

Provided for non-commercial research and educational use.
Not for reproduction, distribution or commercial use.

PLISKA

STUDIA MATHEMATICA
BULGARICA

ПЛИСКА

БЪЛГАРСКИ
МАТЕМАТИЧЕСКИ
СТУДИИ

The attached copy is furnished for non-commercial research and education use only.
Authors are permitted to post this version of the article to their personal websites or institutional repositories and to share with other researchers in the form of electronic reprints.
Other uses, including reproduction and distribution, or selling or licensing copies, or posting to third party websites are prohibited.

For further information on
Pliska Studia Mathematica Bulgarica
visit the website of the journal <http://www.math.bas.bg/~pliska/>
or contact: Editorial Office
Pliska Studia Mathematica Bulgarica
Institute of Mathematics and Informatics
Bulgarian Academy of Sciences
Telephone: (+359-2)9792818, FAX:(+359-2)971-36-49
e-mail: pliska@math.bas.bg

IMPLEMENTATION OF THE EM ALGORITHM FOR MAXIMUM LIKELIHOOD ESTIMATION OF A RANDOM EFFECTS MODEL FOR ONE LONGITUDINAL ORDINAL OUTCOME*

Denitsa Grigorova
Ralitza Gueorguieva

ABSTRACT. Longitudinal data arise when we have repeated measures on subjects over time. The correlated probit model is frequently used for ordered longitudinal data since it allows to seamlessly incorporate different correlation structures. The estimation of the probit model parameters based on direct maximization of the limited information maximum likelihood is a numerically intensive procedure especially when we have repeated measures on subjects. We propose an extension of the EM algorithm for obtaining maximum likelihood estimates for one ordinal longitudinal outcome. The algorithm is implemented in the free software environment for statistical computing and graphics **R**. We use simulations to examine the performance of the developed algorithm and apply the model to data from the Health and Retirement Study (HRS). We apply a bootstrap approach for standard error approximation. Advantages of the presented algorithm include the potential of dealing with high-dimensional random effects and of extending the algorithm to combinations of ordinal and continuous longitudinal outcomes.

*The research was partially supported by appropriated state funds for research allocated to Sofia University (contract No 125/2012), Bulgaria

2010 *Mathematics Subject Classification*: 62J99.

Key words: correlated probit model, EM algorithm, free software environment for statistical computing and graphics **R**, ordinal longitudinal data.

1. Introduction. Longitudinal surveys follow up subjects over time. Modeling such data requires taking into account the correlation of measurements within subject. There are three main classes of models for longitudinal data: random effects models, transition models and marginal models. Marginal models are used when we are primarily interested in inferences for the population mean over time and the correlation structure is of secondary interest. Transition models are used when we are interested in modeling the response as a function of preceding outcomes and covariates. Random effects models are used when we are interested in inferences about individual change rather than average population change over time. Diggle et al. [10] describe the theoretical details of modeling longitudinal data while Weiss [29] provides a more applied overview.

Individuals often drop out of longitudinal studies. The mechanism of missingness needs to be taken into account when we model an incomplete data set. While marginal and transition models need to be extended when there are missing data, random effects models deal seamlessly with missing at random observations, that is, when the missingness depends only on observed outcomes and covariates. Another advantage of random effects models is that they simultaneously describe the mean structure and the correlation structure of the data.

In many longitudinal studies the variable of interest is ordinal. For example, in the Health and Retirement Study (HRS, <http://hrsonline.isr.umich.edu/>) self-rated health is categorical with five levels: excellent (coded as 1), very good (2), good (3), fair (4) and poor (5). The survey follows American citizens born between years 1931 and 1941 and their spouses over 14 years. There are seven waves of data collection at intervals of two years. At each interview the participants provided information about their self-rated health and multiple other variables. Our interest is in modeling how self-assessment of health varies over time. We use the correlated probit model.

Probit models were first proposed by Gaddum [14] and Bliss [3, 4] for binary data. Ashford and Sowden [2] introduced a multivariate extension of the probit model based on an underlying multivariate normal distribution. Aitchison and Silvey [1] proposed a probit model for ordinal data. Ochi and Prentice [27] first introduced a correlated probit model but only for exchangeable binary data. Extensions of this model were proposed by Hedeker and Gibbons [15], Catalano [7], Grilli and Rampichini [16], Gueorguieva and Sanacora [19] among others. The correlated probit model has been extensively used because it is easy to interpret and allows different correlation structures within subject. However, computational problems have always been a challenge and the issue of a general approach

to parameter estimation is still open. Gueorguieva [17] has a detailed overview on correlated probit models.

In the present article we consider a correlated probit model which is suitable for ordinal longitudinal data. A special feature of the model is the assumption of latent normal variables with thresholds that generate the observed ordinal responses. The latent variables can be interpreted as unobserved continuous measures that generate the observed responses. For example, in HRS we may consider health to be an underlying continuous measure. We can not directly observe it, rather we know only the level (categorical variable) at which the subjects rate their health at a particular time point. Conceptually, the unobserved continuous measure depends on measured covariates (both time-independent or time-dependent) via fixed effects and unmeasured covariates via random effects.

The correlated probit model does not have closed form expression for the likelihood function and hence approximations need to be used. There are several methods of statistical inference based on numerical, stochastic or analytical approximations. Most popular appear to be extensions of numerical approximations such as Gauss-Hermite quadrature [13] pp. 306–307 or adaptive Gaussian Quadrature [22]. However, these approaches become too computationally intensive and are not feasible for models with many random effects. Another approach is based on analytical approximations (Breslow and Clayton [5], Wolfinger and O’Connell [31]) but it has been shown to produce bias in the parameter estimates especially for binary data or ordinal data with few categories. Stochastic approximations appear most suitable for correlated probit models for longitudinal data since the computational complexity does not increase exponentially with the increase of number of random effects and they provide unbiased results as the number of generated samples increases.

We consider a stochastic extension of the Expectation Maximization (EM) algorithm [9] or more precisely of the Expectation Conditional Maximization (ECM) algorithm [26]. Ruud [28] is the first to apply the EM algorithm for the estimation of the parameters of probit models. Kawakatsu and Largey [20] extend Ruud’s work to a joint model of a single ordinal and multivariate normal outcomes. Chan and Kuk [8] consider a correlated model for a clustered binary variable and propose an ECM algorithm for parameter estimation. Gueorguieva and Agresti [18] extend their approach to correlated binary and continuous outcomes. Our algorithm extends the approach of Chan and Kuk [8] and Gueorguieva and Agresti [18] to ordinal data by using the parameter transformation proposed by Kawakatsu and Largey [20] for estimation of the threshold parameters.

Since the EM algorithm does not provide direct estimates for the standard errors of the parameters, we use bootstrap for standard error estimation. Bootstrap methods were first introduced by Efron [11] and are resampling methods. They are very useful when the theoretical distribution of a statistic of interest is complicated or unknown. Bootstrap methods can be applied to a broad class of problems (e.g. standard error estimation, hypothesis testing, confidence intervals construction).

The paper is organized as follows. Section 2 defines the correlated probit model and outlines the estimation of the parameters and of their standard errors. Section 3 describes the simulation studies that were performed in order to examine the performance of the algorithm. An application of the model to the HRS data is included in Section 4. Section 5 contains concluding remarks and discussion about possible extensions of the algorithm.

2. Correlated probit model for ordinal longitudinal data. Let y_{ij}^* denote the observed ordinal variable with m levels on the i th subject at time j . We assume that there is a latent normal variable y_{ij} that generated the observed variable. We consider the following random effects model:

$$(1) \quad y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + \mathbf{z}'_{ij}\mathbf{b}_i + \epsilon_{ij}.$$

The rule that relates the latent variable to the observed ordinal variable is:

$$(2) \quad y_{ij}^* = \begin{cases} 1, & y_{ij} \leq \alpha_1; \\ j, & \alpha_{j-1} < y_{ij} \leq \alpha_j, \quad j = 2, \dots, m-1; \\ m, & y_{ij} > \alpha_{m-1}; \end{cases}$$

for some unknown thresholds $\alpha_1, \dots, \alpha_{m-1}$.

The vector of random effects is assumed to be normally distributed q -dimensional and is denoted by $\mathbf{b}_i \sim N(\mathbf{0}, \boldsymbol{\Sigma})$. The error term is normally distributed $\epsilon_{ij} \sim N(0, \sigma^2)$ and is independent of the random effects.

The regression parameters for the fixed effects in model (1) are denoted by the p -dimensional vector $\boldsymbol{\beta}$. The vector of predictors for the fixed effects is \mathbf{x}_{ij} and the vector of predictors for the random effects is \mathbf{z}_{ij} . The covariance matrix $\boldsymbol{\Sigma}$ is a quadratic $q \times q$ positive semi-definite matrix.

From the observed data we can not estimate all of the unknown parameters, so we impose the following identifiability restrictions: the first threshold α_1 is set to zero and the variance of the normal error term σ^2 is set to 1. Other restrictions and re-parameterizations are possible.

The correlated model formulation is very appealing because the underlying normal distribution allows for very rich correlation structure. Also, the model has intuitive interpretation and can be easily extended to multiple ordinal and continuous outcomes. Furthermore, estimating the parameters using maximum likelihood allows the use of all results concerning maximum likelihood estimates. Hence Wald, score and likelihood ratio tests can be used for hypothesis testing and confidence interval construction. Likelihood ratio tests can also be used to compare nested models while information criteria such as Akaike Information Criterion or Schwartz-Bayesian Criterion can be used for model selection.

2.1. Maximum likelihood estimation via the EM algorithm. We extend the stochastic ECM algorithm of Chan and Kuk [8] to estimate the unknown parameters in model (1). The first step is to re-parameterize the thresholds so that they can be explicitly included in the complete data log-likelihood. For this we use the approach of Kawakatsu and Largey [20]. We define the differences between consecutive thresholds with $\delta_i = \alpha_i - \alpha_{i-1}$, $i = 2, \dots, m-1$. We also define $\delta_1 = \delta_m = 1$ for completeness and future use. Then we consider a new variable which is a linear transformation of the latent variable: $y_{ij_{new}} = (y_{ij} - \alpha_{y_{ij}^* - 1}) / \delta_{y_{ij}^*}$. For completeness we denote $\alpha_0 = 0$. For example, if $y_{ij}^* = u$, $u = 1, \dots, m$ then $y_{ij_{new}} = (y_{ij} - \alpha_{u-1}) / \delta_u$. Because the new variable is a linear transformation of the latent variable then it is also normally distributed. However, conditional on the observed categorical variable, it is a truncated normal variable. If we observe the first level of y^* the new variable is truncated at $(-\infty, 0]$, if y^* is between the first and the last level the new variable is truncated at $(0, 1]$, and if we observe the last level of y^* the new variable is truncated at $(0, \infty)$.

2.1.1. Complete data log-likelihood. Complete data log-likelihood $\ln L$ is:

$$\ln L = \ln f(\mathbf{b}, \mathbf{y}_{new}) = \sum_{i=1}^n \ln f(\mathbf{b}_i) f(\mathbf{y}_{i_{new}} | \mathbf{b}_i).$$

Apart from the constants the log likelihood has the following closed form:

$$\begin{aligned} \ln L &= -0.5 \sum_{i=1}^n \ln |\Sigma| - 0.5 \sum_{i=1}^n \mathbf{b}'_i \Sigma^{-1} \mathbf{b}_i + \\ &+ \sum_{i=1}^n \sum_{j=1}^{n_i} \ln \delta_{y_{ij}^*} - 0.5 \sum_{i=1}^n \sum_{j=1}^{n_i} [\delta_{y_{ij}^*} y_{ij_{new}} - (\mathbf{x}'_{ij} \boldsymbol{\beta} + \mathbf{z}'_{ij} \mathbf{b}_i - \alpha_{y_{ij}^* - 1})]^2. \end{aligned}$$

Thus we obtain closed form expressions for the estimators of the unknown parameters $\Gamma = (\beta, \Sigma, \delta_2, \dots, \delta_{m-1})$ by setting the first derivatives of the complete data log-likelihood to zero.

2.1.2. Closed form expressions for the estimators. The estimator for the covariance matrix Σ of the random effects is:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=0}^n \mathbf{b}_i \mathbf{b}'_i.$$

Regression parameters for the fixed effects satisfy the following equation:

$$\sum_{i=1}^n \sum_{j=1}^{n_i} \mathbf{x}_{ij} \mathbf{x}'_{ij} \beta = \sum_{i=1}^n \sum_{j=1}^{n_i} [\delta_{y_{ij}^*} y_{ijnew} - \mathbf{z}'_{ij} \mathbf{b}_i + \alpha_{y_{ij}^*-1}] \mathbf{x}_{ij}.$$

It follows that the regression parameters β are a least squares solution of the regression of \tilde{y}_{ij} on \mathbf{x}_{ij} , where $\tilde{y}_{ij} = \delta_{y_{ij}^*} y_{ijnew} - \mathbf{z}'_{ij} \mathbf{b}_i + \alpha_{y_{ij}^*-1}$.

The equations for δ_k , $k = 2, \dots, m-1$ are quadratic equations of the form $a\delta_k^2 + b\delta_k + c = 0$ which always have real roots and the bigger root is always positive. The constants a, b, c are as follows:

$$a = \sum_{i,j} \sum_{y_{ij}^*=k} (y_{ijnew}^2) + n_{k+1} + \dots + n_m,$$

$$b = - \sum_{i,j} \sum_{y_{ij}^*=k} y_{ijnew} (\mathbf{x}'_{ij} \beta + \mathbf{z}'_{ij} \mathbf{b}_i - \alpha_{k-1}) +$$

$$\sum_{i,j} \sum_{y_{ij}^*>k} (\delta_{y_{ij}^*} y_{ijnew} - \mathbf{x}'_{ij} \beta - \mathbf{z}'_{ij} \mathbf{b}_i + \delta_1 + \dots + \delta_{k-1} + \delta_{k+1} + \dots + \delta_{y_{ij}^*-1}),$$

$$c = -n_k,$$

where n_k is the number of observations categorical variable at k -th level.

All we need to do in order to update the parameter estimates at each step of the algorithm is to find the conditional expectations in the closed form expressions of the estimators. We will show that all of the conditional expectations depend only on the first two moments of the truncated multivariate normal distribution.

2.1.3. Conditional expectations. Let us introduce the following notation:

$$\mathbf{X}_i = \begin{pmatrix} x'_{i1} \\ x'_{i2} \\ \vdots \\ x'_{in_i} \end{pmatrix}, \mathbf{Z}_i = \begin{pmatrix} z'_{i1} \\ z'_{i2} \\ \vdots \\ z'_{in_i} \end{pmatrix}, \boldsymbol{\alpha}_i = \begin{pmatrix} \alpha_{y_{i1}^* - 1} \\ \alpha_{y_{i2}^* - 1} \\ \vdots \\ \alpha_{y_{in_i}^* - 1} \end{pmatrix}, \boldsymbol{\delta}_i^{-1} = \begin{pmatrix} 1/\delta_{y_{i1}} \\ 1/\delta_{y_{i2}} \\ \vdots \\ 1/\delta_{y_{in_i}} \end{pmatrix}.$$

Then the joint distribution of $\mathbf{y}_{i_{new}}$ and \mathbf{b}_i is multivariate normal:

$$\begin{pmatrix} \mathbf{y}_{i_{new}} \\ \mathbf{b}_i \end{pmatrix} \sim N \left[\begin{pmatrix} (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1} \\ \mathbf{0} \end{pmatrix}, \mathbf{V} \right],$$

where $\mathbf{J}_{n_i \times q} \boldsymbol{\delta}_i^{-1}$ is a $n_i \times q$ matrix with columns $\boldsymbol{\delta}_i^{-1}$ and \circ is the Hadamard (elementwise) product and

$$\mathbf{V} = \begin{pmatrix} (\mathbf{Z}_i \boldsymbol{\Sigma} \mathbf{Z}_i' + \mathbf{I}_{n_i}) \circ \boldsymbol{\delta}_i^{-1} \boldsymbol{\delta}_i^{-1'} & \mathbf{Z}_i \boldsymbol{\Sigma} \circ (\mathbf{J}_{n_i \times q} \boldsymbol{\delta}_i^{-1}) \\ \boldsymbol{\Sigma} \mathbf{Z}_i' \circ (\mathbf{J}_{n_i \times q} \boldsymbol{\delta}_i^{-1})' & \boldsymbol{\Sigma} \end{pmatrix}.$$

Let us denote $\boldsymbol{\Sigma}_{\mathbf{B}_i} = (\boldsymbol{\Sigma} \mathbf{Z}_i' \circ (\mathbf{J}_{n_i \times q} \boldsymbol{\delta}_i^{-1})') [(\mathbf{Z}_i \boldsymbol{\Sigma} \mathbf{Z}_i' + \mathbf{I}_{n_i}) \circ \boldsymbol{\delta}_i^{-1} \boldsymbol{\delta}_i^{-1'}]^{-1}$.

Then the conditional distribution of \mathbf{b}_i given $\mathbf{y}_{i_{new}}$ is again normal:

$$\mathbf{b}_i | \mathbf{y}_{i_{new}} \sim N[\boldsymbol{\Sigma}_{\mathbf{B}_i} (\mathbf{y}_{i_{new}} - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1}), \boldsymbol{\Sigma} - \boldsymbol{\Sigma}_{\mathbf{B}_i} (\mathbf{Z}_i \boldsymbol{\Sigma} \circ (\mathbf{J}_{n_i \times q} \boldsymbol{\delta}_i^{-1}))].$$

In the expressions for the estimators we have to find the following conditional expectations: $E(\mathbf{b}_i | \mathbf{y}_i^*)$, $E(\mathbf{b}_i \mathbf{b}_i' | \mathbf{y}_i^*)$, $E(y_{ij_{new}} \mathbf{b}_i | \mathbf{y}_i^*)$. We will show that they depend only on the first two moments of $\mathbf{y}_{i_{new}} | \mathbf{y}_i^*$.

The expectation of the random effects conditional on the observed variable is:

$$\begin{aligned} E(\mathbf{b}_i | \mathbf{y}_i^*) &= E[E(\mathbf{b}_i | \mathbf{y}_{i_{new}}) | \mathbf{y}_i^*] \\ &= E[\boldsymbol{\Sigma}_{\mathbf{B}_i} (\mathbf{y}_{i_{new}} - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1}) | \mathbf{y}_i^*] \\ &= \boldsymbol{\Sigma}_{\mathbf{B}_i} [E(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1}]. \end{aligned}$$

Let us denote $\mathbf{M}_i = \mathbf{y}_{i_{new}} - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1}$. For the estimator of the covariance matrix of the random effects we need:

$$\begin{aligned} E(\mathbf{b}_i \mathbf{b}_i' | \mathbf{y}_i^*) &= E[E(\mathbf{b}_i \mathbf{b}_i' | \mathbf{y}_{i_{new}}) | \mathbf{y}_i^*] \\ &= E[\text{Var}(\mathbf{b}_i | \mathbf{y}_{i_{new}}) + E(\mathbf{b}_i | \mathbf{y}_{i_{new}}) E(\mathbf{b}_i' | \mathbf{y}_{i_{new}}) | \mathbf{y}_i^*] \end{aligned}$$

$$\begin{aligned}
&= \Sigma - \Sigma_{\mathbf{B}_i}(\mathbf{Z}_i \Sigma \circ (\mathbf{J}_{n_i \times q} \delta_i^{-1})) + \Sigma_{\mathbf{B}_i} E[\mathbf{M}_i \mathbf{M}_i' | \mathbf{y}_i^*] \Sigma'_{\mathbf{B}_i} \\
&= \Sigma - \Sigma_{\mathbf{B}_i}(\mathbf{Z}_i \Sigma \circ (\mathbf{J}_{n_i \times q} \delta_i^{-1})) + \\
&\quad \Sigma_{\mathbf{B}_i} [\text{Var}(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) + E(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) E(\mathbf{y}'_{i_{new}} | \mathbf{y}_i^*) \\
&\quad \quad - E(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) [(\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}]' \\
&\quad \quad - [(\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}] E(\mathbf{y}'_{i_{new}} | \mathbf{y}_i^*) \\
&\quad \quad + [(\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}] [(\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}]'] \Sigma'_{\mathbf{B}_i}.
\end{aligned}$$

In the expression for the conditional expectation of b we need:

$$\begin{aligned}
E(y_{ij_{new}} \mathbf{b}_i | \mathbf{y}_i^*) &= E[E(y_{ij_{new}} \mathbf{b}_i | \mathbf{y}_{i_{new}}) | \mathbf{y}_i^*] \\
&= E[y_{ij_{new}} \Sigma_{\mathbf{B}_i} (\mathbf{y}_{i_{new}} - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}) | \mathbf{y}_i^*] \\
&= \Sigma_{\mathbf{B}_i} E[y_{ij_{new}} \mathbf{y}_{i_{new}} - y_{ij_{new}} (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1} | \mathbf{y}_i^*] \\
&= \Sigma_{\mathbf{B}_i} [\text{Cov}(y_{ij_{new}} \mathbf{y}_{i_{new}} | \mathbf{y}_i^*) + E(y_{ij_{new}} | \mathbf{y}_i^*) E(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) \\
&\quad - E(y_{ij_{new}} | \mathbf{y}_i^*) (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \delta_i^{-1}].
\end{aligned}$$

The conditional expectations above are available in closed forms [24] but their calculation is computationally very intensive and proved to be inefficient especially when the dimension of the truncated multivariate distribution d is bigger than 2. We describe a stochastic approximation of the conditional expectations which has good practical properties. In order to find the above conditional expectations which we showed to depend only on the first two moments of $\mathbf{y}_{i_{new}} | \mathbf{y}_i^*$ we use a Monte Carlo method. We generate values from the truncated normal distribution given observed data using Gibbs sampling with the help of the **rtmvnorm** function in the **R** package **tmvtnorm**. The algorithm for the generation of random numbers is described in detail in [30]. We use the sample mean and the sample variance of simulated values to approximate $E(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*)$ and $\text{Var}(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*)$. In mathematical notation this is expressed as:

$$\hat{E}(\mathbf{y}_{i_{new}} | \mathbf{y}_i^*) = \frac{1}{m} \sum_{k=1}^m \mathbf{y}_{i_{new}}^{(k)},$$

$$\widehat{Var}(\mathbf{y}_{i_{new}}|\mathbf{y}_i^*) = \frac{1}{m-1} \sum_{k=1}^m (\mathbf{y}_{i_{new}}^{(k)} - \widehat{E}(\mathbf{y}_{i_{new}}|\mathbf{y}_i^*))(\mathbf{y}_{i_{new}}^{(k)} - \widehat{E}(\mathbf{y}_{i_{new}}|\mathbf{y}_i^*))',$$

where $\mathbf{y}_{i_{new}}^{(k)}$ is the k -th realisation of $\mathbf{y}_{i_{new}}|\mathbf{y}_i^*$ in the generated sample of m random numbers. Our experience shows that even for small m (e.g. 150 or 200) we get adequate results.

2.1.4. ($k+1$)-st iteration of the ECM algorithm. The estimates of the unknown parameters at the $k+1$ -st step of the proposed ECM algorithm are updated as follows:

- The $(k+1)$ -st estimates of the regression parameters $\boldsymbol{\beta}^{k+1}$ are the least square solution of the regression of $E(\tilde{y}_{ij}|\mathbf{y}_i^*; \boldsymbol{\Gamma}^k)$ on \mathbf{x}_{ij} .
- The $(k+1)$ -st estimate of δ_u , $u = 2, \dots, m-1$ is: $\delta_u^{k+1} = \frac{-E[b|\mathbf{y}^*; \boldsymbol{\Gamma}^k] + \sqrt{(E[b|\mathbf{y}^*; \boldsymbol{\Gamma}^k]^2 - 4E[a|\mathbf{y}^*; \boldsymbol{\Gamma}^k]E[c|\mathbf{y}^*; \boldsymbol{\Gamma}^k])}}{2E[a|\mathbf{y}^*; \boldsymbol{\Gamma}^k]}$. In the expression for the expectations of a, b, c we use the already updated estimates $\boldsymbol{\beta}^{k+1}$, δ_i^{k+1} , $i = 2, \dots, u-1$.
- The $(k+1)$ -st estimate of the covariance matrix of random effects is $\widehat{\boldsymbol{\Sigma}}^{k+1} = \frac{1}{n} \sum_{i=0}^n E(\mathbf{b}_i \mathbf{b}_i' | \mathbf{y}_i^*; \boldsymbol{\Gamma}^k)$.

In order to update the estimates we use the approximations of the expectations and the variances described in the previous section.

2.2. Standard error estimation. We use the bootstrap method for standard errors approximation described in [25] pp. 130 – 131. The steps are as follows:

1. We fit model (1) to the observed data set consisting of n individuals using the proposed EM algorithm and obtain the estimates of the unknown parameters denoted by $\widehat{\boldsymbol{\Gamma}} = (\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\Sigma}}, \widehat{\boldsymbol{\delta}})$. To generate a bootstrap sample first we generate n random effects \mathbf{b}_k^b from $N(\mathbf{0}, \widehat{\boldsymbol{\Sigma}})$, $k = 1, \dots, n$. Next we simulate normal values \mathbf{y}_k^b of dimension n_k according to model (1) for every random effect \mathbf{b}_k^b . We use the estimated thresholds via $\widehat{\boldsymbol{\delta}}$ to determine in which interval the normal data \mathbf{y}_k^b , $k = 1, \dots, n$ fall and thus using (2) determine the level of the bootstrap categorical variable \mathbf{y}_k^{b*} . The bootstrap sample consists of the categorical variables \mathbf{y}_k^{b*} , $k = 1, \dots, n$.

2. We apply the EM algorithm to the bootstrap data \mathbf{y}_k^{b*} , $k = 1, \dots, n$ to get estimates for the generated data set $\mathbf{\Gamma}^b$.
3. We use Monte Carlo method to approximate the bootstrap covariance matrix. That means that we repeat step 1 and step 2 B times and calculate the covariance matrix of the B estimated parameters $\mathbf{\Gamma}^b$, $b = 1, \dots, B$:

$$Cov(\hat{\mathbf{\Gamma}}) \approx \sum_{b=1}^B \frac{(\mathbf{\Gamma}^b - \bar{\mathbf{\Gamma}})(\mathbf{\Gamma}^b - \bar{\mathbf{\Gamma}})'}{B - 1},$$

$$\text{where } \bar{\mathbf{\Gamma}} = \sum_{b=1}^B \mathbf{\Gamma}^b / B.$$

3. Simulations. For the implementation of the algorithm we used the free software environment for statistical computing and graphics **R**. The **R** code for fitting the presented models is available from the authors.

We simulated values from the following random intercept model:

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + b_i + \epsilon_{ij},$$

where $\beta_0 = -0.5$, $\beta_1 = 1$, $Var(b_i) = \sigma^2 = 0.01$, $Var(\epsilon_{ij}) = 1$ with thresholds $\alpha_1 = 0$, $\alpha_2 = 1.5$, $\alpha_3 = 3$, $\alpha_4 = 4$.

We simulated 100 samples for two different sample sizes ($n = 100$ and $n = 500$) with 5 repeated measures on each individual. For each approximation of the standard errors we used 75 bootstrap samples which is within the recommended range of 50 to 100 bootstrap replications (Efron and Tibshirani [12]). The results are presented in Table 1.

Note that due to the re-parametrization we estimate the differences in thresholds rather than the thresholds themselves. In both simulations the averages of the estimated parameters are equal within two significant digits after the decimal point to the parameter values from which the samples were generated. The only exception is the estimate of the last threshold difference at the smaller sample size setting but even this estimate is close to the true value. Thus we empirically confirm the unbiasedness of the algorithm.

As expected the estimates get closer to the real values and the standard errors get smaller when we increase the sample size. All of the estimates are statistically significantly different from zero except the variance of the random intercept for

Table 1. Table of estimates and standard errors of both simulation studies

real values	$\beta_0 = -0.5$	$\beta_1 = 1$	$\sigma^2 = 0.01$	$\delta_2 = 1.5$	$\delta_3 = 1.5$	$\delta_4 = 1$
Simulation 1: number of subjects = 500, $n_i = 5$						
mean of estimates	-0.503	1.001	0.010	1.50	1.50	0.997
stand. dev. of estimates	0.06	0.023	0.0006	0.053	0.052	0.036
mean of bootstrap stand. errors	0.059	0.023	0.0006	0.052	0.049	0.040
Simulation 2: number of subjects = 100, $n_i = 5$						
mean of estimates	-0.498	1.007	0.011	1.512	1.498	1.009
stand. dev. of estimates	0.114	0.051	0.010	0.107	0.108	0.090
mean of bootstrap stand. errors	0.138	0.052	0.012	0.122	0.112	0.094

the smaller sample size. This is not surprising since this variance is small and the estimates at the smaller sample size are not as efficient as at the larger sample size.

Finally, the approximate equality of the standard deviations of the estimates and the bootstrap standard errors confirms that the algorithm is converging as expected. However, larger simulation study that varies the parameter settings is necessary to confirm the above observations.

4. Application of the model. We apply the proposed model to the HRS data. The variable of main interest in the study (self-rated health) takes values from excellent (1) to poor (5). Categories (2), (3) and (4) mean very good, good and fair self-rated health respectively. We examine how self-rated health changes over time. We fit the following probit model to the data:

$$\text{Self rated health}_{ij} = \beta_0 + \beta_1 \text{Wave}_{ij} + b_i + \epsilon_{ij}$$

In the analysis we include 7550 individuals in the study who have complete set of observations. The results are presented in Table 2:

Table 2. Table of estimates and standard errors of the model fitted to HRS data

	β_0	β_1	σ^2	δ_2	δ_3	δ_4
estimates	1.23	0.12	2.15	1.58	1.50	1.42
standard errors	0.015	0.0026	0.049	0.011	0.011	0.017

Table 2 shows that all of the parameters in the model are statistically significantly different from zero. The parameter of most interest is the regression coefficient β_1 . It is positive and the z-test statistic for this parameter is $z = 0.12/0.0026 = 46.15$, $p\text{-value} < 0.0001$ and thus we conclude that self-rated health deteriorates significantly over time. Further study including additional covariates may reveal whether this change is associated with particular subject characteristics.

We also note that the variance of the random intercept is significantly different from 0. This implies that the between-subject variability of the self-reported health measurements is large and that there is strong correlation between the repeated measurements on a particular individual.

5. Conclusions. In this paper we considered a correlated probit model for the analysis of repeatedly measured ordinal outcomes. We proposed an extension of the EM algorithm of Chan and Kuk [8] for obtaining maximum likelihood estimates, implemented it in the free software environment for statistical computing and graphics **R** and studied its performance using simulations. We also illustrated the approach on self-reported health data from the Health and Retirement Study (HRS). Our approach has advantages over alternative estimation methods in that it can handle a large number of random effects, it can be easily extended to any combination of binary, ordinal and continuous outcomes and it provides unbiased estimates. It is also easily implemented in the open-source software environment **R**. Using free software is a premise for wider usage and quicker improvement of the code.

There are several possible directions in which the algorithm implementation can be improved. There is a possible extension of the algorithm, called parameter expanded ECM algorithm [21] that can accelerate the speed of convergence of the algorithm. Rather than restrict some parameters (e.g. the variance of the error term) for parameter identifiability up front, this extension allows estimation of all parameters free of restrictions and at the last iteration calculates fully iden-

tifiable functions of the parameter (e.g. the ratios of the regression parameters and the squared root of the variance of the errors estimate). An example of implementation of this algorithm can be found in Gueorguieva and Agresti [18].

It is also possible to improve the implementation of the algorithm by choosing different \mathbf{R} functions or improving the efficiency of the code. We already applied one such optimization. Although there are functions for finding the first two moments of the multivariate truncated normal distribution in the package **mvtnorm** based on the work by Manjunath and Wilhelm [24] they are rather slow. Generating random numbers using Gibbs sampling [6] and finding the first two moments based on that sample proved to be quicker.

Standard error estimation is computationally very intensive. While the bootstrap algorithm can always be applied, it is not efficient. Other approaches may be possible. For example, one might consider the Louis's approximation method [23].

Further research is also needed to extend the algorithm to combinations of ordinal and continuous longitudinal outcomes. Model selection and model diagnostics are also open areas of research.

REFERENCES

- [1] J. AITCHISON, S. D. SILVEY. The generalization of probit analysis to the case of multiple responses. *Biometrika* **44** (1957), No 1–2, 131–140.
- [2] J. R. ASHFORD, R. R. SOWDEN. Multi-variate probit analysis. *Biometrics*, **26**(1970), No 3. 535–546.
- [3] C. I. BLISS. The method of probits. *Science*, 1934, 38–39.
- [4] C. I. BLISS. The method of probits. *Science*, 1934, 409–410.
- [5] N. E. BRESLOW, D. G. CLAYTON. Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, **88** (1993), 925.
- [6] G. CASELLA, E. I. GEORGE. Explaining the Gibbs sampler. *The American Statistician* **46** (1992), No 3, 167–174, .
- [7] P. J. CATALANO. Bivariate modelling of clustered continuous and ordered categorical outcomes. *Statistics in Medicine* **16** (1997), No 8, 883–900.

- [8] J. CHAN, A. KUK. Maximum likelihood estimation for probit-linear mixed models with correlated random effects. *Biometrics* 53 (1997), 86–97.
- [9] A. P. DEMPSTER, N. M. LAIRD, D. B. RUBIN. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, **39** (1977) No, 1–22
- [10] P. J. DIGGLE, K.-L. LIANG, S. L. ZEGER. Analysis of Longitudinal Data New York, Oxford University Press, 1996.
- [11] B. EFRON. Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7**(1979), No 1, 1–26.
- [12] B. EFRON, R. J. TIBSHIRANI. An Introduction to the Bootstrap. Monographs on Statistics & Applied Probability, vol. **57**, New York, Chapman & Hall, 1994.
- [13] L. FAHRMEIR, G. TUTZ. Multivariate Statistical Modelling Based on Generalized Linear Models, 2nd edition. New York, Springer-Verlag, 2001.
- [14] J. H. GADDUM. Reports on biological standards. III. Methods of biological assay depending on a quantal response, 1933.
- [15] R. D. GIBBONS, D. HEDEKER. Application of random-effects probit regression models. *Journal of Consulting and Clinical Psychology* **62** (1994) No 2, 285–296.
- [16] L. GRILLI, C. RAMPICHINI. Alternative specifications of multivariate multi-level probit ordinal response models. *Journal of Educational and Behavioral Statistics* **28** (2003), 3144.
- [17] R. V. GUEORGUEVA. Correlated probit model. In: Encyclopedia of Biopharmaceutical Statistics, chapter 59, 2006, 355–362.
- [18] R. V. GUEORGUEVA, A. AGRESTI. A correlated probit model for joint modeling of clustered binary and continuous responses. *Journal of the American Statistical Association* **96** (2001), 1102–1112.
- [19] R. V. GUEORGUEVA, G. SANACORA. Joint analysis of repeatedly observed continuous and ordinal measures of disease severity. *Statistics in Medicine* **25** (2006), 1307–1322.

- [20] H. KAWAKATSU, A. G. LARGEY. EM algorithms for ordered probit models with endogenous regressors. *Econometrics Journal* **12** (2009), 164–186.
- [21] C. LIU, D. RUBIN, Y. WU. Parameter expansion to accelerate EM: The PX-EM algorithm. *Biometrika* **85**(1998), No 4, 755–770.
- [22] Q. LIU, D. A. PIERCE. A note on Gauss-Hermite quadrature. *Biometrika* **81** (1994), No 3, 624–629.
- [23] T. A. LOUIS. Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society. Series B* **44** (1982), No 2, 226–233.
- [24] B. G. MANJUNATH, S. WILHELM. Moments calculation for the double truncated multivariate normal density. <http://ssrn.com/abstract=1472153>, September 11, 2009.
- [25] G. J. MCLACHLAN, T. KRISHNAN. The EM Algorithm and Extensions, 2 edition. New York, Wiley, Mar. 2008.
- [26] X.-L. MENG, D. B. RUBIN. Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika* **80** (1993), No 2, 267–278.
- [27] Y. OCHI, R. L. PRENTICE. Likelihood inference in a correlated probit regression model. *Biometrika* **71** (1984), No 3, 531–543.
- [28] P. A. RUUD. Extensions of estimation methods using the EM algorithm. *Journal of Econometrics* **49** (September 1991), No 3, 305–341.
- [29] R. E. WEISS. *Modeling Longitudinal Data*. New York, Springer Science+Business Media, 2005.
- [30] S. WILHELM. Gibbs sampler for the truncated multivariate normal distribution. Electronic, April 6 2012. <http://cran.r-project.org/web/packages/tmvtnorm/vignettes/GibbsSampler.pdf>
- [31] R. WOLFINGER, M. O’CONNELL. Generalized linear mixed models: A pseudo-likelihood approach. *Journal of Statistical Computation and Simulation* **48** (1993), 233–243.

Denitsa Grigorova
Sofia University "St. Kliment Ohridski"
Faculty of Mathematics and Informatics
James Bouchier Blvd
Sofia 1164
Bulgaria
e-mail: dgrigorova@fmi.uni-sofia.bg

Ralitza Gueorguieva
Department of Biostatistics
Yale School of Public Health
60 College St
New Haven, CT 06520
USA
e-mail: ralitza.gueorguieva@yale.edu