

A Latent Variable Transition Model for Binary Longitudinal Data with Informative Dropout

M. Ganjali* and Z. Rezaee Ghahrodee

**Department of Statistics, Faculty of Mathematical Sciences, Shahid Beheshti university, Evin, Tehran 19838, Iran*

(received: 18/5/2003 ; accepted: 4/11/2003)

Abstract

A latent variable Markovian model is proposed for longitudinal binary responses with dropout. In this model responses can easily be modeled using probit, logit or any other link. Dropout model is parameterized in such a way that parameters can be dynamically changed on time. Some residuals are also presented. These residuals can be used in the presence of informative dropout. The model is also used in an application where the existence of side effects of using fluvoxamine (a treatment for deregulation of serotonin in the brain) is the response of interest.

Keywords: *Longitudinal data with binary responses; Informative dropout; Transition model; Pearson residuals; Non-stationary process.*

1. Introduction

In a panel or longitudinal study each subject is measured at several occasions. So, in these studies we have the vector of $Y_i = (Y_{i1}, \dots, Y_{iT})$ which its elements are responses for the i th subject ($i=1, \dots, n$) on occasions $t=1, \dots, T$. In these studies two aspects of analysis are important (i) the effect of some covariates, which may change on time, on the responses and (ii) taking into account the correlations between the responses of the same subject.

Often some of the subjects withdraw from the survey before the study is complete and do not return. Subjects that do not stay in the study are said to have dropped out. In dropout, variables of the study can be arranged so that Y_{ij+1}, \dots, Y_{iT} are missing when Y_{ij} is missing, for $j=1 \dots T-1$. For example in a clinical trial, some subjects dropout of the

study for different reasons, for examples side effects of drugs or curing of disease. If this dropout is related in some way to their current response then it is inappropriate to restrict the analysis to the complete sequences or to just the observed components of the sequences. In order to establish whether dropout is informative in this way an attempt must be made to jointly model the dropout and response processes.

Little and Rubin (2002) and Diggle and Kenward (1994, hereafter DK) make important distinctions between the various types of dropout mechanism. DK (1994) defined a dropout process as completely random (CRD) if the dropout mechanism is dependent neither on the current value of the response (Y_t) nor on the previous value of the response (Y_{t-1}), and as random dropout (RD) if it is dependent on Y_{t-1} , but not on Y_t . Dropout is defined as informative (ID) if it is dependent on the current value of response (see also Little, 1995).

In this paper a transition model is presented by using the concept of latent variable for longitudinal binary response with dropout. The non-response model is similar in nature to the model of the DK (1994) and the Molenberghs et al. (1997) model in which the dropout probabilities are functions of the previous and current responses. We extend this approach by allowing the coefficients of parameters in dropout probabilities to change on time. If the data contain enough information to enable us to estimate a model with these additional features, then we will be able to provide more insight into the dropout mechanism.

In the next Section the 4 waves of the Fluvoxamine data (Molenberghs et al., 1997) as motivation are discussed. In Section 3 the transition model is presented. In this section we also give the likelihood, noting that the joint likelihood of the proposed model does not have a closed form. In section 4 we discuss how to look at the residuals to find any inconsistency between data and the model. In section 5, we illustrate our extended model on the 3 waves of the Fluvoxamine data (Molenberghs et al. 1997) where we find a form of non-stationary uninformative dropout. In Section 6 we give a brief conclusion.

2. Empirical illustration: Fluvoxamine data

The example is a 4-wave study of the side effects of using Fluvoxamine (a psychiatric drug). In the original data (discussed in Molenberghs and Lesaffre, 1994, Lesaffre et al., 1996, Molenberghs et al., 1997, and Michiels and Molenberghs, 1997) severity of side effects is an ordinal response with: (0) no side effect, (1) no significant interference with functionality of patient, (2) significant interference with functionality of patient and (3) side-effect surpasses therapeutic effect. We shall focus on a dichotomized version (present/absent) of side effects at four periods. A total of 315 subjects were initially recruited into the study. The extent of the side effects was obtained at weeks 2, 4, 8 and 12 after starting the trial. Also obtained were the sex, age, initial severity (scale 1 to 7), and duration of actual mental disease for each subject.

The data were previously analyzed by Lesaffre et al. (1996) and Molenberghs et al. (1997), who used the data for weeks 2, 4 and 12. Molenberghs et al. (1997) used a marginal model for the responses and a multivariate Dale distribution to take into account the correlation between the responses of the same individual. Their analysis finds that non-response is dependent on the previous, but not on the current value of the response, i.e. dropout is at random (RD).

We use the data for all four visits, but, as we use a transition model, to avoid initial condition problem (Heckman, 1981) we fix the initial response at week 2 (Y_0) and use it as an explanatory variable. The transition model can easily be used to an arbitrary number of periods. Table 1 shows the different patterns of missing responses (responses at weeks 4, 8 and 12) and their frequencies for the 259 individuals without any missing covariate data at weeks 2, 4, 8 and 12.

Table 1. Different patterns of missing data for Fluvoxamine data ('O'=Observed, 'M'=Missed)

Pattern	Frequency
OOO	216
OOM	17
OMM	26

3. The general transition latent variable model with dropout

Let $Y_i = (Y_{i1}, \dots, Y_{iT})$ be the vector of binary responses for the i th individual ($i=1, \dots, n$). The joint general transition (Markovian latent variable) model of order q (the number of previous responses in the model for response at time t) for the responses and the dropout latent variable model where the first response is observed for all individuals are:

$$y_{it}^* = \beta_t' X_{it} + \sum_{j=1}^q \alpha_{t,t-j} y_{it-j} + \varepsilon_{it} \quad \text{for } t = 1, \dots, T \quad (1a)$$

$$R_{it}^* = \xi_t' W_{it} + \gamma_{1t} y_{it-j} + \gamma_{2t} y_{it} + \varepsilon_{iT+t-1} \quad \text{for } t = 2, \dots, T \quad (1b)$$

where ε_t for $t=1 \dots 2T-1$ are uncorrelated errors, which are distributed with means 0. In the following, we shall call the equations (1a) and (1b) as system 1. When $t=1$ we assume that y_{i0} is a predetermined observed value. If we assume normal distributions for these errors we have probit models for response and dropout mechanisms and if we assume logistic distributions for errors we have logit models. So sensitivity of the model parameters to different distributions for errors can easily be obtained.

The vectors X_{it} and $W_{it'}$ for $t=1, 2, \dots, T$ and $t'=2, \dots, T$ are explanatory variables and the vectors of parameters are β and ξ which include a constant value. α 's are the parameters for feedback effects. We suppose that response t is missing i.e. $R_t = 0$ if $R_{it}^* < 0$ for $t=2, 3, \dots, T$ and that if $R_t = 0$ then $R_{t'} = 0$ for $t'=t+1, \dots, T$. The binary responses y_{it} , for $t=1, 2, \dots, T$ are given by

$$y_{it} = \begin{cases} 1 & \text{if } y_{it}^* < 0 \\ 0 & \text{if } y_{it}^* \geq 0 \end{cases}$$

For simplicity dropout models include the effect of the current and previous values of the responses (y_t and y_{t-1}) rather than the complete history and the current values of responses (y_1, \dots, y_t).

The response indicators are

$$R_{it} = \begin{cases} 1 & \text{if } R_{it}^* \geq 0 \\ 0 & \text{Otherwise.} \end{cases}$$

In this model we are able to identify informative dropout process. Informative dropout (ID) occurs when one of γ_{2t} for $t=2\dots T$ is non-zero. We have random dropout (RD) if all of γ_{2t} for $t=2\dots T$ are zero and at least one of γ_{1t} for different values of t is non-zero. We have completely random dropout (CRD) if all of γ_{1t} and γ_{2t} for $t=2\dots T$ are zero.

The dropout model in system (1) similar to the DK (1994) model depends on previous and current responses, but it extends the DK (1994) and Molenberghs et al. (1997) models as this model let the dropout models have different parameters (γ 's are depend on time). So, regressor parameterization is extended to improve the model's flexibility. This makes the dropout processes to be non-stationary.

Model of DK and system (1) can be compared for sequences of length 3. The first order transition model which has the same dropout model as DK and Molenberghs et al. (1997) is:

$$y^*_{it} = \beta'_t X_{it} + \alpha y_{it-1} + \varepsilon_{it}, \quad \text{for } t = 1, 2, 3, \tag{2a, 2b, 2c}$$

$$R^*_{i2} = \xi'_2 W_{i2} + \gamma_1 y_{i1} + \gamma_2 y_{i2} + \varepsilon_{i4}, \tag{2d}$$

$$R^*_{i3} = \xi'_3 W_{i3} + \gamma_1 y_{i2} + \gamma_2 y_{i3} + \varepsilon_{i5}, \tag{2e}$$

where ε_j for $j=1, \dots, 5$ are uncorrelated errors, which are distributed with means 0.

The extended model is

$$y^*_{it} = \beta'_t X_{it} + \alpha_t y_{it-1} + \varepsilon_{it}, \quad \text{for } t = 1, 2, 3, \tag{3a, 3b, 3c}$$

$$R^*_{i2} = \xi'_2 W_{i2} + \gamma_{12} y_{i1} + \gamma_{22} y_{i2} + \varepsilon_{i4}, \tag{3d}$$

$$R^*_{i3} = \xi'_3 W_{i3} + \gamma_{13} y_{i2} + \gamma_{23} y_{i3} + \varepsilon_{i5}, \tag{3e}$$

where the dropout models have different parameters and the feedback effect let to be changed with time. Note also that responses in DK (1994) model are continuous (y^*_{it} are continuous observed variables) and DK used a multivariate normal distribution for the responses. So, they used a marginal model for responses not a transition model (see, Diggle et al, 1994).

3.1 The likelihood function

If we just observe the first response, the likelihood for the general transition model of order 1 is:

$$L_i = f(y_{i1}, R_{i2} = 0 | y_{i0}) \\ = \sum_{y_{i2}=0}^1 f(y_{i2} | y_{i1}) f(y_{i1} | y_{i0}) p(R_{i2} = 0 | y_{i1}, y_{i2}), \quad (4)$$

where for simplicity we suppress dependence on covariates in this equation. If responses of the i th individual are observed until time T_i ($T_i < T$) and individual dropout at this time, the likelihood is:

$$L_i = f(y_{i1}, y_{i2}, \dots, y_{iT_i-1} | y_{i0}) p(R_{i2} = 1, \dots, R_{iT_i-1} = 1, R_{iT_i} = 0 | y_{i0}, y_{i1}, \dots, y_{iT_i-1}) \\ = \sum_{y_{iT_i}=0}^1 \{ [\prod_{t=1}^{T_i} f(y_{it} | y_{it-1})] [\prod_{j=2}^{T_i-1} p(R_{ij} = 1 | y_{ij-1}, y_{ij})] p(R_{iT_i} = 0 | y_{iT_i-1}, y_{iT_i}) \}, \quad (5)$$

and the likelihood for an individual with complete responses is:

$$L_i = f(y_{i1}, y_{i2}, \dots, y_{iT} | y_{i0}) p(R_{i2} = 1, \dots, R_{iT-1} = 1, R_{iT} = 1 | y_{i0}, y_{i1}, \dots, y_{iT}) \\ = [\prod_{t=1}^T f(y_{it} | y_{it-1})] [\prod_{j=2}^T p(R_{ij} = 1 | y_{ij-1}, y_{ij})]. \quad (6)$$

The overall likelihood is the product of the individuals likelihood, i.e.,

$L = \prod_{i=1}^n L_i$. This likelihood can easily be extended and written for any value of q .

This likelihood does not have a closed form for parameter estimates and should be solved by a numerical algorithm. We use NAG (1996) routine E04UCF to find the parameter estimates. E04UCF is a FORTRAN routine to minimize a smooth function (minus logarithm of likelihood in our work) subject to constraints using a sequential quadratic programming (SQP, Fletcher, 2000) method. In this routine all unspecified derivatives are approximated by finite differences.

4. Residuals

For the transition model of order 1 the Pearson residuals may be found by:

$$r_{it} = \frac{y_{it} - E(Y_{it} | y_{i,t-1}, X_{it})}{\sqrt{\text{Var}(Y_{it} | y_{i,t-1}, X_{it})}}. \quad (7)$$

This can not be evaluated for missing responses and does not use the information coming from dropout mechanism. These residuals (r_{it} in equation 7) can be modified to consider the dropout mechanism. For this, predicted value at time t should be evaluated given the condition that response at time t is observed ($R_{it} = 1$). So residuals could be found by

$$\begin{aligned}
 r_{it} &= \frac{y_{it} - E(Y_{it} | y_{i,t-1}, R_{it} = 1, X_{it})}{\sqrt{Var(Y_{it} | y_{i,t-1}, R_{it} = 1, X_{it})}} \\
 &= \frac{y_{it} - \Pi_{it}}{[\Pi_{it}(1 - \Pi_{it})]^{1/2}}
 \end{aligned} \tag{8}$$

where

$$\begin{aligned}
 \Pi_{it} &= pr(Y_{it} = 1 | R_{it} = 1, y_{i,t-1}, X_{it}) \\
 &= \frac{pr(Y_{it} = 1, R_{it} = 1 | y_{i,t-1}, X_{it})}{pr(R_{it} = 1 | y_{i,t-1}, X_{it})} \\
 &= \frac{pr(R_{it} = 1 | y_{it} = 1, y_{i,t-1}, X_{it})pr(Y_{it} = 1 | y_{i,t-1}, X_{it})}{\sum_{j=0}^1 pr(R_{it} = 1 | y_{it} = j, y_{i,t-1}, X_{it})pr(Y_{it} = j | y_{i,t-1}, X_{it})}.
 \end{aligned}$$

The estimated Pearson residuals (\hat{r}_{it}) can be found by using the maximum likelihood estimates of the parameters in system (1).

5. Model and results for Fluvoxamine data

In this section, we give some model fits to the Fluvoxamine data of section 2. Let $Y_i = (Y_{i1}, Y_{i2}, Y_{i3})$ be the vector of binary responses for the i th individual ($i=1\dots 259$) on weeks 4, 8 and 12 and Y_{i0} be the response on week 2. Also, let $Y^*_i = (Y^*_{i1}, Y^*_{i2}, Y^*_{i3})$ be the corresponding vector of latent variables for Y_i . At first we use following model for Fluvoxamine data (Model I):

$$y_{i1}^* = \beta_1' X_{i1} + \alpha_{10} y_{i0} + \varepsilon_{i1}, \quad (9a)$$

$$y_{i2}^* = \beta_2' X_{i2} + \alpha_{20} y_{i0} + \alpha_{21} y_{i1} + \varepsilon_{i2}, \quad (9b)$$

$$y_{i3}^* = \beta_3' X_{i3} + \alpha_{30} y_{i0} + \alpha_{31} y_{i1} + \alpha_{32} y_{i2} + \varepsilon_{i3}, \quad (9c)$$

$$R_{i2}^* = \xi_2' W_{i2} + \gamma_1 y_{i1} + \gamma_2 y_{i2} + \varepsilon_{i4}, \quad (9d)$$

$$R_{i3}^* = \xi_2' W_{i3} + \gamma_1 y_{i2} + \gamma_2 y_{i3} + \varepsilon_{i5}, \quad (9e)$$

The vector of X on time t includes 1, sex, age, initial severity (SEVE) and logarithm of duration (lnDUR) and the vector of W include 1 and logarithm of duration. We use probit link for all responses and dropout models. We also fit RD model ($\gamma_2 = 0$) and CRD model ($\gamma_1 = 0, \gamma_2 = 0$). These results are given in table 2.

The values of minus logarithm of likelihood (-logL) in the table 2 show that dropout is at random (RD). Dropout depends on the previous response, but not on the current response. Table 2 also shows that some of the explanatory variables have no significant effect on the responses.

With considering Model I as full model and using a backward approach, we remove some of these variables (see table 3) to find a more parsimonious model. We compare two models M_0 and M_1 (with L_0 and L_1 as likelihoods evaluated at the ML estimates for each model) with M_0 a special case of M_1 , using:

$$G^2 = -2(\log L_0 - \log L_1),$$

which has an approximate chi-square null distribution (i. e. under the assumption that model M_0 is the true model) with d.f. equal to the difference between number of parameters in two models. For example in table 3, (a) compares a model with assumption of no sex effects on all responses (a = M_0) with model I (M_1) and (b) compares a model with assumption of no sex and severity effects on all responses (b = M_0) with a model with assumption of no sex effects on all responses (now a = M_1) and so on.

Table 2: Results for the responses and dropout models

Par	Model I		RD Model		CRD Model	
	Est.	Se.	Est.	Se.	Est.	Se.
β_{01}	-1.474	0.716	-1.474	0.700	-1.474	0.697
$\beta_{11}(AGE)$	0.007	0.007	0.007	0.007	0.007	0.007
$\beta_{21}(SEX)$	0.161	0.194	0.161	0.194	0.161	0.194
$\beta_{31}(\ln DUR)$	-0.034	0.076	-0.034	0.076	-0.034	0.076
$\beta_{41}(SEVE)$	0.003	0.122	0.003	0.120	0.003	0.120
$\alpha_{10}(Y_0)$	1.826	0.193	1.826	0.193	1.826	0.193
β_{02}	-1.220	0.812	-1.181	0.800	-1.181	0.814
$\beta_{12}(AGE)$	0.016	0.008	0.015	0.008	0.015	0.008
$\beta_{22}(SEX)$	0.145	0.283	0.222	0.227	0.222	0.227
$\beta_{32}(\ln DUR)$	0.096	0.111	0.125	0.090	0.125	0.090
$\beta_{42}(SEVE)$	-0.222	0.151	-0.238	0.140	-0.238	0.143
$\alpha_{20}(Y_0)$	0.283	0.263	0.302	0.263	0.302	0.264
$\alpha_{21}(Y_1)$	1.715	0.301	1.791	0.254	1.791	0.255
β_{03}	-1.985	1.058	-2.096	1.064	-2.096	1.068
$\beta_{13}(AGE)$	-0.003	0.011	0.000	0.010	0.000	0.010
$\beta_{23}(SEX)$	0.173	0.272	0.175	0.273	0.175	0.273
$\beta_{33}(\ln DUR)$	0.201	0.119	0.212	0.114	0.212	0.114
$\beta_{43}(SEVE)$	-0.001	0.175	0.002	0.177	0.002	0.178
$\alpha_{30}(Y_0)$	0.277	0.321	0.229	0.309	0.229	0.309
$\alpha_{31}(Y_1)$	0.388	0.332	0.426	0.334	0.426	0.334
$\alpha_{32}(Y_2)$	2.279	0.331	2.310	0.309	2.310	0.309
ξ_0	1.846	0.172	1.839	0.165	1.669	0.144
$\xi_1(\ln DUR)$	-0.186	0.075	-0.166	0.064	-0.177	0.064
γ_1	-0.746	0.607	-0.391	0.165	-	-
γ_2	0.549	1.006	-	-	-	-
-logL	418.893		419.040		421.877	

Which has an approximate chi-square null distribution (i. e. under the assumption that model M_0 is the true model) with d.f. equal to the difference between number of parameters in two models. For example in table 3, (a) compares a model with assumption of no sex effects on all responses ($a = M_0$) with model I (M_1) and (b) compares a model with assumption of no sex and severity effects on all responses ($b = M_0$) with a model with assumption of no sex effects on all responses (now $a = M_1$) and so on.

Table 3: A backward selection approach to obtain a parsimonious model

Removed par. (var.)	-2logL	G^2	d. f.	P-value
(a)Sex from Model I	839.154	2.368	3	0.500
(b)Seve from (a)	841.143	1.989	3	0.574
(c)Second and Higher order effects from (b)	846.050	4.907	3	0.179
(d)Age at Periods 1 and 3 from (c)	846.113	0.063	2	0.969
(e)lnDur at period 1 and 2 from (d)	847.761	1.648	1	0.199

Table 3 shows that age, sex, lnDUR and SEVE have no significant effect on Y_1 , sex, lnDUR, SEVE and Y_0 have no significant effect on Y_2 , and Y_0 , age, sex, SEVE and Y_1 also have no significant effect on Y_3 .

To see whether dropout model parameters have different effect on time, we let the parameters in dropout models change on time (see model in system 3, extended model). Results (with removing non significant covariates) are given in table 4. We find a significant time varying effect of previous responses on dropout models. Previous outcome shows no effect on dropout model in period three, but it is strongly significant on the dropout model in period two. We also find a time varying duration of disease effect on dropout models. Duration of disease has a negative effect on the dropout model in period 2, but it has no significant effect on the dropout model in period 3.

Table 4: Results for the extended model

Par .	Extended model				
	Est.	Se.	Par.	Est.	Se.
β_{01}	1.166	0.151	ξ_{02}	2.166	0.291
$\alpha_{10}(Y_0)$	1.819	0.189	$\xi_{12}(\ln DUR)$	0.331	0.107
β_{02}	2.066	0.345	$\gamma_{12}(Y_1)$	1.156	0.325
$\beta_{12}(AGE)$	0.015	0.007	$\gamma_{22}(Y_2)$	1.463	1.003
$\alpha_{21}(Y_1)$	1.717	0.222	ξ_{03}	1.585	0.229
β_{03}	1.820	0.269	ξ_{13}	0.016	0.111
$\beta_{43}(\ln DUR)$	0.216	0.109	$\gamma_{13}(Y_2)$	0.131	0.994
$\alpha_{32}(Y_2)$	2.652	0.253	$\gamma_{23}(Y_3)$	0.131	1.200
$-\log L$	421.660				

Results in table 4 also show a first order effect of responses. Individuals who have side effect on previous response are more likely to have side effect on current response. There is also age effect on response in period two. The older the individual is the more likely is her/him to have the side effects on period two. Response in period 1 have negative effect on dropout on time two which means individuals who have side effects on time 1 are more likely to dropout on time 2. Using dropout model on time 3, we find that the probability of dropout on time 3 (given observing the response on time 2) is independent of, Y_2 , Y_3 and $\ln DUR$ and it is 0.073.

Residuals for the first response do not show any outliers. Residuals for the second response (using equation 7) show 4 individuals as outliers and residuals for the third response show 7 individual as outliers. To check for sensitivity of the results to different distributions for the errors in dropout and response models, we use logit link for all dropout and response models. However, our main interpretations are the same as what we find using probit link.

6. Conclusion

We use a transition model for longitudinal binary response with informative dropout. The model is so flexible to be used with different distributions for measurement errors in the model. In this model

dropout model parameters, in contrast with the DK model, let change on time. For Fluvoxamine data (Molenberghs et al. 1997) we find a form of non-stationary uninformative dropout. We obtain that dropout model on time 2 is dependent on previous outcome, but, on time 3, it is not dependent on this variable. In these data, response on time t is strongly dependent on the response on time $t-1$, but not on the other responses (a first order model is sufficient). Some residuals are also presented which, in the case of informative dropout, can provide a better performance than the usual Pearson residuals. For further work the model can be extended to be used for longitudinal ordinal response with dropout.

Acknowledgement

We are grateful to the referees for improving comments.

References

- Diggle, P.J. and Kenward, M. G. (1994) *Informative Drop-out in longitudinal data analysis*. Journal of Applied Statistics **43**, 49-93.
- Diggle, P.J., Liang, K.Y., Zeger, S.L. (1994) *Analysis of longitudinal data*. London: Chapman and Hall.
- Fletcher, R. (2000) *Practical Methods of Optimization*. John Wiley.
- Heckman, J. (1981) *Statistical models for discrete panel data*, in Manski, C. & McFadden D., *Structural Analysis of discrete data with econometric applications*, 114-195, Cambridge, Mass: MIT press.
- Lesaffre, E., Molenberghs, G., and Dewulf, L. (1996) *Effect of dropouts in a longitudinal study: An application of a repeated ordinal model*, Statistics in Medicine, **15**, 1123-1141.
- Little, R.J.A. (1995) *Modeling the drop-out mechanism in repeated-measures studies*. Journal of the American Statistical Association, **90**, 1112-1121.
- Little, R.J. A, Rubin, D.B. (2002) *Statistical Analysis With Missing Data*. New York. John Wiley.
- Michiels B., Molenberghs G. (1997) *Protective estimation of longitudinal categorical data with nonrandom dropout*. Communications in statistics: Theory and methods, 26, p. 65-94.

-
- Molenberghs, G., Kenward, M.G. and Lesaffre, E. (1997) *The analysis of longitudinal ordinal data with nonrandom drop-out*. *Biometrika*, **84**, 33-44.
- Molenberghs G., Lesaffre, E. (1994) *Marginal Modelling of Correlated Ordinal Data Using a Multivariate Plackett Distribution*. *Journal of the American Statistical Association*, **89**, 633-644.
- NAG. (1996). *Numerical Algorithms Group Manual*. Mark 16. Oxford, U.K.