

## Analysis of longitudinal binary outcomes in clinical trials with low percentage of missing values

Alireza Akbarzadeh Baghban<sup>1,\*</sup>, Erfan Ghasemi<sup>2</sup>, Farid Zayeri<sup>3</sup>, Saeed Asgary<sup>4</sup>, Mahshid Namdari<sup>2</sup>

<sup>1</sup> Department of Basic Sciences, School of Rehabilitation Sciences, Shahid Beheshti University of Medical Sciences, Tehran, Iran

<sup>2</sup> Department of Biostatistics, Faculty of Paramedical Sciences, Shahid Beheshti University of Medical Sciences, Tehran, Iran

<sup>3</sup> Proteomics Research Center, Faculty of Paramedical Sciences, Shahid Beheshti University of Medical Sciences, Tehran, Iran

<sup>4</sup> Iranian Center of Endodontic Research, Dental Research Center, Shahid Beheshti University of Medical Sciences, Tehran, Iran

\*Corresponding author: e-mail address: [akbarzad@sbmu.ac.ir](mailto:akbarzad@sbmu.ac.ir) (Alireza Akbarzadeh Baghban)

### ABSTRACT

In interventional or observational longitudinal studies, the issue of missing values is one of the main concepts that should be investigated. The researcher's main concerns are the impact of missing data on the final results of the study and the appropriate methods that missing values should be handled. Regarding the role and the scale of the variable that missing values have been occurred and the structure of missing values, different methods for analysis have been presented. In this article, the impact of missing values on a binary response variable, in a longitudinal clinical trial with three follow up sessions has been investigated. Propensity Score, Predictive Model Based and Mahalanobis imputation strategies with complete case and available data methods have been used for dealing with missing values in the mentioned study. Three models; Random intercept, Marginal GEE and Marginalized Random effects models were implemented to evaluate the effect of covariates. The percentage of missing responses in each of the treatment groups, throughout the course of the study, differs from 6.8 to 14.1. Although, the estimate of variance component in random intercept and marginalized random effect models were highly significant ( $p < 0.001$ ) the same results were obtained for the effect of independent variables on the response variable with different imputation strategies. In our study according to the low missing percentage, there were no considerable differences between different methods that were used for handling missing data.

**Keywords:** Missing data; Longitudinal study; Binary response; Multiple imputation

### INTRODUCTION

Some of clinical trials use longitudinal designs because treatments are expected to change the response over time [1]. Missing response measurements are a common problem in longitudinal trials [2]. Missingness can substantially reduce the number of cases in a data set and subsequently decrease the power. On the other hand, completers may be a non-random sample from the original group that was assessed. So, missingness can result in biased treatment comparisons [3]. In a review of 331 articles during a 18 month period of 2004–2005 in The New England Journal of Medicine, only 26 (8%) have reported some form of missing data methods [4]. It seems that few authors in the field of medicine considered the impact of missing data. Usually, for analyzing longitudinal studies with missingness, complete

case or available data methods are used. In fact, there is a considerable difference between optimum statistical methodologies and methods that are commonly used in practice. The use of complete case methods without checking the mechanism of missingness can cause inefficient and potentially biased estimates [2].

According to Little & Rubin [5], a missing data mechanism is said to be missing completely at random (MCAR) if missingness is not related to any observed or unobserved factor. On the other hand, missing at random (MAR) assumes that conditional on observed factors, the missingness is independent of the unobserved data. Otherwise, when the missingness depends on unobserved quantities, it will be termed as non-ignorable or not missing at random and its abbreviation is NMAR. Since different statistical methods are valid only under certain missing mechanisms; the appropriate method

for the analysis should be selected according to the mechanism of missingness. However, it should be noted that the true missing data mechanism is not testable from the data. It is only possible to suggest that the data are not consistent with the MCAR assumption and no amount of clever modeling can overcome this issue [6]. Several models have been proposed for the analysis of longitudinal binary responses. Some of them are marginal models with the GEE approach, random-effects models, marginalized random effect and transitional models. They are extensions of the well-known logistic regression for correlated data; that is a particular case of the generalized linear models with a logistic link function [7].

In marginal models, the population-averaged effect of covariates on the longitudinal response is directly specified and the regression coefficients have interpretation for the population. Whereas the random-effects model gives relationships conditionally on having certain individual characteristics modeled by the random effects. On the other hand, in marginalized models, the population averaged effect of covariates on responses is specified conditionally on random effects or previous history of responses [8]. In the case of missing values, parameter estimates based on marginal GEE models can be biased under MAR [9, 10]. To overcome this problem under MAR mechanism, multiple imputations based on generalized estimating equation (MI-GEE) or random effect models can be used [11, 12]. By imputation, missing values are filled in with particular values by specific procedures. Then standard methods for data analysis can be held on the complete (imputed) data set. The goal of missing data imputation is to preserve important characteristics, such as mean and variance, of the whole data set, so the results would be efficient and unbiased.

The multiple imputation (MI) method produces more than one imputed data set. Each imputed data set contains slightly different imputed values. The data analysis procedure is then conducted on multiple imputed data sets and the results from different data sets are combined using Rubin's rules [13]. In this paper marginal model with the GEE approach, random intercept model and marginalized random effect model with different imputation strategies for dichotomized outcome is utilized in a longitudinal dental clinical trial study.

## MATERIAL AND METHODS

### Statistical approaches:

Marginal models with Generalized estimating equations (GEEs) are formalized by Liang and Zeger [9] to extend generalized linear models (GLMs) to a regression setting with correlated observations. GEEs are used to characterize the marginal expectation of a set of outcomes as a function of a set of study covariates [14]. Otherwise, Random-effects models which are also called generalized linear mixed models are an alternative but closely related approach to GEE marginal models. The underlying premise of linear mixed effects models is that some subsets of the regression parameters vary randomly from one individual to another. In mixed effects models the mean response is modeled as a combination of population characteristics that are assumed to be shared by all individuals, and subject-specific effects that are unique to a particular individual. However, in the simplest case only a random intercept model can be introduced [15].

The marginal model can be written as

$$\text{logit}[P(Y_{ij} = 1|X_{ij})] = X_{ij}\beta \quad (1)$$

and the random intercept model as

$$\text{logit}[P(Y_{ij} = 1|X_{ij}, U_i)] = X_{ij}\beta^* + U_{i0} \quad (2)$$

Where  $U_{i0}$  is the random intercept [16]. It is important to note that the vectors  $\beta$  in model 1 and  $\beta$  in model 2 are not equal and the estimators estimate different things. The functional form of marginalized random effect model is

$$\text{logit} E(Y_{ij}|X_{ij}) = X_{ij}\beta \quad (3)$$

and

$$\text{logit} E(Y_{ij}|b_i, X_{ij}) = \Delta_{ij} + b_{ij} \quad (4)$$

where equation 3 is a marginal logistic regression for the average response as a function of covariates and equation 4 can describe the dependence among the longitudinal measurements (17).

Different assumptions are required for these models. Unfortunately, marginal model using the GEE can be applied either to complete datasets or the mechanism of missingness is completely at random. Parameter estimates based on GEE may be biased and conclusions may be misleading with other kinds of missing mechanisms. Unlike the marginal models, random effect and marginalized models only need MAR assumption [16, 17].

### *Approaches for handling missing values in dichotomized response variable:*

A series of approaches for the situation of missing data in the dependent variable has been proposed. One of them is imputation methods, which replace missing data with estimated values; thereby a complete data set emerges. Multiple imputation (MI) can be used in situations which the MAR assumption is accepted. Different strategies can be used for imputation. The strategies that were used in the current study are described in the following subsection.

#### *1: Complete case analysis*

In complete case analysis, only the cases with completed data are included for analysis, while cases with missing data are excluded. When the data are MCAR, the complete case analysis approach, using either random effect or the marginal model such as GEE approach, is valid for analyzing binary outcomes [18].

#### *2: Available data Methods*

Available data methods are a collection of techniques that can incorporate vectors of repeated measure of unequal length in the analysis. One of the popular available data methods is generalized estimating equations methods. These methods can be used under MCAR mechanism [15].

#### *3: Mahalanobis Distance Method*

In general statistical analysis, the Mahalanobis Distance is a metric that can be used to measure the similarity/dissimilarity between two vectors. The Mahalanobis distance is used to identify cases that have similar characteristics to cases that have missing values. Missing data are filled in by sampling from the closest cases. The multiple imputations are independent repetitions drawn from the range of closest cases. For each case containing a missing value, the Mahalanobis distance between that case and all other cases within the dataset, is calculated using equation 5. The distance is calculated using covariates specified where;  $y$  is the vector of the covariates for the case with the missing value,  $x_i$  is the vector for the  $i_{th}$  fully observed case in the dataset and  $S$  is the covariance matrix for the set of covariates being used in the calculation of the Mahalanobis distance.

$$d(\vec{x}_i, \vec{y}) = \sqrt{(\vec{x}_i - \vec{y})^T S^{-1} (\vec{x}_i - \vec{y})} \quad (5)$$

Each missing value from the imputation variable  $y$  is imputed by values randomly drawn from a subset of observed values, i.e. its donor

pool, with the shortest Mahalanobis distance to the missing data entry that is to be imputed. The Donor Pool defines a set of cases with observed values for that imputation variable [19]. The imputed values are real numbers in  $(-\infty, \infty)$  interval. They can be dichotomized with 0.5 cut point. The imputation based on Mahalanobis Method can be performed with SOLAS 4 software.

#### *4: Predictive Model Based Method*

With the categorical data, the discriminate method is applied in Predictive Model Based imputation. Multiple imputations are generated using a regression model of the imputation variable on a set of user-specified covariates. The imputations are generated via randomly drawn regression model parameters from the Bayesian posterior distribution based on the cases for which the imputation variable is observed plus a randomly drawn error-term. The randomly drawn error-term is added to the imputations to prevent over-smoothing of the imputed data. The regression model parameters are drawn from a Bayesian posterior distribution in order to reflect the extra uncertainty due to the fact that the regression parameters can be estimated, but not determined, from the observed data [12, 20].

#### *5: Propensity Score Method*

The propensity score is the conditional probability of being missing given the observed data. It can be estimated by the means of logistic regression model with a binary outcome indicating whether the data are missing or not [18]. Imputation based on propensity score method is based on the following steps. First, for the variable with missing values, a logistic model is fitted for the probability of missingness (the "propensity score") as a function of all previous variables in the data set. The observations are then grouped based on these propensity scores, and an approximate Bayesian bootstrap imputation is applied to each group. This is done first by drawing a sample with replacement from the set of nonmissing observations, and then assigning the missing observations by sampling from this subset of nonmissing values [21, 22].

### **APPLIED EXAMPLE**

Our applied example is about a longitudinal randomized clinical trial (RCT) in the field of dentistry. The trial was approved and evaluated by the Iranian Ministry of Health as well as by

the Ethics Committee of the Dental Research Center of Shahid Beheshti University of Medical Sciences, Tehran, Iran. The purpose of this RCT was to compare the existence of periapical lesion after one-visit endodontic therapy (OET) and a pulpotomy performed with a new endodontic biomaterial (CEM cement; PCC) in human permanent molars with irreversible pulpitis. A total of 383 selected patients who were met the inclusion criteria were randomly allocated into the OET group (n =190) and the PCC group (n =193). These patients were recruited from 23 healthcare centers in four states and five medical universities of Iran between April and September 2008.

All treated teeth were followed-up clinically and radiographically at baseline and during 6th, 12th and 24th month after pulpotomy or endodontic therapy. The outcome in our model is existence of periapical lesion vs. no periapical lesion. Predictors in the logistic regression included the time of follow up, an indicator of gender (0=female, 1=male), age (in years), an indicator of type of treatment for PCC or OET, an indicator of education levels for academic degree, diploma or under diploma and an

indicator of marital status (0=married, 1=single).

In our analyses, three imputation methods (Mahalanobis Distance, Predictive Model Based and propensity score) were used to address the problem of potentially informative missingness and five iteration was selected for multiple imputation [23]. The imputation was implemented using SOLAS 4. The marginal model with exchangeable logOR structure [24] and random intercept model were used to identify the impact of covariates on the existence of periapical lesion.

The marginal and random intercept models are implemented using GENMOD and NLMIXED procedures. For combining the results of multiple imputations, MIANALYZE procedure in SAS 9.2 is used. The marginalized random effects model is run with lnMLE package in R 2.9.0.

## RESULTS

Table 1 presents the marginal percentages of missing and available responses in each of the treatment groups, throughout the course of the study and an overview of missing patterns is illustrated in Table 2.

**Table 1.** distribution of responses in each of the treatment groups throughout the follow up time

Group	Time of Follow up	missing Number (%)	Lesion Number (%)	No lesion Number (%)	Total Number (%)
OET	baseline	-	61(32.1)	129(67.9)	190(100)
	6 <sup>th</sup> month	13(6.8)	49(25.8)	128(67.4)	190(100)
	12 <sup>th</sup> month	22(11.6)	35(18.4)	133(70)	190(100)
	24 <sup>th</sup> month	26(13.7)	32(16.8)	132(69.5)	190(100)
PCC	baseline	-	54(28)	139(72)	193(100)
	6 <sup>th</sup> month	16(8.3)	17(8.8)	160(82.9)	193(100)
	12 <sup>th</sup> month	26(13.5)	13(6.7)	154(79.8)	193(100)
	24 <sup>th</sup> month	27(14.1)	23(11.9)	143(74.1)	193(100)

OET= one-visit endodontic therapy

PCC= pulpotomy performed with CEM cement

**Table 2.** Overview of missingness patterns and the frequencies with which they occur (O: observed and M:missing)

Follow up time			PCC		OET	
			Number	%	Number	%
6	12	24				
O	O	O	137	71.0	140	73.7
O	M	O	13	6.7	14	7.4
O	O	M	19	9.8	15	7.9
M	O	O	11	5.7	10	5.3
O	M	M	8	4.1	8	4.2
M	M	O	5	2.6	-	
M	O	M	-		3	1.6

OET= one-visit endodontic therapy

PCC= pulpotomy performed with CEM cement

It seems to be an equal percentage of missing values in the treatment groups over the time. The mean age of patients in OET was 26.12 (SD=7.8) and 26.7 (SD=8.4) in PCC. There was no significant difference between the mean age of patients in groups ( $p=0.48$ ); 115 (60.5%) of OET and 126 (65.3%) of PCC were female ( $p=0.34$ ).

The estimated Parameters of the marginal, random intercept and marginalized random effects models with different multiple imputation strategies are displayed in Tables 3-5. It seems that different imputation strategies lead to the same results.

The general conclusion from the comparison between the 3 modeling strategies is the same for the 3 models but the magnitude of estimated parameters are different. The estimates from the marginal model are lower than those from the random intercept model. The differences

between estimates of different modeling approaches are expected, since, the interpretations of their estimates are different. It should be noted that the differences between the estimates of marginal and random effect approaches are largely dependent on the inter-individual heterogeneity. This heterogeneity can be measured by the intercept variance in the random models.

In the random intercept model, the estimate of random intercept variance with different imputation strategies differ from 4.70 to 6.27, which is highly significant and in marginalized random effect model the estimate of random intercept variance with different imputation strategies differ from 0.77 to 0.92. In three models PCC has lower odds of lesion compared to OET and time has a positive effect on treatment success.

**Table 3:** parameter estimates for the marginal models: with GEE with different multiple imputation strategies

Effect		Method					
		Available Data	Mahalanobis	Predictive Model Based	Propensity Score	Complete Case	
Intercept	$\beta$	0.4316	0.4813	0.5442	0.3875	0.7565	
	SE	0.5040	0.5041	0.5126	0.5060	0.5948	
	P value	0.3918	0.3397	0.2884	0.4438	0.2034	
Age	$\beta$	0.0240	0.0243	0.0228	0.0259	0.0190	
	SE	0.0135	0.0135	0.0135	0.0138	0.0159	
	P value	0.0760	0.0724	0.0921	0.0604	0.2326	
Gender (Male vs. Female)	$\beta$	-0.3159	-0.3367	-0.3027	-0.3149	-0.5747	
	SE	0.2058	0.2067	0.2083	0.2061	0.2476	
	P value	0.1249	0.1034	0.1462	0.1266	0.0203*	
Marital (Single vs. Married)	$\beta$	-0.0179	-0.0131	-0.0443	0.0319	0.0329	
	SE	0.2572	0.2576	0.2573	0.2565	0.3181	
	P value	0.9446	0.9593	0.8630	0.9010	0.9176	
Education	Under Diploma vs. Academic Degree	$\beta$	-0.2899	-0.3150	-0.3266	-0.2613	-0.2860
		SE	0.2729	0.2738	0.2793	0.2704	0.3320
		P value	0.2880	0.2498	0.2422	0.3338	0.3890
	Diploma vs. Academic Degree	$\beta$	-0.0717	-0.1241	-0.1605	-0.0651	0.1092
		SE	0.2806	0.2815	0.2873	0.2772	0.3479
		P value	0.7983	0.6594	0.5764	0.8143	0.7536
Group (PCC vs.OET)	$\beta$	0.6209	0.6639	0.6550	0.5962	0.6232	
	SE	0.1919	0.1931	0.1942	0.1897	0.2367	
	P value	0.0012*	0.0006*	0.0007*	0.0017*	0.0085*	
Time	$\beta$	0.3722	0.3186	0.2542	0.3422	0.3550	
	SE	0.0898	0.0874	0.0810	0.0910	0.1085	
	P value	<.0001*	0.0003*	0.0017*	0.0002*	0.0011*	

OET= one-visit endodontic therapy

PCC= pulpotomy performed with CEM cement

**Table 4:** parameter estimates for the Random effects model with different multiple imputation strategies

Effect			Method				
			Available Data	Mahalanobis	Predictive Model Based	Propensity Score	Complete Case
Intercept		$\beta$	0.9599	1.0766	1.1954	0.7492	1.6211
		SE	0.8776	0.8990	0.9497	0.8649	1.1218
		P value	0.2747	0.2311	0.2081	0.3866	0.1496
Age		$\beta$	0.0367	0.0382	0.0383	0.0420	0.0272
		SE	0.0247	0.0254	0.0266	0.0243	0.0293
		P value	0.1374	0.1338	0.1497	0.0846	0.3543
Gender (Male vs. Female)		$\beta$	-0.5276	-0.5777	-0.5447	-0.4944	-0.9555
		SE	0.3321	0.3394	0.3598	0.3207	0.3946
		P value	0.1129	0.0887	0.1301	0.1232	0.0161*
Marital (Single vs. Married)		$\beta$	-0.0948	-0.0766	-0.1232	0.0515	-0.0791
		SE	0.3993	0.4092	0.4301	0.3878	0.4812
		P value	0.8126	0.8514	0.7745	0.8943	0.8696
Education	Under Diploma vs. Academic Degree	$\beta$	-0.5110	-0.5799	-0.6292	-0.4111	-0.5107
		SE	0.4536	0.4635	0.4910	0.4421	0.5751
		P value	0.2606	0.2110	0.2000	0.3526	0.3754
	Diploma vs. Academic Degree	$\beta$	-0.1524	-0.2605	-0.3194	-0.1171	0.1246
		SE	0.4578	0.4704	0.4980	0.4478	0.5753
		P value	0.7393	0.5797	0.5213	0.7938	0.8287
Group (PCC vs.OET)		$\beta$	0.9598	1.0303	1.0764	0.8522	0.9043
		SE	0.3187	0.3274	0.3439	0.3066	0.3820
		P value	0.0028*	0.0017*	0.0017*	0.0055*	0.0186*
Time		$\beta$	0.6078	0.5358	0.4530	0.5364	0.5554
		SE	0.1289	0.1289	0.1221	0.1278	0.1474
		P value	<.0001*	<.0001*	0.0002*	<.0001*	0.0002*
Variance component: Intercept (Random intercept variance)		Estimate	5.0636	5.4491	6.2726	4.6983	5.1719
		SE	0.9734	1.0130	1.1021	0.8671	1.1329
		P value	<.0001*	<.0001*	<.0001*	<.0001*	<.0001*

OET= one-visit endodontic therapy

PCC= pulpotomy performed with CEM cement

## DISCUSSION

Three models were used to characterize the changes in clinical success of two endodontic treatments over time: the marginal model with the GEE approach, the random intercept model and the marginalized random effect model. In addition to the complete case and available data methods, three different imputation strategies were implemented.

The models can be compared to show that the estimated parameters can be different, and to explain these differences. Different imputation strategies can be evaluated for their impact in filling missing values and the final result.

The general conclusion is the same for the 3 models but the estimated parameters are considerably different in marginal and marginalized random effect models with respect to the random intercept model. For instance, the endodontic treatment is very significant in all of the models; but parameter estimates are 0.62, 0.58 and 0.96 for PCC. Differences between the estimators of the marginal and random intercept models are expected. The marginal model expresses averaged relationships without taking into account the fact that the same subjects are considered at each time interval, whereas the random-effects model gives relationships conditionally on having certain individual characteristics modeled by the random effects.

**Table 5:** parameter estimates for the Marginalized Random effects model with different multiple imputation strategies

Effect		Method					
		Available Data	Mahalanobis	Predictive Model Based	Propensity Score	Complete Case	
Intercept	$\beta$	0.4930	0.5556	0.5888	0.3996	0.8302	
	SE	0.5543	0.5576	0.5405	0.5594	0.6797	
	P value	0.3738	0.3193	0.2761	0.4771	0.2220	
Age	$\beta$	0.0223	0.0225	0.0219	0.0257	0.0177	
	SE	0.0164	0.0164	0.0153	0.0165	0.0180	
	P value	0.1731	0.1716	0.1533	0.1231	0.3242	
Gender (Male vs. Female)	$\beta$	-0.2987	-0.3170	-0.2874	-0.2939	-0.5633	
	SE	0.2010	0.2002	0.2025	0.1969	0.2390	
	P value	0.1368	0.1134	0.1559	0.1358	0.0184*	
Marital (Single vs. Married)	$\beta$	-0.0407	-0.0303	-0.0472	0.0364	-0.0110	
	SE	0.2368	0.2409	0.2475	0.2356	0.2971	
	P value	0.8633	0.8998	0.8487	0.8775	0.9704	
Education	Under Diploma vs. Academic Degree	$\beta$	-0.3074	-0.3414	-0.3531	-0.2629	-0.2909
		SE	0.2794	0.2804	0.2788	0.2781	0.3495
		P value	0.2711	0.2235	0.2054	0.3467	0.4053
	Diploma vs. Academic Degree	$\beta$	-0.0730	-0.1374	-0.1684	-0.0620	0.1166
		SE	0.2824	0.2848	0.2829	0.2814	0.3560
		P value	0.7961	0.6297	0.5519	0.8262	0.7432
Group (PCC vs.OET)	$\beta$	0.5812	0.6094	0.6032	0.5354	0.5447	
	SE	0.1974	0.1997	0.1936	0.1950	0.2351	
	P value	0.0032*	0.0025*	0.0018*	0.0063*	0.0205*	
Time	$\beta$	0.3783	0.3232	0.2581	0.3438	0.3541	
	SE	0.0690	0.0685	0.0694	0.0691	0.0922	
	P value	<.0001*	0.0001*	0.0002*	0.0001*	0.0001*	
Variance components: Intercept (Random intercept variance)	Estimate	0.8121	0.8286	0.9158	0.7737	0.8243	
	SE	0.1019	0.0988	0.0908	0.0938	0.1114	
	P value	<.0001*	<.0001*	<.0001*	<.0001*	<.0001*	

OET= one-visit endodontic therapy

PCC= pulpotomy performed with CEM cement

In other words, it can be said that in the marginal model, the exponential of an estimate is a population-averaged odds ratio for clinical success and concerns the sub-population that shares a characteristic relative to the sub-population not sharing the mentioned characteristic.

In the random model, the exponential of an estimate is an odds ratio for a specific person that has a characteristic relative to the same person if she/he were free of that characteristic [16].

In accordance with the results of our study; Nehaus and et al noted that the estimates from the marginal model are systematically lower than those from the random effect model [25].

As it was noted the estimates from the marginal model are systematically lower than those from the random effect models. In addition, the estimate of standard errors in marginal model is smaller compared with random effect models too.

As a result, the same conclusion will be obtained for three classes of models in checking the null hypothesis. For selecting the best model, it should be noted that the missing data mechanism must be examined. The marginal model with GEE approach assumes that the sample is representative of the whole population at each time point and the missing mechanism is MCAR. In contrast, with MAR mechanism the random-effects and marginalized random effect

models are appropriate. Three different imputation strategies were considered in this article lead to similar results, maybe this could be due to the relatively low number of missing values in the data set. Imputation methods resulted in estimates that were more similar to available data method compared with complete case estimates.

In general available data methods are more efficient than complete case methods, because they use partial information obtained from those who dropout [15].

## REFERENCES

1. Yang X, Shoptaw S. Assessing missing data assumptions in longitudinal studies: an example using a smoking cessation trial. *Drug Alcohol Depend.* 2005;77:213–25.
2. Horton NJ, Kleinman KP. Much ado about nothing: A comparison of missing data methods and software to fit incomplete data regression models. *Am Stat.* 2007;61(1):79-90.
3. Myers WR. Handling missing data in clinical trials: An overview. *Drug Info J.* 2000;34(2):525-33.
4. Horton NJ, Switzer SS. Statistical Methods in the Journal (research letter). *N Engl J Med.* 2005;353:1977–79.
5. Little RA, Rubin DB. *Statistical Analysis with Missing Data.* 2 ed. New York: Wiley; 2002.
6. Burzykowski T, Carpenter J, Coens C, Evans D, France L, Kenward M, et al. Missing data: discussion points from the PSI missing data expert group. *Pharm Stat.* 2010;9(4):288-97.
7. McCullagh P, Nelder JA. *Generalized Linear Models.* 2 ed. London: Chapman and Hall; 1989.
8. Lee K, Daniels MJ. Marginalized models for longitudinal ordinal data with application to quality of life studies. *Stat Med.* 2008;27(21):4359-80.
9. Liang KY, Zeger SL. Longitudinal data analysis using generalized linear models. *Biometrika.* 1986;73:12–22.
10. Laird NM. Missing data in longitudinal studies. *Stat Med.* 1988;7:305–15.
11. Yoo B. The impact of dichotomization in longitudinal data analysis: a simulation study. *Pharm Stat.* 2010;9(4):298-312.
12. Rubin DB. *Multiple imputation for Non response in Surveys.* New York: Wiley; 1987.

## CONCLUSION

Missing data can cause a reduction in efficiency or precision of the results of the trial. However, the amount of decrease in precision is highly related to the amount of missing data. Although analyses of complete data can be less efficient than methods which use all available data or data sets that their missing values have been imputed, the results of our study show that: when the percentage of missing data is low, different imputation strategies or available data analysis approaches lead to quite similar results.

13. Schafer JL, Graham JW. Missing data: Our view of the state of the art. *Psychol Methods.* 2007;7(2):147-77.
14. Horton NJ, Lipsitz SR. Review of Software to Fit Generalized Estimating Equation Regression Models. *Am Stat.* 1999;53:160-69.
15. Fitzmaurice GM, Laird NM, Ware JH. *Applied Longitudinal Analysis.* New York: Wiley; 2004.
16. Carriere I, Bouyer J. Choosing marginal or random effects models for longitudinal binary responses: application to self-reported disability among older persons. *BMC Med Res Methodol.* 2002;2(1):15.
17. Heagerty PJ. Marginally Specified Logistic-Normal Models for Longitudinal Binary Data. *Biometrics.* 1999;55(3):688-98.
18. Ma J, Akhtar-Danesh N, Dolovich L, Thabane L, and the CHAT investigators. Imputation strategies for missing binary outcomes in cluster randomized trials. *BMC Med Res Methodol.* 2011;11:18.
19. Mahalanobis Distance Matching Method. Available from: <http://www.statistical-solutions-software.com/solas-for-missing-data-analysis/solas-4-0-new-features/mahalanobis-distance-matching-method/>. Accessed September 25, 2012.
20. Gelman A, Carlin J, Stern H, Rubin DB. *Bayesian Data Analysis.* New York: Chapman and Hall; 1995.
21. Lavori PW, Dawson R, Shera D. A multiple imputation strategy for clinical trials with truncation of patient data. *Stat Med.* 1995;14:1913-25.
22. Kang T, Kraft P, Gauderman WJ, Thomas D. Multiple imputation methods for longitudinal blood pressure measurements from the Framingham Heart Study. *BMC Genet.* 2003;4(Suppl1):S43.



23. Fitzmaurice GM, Davidian M, Molenberghs G, Verbeke G. *Longitudinal Data Analysis*. New York: Chapman & Hall/CRC; 2009.

24. Diggle PJ, Heagerty P, Liang KY, Zeger SL. *Analysis of Longitudinal Data*. 2nd ed. Oxford University Press; 2002.

25. Neuhaus JM, Kalbfleisch JD, Hauck WW. A comparison of cluster-specific and population-averaged approaches for analyzing correlated binary data. *Int Stat Rev*. 1991;59:25-36.