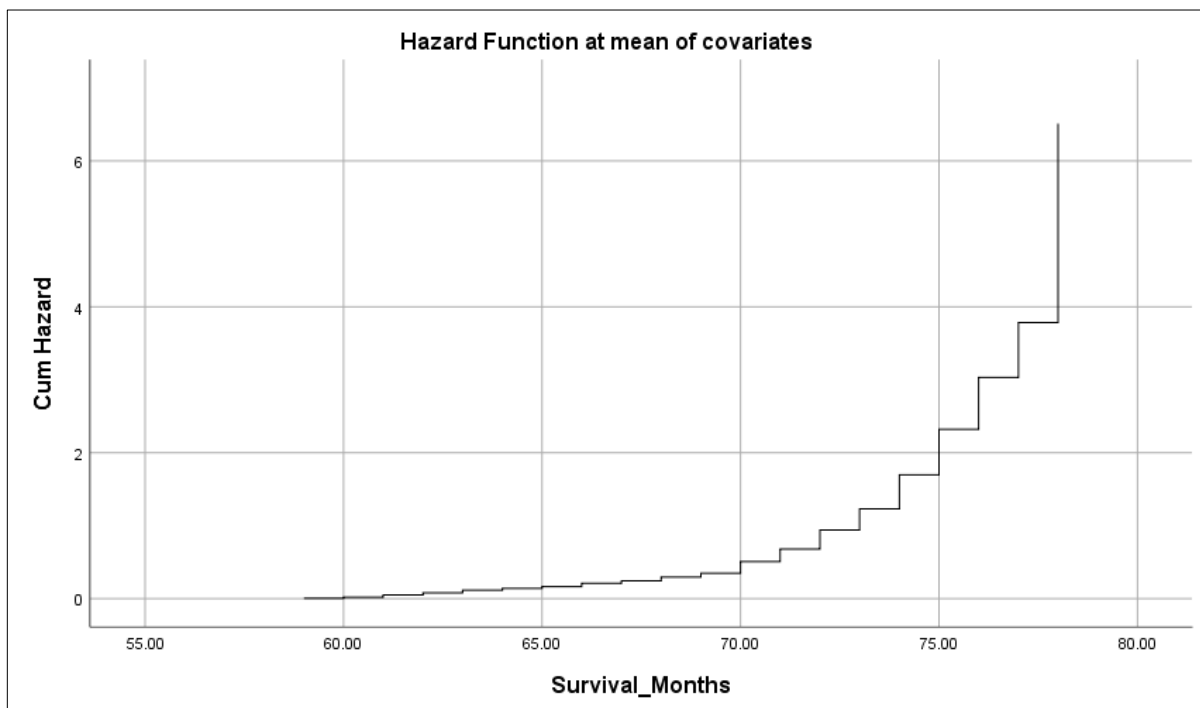


Fig 1: Hazard Function at mean of Covariate

The estimated variables by Cox regression are given in following: Age Group (X_1), Rural Urban (X_2), Medical History (X_3), Laterality (X_4), Stage (X_5), Recurrence and Metastatic (X_6), Surgery (X_7), Chemo Therapy (X_8), Radio Therapy (X_9), Hormonal Therapy (X_{10}). In step-1, the partial test shows that only Pg variables are statistically significant (P-value <5%). The backward stepwise method is used to extract the least influencing factors so that the final model is obtained in step-4 and the same method of Cox PH analysis applied all the three data set.

The β regression coefficient of the obtained models are all positive ($\beta > 0$) with the value of $\exp(\beta) > 0$, meaning that all factors included in the model influence the event speed (death). That is, the risk of failure of depending on advanced stage of BC is 1.617 times greater than those lower stages. The risk of death of BC patients with recurrent and metastatic is 0.623 times greater than those that do not have recurrent and metastatic.

Table 1: Partial Test with Backward Stepwise Method for deletion and imputation data set

	Indepented Variables	Pairwise Deletion				Listwise Deletion				Imputation Method			
		β	Wald	Sig.	$\exp(\beta)$	β	Wald	Sig.	EXP(β)	β	Wald	Sig.	EXP(β)
Step-1	X_1	-.14	2.61	.106	.869	-.11	1.292	.256	.897	-.15	2.850	.091	.864
	X_2	-.12	.506	.477	.886	-.13	.468	.494	.882	-.12	.481	.488	.889
	X_3	-.05	.097	.756	.946	.04	.048	.827	1.043	-.06	.096	.757	.947
	X_4	-.04	.067	.796	.960	.02	.008	.930	1.015	-.07	.203	.652	.933
	X_5	.52	10.11	.001	1.68	.53	8.637	.003	1.693	.55	12.110	.001	1.734
	X_6	-.45	3.23	.073	.633	-.33	1.361	.243	.720	-.48	3.583	.058	.619
	X_7	.12	.045	.833	1.12	.02	.001	.973	1.022	.31	.529	.467	1.358
	X_8	-.16	.285	.593	.848	-.36	.998	.318	.694	-.13	.205	.651	.878
	X_9	.21	1.02	.314	1.23	.23	1.040	.308	1.259	.25	1.457	.227	1.277
	X_{10}	-.07	.152	.697	.929	-.01	.002	.963	.990	-.05	.065	.798	.953
Step-2	X_1	-.11	1.948	.163	.896	-.08	.981	.322	.918	-.11	1.864	.172	.899
	X_3	-.04	.062	.804	.958	.08	.201	.654	1.087	-.06	.146	.703	.938
	X_5	.46	10.992	.001	1.580	.47	9.764	.002	1.599	.48	12.411	.000	1.616
	X_6	-.44	3.078	.079	.644	-.30	1.187	.276	.739	-.45	3.227	.072	.639
Step-3	X_1	-.10	2.060	.151	.904	-.10	1.813	.178	.902	-.09	1.783	.182	.911
	X_5	.46	11.162	.001	1.584	.46	9.706	.002	1.596	.48	12.625	.000	1.622
	X_6	-.44	3.07	.080	.644	-.30	1.204	.273	.738	-.45	3.226	.072	.639
Step-4	X_5	.45	10.950	.001	1.576	.46	9.560	.002	1.588	.48	12.497	.000	1.617
	X_6	-.46	3.484	.062	.627	-.33	1.376	.241	.723	-.47	3.629	.057	.623

Table 2: Overall Score in Backward Stepwise Method for deletion and imputation data set

Handling of Missing Data	Stepwise Method	-2 Log Likelihood	Overall (score)		
			Chi-square	DF	Sig.
Pairwise Deletion	Step-1	1461.524	22.994	10	.014
	Step-2	1465.476	18.326	4	.001
	Step-3	1465.569	18.151	3	.000
	Step-4	1466.733	17.023	2	.000
Listwise Deletion	Step-1	1268.868	17.071	10	.073
	Step-2	1271.815	14.268	4	.006
	Step-3	1272.017	14.119	3	.003
	Step-4	1273.827	12.369	2	.002
Imputation Method	Step-1	1566.067	25.330	10	.005
	Step-2	1569.285	21.383	4	.000
	Step-3	1569.431	21.134	3	.000
	Step-4	1571.208	19.399	2	.000

In Table-1: The result of cox proportional hazard analysis showed that the most significant pg variable to the probability of death was the presence of advanced stage and tumor recurrence with metastatic. Table-2 indicated the overall score for the data set of pairwise deletion and imputation methods are more appropriate in missing data analysis. Meanwhile, the Listwise deletion leads reduced the efficiency and lower the precision of the model estimates.

4.2. Kaplan-Meier Analysis

The K-M estimated the probability of survival curve for missing data technique of pairwise deletion, listwise deletion and imputation method. The following figures are according to the two impact variables of stages and recurrence with metastasis.

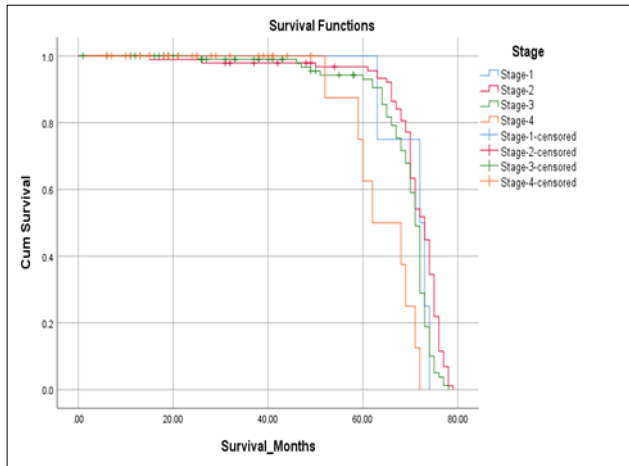


Fig 2: Tumor stage for Pairwise deletion

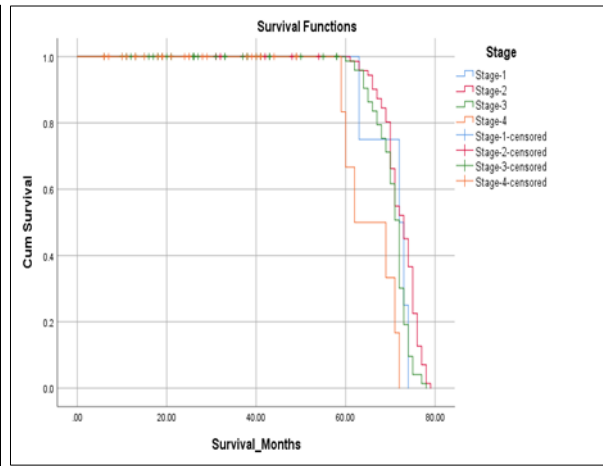


Fig 3: Tumor stage for Listwise Deletion

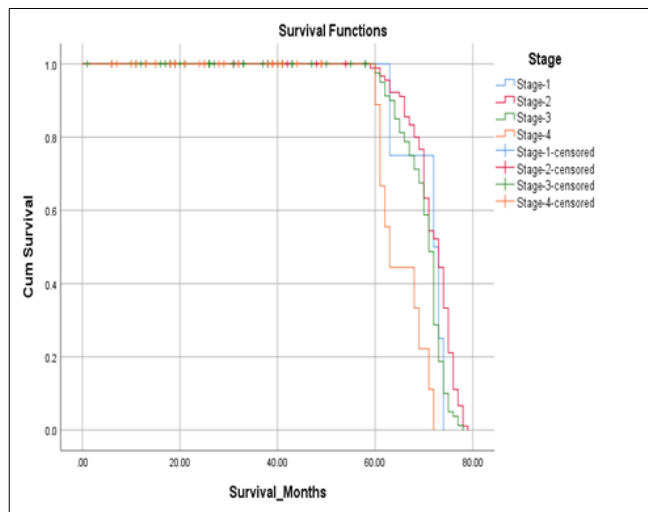


Fig 4: Tumor stage for Imputation Method

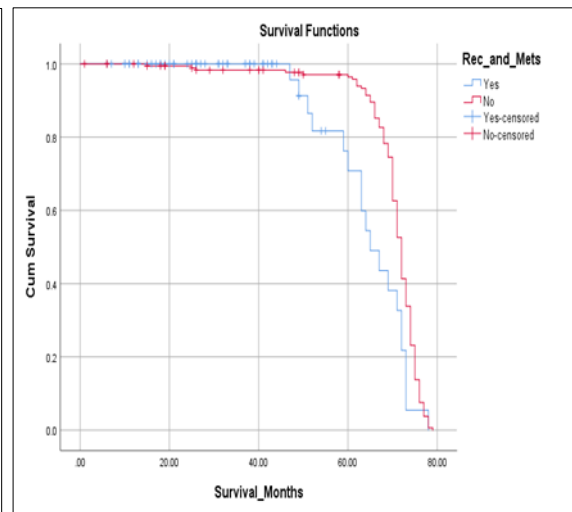


Fig 5: Recu and Mets for Pairwise Deletion

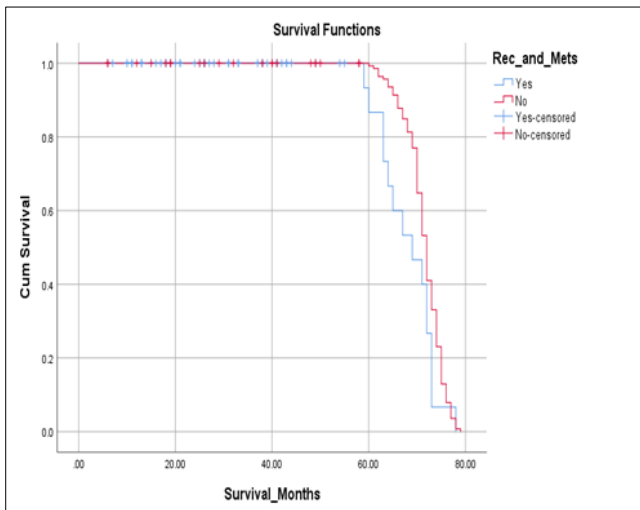


Fig 6: Recu and Mets for Listwise Deletion

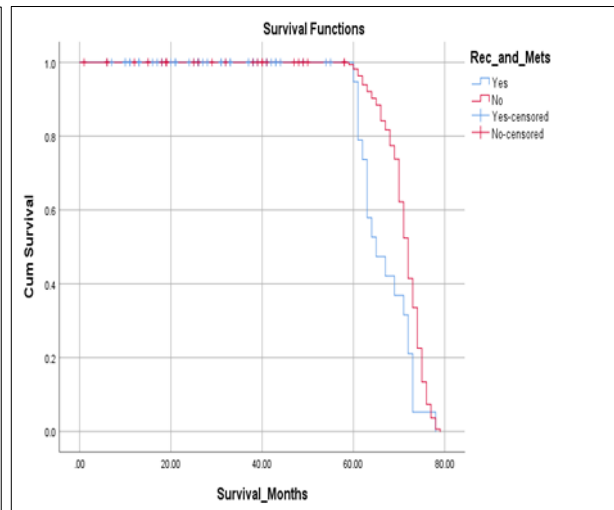


Fig 7: Recu and Mets for Imputation Method

The stages and recurrence with metastasis are the most important factor variable in the BC and these can determine the conditions of the cancer patients. In figure-2, 3 & 4 survival probability shown that the cancer Stages of patient with BC and clearly seen which stage is mostly survived or died in these BC. The survival rate of stage-1 and stage-2 patients was very high compared to stage-3 and stage-4 and risk rate of BC patients in stage-1 and stage-2 was very low when compared to other stages. In figure-5, 6 & 7 survival

probability shown that the recurrence with metastatic cancer patient and clearly seen who have not spread the cancer they are only mostly survived in these BC.

4.3. Log Rank Test

The log rank test to determine if there is a difference between the survivals curves. The log rank test of significant or not significant in Pg variables are given in following table.

Table 3: Log Rank test used for Pg variable affecting Survival of BC

Log Rank Test / Pg Variable	df	Pairwise Deletion		Listwise Deletion		Imputation	
		χ^2	Sig.	χ^2	Sig.	χ^2	Sig.
Age Group	4	2.927	.570	.993	.911	4.851	.303
Area	1	.227	.634	.064	.800	.172	.678
Medical History	1	.070	.791	.147	.702	.026	.872
Laterality	1	.016	.898	.436	.509	.577	.447
Stage	3	26.162	.000	20.801	.000	31.293	.000
Recurrence & Metastatic	1	7.841	.005	15.173	.000	9.478	.004
Surgery	1	1.025	.311	1.899	.168	1.032	.310
Chemo Therapy	1	2.133	.144	1.486	.223	2.246	.134
Radiation Therapy	1	.780	.377	.706	.401	3.120	.077
Hormonal Therapy	1	.004	.952	.000	.994	.140	.708

Based on the Log Rank Test in Table-3, the equality of survival distribution of the BC variables Cancer Stages and Recurrence with metastatic were statistically recorded a p-value (0.000 and 0.004) makes a significant difference and other variables have statistically no significant difference. Meanwhile, beauty of this study is in the handling the missing

data technique, the imputation and pairwise deletion methods had outperforming when compared to list wise deletion method.

4.4. Comparison of survival probabilities

Table 4: Survival Probability for Deletion and Imputation Methods

Handling Of Missing Data Technique	Method Of Probability	BC Five Year Survival Probability by Percentage				
		1 st Year	2 nd Year	3 rd Year	4 th Year	5 th Year
Pairwise Deletion	MISP	95%	84%	76%	68%	64%
	MASP	96%	87%	80%	71%	66%
	K-M	96%	87%	81%	70%	65%
List wise Deletion	MISP	95%	82%	74%	66%	62%
	MASP	95%	82%	74%	66%	62%
	K-M	94%	81%	73%	65%	61%
Imputation method	MISP	95%	85%	78%	70%	65%
	MASP	96%	87%	80%	72%	67%
	K-M	96%	86%	79%	71%	69%

Table-4: Shows the cumulative survival probabilities at the end of each year from the date of completion of treatment through different methods. These estimates are obtained by using MISP, MASP and K-M methods. In general, by all the methods estimates of the cumulative probabilities have been decreased as the survival period has increased. The higher probabilities have been estimated by MASP. i.e., the estimates of MISP and MASP provide the two extreme values of the survival band within which the true survival probability lies. The three estimates are similar but not identical. The overall five-year survival probability (%) for the BC patients has been found to be 67%, which is very much similar to other methods. However, this overall survival probability may not be an appropriate one, since the stage of the disease at diagnosis is one of the significant factors associated with the number of deaths occurred.

5. Conclusion

The K-M, Cox PH, MISP and MASP survival results of the study showed that age, medical history, resident, laterality of breast, stage, recurrence, metastasis, surgery, chemo therapy, radiation therapy and hormone therapy affected the time to death of BC patients 2013 at Adyar Cancer Hospital. The K-M estimated the survival month of the BC is 69 months. The analyses Cox PH found main factor behind the poor survival time is that the treated patients is already in the advanced stage and recurrent with metastatic. The comparison between the MISP, MASP and K-M analysis the MASP and K-M showed similar together and most useful to survival analysis. The beauty of this survival analysis with missing data studies, here we have clearly outlined how to handling missing data in survival analysis. We reports that among the three methods

used in this survival analysis, the pairwise deletion and imputation methods are more suitable for all type of analysis. The information loss is high and the model accuracy is very low for the listwise deletion method when compared to the other two methods. So it's best to avoid using the listwise deletion method when dealing with missing data.

6. Recommendation

Health professionals, governments and NGO should raise awareness of early cancer screening and should also encourage women to be diagnosed at an early stage to improve mortality risk, and cancer screening facilitation and scheduling should be planned and scheduled in rural areas of this region to elucidated mortality risk.

7. References

1. Abay M, Tuke G, Zewdie E, Abraha TH, Grum T, Brhane E. Breast self-examination practice and associated factors among women aged 20-70 years attending public health institutions of Adwa town, North Ethiopia. BMC Research Notes, 2018;11(1):622,1-7.
2. Amin MB, Edge SB, Greene FL, Stephen B, Compton CC, Gershenwald JE, *et al.* American Joint Committee On Cancer Staging manual. 8th Edition, New York: Springer, 2017.
3. American Cancer Society. Stages of Breast Cancer Internet, 2017-2018.
4. American Cancer Society. Breast Cancer Fact & Figures Internet, 2019-2020.
5. Desantis C, Siegel R, Jemal A. Breast cancer facts and figures, 2013-2014. American Cancer Society, 2013, pp 1-38.

6. Diabate M, Coquille L, Samson A. Parameter estimation and treatment optimization in a stochastic model for immunotherapy of cancer. arXiv Preprint ArXiv, 1806.01915., 2018.
7. Becker WE, Powers JR. Student performance, attrition, and class size given missing student data. *Economics of Education Review*. 2001;20:377-388.
8. Becker WE, Walstad WB. Data loss from pretest to posttest as a sample selection problem. *The Review of Economics and Statistics*. 1990;72(1):184-188.
9. Diggle PJ, Heagerty P, Liang KY, Zeger SL. *Analysis of longitudinal data* (second edition). Oxford University Press. New York, 2002.
10. Fayers PM, Machin D. *Quality of life: the assessment, analysis and interpretation of patient-reported outcomes* (Second Edition). West Sussex: John Wiley & Sons, 2007.
11. Felix AJ, Kannan R. *Statistical models in survival analysis*, Chap-3, 2007, pp 50-51.
12. Fitzmaurice GM. Methods for handling dropouts in longitudinal clinical trials. *Statistica Neerlandica*, 2003;57:75-99.
13. Fitzmaurice GM, Laird NM, Ware JH. *Applied Longitudinal Analysis*. John Wiley and Sons. New York, 2004.
14. Hedeker D, Gibbons RD. *Longitudinal Data Analysis*. John Wiley & sons. New Jersey, 2006.
15. Hesketh SR, Skrondal A. *Multilevel and Longitudinal Modeling Using Stata* (2nd ed.). College Station, TX: Stata Press, 2008.
16. Holt D. Missing data and nonresponse. In J.P. Keeve (Ed.) *Educational research, methodology, and measurement: An international handbook* (2ed.). New York: Elsevier Science Ltd. 1997.
17. Kantelhardt E, Zerche P, Mathewos A, Trocchi P, Addissie A, Aynalem A, *et al.* Breast cancer survival in Ethiopia: A cohort study of 1,070 women. *International Journal of Cancer*. 2014;135(3):702-709.
18. Kim JO, Curry J. The treatment of missing data in multivariate analysis. *Sociological Methods & Research*. 1977;6(2):215-240.
19. Lee ET, Wang J. *Statistical methods for survival data analysis*. Edition-3, John Wiley & Sons. New York, 2003.
20. Little RJA, Rubin DB. *Statistical Analysis with Missing Data*. John Wiley & Sons, New York, 1987.
21. Mohanraj J, Srinivasan MR. Missing longitudinal data analysis with covariance structure. *Aligarh Journal Of Statistics*. 2018;38:83-102.
22. Parkin DM, Bray F, Ferlay J, Jemal A. *Cancer in Africa 2012*. *Cancer Epidemiology and Prevention Biomarkers*. 2014;23(6):953-966.
23. Pavithra V, Kannan. Impact of Missing Data in Survival Analysis. *High Technology Letters*. 2022;28(9):374-384.
24. Rezaianzadeh A, Peacock J, Reidpath D, Talei A, Hosseini SV, Mehrabani D. Survival analysis of 1148 women diagnosed with breast cancer in Southern Iran. *BMC Cancer*. 2009;9(1):168.
25. Rubin DB. Inference and missing data. *Biometrika*. 1976;63:581-592.
26. Rubin DB. *Multiple Imputation for Non-response in Surveys*. John Wiley & Sons. New York, 1987.
27. Scheffer J. Dealing with missing data. *Research Letters in the Information and Mathematical Sciences*. 2002;3:153-160.
28. Sinaga ES, Ahmad RA, Hutajulu SH. *Berita kedokteran masyarakat*. Fakultas Kedokteran, Universitas Gadjah Mada, 2017, 33.
29. Torre LA, Bray F, Siegel RL, Ferlay J, Lortet-Tieulent J, Jemal A. *Global cancer statistics, 2012*. CA: A Cancer Journal for Clinicians. 2015;65(2):87-108.