



## A selection modelling approach to analysing missing data of liver Cirrhosis patients

Dilip C. Nath<sup>1</sup>, Ramesh K. Vishwakarma<sup>2</sup>, Atanu Bhattacharjee<sup>3</sup>

<sup>1</sup>Assam University, Silchar, India, e-mail: dilipc.nath@gmail.com,

<sup>2</sup>Department of Statistics, Gauhati University, India, e-mail:  
ramesh.biostats@gmail.com,

<sup>3</sup>Department of Biometrics, Chiltern Clinical Research Ltd, India, corresponding author  
e-mail: atanustat@gmail.com

### SUMMARY

Methods for dealing with missing data in clinical trials have received increased attention from the regulators and practitioners in the pharmaceutical industry over the last few years. Consideration of missing data in a study is important as they can lead to substantial biases and have an impact on overall statistical power. This problem may be caused by patients dropping before completion of the study. The new guidelines of the International Conference on Harmonization place great emphasis on the importance of carefully choosing primary analysis methods based on clearly formulated assumptions regarding the missingness mechanism. The reason for dropout or withdrawal would be either related to the trial (e.g. adverse event, death, unpleasant study procedures, lack of improvement) or unrelated to the trial (e.g. moving away, unrelated disease). We applied selection models on liver cirrhosis patient data to analyse the treatment efficiency comparing the surgery of liver cirrhosis patients with consenting for participation HFLPC (Human Fatal Liver Progenitor Cells) infusion with surgery alone. It was found that comparison between treatment conditions when missing values are ignored potentially leads to biased conclusions.

**Key words:** selection model; model for end-stage liver disease; missing not at random

### 1. Introduction

Cirrhosis is a state of the liver in which the tissue becomes trapped within a sea of scar and struggle to regenerate. It causes gradual shrinkage of the size of the liver. In liver cirrhosis cases, the duration between transplantation and recovery is a crucial period for patients. Patients are generally monitored through follow up periods with liver functioning effects. The

severity of the disease in a patient with liver cirrhosis is measured using the model for end stage liver disease (MELD). The MELD Score (UNOS Modification)(Kamath and Kim, 2007) is calculated as follows

$$\begin{aligned} \text{MELD Score} = & 6.57 \times \log_e(\text{serum creatinine}) + \\ & + 3.78 \times \log_e(\text{bilirubin}) + 11.2 \times \log_e(\text{INR}) \end{aligned} \quad (1)$$

The MELD is a useful tool for determining liver status in patients. It is a widely used method for organ allocation in liver transplantation. The score function is formulated with consideration of liver and renal functions. Several studies have concluded that the low MELD scores can indicate a higher risk of mortality in liver transplanted patients. Different biochemical parameters are used to calculate the MELD score (Kamath and Kim, 2007). The MELD score is used as a primary end point to perform the analysis. However, to better investigate the effect of treatment, one is often interested in evaluating how the parameters of interest change over time. These kinds of follow-up studies are also known as longitudinal studies, and are very common for clinical trials. Every study has several variables among which one is the target or primary response variable and the others are potentially explanatory variables or covariates. Anyone of them can change from visit to visit and could be measured over time but usually the response variable is obtained by repeated measures. We are interested in how the change or variation of the response variable can be explained by the explanatory variables, and statistical modeling is exploited to approximate the relationship between the two kinds of variables.

Diggle, Liang and Zeger (2002) presented a detailed summary of models used for longitudinal data. Popular choices are linear models, generalized linear models, transition models and mixed effects models. We can obtain the likelihood function (or conditional likelihood function) based on the models used as the joint probability distribution function overall subjects and all time points. MLE (Maximum Likelihood Estimation) methods are then used to make the estimation and inference. With the same kinds of models, the computation of estimates becomes complicated when a few data points are missing. But missing data are very common for clinical trials with longitudinal studies, reflecting the problematic nature of the phenomenon under study.

A wide range of statistical models for analysing outcomes with missing data are available, but their validity depends on the nature of the missing

data mechanism as well as on the assumptions used. The proportion of missing data is sometimes noticeably large enough to cause the results to be significantly biased. Although investigators and clinical research co-ordinators may devote substantial efforts to minimising the number of missing values, some amount of missing data is inevitable in the practice of randomised clinical trials. The data with missing observations can be classified into two categories: intermittent missing values (i.e., missing values due to occasionally absence), with observed values afterwards, and dropout values (i.e., missing values due to earlier withdrawal), with no more observed values. Reasons for dropout are often study-related, e.g., negative side-effects of the tested medicines, ineffectiveness of the intervention, and inappropriate conduction of the therapy (Dragset, 2009; Kaciroti and Raghunathan, 2014). Without careful handling of dropouts, either biased parameter estimates or invalid inferences would result. This may also be true for intermittent missing values. It is not possible for us to compute the original joint probability function of all the repeated measures directly with missing data ((Kaciroti et al., 2012)). There are at least three ways to deal with this situation: using complete case only, analysing as incomplete, and imputation. The complete cases only approach means just keeping the complete cases with specified weight and dropping all other cases with missing values. The imputation method is to fill in the missing values with randomly generated ones; but generating appropriate values to approximate the original ones is still challenging (Kim and Yu, 2012). Analysis as an incomplete model means to integrate out the intermittent missing values and to ignore the dropout missing values, which seems a straightforward choice, but is computationally expensive and not always stable. All of these three approaches focus on the response variable since we are mainly concerned about estimating the effect of the covariates on the response variable. If the missing data do not affect the estimation at all, we can just use the available data directly and ignore the missing data. Hence the problem arises of the ignorability of the missing values. Usually this is unknown and so a safe and natural method is to use the repeated measures model and the missingness indicators jointly (Kaciroti and Raghunathan, 2014; Kaciroti et al., 2009). There exist three ways to factor the joint distribution of the complete data and missingness indicators: outcome-dependent factorisation, pattern dependent factorisation, and parameter-dependent factorisation. Correspondingly there are three kinds of models: selection models, pattern-mixture models and shared-parameter models (Dragset, 2009). For analysing data with missing values that are po-

tentially missing not at random (MNAR), two widely used approaches are available: pattern-mixture models and selection models. Both approaches derive their inferences based on the joint distribution of outcome and the response data indicator (Kaciroti et al., 2008; Daniels and Hogan, 2008). Selection models partition the function of outcome and response data indicator as the product of a function of outcome and a function of outcome given response data indicator (Nath and Bhattacharje, 2012; Fitzmaurice et al., 2008). This requires explicit modeling of the missing-data mechanism where the probability that a subject is missing may depend on observed and unobserved values (Satty and Mwambi, 2013).

In this study, we primarily focused on selection modeling frameworks to account for non-random dropout. We demonstrate the application of selection models for handling dropout in longitudinal data where the dependent variable is missing across time. We consider the construction of a selection model that uses mixed models and where the outcomes are continuous, to describe the dependency of dropout indicators on the longitudinal measurement. The primary objectives are to investigate the potential influence that dropout might exert on the dependent measurement on the considered data, as well as how to deal with incomplete sequences. We apply this method to a data set arising from a liver cirrhosis study. Section 2 describes the data methodology, including also notation, general concepts and discussion related to the selection model, the model used in the analysis. In Section 3, an application of the selection model to liver cirrhosis data is described. The results obtained are elaborated in Section 4. Section 5 contains a discussion of this application, and finally the findings and drawbacks are presented in Section 6.

## 2. Data and methodology

### 2.1. The data

The data considered for this study were taken from the path [www.mayo.edu/int-med/gi/model/mayomodl-5-unos.htm](http://www.mayo.edu/int-med/gi/model/mayomodl-5-unos.htm), accessed on Jan 28, 2013. Patients with MELD score of 12 to 24 were considered not to qualify for a solid liver transplant because of standard co-morbidities and should have an estimated life expectancy of approximately 6 to 18 months. Patients with liver cirrhosis aged between 18 and 70 years, with no gender restriction, were eligible to take part in the study. Assessments for early determination

of patient eligibility (14 to 10 days prior to the day of planned cell transplantation) included a full physical examination and medical history, and grading of encephalopathy. The patients were brought to the radiology suite approximately 30 to 60 minutes before splenic artery catheterisation. After initiation of conscious sedation, a catheter was inserted into the femoral artery and, under fluoroscopic guidance, passed into the splenic artery. The final position of the catheter was confirmed with a small volume of contrast dye. Blood pressure, heart rate, respiratory rate and O<sub>2</sub> saturation were monitored frequently during cell infusion and until the catheter was removed. Confirmation of splenic artery and splenic vein patency with contrast media was performed just before cell infusion and again just after cell infusion. Post-infusion patients were kept under close surveillance in the intensive liver care unit (ILCU) for 5 days (3 days is normal in case of post-catheterisation). In case of required medical attention the stay was prolonged. All patients (the new therapy group and conventional therapy group) were provided with a schedule of follow up visits and were instructed to return for next follow-up visit. All relevant biochemical, pathological and radiological examinations were performed at each visit along with abdominal ultrasound examination to analyse the portal vein caliber and presence of ascites. The principal investigator was asked to review the subject safety data as it was generated. Study outcome and adverse event data were reviewed to determine whether there was any change to the anticipated benefit-to-risk ratio of study participation and whether the study should continue as originally designed, should be changed, or should be terminated. Any changes to the anticipated benefit-to-risk ratio of study participation or recommendations related to continuing, changing or terminating the study were to be reported to the ethics committee. In this study a total of 114 patients with cirrhosis were randomised in two different treatment groups. The patients were asked to read the informed consent form if they were interested in participating in the study. After agreeing to participate in the study, the selection tests and procedures were performed to determine the eligibility of patients to participate. The study doctor reviewed the results of the screening tests to determine if a patient could participate in the study. If so, the patient was assigned randomly to the treatment or control group based on a randomization procedure. In the treatment group, patients received human fetal liver progenitor cell (HFPLC) infusion as background conventional medical treatment, while in the control group patients received only the conventional medical treatment without the infusion. The duration

of the study was 36 months, with 7 visits (baseline, 3 months, 6 months, 9 months, 12 months, 24 months and 36 months). The parameters under consideration were MELD Scores, therapy and visit. Each patient's MELD score was taken at seven different time points.

## 2.2. Notations

Let us suppose  $y_i$  is the value of MELD scores for the  $i^{th}$  individual and  $x_i$  is the value of the covariate for the corresponding  $i^{th}$  individual. We divide  $y_i$  into two parts  $(y_i^{obs}, y_i^{miss})$ , with  $y_i^{obs}$  representing the observed values, and  $y_i^{miss}$  representing values that would be observed if they were not missing  $\theta$  represents parameters of the model for repeated measures, and  $\phi$  represents the parameters of the missingness mechanism. Let  $r_i$  be the missingness indicator for repeated observations on subject  $i$ . Suppose the missing data are due to dropout; then the measurements for each subject can be considered up to a defined time point, after which all data are unrecorded. In this case, a dropout indicator can then be defined as  $D_i$ , given by  $D_i = 1 + \sum_{j=1}^n r_{ij}$ , denoting the instance at which dropout first occurred.

## 2.3. Methods

Data with missing observations can make parameter estimation and inference much more complicated. Evaluation of the likelihood function requires integration. Missing observations can apply to both the dependent variables and the covariates. In modeling missing data, it is frequently necessary to adopt a joint model for the measurement process together with the dropout process. Therefore, the full data density is given by

$$f(y_i, r_i | X_i, Z_i, \theta, \phi), \quad (2)$$

where  $X_i$  denotes the design matrix for fixed effects, and  $Z_i$  denotes the design matrix for random effects. When missing values are introduced by dropout, the pattern of missingness can be signified by a scalar  $d_i$ . Considering the above model in equation (2), we can factorise this joint density function in the form of a selection model defined by the conditional factorisations of the joint distribution of Y and R; both are discussed in more detail in Little (1995) and stated briefly below. A selection model is based on the following factorisation

$$f(y_i, r_i | X_i, Z_i, \theta, \phi) = f(y_i | X_i, Z_i, \theta) \times f(r_i | y_i, X_i, \phi), \quad (3)$$

where the first factor in the above factorisation represents the marginal density of the measurement process, while the second factor represents the density of the dropout process, conditional on the measurements. The missing data processes have been developed by Rubin (1976) and Little and Rubin (2014) through the selection model framework. They make distinctions among different missing data processes. These processes can be formulated based on the second factor of equation-(6), i.e.,

$$f(r_i|y_i, X_i, \phi) = f(r_i|y_i^{obs}, y_i^{miss}, X_i, \phi). \quad (4)$$

Thus, if the distribution of the missingness process is reduced to  $f(r_i|y_i, X_i, \phi) = f(r_i, X_i, \phi)$  i.e., the process is independent of the measurements, then the process is defined as missing completely at random (MCAR). If the missingness probability depends on the observed measurement  $y_i^{obs}$ , but not on the missing measurements  $y_i^{miss}$ , i.e.,  $f(r_i|y_i, X_i, \phi) = f(r_i|y_i^{obs}, X_i, \phi)$ , then the process is termed missing at random (MAR). Finally, data are missing not at random (MNAR) or exhibit an informative process, when the missingness probability depends on the unobserved measurement,  $y_i^{miss}$ , and possibly on the observed measurement,  $y_i^{obs}$ , i.e.,  $f(r_i|y_i, X_i, \phi) = f(r_i|y_i^{obs}, y_i^{miss}, X_i, \phi)$ .

#### 2.4. Modelling with missing values

As mentioned above, a selection model factors the joint distribution into two parts: the marginal measurement model that describes the distribution of the complete measurements, and the missingness model that describes the conditional distribution of the response indicators given the observed and unobserved measurements. In other words, in a selection model, we first specify a distribution for the measurement, and then suggest a manner in which the probability of being observed depends on the data. For continuous responses, using a selection model formulation as in equation (6), Diggle and Kenward (1994) combine a multivariate Gaussian linear model together with the dropout model. In the same way, we consider the measurement model to be of the linear mixed-effects model type (Laird and Ware, 1982). Recall that  $y_{ij}$  is the response of interest for the  $i$ th subject in the study, where  $i = 1, \dots, n$ , at time point  $j$ , where  $j = 1, \dots, n_i$ . More generally, the model for  $y_i$  the  $(n_i \times 1)$  vector of responses for the  $i^{th}$  subject can be written as

$$y_i = X_i\beta + Z_ib_i + \epsilon_i, \quad (5)$$

where  $X_i$  and  $Z_i$  are known  $(n_i \times p)$  and  $(n_i \times q)$  design matrices for fixed and random effects respectively,  $\beta$  is the  $(p \times 1)$  vector of fixed effects,  $\mathbf{b}_i$  is the  $(q \times 1)$  vector of the random effects distributed as  $N(0, K)$ ,  $\epsilon_i$  is the  $(n_i \times 1)$  vector of the residual components distributed independently as  $N(0, \Sigma_i)$ ,  $K$  is the general  $(q \times q)$  covariance matrix with  $(i, j)^{th}$  element and  $\Sigma_i$  is the  $(n_i \times n_i)$  error covariance matrix.

As noted previously, we focus only on incompleteness due to dropout, and thus we assume that the first measurement  $Y_{i1}$  is measured for all subjects in the study. In accordance with the notation introduced in section 2.2, the selection model arises when the joint likelihood of the measurement process and the dropout process is factorised as follows

$$f(y_i r_i | X_i, Z_i, \theta, \phi) = f(y_i | X_i, \theta) \times f(r_i | y_i, Z_i, \phi). \quad (6)$$

We use the linear mixed-effects model introduced in equation (5) to model the measurement process, together with a logistic regression to describe the dropout process. According to Diggle and Kenward (1994), the model for the dropout process is based on a logistics regression for the conditional probability of dropout at occasion  $j$ , given that the subject is still in the study. Again, the  $m_i(y_{ij}, h_{ij})$  denotes this probability of dropout at time  $j$ , where  $h_{ij} = (Y_{i1}, Y_{i2}, \dots, Y_{i,j-1})$  is a vector possibly containing all observed measurements up to and including occasion  $j - 1$  as well as relevant covariates in the conditional probability of dropout model. Theoretically, the dependence on future unobserved measurements is possible to justify but not straight forward; for simplicity, we model dependence only on the first-order history. Therefore, modeling the dropout mechanism may be simplified by allowing dropout to depend upon the current measurement and immediately preceding measurement only, with corresponding regression coefficients, i.e.,  $\phi_1$  and  $\phi_2$ . In particular, for subjects with observed measurements, dropout depends on the measurement prior to the last measurement ( $y_{i,j-1}$ ) and the current unobserved measurement ( $y_{ij}$ ). A commonly used version of such a logistic dropout model is

$$\text{logit } \Pr(D_i = j | D_i \leq j, h_{ij}, y_{ij} \phi) = \phi_0 + \phi_k z_i + \phi_1 y_{i,j-1} + \phi_2 y_{ij}, \quad (7)$$

where  $\phi_0$  and  $\phi_k$  denote respectively the intercept and the vector of parameters for covariates  $Z_i$ , respectively. The model in equation (7) contains special cases corresponding to MAR and MCAR mechanisms that can be obtained from  $(\phi_2 = 0, \phi_1 > 0)$  and  $(\phi_2 = \phi_1 = 0)$ , respectively.



The selection models originated from the Tobit models for analysing missing data (Heckman, 1976). Verbeke and Molenburghs, (2009) addressed the theoretical translation from the Tobit model to Diggle and Kenward's selection model. Subsequently, Troxel, Harrington, and Lipsitz, (1998) extended it to the non-monotone setting. Selection models for categorical and other type of measures were also developed; see Fitzmaurice, et al. (1995), Molenberghs et al. (1997), Nordheim (1984), and Kenward and Molenburghs (1999).

### 3. Application to the liver cirrhosis data

In this section, we describe the application of the selection model to data from liver cirrhosis patients. More background details of this data are given in Nath et. al (2015).

#### 3.1. Fitting of selection model

In line with Diggle and Kenward (1994), we fit the selection models to the liver cirrhosis data by combining the measurement model with the logistic regression for the dropout model. The combined model for measurement dropout will be fitted to the MELD by maximum likelihood using a generic function maximisation routine. We use the linear mixed-effects model in the form

$$y_i = X_i\beta + Z_i\beta_i + \epsilon_i \quad (8)$$

in order to obtain initial values for the parameters of the measurement model. In the fitted model, we assume different intercepts and treatment effects for each of the seven time points, with a  $(7 \times 7)$  unstructured variance covariance matrix. Specifically, we consider a multivariate normal model, with unconstrained time trend under the conventional and new therapy groups. Since the liver cirrhosis data cover 114 subjects ( $i = 1, \dots, 114$ ) on seven time points ( $j = 3, 6, 9, 12, 24, 36$ ), the model can be written as

$$y_{ij} = \text{Visit}_j\beta + \beta\text{TRT}_i + \epsilon_i, \quad (9)$$

where  $\text{TRT}_i = 0$  for conventional therapy and  $\text{TRT}_i = 1$  for new therapy. In this way, the parameter estimates and standard errors as well as p-values for the eight mean model parameters can be obtained. To fit this model, we use the lme4 package for R software. Next, we consider the dropout

model. The dropout will be allowed to be independent of covariates. We fit the model with an intercept, an effect for previous outcome and an effect for the current unobserved measurement, corresponding to MCAR, MAR and MNAR, respectively. Dependence on future unobserved measurements is theoretically possible; we model dependence on the current missing observation through the dropout model. The probability of MELD scores is assumed to follow the logistic regression model (a commonly used model for dropout processes; see, Molenberghs and Kenward 2007) in equation (7). Estimation of a selection model for MNAR can be seen as a major complication as the dropout indicators depend on the unobserved measurement. For example, in the selection model mentioned above, the dropout indicators depend in part on the unobserved longitudinal measurements at the time of dropout. This leads to complexity in assessing the likelihood function (Diggle and Kenward, 1994). Generally, the parameters were estimated using code written in SAS provided by Dmitrienko et al. (2005) that maximises the log-likelihood for the model using PROC IML. However, a recent development suggests that the MCMC approach can be used to handle this situation. The MCMC approach was used to find the response variable and the selection model that fit the missing data indicators, given the responses are defined as follows. Let  $i$  be the subject index,  $j = (0, 3, 6, 9, 12, 24, 36)$  be the visit index, and  $m = (1, 2)$  (1=new therapy, 2=conventional therapy) be the treatment index. We modelled the response variable  $MELDScores_i = change_{ji}$  using a multivariate normal distribution with mean  $\mu_i = (\mu_{0i}, \mu_{3i}, \mu_{6i}, \mu_{9i}, \mu_{12i}, \mu_{24i}, \mu_{36i})$  and covariance matrix  $\Sigma$ . The mean parameter for the  $i^{th}$  subject at the  $j^{th}$  visit has the regression model,

$$\mu_{ij} = \mu_{mj} \times (\text{Baseline MELD Score} - \text{constant}) + \epsilon_i, \quad (10)$$

where  $\mu_{mj}$  is the treatment effect for the  $j^{th}$  visit. All parameters are given a flat prior:  $\pi(m_{kj}), \pi(\beta_j), \pi(\epsilon_i)$  equivalent to 1.

The group intercept  $\epsilon_i$  with a flat prior makes the design matrix rank deficient. A common strategy is to use a constraint by setting one of the redundant  $\epsilon_i$  parameters to 0 and reducing the total number of parameters  $\epsilon$  by 1. We did not place constraints on the  $m_k$  parameters because the model does not have overall (multidimensional) intercepts over the visits. The logistic function was used to model the dropout probabilities.

$$\text{logit}(\Pr(D_i = j | D_i \leq j, h_{ij}, y_{ij} \phi)) = \phi_0 + \phi Z_i + \phi_1 y_{ij-1} + \phi_2 y_{ij}, \quad (11)$$

where  $\phi_0$  and  $\phi_k$  denote the logistic regression coefficients for the current and immediately previous observations respectively, and  $j$  denotes the time points. The covariates include previous and current (possibly missing) response variables, because we want to infer conditionally on not having withdrawn earlier if either of the following might affect a patient's willingness to continue the trial: MELD score improvement from previous visit or the current state of wellness. This type of selection model is referred to as the Diggle-Kenward selection model (Diggle and Kenward 1994; Daniels and Hogan 2007). Because there is an expected treatment effect, we want the logistic regression to include separate intercepts and regression coefficients for each treatment group. In the model, the covariance matrix takes on an inverse Wishart prior distribution, and the remainder of the parameters are assigned flat priors.

#### 4. Results

In this study, many MELD measurements contain missing values because laboratory specimens were lost or inadequate, or patient follow-up was terminated. In addition, all subjects have observed values at baseline and at 3, 6, 9, 12, 24 and 36 months of follow-up. Some individuals underwent liver transplantation with either new or conventional therapy, but dropped out of the study before the scheduled post-baseline time. Most of the individuals began dropping of the study from 12 months onwards. Therefore, the data presents three possible dropouts patterns (dropout at time points 12, 24, or 36). All 114 patients are observed at the first occasion (baseline), whereas there are a total of 103, 94, 86, 67, 14 and 3 patients seen at the month 3, month 6, month 9, month 12, month 24 and month 36 respectively. The percentage of patients that are still in the study and dropped out from the study after each follow-up visit are summarised by treatment group in Table 1.

Table 2 describes the baseline and demography characteristics of liver cirrhosis patients in the new therapy and conventional therapy groups. The mean age of liver cirrhosis patient in the new therapy group is 48.6 with standard deviation (SD) 9.38, whereas in the conventional therapy group, the mean age is 49.85 with SD 11.06. The distribution of males in the new therapy and conventional therapy groups is 10(19.2%) and 40(64.5%) respectively. Accordingly, the distribution of females in the new therapy and conventional therapy groups is 42(88.8%) and 22(35.5%) respectively. The mean height of liver cirrhosis patients in the treatment group is 164.9 with

**Table 1.** Classification of dropouts in MELD scores obtained from liver cirrhosis patients

visit	dropout(1=Available 0=Missing)	Conventional Therapy	New Therapy
Baseline	1	62(54.4%)	52(45.6%)
Month 3	0	10(90.9%)	01(09.1%)
Month 3	1	52(50.5%)	51(49.5%)
Month 6	0	17(85%)	3(15%)
Month 6	1	45(47.9%)	49(52.1%)
Month 9	0	24(85.7%)	4(14.3%)
Month 9	1	38(44.2%)	48(55.8%)
Month 12	0	34(72.3%)	13(27.7%)
Month 12	1	28(41.8%)	39(58.2%)
Month 24	0	55(55.6%)	44(44.4%)
Month 24	1	7(46.7%)	8(53.3%)
Month 36	0	62(55.9%)	49(44.1%)
Month 36	1	0	3(100%)

**Table 2.** Descriptive statistics on baseline observations of liver cirrhosis patients

Parameters	Treatment Group Mean (SD)	Control Group Mean (SD)
Age	48.6(9.38)	49.9(11.06)
Gender		
Male	10(19.2%)	40(64.5%)
Female	42(88.8%)	22(35.5%)
Height	164.9(4.46)	165.4(5.84)
Weight	65.7(5.24)	69.5(8.78)
Respiratory rate (RR)	25.9(16.59)	21.5(1.62)
Heart Rate	72.9(15.61)	77.2(2.08)

SD 4.46, whereas in the conventional therapy group the mean height is 165.4 with SD 5.84. The mean weight of liver cirrhosis patients in the treatment group is 65.7 with SD 5.24, whereas in the conventional therapy group, the mean weight is 69.5 with SD 8.87. The mean respiratory rate (RR) of liver cirrhosis patients in the treatment group is 25.9 with SD 16.59, whereas in the conventional therapy group the mean RR is 21.4 with SD 1.62. The mean heart rate (HR) of liver cirrhosis patients in the treatment group is 72.9 with SD 15.61, whereas in the conventional therapy group the mean HR is 77.2 with SD 2.08. The parameter estimates obtained from mixed models for the MCAR, MAR, complete cases and all cases approaches are summarised in Table 3. The all cases and complete case approaches have the same estimates for parameters, because the mixed-effect model by default ignores missing values at the time of analysis. In the MCAR approach,

**Table 3.** Parameter estimates for liver cirrhosis patients from mixed-effects models for various techniques

Analysis for all cases				
Effect	Estimate	SE	t-Value	P-Value
Intercept	-0.9297	0.7385	-1.26	0.2129
TRT (1)	0.5343	0.4658	1.15	0.2578
baseline	0.05143	0.02605	1.97	0.0492
visit*TRT (0)	0.8324	0.1286	6.47	<.0001
visit*TRT (1)	1.3233	0.1125	11.77	<.0001
Analysis for Complete cases				
Intercept	-0.9297	0.7385	-1.26	0.2129
TRT (1)	0.5343	0.4658	1.15	0.2578
baseline	0.05143	0.02605	1.97	0.0492
visit*TRT (0)	0.8324	0.1286	6.47	<.0001
visit*TRT (1)	1.3233	0.1125	11.77	<.0001
Analysis with MCAR				
Intercept	4.8597	0.5906	8.23	<.0001
TRT (1)	-0.2906	0.4166	-0.7	0.4886
baseline	-0.1337	0.02021	-6.62	<.0001
visit*TRT (0)	0.4977	0.07222	6.89	<.0001
visit*TRT (1)	0.6769	0.07886	8.58	<.0001
Analysis with MAR				
Intercept	6.4382	0.4637	13.89	<.0001
TRT (1)	-0.5767	0.3425	-1.69	0.3386
baseline	-0.2375	0.0192	-12.36	<.0001
visit*TRT (0)	0.8764	0.0546	16.05	<.0001
visit*TRT (1)	0.78653	0.0566	13.89	<.0001

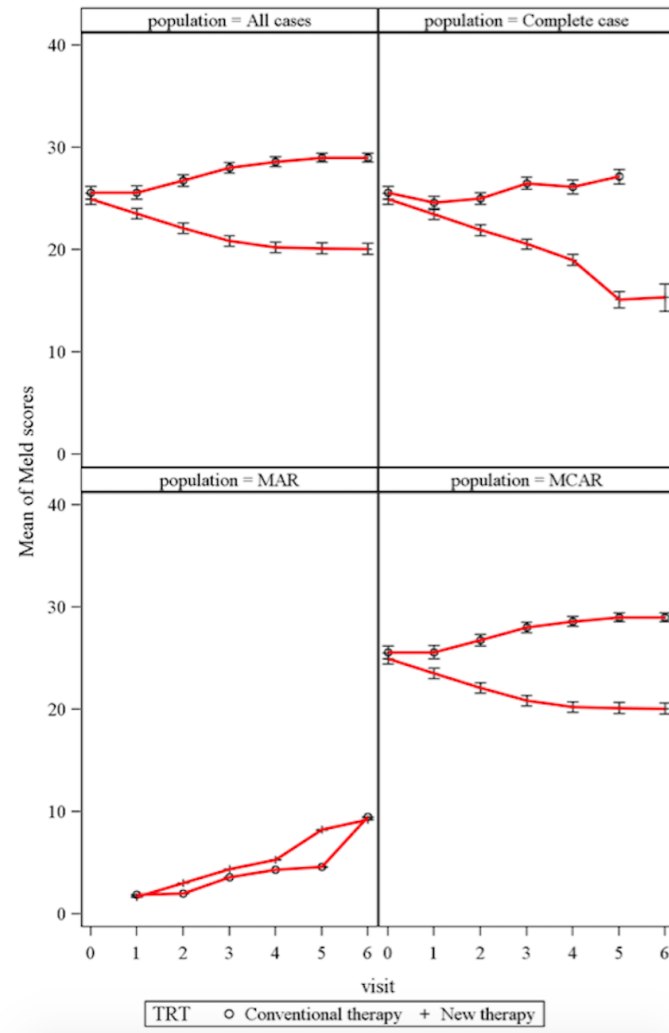
the effects of estimates are reduced as compared to the all cases and complete case approaches. The MAR approach also gives a lower effect of estimates as compared to the all cases, complete case and MCAR approaches. The intercept terms are not significant in the all cases and complete case approaches. Table 4 summarises the parameter estimates obtained from selection modelling. Figure 2 describes the distribution of residuals for change from baseline when data are considered for all cases.

The  $\beta$  parameters indicate approximately the same declining rate from the base MELD score values up to month 12. The  $\beta$  parameters are showing increasing trend after month 12 for MELD scores from the base value. The treatment effects are declining from month 3 to month 12. However, they show a positive trend from month 24 onwards, indicating that the therapy was effective in the long term. The posterior mean estimates for  $(\phi_{\text{month}3})$  and  $(\phi_{\text{month}36})$  are -1.32 and -4.81, respectively. The negative

**Table 4.** Parameter estimates of liver cirrhosis patients from the selection model

Parameter	Mean	Standard Deviation	Percentiles 25%,50%,75%	HPD Interval
$\beta_{month3}$	0.3371	0.2269	0.0737 ,0.3536 ,0.5552	0.0119; 0.6427
$\beta_{month6}$	0.2670	0.1656	0.1255,0.2405,0.4309	-0.00412; 0.5340
$\beta_{month9}$	-0.0593	0.1265	-0.1631,-0.0319,0.0303	-0.2880; 0.1540
$\beta_{month12}$	-0.1142	0.1463	-0.2152,-0.0557,-0.0109	-0.3955; 0.0764
$\beta_{month24}$	-0.0026	0.0291	-0.0248,-0.00879,0.0122	-0.0411; 0.0602
$\beta_{month36}$	0.3513	0.2226	0.1004,0.4245,0.5706	0.0130; 0.6266
Treatment effect:				
month 3	5.7152	4.5286	0.9329 ,6.0906 ,10.0488	-1.2996; 11.7930
month 6	4.7060	3.4095	2.0542,4.2195,8.0191	-1.5188; 9.7635
month 9	-0.9578	2.5164	-3.0633, -0.7817,0.8800	-5.2739; 3.2363
month 12	-1.3984	2.6968	-3.1089 ,-0.3197,0.5414	-6.9806; 1.9360
month 24	3.3286	1.1697	2.5464 ,3.2672 ,4.0799	0.6804; 5.5004
month 36	14.7254	4.6324	11.0381,14.9410 ,18.3104	6.8502; 23.2434
Dropout Model Process				
$\phi_{month3}$	-1.3169	0.4087	-1.5806 -1.3014,-1.0353	-2.1341; -0.5361
$\phi_{month6}$	21.483	25.808	5.48441,1.87045,3.44980	-2.858; -0.81183
$\phi_{month9}$	40.038	75.886	-7.1619,3.29607,8.52209	-1.089; -0.20873
$\phi_{month9}$	50.180	12.344	-2.6232,3.60305,1.32451	-1.944; -0.29778
$\phi_{month12}$	3.7124	0.5388	3.3332,3.6766, 4.0474	2.7117; 4.7995
$\phi_{month24}$	-4.8126	0.6746	-5.2229 -4.7514,-4.3364	-6.1498; -3.5620

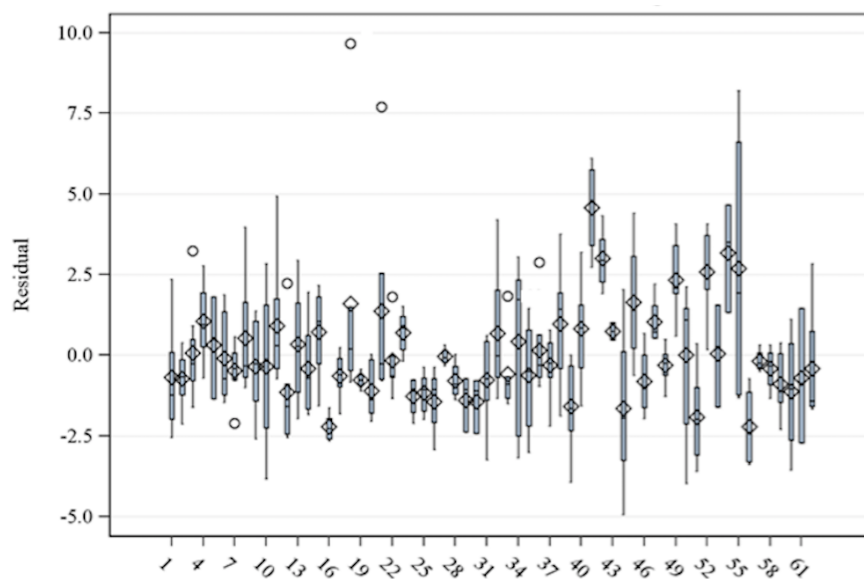
values suggest that patients felt worse (increase in MELD score) at their previous visit because they were more likely to drop out. Figure 1 shows line plots for mean with SE from the calculated data for various approaches. Figure 3 shows the distribution of residuals for change from baseline with the complete case approach. Figure 4 provides the distribution of residuals for change from baseline with the MCAR approach. It is clear from Figure 3 and Figure 4 that the residual distributions are similar for the complete case and all cases approaches. The plots of the posterior densities of treatment and dropout variables are displayed in Figure 5 and Figure 6. The trace plot shows the values of parameters considered during a runtime of 500000 iterations. The marginal density plot is the (smoothened) histogram of the values in the trace-plot, i.e. the distribution of the values of the treatment effect and dropout variable in the chain. The marginal density plot in Figure 5 indicates that the higher posterior densities of treatment effect lie between the differences of treatment effects from new therapy and conventional therapy, of -6 and -4. However, in Figure 6, the plot indicates that the higher posterior densities of the dropout variable lie between visits 3 and 4.



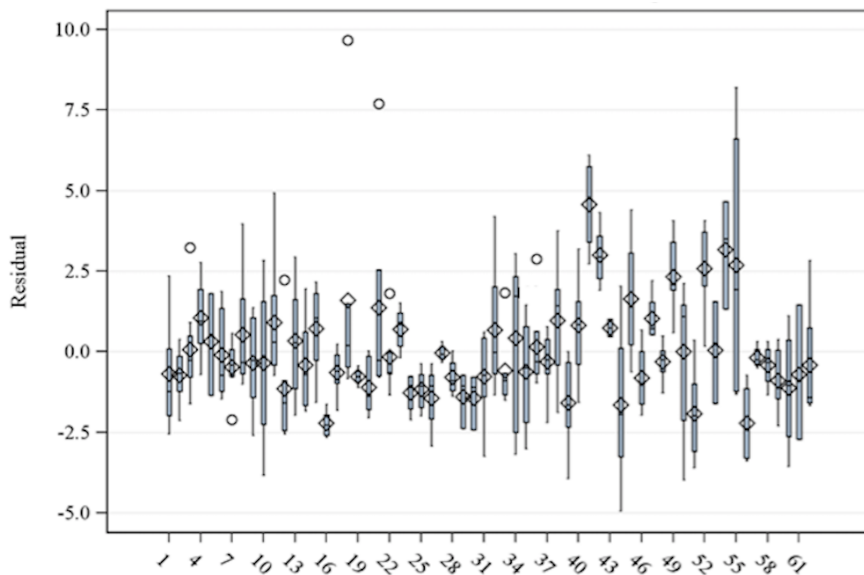
**Figure 1.** Line plot for mean with SE of MELD scores obtained from liver cirrhosis patients' data, for various approaches.

## 5. Discussion

The limitation of missing data analysis is that the true model and mechanism for measurements and missingness are usually unverifiable (Graham, 2012). Thus, in many settings, the selection model should be viewed as

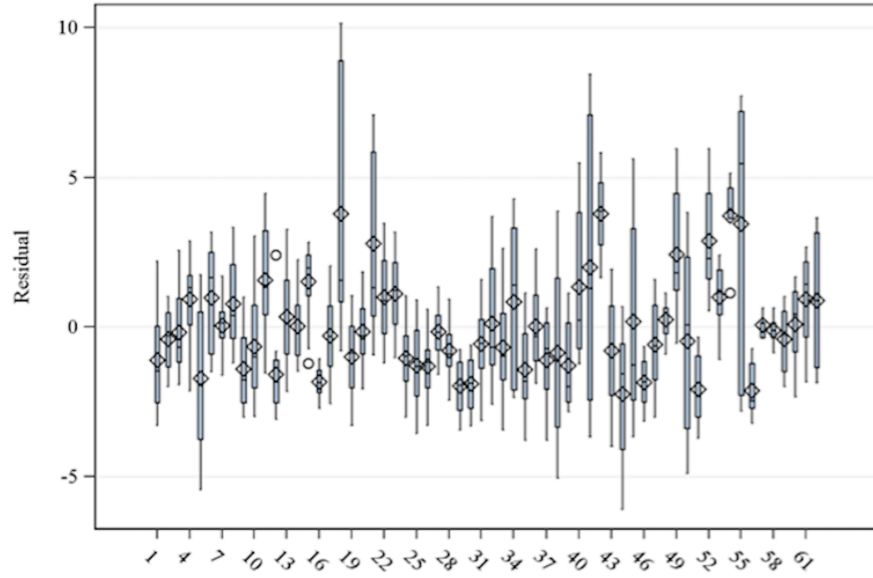


**Figure 2.** Distribution of residuals for change from baseline for MELD score when liver cirrhosis data are considered for all cases.



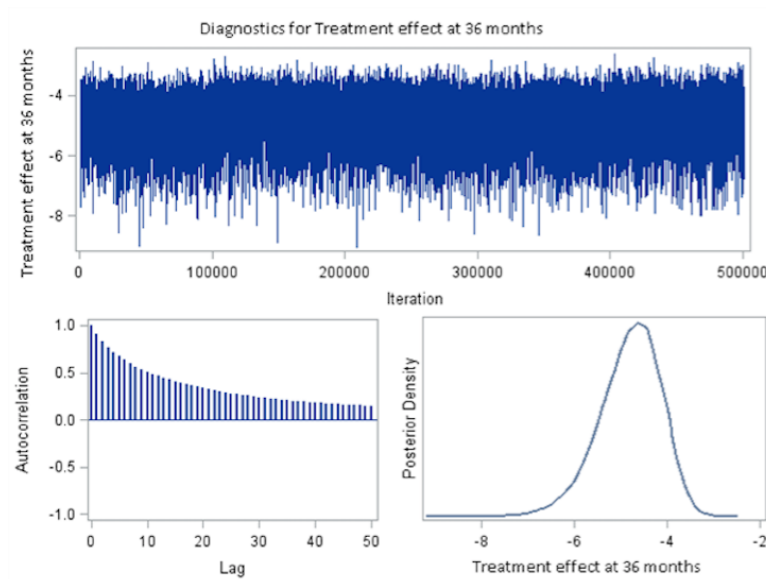
**Figure 3.** Distribution of residuals for change from baseline for MELD score when liver cirrhosis data are considered under the complete case approach.





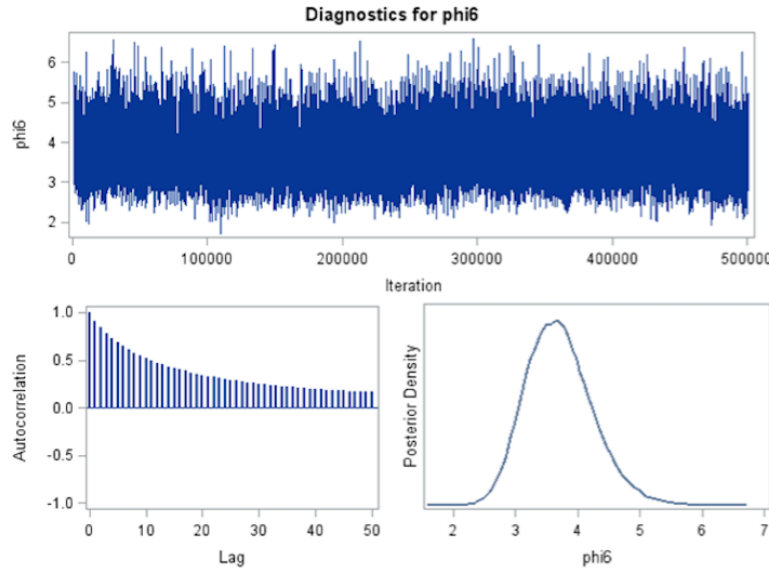
**Figure 4.** Distribution of residuals for change from baseline for MELD score when liver cirrhosis data are considered under the MCAR approach

a model with rich assumptions. The guideline is to always to investigate the sensitivity of the inferences on fixed parameters to varying assumptions. However, extending the use of a specific model beyond its assumptions might not be supported by the practical data. This is especially true in phase I, II, and III clinical trials, where sample sizes are usually not large enough to support overfitting of the model. Selection models are generalised versions of standard longitudinal models (marginal models using GEE, mixed-effects models, and transition models). For example, the mixed-effects model ignoring missing values is nothing but a selection model with ignorable missingness mechanisms (Daniels and Hogan, 2008). We applied selection models to liver cirrhosis patient data to analyse the effectiveness of treatment, comparing the surgery of liver cirrhosis patients receiving HFLPC (Human Fatal Liver Progenitor Cells) infusion against a group receiving surgery alone. In this study, we have illustrated an application to analysing incomplete follow up data, where the response variable is missing throughout visits. We gave attention to the situation in which responses are continuous. The model considered was the selection model. The study focused on specific cases of



**Figure 5.** Plots for diagnostics of the treatment effect from liver cirrhosis patients' data.

selection model; that is, a Diggle and Kenward (1994) model. In the context of the selection model, we used logistic regression for modelling dropout. However, a number of various probabilities can be used, for instance, using the survival analysis approach, the span of duration of treatment or placebo before dropout can also be modelled. However, in our study, the survival model for dropout cannot be used because the time to event is not exactly decided by design. For example, if any patient is not seen at month 12, the exact time to dropout could hypothetically be any time between month 6 and 12. The aim was to investigate the possible influence of dropout on the response measurement on the liver cirrhosis data and also to deal with incomplete sequences. The selection model implied that the dropout mechanisms were not completely at random. In other words, in the context of the implicit model, there was much indication of the prevalence of an MAR rather than an MCAR dropout process. However, many authors (Diggle and Kenward 1994; Molenbergh and Verbeke 2005; Williamson 2006) have stated that caution is necessary when drawing such a conclusion only from the data under analysis. A problem arises when dealing with dropout that is MNAR. Given this difficulty, in a longitudinal study, it is important to



**Figure 6.** Plots for diagnostics of the dropout variables from liver cirrhosis patients' data.

understand that this assumption gives rise to dropout that is not likely to be known in the application setting. Therefore, the different proposed application methods to address dropout that are MNAR cannot easily be verified. For example, one often does not know if the dropout process is precisely captured by a particular method used. Molenbergh and Kenward (2007) suggested that one should apply several approaches to the same data problem. According to Xu and Blozi (2011), if parameter estimates are comparable under different methods, this can indicate that the dropout process may be ignored. However, if different methods give different estimates of the parameters of the longitudinal model, this can indicate that the dropout process can be considered as an important element for the description of the data in the analysis.

## 6. Conclusion

The measurement process and dropout process are often unverifiable. Support has been found for the recommendation that in many settings, multiple strategies or models such as selection and other models, e.g. pattern-mixture models, should be applied to the same data set in order to investigate the

impact of the assumption on dropout or missingness. The most notable limitation with practical data analysis with missing values is that the true model and mechanism for measurements and missingness are usually unverifiable. Thus, in many settings, selection model should be viewed as models with rich assumptions. In the future, we also need to study the models for categorical data, which are also very common in practice. We should extend our methods to be able to also handle missingness in covariates. We should also try to handle functional analysis on the response variables (e.g. power spectra density responses) which is part of multivariate response longitudinal data analysis. Finally, in addition to parametric models, we shall also study non-parametric and semi-parametric models for incomplete longitudinal data analysis. In our study the selection model implied that the dropout mechanisms were not completely at random. In other words, in the context of the implicit model, there was much indication of the prevalence of an MAR rather than an MCAR dropout process.

### Declaration of interest

The authors declare that there is no conflict of interest that could be perceived as prejudicing the impartiality of the research reported.

### Acknowledgements

The authors thank the two anonymous referees for their cautious reading and constructive suggestions which have led to improvement on earlier versions of the manuscript.

### REFERENCES

- Daniels M.J., Hogan J.W. (2008): *Missing data in longitudinal studies: Strategies for Bayesian modeling and sensitivity analysis*. CRC Press.
- Diggle P., Kenward M.G. (1994): Informative drop-out in longitudinal data analysis. *Applied statistics* pages 49–93.
- Dragset I.G. (2009): Analysis of longitudinal data with missing values. Project Thesis. Norwegian University of Science and Technology.
- Fitzmaurice G., Davidian M., Verbeke G., Molenberghs G. (2008): *Longitudinal data analysis*. CRC Press.
- Graham J.W. (2012): *Missing data: Analysis and design*. Springer Science & Business Media.

- Heckman J.J. (1976): The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. In *Annals of Economic and Social Measurement* volume 5 pages 475–492. NBER.
- Kaciroti N.A., Raghunathan T. (2014): Bayesian sensitivity analysis of incomplete data: bridging pattern-mixture and selection models. *Statistics in medicine* 33(27): 4841–4857.
- Kaciroti N.A., Raghunathan T.E., Anthony Schork M., Clark N.M. (2008): A bayesian model for longitudinal count data with non-ignorable dropout. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 57(5): 521–534.
- Kaciroti N.A., Raghunathan T.E., Taylor J.M., Julius S. (2012): A bayesian model for time-to-event data with informative censoring. *Biostatistics* 13(2): 341–354.
- Kaciroti N.A., Schork M.A., Raghunathan T., Julius S. (2009): A bayesian sensitivity model for intention-to-treat analysis on binary outcomes with dropouts. *Statistics in medicine* 28(4): 572–585.
- Kamath P.S., Kim W. (2007): The model for end-stage liver disease (meld). *Hepatology* 45(3): 797–805.
- Kim J.K., Yu C.L. (2012): A semiparametric estimation of mean functionals with nonignorable missing data. *Journal of the American Statistical Association*.
- Laird N.M., Ware J.H. (1982): Random-effects models for longitudinal data. *Biometrics* 38(4): 963–974.
- Nath D.C., Bhattacharje A. (2012): Pattern mixture modeling: An application in anti diabetes drug therapy on serum creatinine in type 2 diabetes patients. *Asian Journal of Mathematics & Statistics* 5(3): 71.
- Satty A., Mwambi H. (2013): Selection and pattern mixture models for modelling longitudinal data with dropout: An application study. *SORT-Statistics and Operations Research Transactions* 1(2): 131–152.