

On statistical Computing via EM Algorithm in Logistic Linear Models Involving Non-ignorable Missing data *

Yuna Jun[†] Guoqi Qian[‡] Jeong-soo Park[§]

Abstract

Many data sets obtained from surveys or medical trials often include missing observations. When these data sets are analyzed, it is general to use only complete cases. However, it is possible to have big biases or involve inefficiency. In this paper, we consider a method for estimating parameters in logistic linear models involving non-ignorable missing data mechanism. A binomial response and normal exploratory model for the missing data are used. We fit the model using the EM algorithm. The E-step is derived by Metropolis-hastings algorithm to generate a sample for missing data and Monte-carlo technique, and the M-step is by Newton-Raphson to maximize likelihood function. Asymptotic variances of the MLE's are derived and the standard error and estimates of parameters are compared.

Keywords: EM algorithm, logistic linear models, missing data mechanism, maximum likelihood estimation, non-ignorable missing data.

1 Introduction

Most of data sets which are obtained have missing data among observed data. When these data sets are analyzed, it is general to use only complete cases with data all observed after removing missing data. However, there are some problems, if the missing data is related to values of the missing variables. It is possible to have big biases or involve inefficiency. Therefore, we use robust statistical methods to consider these problems, Little and Rubin(1987) proposes the missing data mechanism and Baker and Laird(1988) used the EM algorithm to obtain maximum likelihood estimates.

In this paper, we have a method for estimating parameters in logistic linear models involving non-ignorable missing data mechanism. We use a binomial response model and normal model for the missing data. We fit the model using the EM algorithm. The E-step is derived by Metropolis-hastings algorithm to generate a sample for missing data and Monte-carlo technique, and the M-step is by Newton-Raphson to maximize likelihood function. We estimate the logistic regression coefficients when the response and a covariate are missing using the EM algorithm.

The rest of this paper is organized as follows. In section 2 we state notation and model. In section 3 we derive the E and M steps of the EM algorithm. In section 4 we describe the method how to estimate parameters. In section 5 we illustrate our results with one example test.

*This research was supported by KOSEF(F01-2004-000-10351-0)

[†]Department of Statistics, Chonnam National University, Gwangju, Korea email: umumu@hanmail.net

[‡]Department of Statistics, La Trobe University, Melbourne, Australia; email: g.qian@latrobe.edu.au

[§]Department of Statistics, Chonnam National University, Gwangju, Korea email: jspark@chonnam.ac.kr

2 Notation and Model

Suppose that y_1, \dots, y_n are independent observations, where each y_i has a binomial distribution with sample size m_i and success probability p_i . Let $X_i = (x_{1i}, x_{2i})^t$ is a 2×1 random vector of covariates. x_{1i} and x_{2i} are independent observations and each covariates have normal distribution with mean μ_1, μ_2 and variance σ_1^2, σ_2^2 . Further, let $\beta^t = (\beta_0, \beta_1, \beta_2)$ is regression coefficients assuming to include intercept coefficient.

$$\text{logit}(\pi_i) = X_i^t \beta$$

$$p(y_i | X_i, \beta) = \frac{\exp\{y_i X_i^t \beta\}}{1 + \exp\{X_i^t \beta\}}. \quad (1)$$

We assume that x_{1i} is completely observed and y_i and x_{2i} partially missing.

The missing data mechanism is defined as $r_i = 0$ if y_i is observed and $s_i = 0$ if x_{2i} is observed so now we assume $p(r_i) = \psi$, $p(s_i) = \phi$.

$$\begin{aligned} \text{logit}(\psi_i) &= X_i^t \delta + y_i \omega \\ \text{logit}(\phi_i) &= X_i^t \alpha + y_i \tau, \quad i = 1, 2, \dots, n. \end{aligned}$$

The conditional probability for r_i and s_i is derived by equations (2) and (3).

$$p(r_i | X_i, y_i, \delta, \omega) = \frac{\exp\{r_i(X_i^t \delta + y_i \omega)\}}{1 + \exp\{X_i^t \delta + y_i \omega\}} \quad (2)$$

$$p(s_i | X_i, y_i, \alpha, \tau) = \frac{\exp\{s_i(X_i^t \alpha + y_i \tau)\}}{1 + \exp\{X_i^t \alpha + y_i \tau\}} \quad (3)$$

We derive the joint probability function as

$$\begin{aligned} p(y_i, x_{2i}, r_i, s_i | x_{1i}) &= p(r_i | y_i, X_i, \delta, \omega) p(s_i | y_i, X_i, \alpha, \tau) p(y_i | X_i, \beta) p(x_{2i} | x_{1i}) \\ &\propto \frac{\exp\{r_i(X_i^t \delta + y_i \omega)\}}{1 + \exp\{X_i^t \delta + y_i \omega\}} \times \frac{\exp\{s_i(X_i^t \alpha + y_i \tau)\}}{1 + \exp\{X_i^t \alpha + y_i \tau\}} \times \exp\{X_i^t \beta y_i\} \\ &\times (1 + \exp\{X_i^t \beta\})^{-m_i} \times (2\pi\sigma_2^2)^{-1/2} \times \exp\left\{-\frac{(x_{2i} - \mu_2)^2}{2\sigma_2^2}\right\}. \quad (4) \end{aligned}$$

Therefore, we can write down the complete-data log-likelihood by

$$\begin{aligned} \log \ell(\theta | y, X_i, r, s) &= \sum_{i=1}^n \log\left(\frac{\exp\{r_i(X_i^t \delta + y_i \omega)\}}{1 + \exp\{X_i^t \delta + y_i \omega\}}\right) + \sum_{i=1}^n \log\left(\frac{\exp\{s_i(X_i^t \alpha + y_i \tau)\}}{1 + \exp\{X_i^t \alpha + y_i \tau\}}\right) \\ &+ \sum_{i=1}^n X_i^t \beta y_i - \sum_{i=1}^n m_i \log(1 + \exp\{X_i^t \beta\}) - \frac{n}{2} \log(2\pi\sigma_2^2) \\ &- \sum_{i=1}^n \frac{(x_{2i} - \mu_2)^2}{2\sigma_2^2}. \quad (5) \end{aligned}$$

where $\theta = (\beta, \delta, \omega, \alpha, \tau, \mu_2, \sigma_2^2)$ is the parameters related to develop EM algorithm.

3 E-step and M-step of the EM algorithm

To compute maximum likelihood estimates in non-ignorable missing data, we use the expected log-likelihood. We consider the expectation in response variable y_i is missing and a covariate x_{2i} is missing and both of them are missing. The expected log-likelihood for E-step can be written by

$$E[l(\theta; X_i, y_i, r_i, s_i)] = \begin{cases} \sum_{y_i=0}^{m_i} l(\theta; X_i, y_i, r_i, s_i) p(y_i | X_i, r_i, s_i) & (\text{if } y_i \text{ has missing components.}) \\ \int l(\theta; X_i, y_i, r_i, s_i) p(x_{2i} | x_{1i}, y_i, r_i, s_i) dx_{2i, \text{mis}} & (\text{if } x_{2i} \text{ has missing components.}) \\ \sum_{y_i=0}^{m_i} \int l(\theta; X_i, y_i, r_i, s_i) p(y_i, x_{2i} | x_{1i}, r_i, s_i) dx_{2i, \text{mis}} & (\text{if } y_i \text{ and } x_{2i} \text{ have missing components.}) \\ l(\theta; X_i, y_i, r_i, s_i) & (\text{if } x_{2i}, y_i \text{ are observed components.}) \end{cases}$$

where $l(\theta; X_i, y_i, r_i, s_i)$ is complete-data log-likelihood $\log p(y_i, X_i, r_i, s_i)$. Equation (5) leads to the E-step for the i th observation as

$$\begin{aligned} Q(\theta, \theta^r) &= \sum_{i=1}^n \sum_{y_i=0}^{m_i} \int w_{i(r)} l(\theta; X_i, y_i, r_i, s_i) dx_{2i, \text{mis}} \\ &= \sum_{i=1}^{n_1} l(\theta; X_i, y_i, r_i, s_i) + \sum_{i=n_1+1}^{n_2} \sum_{y_i=0}^{m_i} l(\theta; X_i, y_i, r_i, s_i) p(y_i | X_i, r_i, s_i, \theta^r) \\ &\quad + \sum_{n_2+1}^{n_3} \int l(\theta; X_i, y_i, r_i, s_i) p(x_{2i} | x_{1i}, y_i, r_i, s_i, \theta^r) dx_{2i, \text{mis}} \\ &\quad + \sum_{n_3+1}^n \sum_{y_i=0}^{m_i} \int l(\theta; X_i, y_i, r_i, s_i) p(y_i, x_{2i} | x_{1i}, r_i, s_i, \theta^r) dx_{2i, \text{mis}}. \end{aligned} \quad (6)$$

Where θ^r is r^{th} iteration estimates, $p(y_{i, \text{mis}} | X_i, y_{i, \text{obs}}, r_i, s_i)$, $p(x_{2i, \text{mis}} | X_{i, \text{obs}}, y_i, r_i, s_i)$ and $p(y_{i, \text{mis}}, x_{2i, \text{mis}} | X_{i, \text{obs}}, y_{i, \text{obs}}, r_i, s_{2i})$ are the conditional probability of the missing data given the observed data and regarded as the weights. The weights have the form as

$$\begin{aligned} & p(y_{i, \text{mis}}, x_{2i, \text{mis}} | X_{i, \text{obs}}, r_i, s_i, \theta^r) \\ &= \frac{p(y_i | X_i, \theta^r) p(x_{2i} | x_{1i}) p(r_i | y_i, X_i, \theta^r) p(s_i | y_i, X_i, \theta^r)}{\sum_{y_i=0}^{m_i} \int p(y_i | X_i, \theta^r) p(x_{2i} | x_{1i}) p(r_i | y_i, X_i, \theta^r) p(s_i | y_i, X_i, \theta^r)} \\ &\propto p(y_i, x_{2i}, r_i, s_i | x_{1i}, \theta^r). \end{aligned} \quad (7)$$

$$\begin{aligned} & p(x_{2i, \text{mis}} | X_{i, \text{obs}}, y_i, r_i, s_i, \theta^r) \\ &= \frac{p(x_{2i} | x_{1i}, \theta^r) p(s_i | y_i, X_i, \theta^r)}{\int p(x_{2i} | x_{1i}, \theta^r) p(s_i | y_i, X_i, \theta^r)} \propto p(x_{2i} | x_{1i}, \theta^r) p(s_i | y_i, X_i, \theta^r) \\ &\propto \frac{\exp\{s_i(X_i^t \alpha + y_i \tau)\}}{1 + \exp\{s_i(X_i^t \alpha + y_i \tau)\}} \times (2\pi\sigma_2^2)^{-1/2} \times \exp\left\{-\frac{(x_{2i} - \mu_2)^2}{2\sigma_2^2}\right\}. \end{aligned} \quad (8)$$

$$\begin{aligned} & p(y_{i, \text{mis}} | X_i, r_i, s_i, \theta^r) \\ &= \frac{p(y_i | X_i, \theta^r) p(r_i | y_i, X_i, \theta^r)}{\sum_{y_i=0}^{m_i} p(y_i | X_i, \theta^r) p(r_i | y_i, X_i, \theta^r)} \propto p(y_i | X_i, \theta^r) p(r_i | y_i, X_i, \theta^r). \end{aligned} \quad (9)$$

To do E-step we need to generate a sample from weights function. For generation of equation (7), (8), (9) we use Metropolis-hastings algorithm. In the M-step, we maximize the log-likelihood and estimate $\theta^{r+1} = (\beta, \delta, \omega, \alpha, \tau, \mu_2, \sigma_2^2)$ converging. The Newton-Raphson algorithm is often used to maximize function for linear model. The equations for parameters $\theta^{r+1} = (\beta, \delta, \omega, \alpha, \tau)$ in the M-step at the $(r+1)^{st}$ EM iteration and the $(t+1)^{st}$ Newton-Raphson iteration take the form.

$$\beta^{r+1} = \beta^r + \left(-\frac{\partial^2 Q(\theta, \theta^r)}{\partial \beta \partial \beta^t}\right)^{-1} \times \left(\frac{\partial Q(\theta, \theta^r)}{\partial \beta}\right) \quad (10)$$

Iterating E-step and M-step, the $(r+1)^{st}$ step estimates of $\theta^{r+1} = (\beta, \delta, \omega, \alpha, \tau)$ can be obtained by these derivation.

$$\begin{aligned} \frac{\partial}{\partial \beta} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t \mathbf{y}_i + \sum_{i=n1+1}^{n2} E(\mathbf{x}_i^t \mathbf{y}_i | \mathbf{x}_{\text{obs}}, \theta^r) + \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t \mathbf{y}_i | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) \\ &+ \sum_{i=n3+1}^{n4} E(\mathbf{x}_i^t \mathbf{y}_i | \mathbf{x}_{\text{obs}}, \theta^r) + \sum_{i=1}^{n1} \mathbf{x}_i^t \pi_i - \sum_{i=n1+1}^{n2} E(\mathbf{x}_i^t \pi_i | \mathbf{x}_{\text{obs}}, \theta^r) \\ &+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t \pi_i | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n E(\mathbf{x}_i^t \pi_i | \mathbf{x}_{\text{obs}}, \theta^r) \\ \frac{\partial^2}{\partial \beta \partial \beta^t} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t \pi_i (\pi_i - 1) \mathbf{x}_i + \sum_{i=n1+1}^{n2} E(\mathbf{x}_i^t \pi_i (\pi_i - 1) \mathbf{x}_i | \mathbf{x}_{\text{obs}}, \theta^r) \\ &+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t \pi_i (\pi_i - 1) \mathbf{x}_i | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n E(\mathbf{x}_i^t \pi_i (\pi_i - 1) \mathbf{x}_i | \mathbf{x}_{\text{obs}}, \theta^r) \end{aligned} \quad (11)$$

where $\pi_i = \exp\{\mathbf{x}_i^t \beta\} / (1 + \exp\{\mathbf{x}_i^t \beta\})$,

$$\begin{aligned} \frac{\partial}{\partial \delta} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t (\mathbf{r}_i - \psi_i) + \sum_{i=n1+1}^{n2} E(\mathbf{x}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \\ &+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n E(\mathbf{x}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \\ \frac{\partial^2}{\partial \delta \partial \delta^t} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) + \sum_{i=n1+1}^{n2} E(\mathbf{x}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \\ &+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n E(\mathbf{x}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \end{aligned} \quad (12)$$

$$\begin{aligned} \frac{\partial}{\partial \omega} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{y}_i^t (\mathbf{r}_i - \psi_i) + \sum_{i=n1+1}^{n2} E(\mathbf{y}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \\ &+ \sum_{i=n2+1}^{n3} E(\mathbf{y}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n E(\mathbf{y}_i^t (\mathbf{r}_i - \psi_i) | \mathbf{x}_{\text{obs}}, \theta^r) \end{aligned}$$

$$\begin{aligned}
\frac{\partial^2}{\partial \omega \partial \omega^t} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{y}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) + \sum_{i=n1+1}^{n2} \mathbf{E}(\mathbf{y}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
&+ \sum_{i=n2+1}^{n3} E(\mathbf{y}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) \mid \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^{n4} \mathbf{E}(\mathbf{y}_i^t \mathbf{r}_i \psi_i (1 - \psi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r)
\end{aligned} \tag{13}$$

where $\psi_i = \exp\{\mathbf{x}_i^t \delta + \mathbf{y}_i \omega\} / (1 + \exp\{\mathbf{x}_i^t \delta + \mathbf{y}_i \omega\})$,

$$\begin{aligned}
\frac{\partial}{\partial \alpha} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t (s_i - \phi_i) + \sum_{i=n1+1}^{n2} \mathbf{E}(\mathbf{x}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
&+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n \mathbf{E}(\mathbf{x}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
\frac{\partial^2}{\partial \alpha \partial \alpha^t} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{x}_i^t s_i \phi_i (1 - \phi_i) + \sum_{i=n1+1}^{n2} \mathbf{E}(\mathbf{x}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
&+ \sum_{i=n2+1}^{n3} E(\mathbf{x}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n \mathbf{E}(\mathbf{x}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r)
\end{aligned} \tag{14}$$

$$\begin{aligned}
\frac{\partial}{\partial \tau} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{y}_i^t (s_i - \phi_i) + \sum_{i=n1+1}^{n2} \mathbf{E}(\mathbf{y}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
&+ \sum_{i=n2+1}^{n3} E(\mathbf{y}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n \mathbf{E}(\mathbf{y}_i^t (s_i - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
\frac{\partial^2}{\partial \tau \partial \tau^t} Q(\theta, \theta^r) &= \sum_{i=1}^{n1} \mathbf{y}_i^t s_i \phi_i (1 - \phi_i) + \sum_{i=n1+1}^{n2} \mathbf{E}(\mathbf{y}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r) \\
&+ \sum_{i=n2+1}^{n3} E(\mathbf{y}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \mathbf{y}_i, \theta^r) + \sum_{i=n3+1}^n \mathbf{E}(\mathbf{y}_i^t s_i \phi_i (1 - \phi_i) \mid \mathbf{x}_{\text{obs}}, \theta^r)
\end{aligned} \tag{15}$$

where $\phi_i = \exp\{\mathbf{x}_i^t \alpha + \mathbf{y}_i \tau\} / (1 + \exp\{\mathbf{x}_i^t \alpha + \mathbf{y}_i \tau\})$,

The $(r+1)^{\text{st}}$ estimates of μ_2 , σ_2^2 are obtained to maximize the log-likelihood by solving the first derivation.

Therefore, we take μ_2^{r+1} , $\sigma_2^{2(r+1)}$ by

$$\mu_2^{r+1} = \frac{1}{n} E(x_{2i} \mid x_{1i}, y_i, r_i, s_i) \tag{16}$$

$$\sigma_2^{2(r+1)} = \frac{1}{n} E((x_{2i} - \mu_2)^2 \mid x_{1i}, y_i, r_i, s_i) \tag{17}$$

4 Summary

We have proposed a method for estimating parameters in logistic linear models when the response is missing or the covariate variable is missing. We used missing data mechanism to avoid serious biases or inefficiency because of model without considering missing data. We generated sample for missing data by Metropolis-hastings algorithm to compute weights. For estimating parameters on the incomplete data set, EM algorithm is general to use. Finding the incomplete data likelihood is quite difficult. Therefore, we compute E-step by the expected log-likelihood and carry out the M-step for a EM iteration. If there is convergence, then we repeat the method until convergence. A small-scale monte-carlo simulation study to evaluate the performance of our proposed method will be presented at the conference.

References

- Baker, S. G. and Laird, N. M. (1988), "Regression analysis for categorical variables with outcome subject to nonignorable nonresponse". *J. Am. Statist. Ass.*, 83, 62-69.
- Donald B. Rubin (Jun., 1974), "Characterizing the Estimation of Parameters in Incomplete-Data Problems", *Journal of the American Statistical Association*, 69, No.346, 467-474.
- Gelman, Andrew (2004), Bayesian data analysis, *Chapman Hall/CRC*.
- Joseph G. Ibrahim (Sep., 1990), "incomplete data in Generalized linear models", *Journal of the American Statistical Association*, 85, No.411, 765-769.
- Joseph G. Ibrahim, Stuart R. Lipsitz (Sep., 1996), "Parameter Estimation from Incomplete Data in Binomial Regression when the missing Data mechanism is nonignorable", *Biometrics*, 52, No.3, 1071-1078.
- Joseph G. Ibrahim, Stuart R. Lipsitz (1999), "Missing covariates in generalized linear models when the missing data mechanism is non-ignorable", *J. R. Statist. Soc.* 61. 173-190.
- Lipsitz, S. R. and Ibrahim, J. G. (1996), "A conditional model for incomplete covariates in parametric regression models. *Biometrika*, 83, 916-922.
- Little, R. J. A., Schluchter, M. (1985), "Maximum likelihood estimation for mixed continuous and categorical data with missing values", *Biometrika*, 72, 497-512.
- Robins, J. M., Rotnitzky, A. and Zhao, L. P. (1994), "Estimation of regression coefficients when some regressors are not always observed" *J. Am. Statist. Ass.*, 89, 846-866.
- Robert, Christian P. (1999), "Monte Carlo statistical methods", *Springer Berlin*.