

EM ALGORITHM FOR MLE OF A PROBIT MODEL FOR MULTIPLE ORDINAL OUTCOMES

Denitsa Grigorova, Elitsa Encheva, Ralitzia Gueorguieva

ABSTRACT. The correlated probit model is frequently used for multiple ordered data since it allows to incorporate seamlessly 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. We propose an extension of the EM algorithm for obtaining maximum likelihood estimates for a correlated probit model for multiple ordinal outcomes. The algorithm is implemented in the free software environment for statistical computing and graphics R. We present two simulation studies to examine the performance of the developed algorithm. We apply the model to data on 121 women with cervical or endometrial cancer. Patients developed normal tissue reactions as a result of post-operative external beam pelvic radiotherapy. In this work we focused on modeling the effects of a genetic factor on early skin and early urogenital tissue reactions and on assessing the strength of association between the two types of reactions. We established that there was an association between skin reactions and polymorphism XRCC3 codon 241 (C>T) (rs861539) and that skin and urogenital reactions were positively correlated.

ACM Computing Classification System (1998): G.3.

Key words: correlated probit model, EM algorithm, ordered data, polymorphism XRCC3 codon 241 (C>T) (rs861539), random effects.

1. Introduction. Probit models were first introduced by Bliss [7, 8] and Gaduum [16] for binary data. The main feature of probit models is the assumption of a latent variable which determines the level of the observed ordinal response through thresholds. The usefulness of the model is not affected when the existence of the latent variable does not seem natural.

Aitchison and Silvey [1] proposed a probit model for ordinal data. Ashford and Sowden [6] introduced a multivariate extension of the probit model based on an underlying multivariate normal distribution. Ochi and Prentice [31] first introduced a correlated probit model but only for exchangeable binary data. Extensions of this model were proposed by Hedeker and Gibbons [17], Catalano [10], Grilli and Rampichini [20], Gueorguieva and Sanacora [23] among others. Gueorguieva [21] has a detailed overview on correlated probit models. Correlated probit models are widely used for modeling of multiple categorical variables or clustered/longitudinal ordinal outcomes for these models have two main advantages. They are easy for interpretation and they allow rich correlation structure of the latent variables via random effects and/or correlated errors. That allows to take into account the natural dependence of the measurements on the same subject or within cluster.

The correlated probit model does not have closed form expression for the likelihood function. Approximations need to be used in order to obtain estimates of the unknown parameters. 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 [15, pp. 306–307] or adaptive Gaussian Quadrature [26]. Another approach is based on analytical approximations (Breslow and Clayton [9], Wolfinger and O’Connell [37]) but it has been shown to produce bias in the parameter estimates especially for binary data or ordinal data with few categories. A third approach is the Expectation-Maximization (EM) algorithm [13]. An extension of the EM algorithm is the Expectation/Conditional Maximization (ECM) algorithm [30] which is used in cases of complicated M-step. Ruud [34] is the first to apply the EM algorithm for the estimation of the parameters of probit models. Kawakatsu and Largey [24] extend Ruud’s work to a joint model of a single ordinal and multivariate normal outcomes. Chan and Kuk [11] consider a correlated model for a clustered binary variable and propose an ECM algorithm for parameter estimation.

Our algorithm is a modification of the algorithm of Chan and Kuk [11] and Grigorova and Gueorguieva [19] to multivariate ordinal data by using the parameter transformation proposed by Kawakatsu and Largey [24] for estimation

of the threshold parameters.

We apply the model to data on 121 women with cervical or endometrial cancer. All of the cancer patients received post-operative external beam pelvic radiotherapy. They were followed at the Medical University — Sofia in the period from 2006 to the beginning of 2008. Skin, gastrointestinal and urogenital side effects in the patients were observed and recorded. In this work we focused on modeling of the early (starting from the first day of the radiotherapy to 3 months after it) skin and the early urogenital normal tissue reactions.

A large number of genes are responsible for the biological response of healthy tissues to ionizing radiation including DNA repair genes. Since the damage of the DNA molecule is a process that occurs immediately after exposure, it is logical to assume that impaired reparative processes of DNA may be responsible for the development of radiation adverse events [32]. Given the importance of DNA repair for cell and tissue response after radiation exposure, SNPs (single nucleotide polymorphisms) in genes responsible for signaling DNA damage and reparative mechanisms are suitable candidates in the search for genetic basis of normal tissues' radio-sensitivity. Some studies revealed that the XRCC3 gene plays a key role in the repair of DNA changes induced by ionizing radiation and oxidative stress. It is involved in double breaks DNA repair by homologues recombination [38]. SNPs in this gene affect the risk for development of various malignancies and are associated with different biological markers of impaired DNA repair [4, 5].

The aim of our investigation was to assess the strength of association between early skin and early urogenital tissue reactions on the one hand and associations between them and a particular SNP of XRCC3 (241 Thr/Met) (rs861539) on the other hand. This polymorphism is located in exon 8 and has a potential functional effect [2].

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 cancer data is included in Section 4. Section 5 contains concluding remarks and discussion about possible extensions of the algorithm.

2. Model. We observe p ordinal outcomes on the same subject i with respectively m_1, m_2, \dots, m_p levels denoted by $\mathbf{y}_i^* = (y_{i1}^*, y_{i2}^*, \dots, y_{ip}^*)'$. We use bold type for vectors and matrices. We assume that latent normal variables y_{ij} , $j = 1, 2, \dots, p$ generate the observed variables. We consider the following

correlated probit model for the latent variables:

$$(1) \quad \begin{aligned} y_{ij} &= \mathbf{x}'_{ij}\boldsymbol{\beta}_j + \mathbf{z}'_{ij}\mathbf{b}_{ij} + \epsilon_{ij}, \quad j = 1, 2, \dots, p, \text{ where we observe} \\ y_{ij}^* &= \begin{cases} 1, & \text{if } y_{ij} \leq \alpha_{j,1}; \\ l, & \text{if } \alpha_{j,l-1} < y_{ij} \leq \alpha_{j,l}, \quad l = 2, \dots, m_j - 1; \\ m_j, & \text{if } y_{ij} > \alpha_{j,m_j-1}; \end{cases} \end{aligned}$$

for some thresholds $\alpha_{j,1}, \dots, \alpha_{j,m_j-1}, j = 1, 2, \dots, p$.

We assume a normal distribution of the q -dimensional vector of the random effects $\mathbf{b}_i = (\mathbf{b}'_{i1}, \dots, \mathbf{b}'_{ip})' \sim N(\mathbf{0}, \boldsymbol{\Sigma})$. The covariance matrix $\boldsymbol{\Sigma}$ is a quadratic $q \times q$ positive semi-definite matrix. The error terms are independent normally distributed $\epsilon_{ij} \sim N(0, \sigma^2)$. We also assume that the random effects and the error terms are independent of each other.

The regression parameters for the fixed effects in model (1) are denoted by q_j -dimensional vectors $\boldsymbol{\beta}_j, j = 1, \dots, p$. The vectors of predictors for the fixed effects are $\mathbf{x}_{ij}, j = 1, \dots, p$ and the predictors for the random effects are $\mathbf{z}_{ij}, j = 1, \dots, p$.

From the observed data it is not possible to uniquely estimate all of the unknown parameters, so we pose the following restrictions: the first thresholds $\alpha_{j,1}, j = 1, \dots, p$ are set to zero and the variance of the normal error terms σ^2 is set to 1. Some other restrictions and reparameterisations are possible.

2.1. EM algorithm for MLE. We propose an EM algorithm [13] for estimation of the unknown parameters and thresholds in model (1).

The EM algorithm is an iterative procedure for obtaining maximum likelihood estimates for models that depend on unobserved data. In our model the unobserved data are the latent variables and the random effects. Each iteration of the EM algorithm consists of two steps: E-step (Expectation step) and M-step (Maximisation step). Let us denote with \mathbf{X} the observed data, with \mathbf{Z} the unobserved data and with $\boldsymbol{\Gamma}$ the unknown parameters of the model. The two steps at the $(k + 1)$ -st iteration of the algorithm are:

- E-step: $Q(\boldsymbol{\Gamma}|\boldsymbol{\Gamma}^{(k)}) = E_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\Gamma}^{(k)}} [\ln L(\boldsymbol{\Gamma}; \mathbf{X}, \mathbf{Z})]$, where the ‘complete data’ likelihood function is $L(\boldsymbol{\Gamma}; \mathbf{X}, \mathbf{Z}) = f(\mathbf{X}, \mathbf{Z}|\boldsymbol{\Gamma})$, where $f(\cdot)$ is a density function,
- M-step: $\boldsymbol{\Gamma}^{(k+1)} = \arg \max_{\boldsymbol{\Gamma}} Q(\boldsymbol{\Gamma}|\boldsymbol{\Gamma}^{(k)})$.

The algorithm starts with initial values for the unknown parameters $\boldsymbol{\Gamma}^{(0)}$, iterates between the E-step and the M-step and stops when a converging criterion is met.

Our choice for converging criterion is when $|\mathbf{\Gamma}^{(k+1)} - \mathbf{\Gamma}^{(k)}| < \epsilon$ for each element of the vector, where ϵ is a preselected small number.

The first difficulty in applying the EM algorithm to our model is the introduction of the thresholds in the complete data likelihood. We adopt the approach by Kawakatsu and Largey [24] who extend Ruud's work [34]. According to their method we define the differences between consecutive thresholds with $\delta_{j,i} = \alpha_{j,i} - \alpha_{j,i-1}$, $i = 2, \dots, m_j - 1$, $j = 1, 2, \dots, p$ (we define additionally $\delta_{j,1} = \delta_{j,m_j} = 1$).

It follows the connection $\alpha_{j,i} = \sum_{k=2}^i \delta_{j,k}$, $j = 1, 2, \dots, p$, $i = 2, \dots, m_j - 1$. Then

we consider new variables which are a linear transformation of the latent variables. The new variables are denoted by $y_{ijnew} = (y_{ij} - \alpha_{j,y_{ij}^* - 1}) / \delta_{j,y_{ij}^*}$, $j = 1, 2, \dots, p$, where $\alpha_{j,0} = 0$, $j = 1, 2, \dots, p$ and $\mathbf{y}_{i new} = (y_{i1new}, y_{i2new}, \dots, y_{ipnew})'$.

Since the new variables are a linear transformation of the latent variables, they also have a normal distribution. But given the observed variables, the transformed variables have truncated multivariate normal distribution with boundaries of truncation independent of the unknown parameters.

If we observe the first level of y_{ij}^* the new variable y_{ijnew} is truncated at $(-\infty, 0]$, if y_{ij}^* 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_{ij}^* the new variable is truncated at $(0, \infty)$.

We use the approach by Chan and Kuk [11] in order to find closed form expressions for the unknown parameters $\mathbf{\Gamma} = (\beta'_1, \beta'_2, \dots, \beta'_p, \Sigma, \delta'_1, \delta'_2, \dots, \delta'_p)$, where $\delta_j = (\delta_{j,2}, \dots, \delta_{j,m_j-1})$, $j = 1, \dots, p$.

2.1.1. Complete data log-likelihood. The complete data log-likelihood has the following form:

$$\ln L = \ln f(\mathbf{b}, \mathbf{y}_{new}) = \sum_{i=1}^n \ln f(\mathbf{b}_i) f(\mathbf{y}_{i new} | \mathbf{b}_i) = \sum_{i=1}^n \ln [f(\mathbf{b}_i) \prod_{j=1}^p f(y_{ijnew} | \mathbf{b}_i)],$$

where $f(\cdot)$ denotes a normal density function.

From the model definition and the assumption for the distribution of the random effects it follows that apart from the constants the log-likelihood is:

$$\begin{aligned} \ln L = & -0.5 \sum_{i=1}^n \ln |\Sigma| - 0.5 \sum_{i=1}^n \mathbf{b}'_i \Sigma^{-1} \mathbf{b}_i + \\ & + \sum_{i=1}^n \ln \delta_{1,y_{i1}^*} - 0.5 \sum_{i=1}^n [\delta_{1,y_{i1}^*} y_{i1new} - (\mathbf{x}'_{i1} \beta_1 + \mathbf{z}'_{i1} \mathbf{b}_{i1} - \alpha_{1,y_{i1}^* - 1})]^2 \end{aligned}$$

$$\begin{aligned}
& + \dots \\
& + \sum_{i=1}^n \ln \delta_{p, y_{ip}^*} - 0.5 \sum_{i=1}^n [\delta_{1, y_{i1}^*} y_{i1_{new}} - (\mathbf{x}'_{i1} \boldsymbol{\beta}_1 + \mathbf{z}'_{ip} \mathbf{b}_{ip} - \alpha_{1, y_{i1}^* - 1})]^2.
\end{aligned}$$

2.1.2. Closed form expressions for the estimators. We obtain closed form expressions for the estimators of the unknown parameters by setting the first derivatives of the complete data log-likelihood to zero.

The estimator for the covariance matrix $\boldsymbol{\Sigma}$ of the random effects is:

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

The regression parameters for the fixed effects $\boldsymbol{\beta}_j, j = 1, 2, \dots, p$ satisfy the following system of equations:

$$\sum_{i=1}^n \mathbf{x}_{ij} \mathbf{x}'_{ij} \boldsymbol{\beta}_j = \sum_{i=1}^n [\delta_{j, y_{ij}^*} y_{ij_{new}} - \mathbf{z}'_{ij} \mathbf{b}_{ij} + \alpha_{j, y_{ij}^* - 1}] \mathbf{x}_{ij}.$$

It follows that the regression parameters $\boldsymbol{\beta}_j$ are a least square solution of regression of \tilde{y}_{ij} on \mathbf{x}_{ij} , where $\tilde{y}_{ij} = \delta_{j, y_{ij}^*} y_{ij_{new}} - \mathbf{z}'_{ij} \mathbf{b}_{ij} + \alpha_{j, y_{ij}^* - 1}, j = 1, 2, \dots, p$.

The equations for $\delta_{j, k}, k = 2, \dots, m_j - 1, j = 1, 2, \dots, p$ are quadratic equations of the form: $a_j \delta_{j, k}^2 + b_j \delta_{j, k} + c_j = 0$, which always have real roots and the larger root is always positive. The constants a_j, b_j, c_j are as follows:

$$\begin{aligned}
a_j &= \sum_i \sum_{y_{ij}^* = k} (y_{ij_{new}}^2) + n_{j, k+1} + \dots + n_{j, m}, \\
b_j &= - \sum_i \sum_{y_{ij}^* = k} y_{ij_{new}} (\mathbf{x}'_{ij} \boldsymbol{\beta}_j + \mathbf{z}'_{ij} \mathbf{b}_{ij} - \alpha_{j, k-1}) + \\
& \sum_i \sum_{y_{ij}^* > k} (\delta_{j, y_{ij}^*} y_{ij_{new}} - \mathbf{x}'_{ij} \boldsymbol{\beta}_j - \mathbf{z}'_{ij} \mathbf{b}_{ij} + \delta_{j, 2} + \dots + \delta_{j, k-1} + \delta_{j, k+1} + \dots + \delta_{j, y_{ij}^* - 1}), \\
c_j &= -n_{j, k},
\end{aligned}$$

where $n_{j, k}$ is the number of the observations of the categorical variable y_j^* at the k th level.

In order to update the new estimates of the parameters we need to express the conditional expectations in the closed form expressions for the estimators. We will show that all of the conditional expectations depend only on the first two moments of the truncated multivariate normal distribution.

Let us have the following notation: $\mathbf{X}_i = \begin{pmatrix} \mathbf{x}'_{i1} & \mathbf{0}' & \dots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{x}'_{i2} & \dots & \mathbf{0}' \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \dots & \mathbf{x}'_{ip} \end{pmatrix}$,

$$\mathbf{Z}_i = \begin{pmatrix} z'_{i1} & \mathbf{0}' & \dots & \mathbf{0}' \\ \mathbf{0}' & z'_{i2} & \dots & \mathbf{0}' \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \dots & z'_{ip} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\alpha}_i = \begin{pmatrix} \alpha_{1,y_{i1}^*-1} \\ \alpha_{2,y_{i2}^*-1} \\ \vdots \\ \alpha_{p,y_{ip}^*-1} \end{pmatrix},$$

$$\boldsymbol{\delta}_i^{-1} = \begin{pmatrix} 1/\delta_{1,y_{i1}} \\ 1/\delta_{2,y_{i2}} \\ \vdots \\ 1/\delta_{p,y_{ip}} \end{pmatrix}.$$

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 \circ is the Hadamard (element-wise) product and the covariance matrix \mathbf{V} is:

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

and $\mathbf{J}_{p \times q} \boldsymbol{\delta}_i^{-1}$ is $p \times q$ matrix with columns $\boldsymbol{\delta}_i^{-1}$.

Let us denote $\mathbf{M}_i = \mathbf{y}_{i_{new}} - (\mathbf{X}_i \boldsymbol{\beta} - \boldsymbol{\alpha}_i) \circ \boldsymbol{\delta}_i^{-1}$ and with $\boldsymbol{\Sigma}_{\mathbf{B}_i} = [\boldsymbol{\Sigma} \mathbf{Z}'_i \circ (\mathbf{J}_{p \times q} \boldsymbol{\delta}_i^{-1})'] [(\mathbf{Z}_i \boldsymbol{\Sigma} \mathbf{Z}'_i + \mathbf{I}_p) \circ \boldsymbol{\delta}_i^{-1} \boldsymbol{\delta}_i^{-1'}]^{-1}$.

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

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

In the expressions for the estimators we have to calculate 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 given the observed variable is:

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

The expectation of the second moment of the random effects given the observed variable is:

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

The last expectation that we need is:

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

2.1.3. $(k+1)$ -st iteration of the EM algorithm. We use an extension of the EM algorithm called Expectation/Conditional Maximization algorithm [30]. The E-step at the $(k+1)$ -st iteration of the proposed algorithm consists of finding of following expectations: $E(\mathbf{b}_i | \mathbf{y}_i^*; \Gamma^k)$, $E(\mathbf{b}_i \mathbf{b}'_i | \mathbf{y}_i^*; \Gamma^k)$, $E(y_{ij_{new}} \mathbf{b}_i | \mathbf{y}_i^*; \Gamma^k)$, where Γ^k are the k th estimates of the unknown parameters Γ . The M-step consists of several computationally simpler CM-steps. In each CM-step we maximise the expectation of the complete data log-likelihood function in respect to some parameters when the other parameters are kept fixed. We will write down the estimates of the unknown parameters at $(k+1)$ st iteration of the EM algorithm:

- The $(k+1)$ st estimate of regression parameters β_j^{k+1} , $j = 1, 2, \dots, p$ is a least square solution of regression of $E(\tilde{y}_{ij} | \mathbf{y}_i^*; \Gamma^k)$ on x_{ij} .
- The $(k+1)$ st estimates of $\delta_{j,u}$, $u = 2, \dots, m_j - 1$, $j = 1, 2, \dots, p$ are $\delta_{j,u}^{k+1} = (-E[b_j | \mathbf{y}^*; \Gamma^k] + \sqrt{(E[b_j | \mathbf{y}^*; \Gamma^k]^2 - 4E[a_j | \mathbf{y}^*; \Gamma^k]E[c_j | \mathbf{y}^*; \Gamma^k])}) / 2E[a_j | \mathbf{y}^*; \Gamma^k]$ and in the expression for a_j, b_j, c_j we use the updated estimates $\beta_j^{k+1}, \delta_{j,i}^{k+1}, i = 2, \dots, u - 1$.

- The $(k+1)$ st estimate of the covariance matrix of random effects is $\hat{\Sigma}^{k+1} = \frac{1}{n} \sum_{i=0}^n E(\mathbf{b}_i \mathbf{b}_i' | \mathbf{y}_i^*; \Gamma^k)$.

The algorithm starts with initial values for the unknown parameters Γ^0 , iterates between E-step and M-step and stops when $|\Gamma^{k+1} - \Gamma^k| < \epsilon$ for each element of the vector, where ϵ is a preselected small number (for example $\epsilon = 0.0001$).

2.2. Standard error estimation. We use the bootstrap method for standard errors approximation described in [29] 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 ECM algorithm and obtain the estimates of the unknown parameters denoted by $\hat{\Gamma} = (\hat{\beta}'_1, \hat{\beta}'_2, \dots, \hat{\beta}'_p, \hat{\Sigma}, \hat{\delta}'_1, \hat{\delta}'_2, \dots, \hat{\delta}'_p)$. To generate a bootstrap sample first we generate n random effects \mathbf{b}_k^b from $N(\mathbf{0}, \hat{\Sigma})$, $k = 1, \dots, n$. Next we simulate normal values \mathbf{y}_k^b of dimension p according to the fitted model for every random effect \mathbf{b}_k^b . We use the estimated via $\hat{\delta}_j, j = 1, \dots, p$ thresholds to determine in which interval the normal data $\mathbf{y}_k^b, k = 1, \dots, n$ fall and determine the levels 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 ECM algorithm to the bootstrap data $\mathbf{y}_k^{b*}, k = 1, \dots, n$ to obtain estimates for the generated bootstrap data set Γ^b .
3. We use the Monte Carlo method to approximate the bootstrap covariance matrix. This means that we repeat step 1 and step 2 B times and calculate the covariance matrix of the B estimated parameters $\Gamma^b, b = 1, \dots, B$:

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

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

3. Simulations. We simulated values from the following correlated probit model for two ordinal outcomes with three levels each:

$$\begin{aligned}y_{i1} &= \beta_{10} + \beta_{11}x_{i1} + b_{i1} + \epsilon_{i1}, \\y_{i2} &= \beta_{20} + \beta_{21}x_{i2} + b_{i2} + \epsilon_{i2},\end{aligned}$$

where $\beta_{10} = -0.5, \beta_{11} = 1, \beta_{20} = 1, \beta_{21} = -0.5, \text{Var}(\epsilon_{ij}) = 1, j = 1, 2$ with thresholds $\alpha_{1,1} = \alpha_{2,1} = 0, \alpha_{1,2} = 1.2, \alpha_{2,2} = 0.7$ and covariance matrix of the random intercepts

$$\text{Var} \begin{pmatrix} b_{i1} \\ b_{i2} \end{pmatrix} = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}.$$

We simulated 100 samples for two sample sizes ($n = 100$ and $n = 500$). For each approximation of the standard errors we used 50 bootstrap samples which is within the recommended range of 50 to 100 bootstrap replications (Efron and Tibshirani [14]). The results are presented in Table 1.

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

parameters	β_{10}	β_{11}	β_{20}	β_{21}	$\delta_{1,2}$	$\delta_{2,2}$	σ_{11}	σ_{12}	σ_{22}
values	-0.5	1	1	-0.5	1.2	0.7	1	-0.8	1
Simulation 1: number of subjects = 100									
mean of estimates	-.550	1.057	1.070	-0.537	1.233	0.765	1.140	-0.940	1.134
stand. dev. of estimates	0.354	0.316	0.299	0.246	0.217	0.151	0.400	0.428	0.405
mean of bootstrap stand. errors	0.371	0.335	0.291	0.237	0.216	0.167	0.452	0.485	0.456
Simulation 2: number of subjects = 500									
mean of estimates	-.494	0.992	1.004	-0.505	1.203	0.703	1.003	-0.802	1.003
stand. dev. of estimates	0.149	0.141	0.116	0.067	0.097	0.067	0.166	0.181	0.170
mean of bootstrap stand. errors	0.166	0.148	0.118	0.084	0.087	0.068	0.160	0.173	0.161

Note that due to the re-parametrization we estimate the differences in the thresholds rather than the thresholds themselves, but they coincide in the case of only three levels of the categorical variables. In the second simulation 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 except for two parameters, but they differ from the true values by < 0.01 . For the first simulation study we obtain biased estimates, which may be explained

with the sample size. From statistical theory we know that maximum likelihood estimates are only asymptotically unbiased. In larger samples the distribution of the estimates will approximate normality even more closely due to the properties of MLE.

As expected the estimates get closer to the real values and the standard errors get smaller when we increase the sample size. All the estimates are statistically significantly different from zero except β_{10} for the smaller sample size.

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

3.1. Implementation of the algorithm. For the implementation of the algorithm we used the free software environment for statistical computing and graphics R [33]. The R code for fitting the presented models is available on the journal's web site or from the authors.

We want to point out several things regarding the implementation of the proposed ECM algorithm. In the package `mvtnorm` [36] there are functions for analytical finding of the first two moments of multivariate truncated normal distribution based on the work by Manjunath and Wilhelm [28]. There are also functions for generating random numbers using Gibbs sampling [35] which allows stochastic approximation of the first two moments of the truncated normal distribution. But when the multiple outcomes are only two, the analytical calculation is more precise and at least as fast as the stochastic, so we recommend it.

A good choice of starting points for the regression parameters and thresholds in model (1) for the proposed ECM algorithm are estimates from a model without random effects. Selecting large values as starting points for the variances of the random effects should be avoided. Problems with performance of the algorithm may occur with starting points corresponding to a multivariate truncated normal distribution for which the truncation area is close to 0. In such cases finding analytical solutions for the moments of the truncated normal distribution may fail. Generating random numbers via Gibbs sampling may also fail.

For data from the first simulation study it takes less than a second on average to perform one iteration of the algorithm, while for a data set from the second simulation study the time for performance of one iteration is less than 4 seconds on average on an Intel(R) Core(TM) i3 CPU @2.27 GHz with 4 GB RAM. In our first simulation study it took 60 iterations on average for the algorithm to converge and in our second simulation study – 50 iterations on average.

4. Application of the model. We apply the proposed model to a data set from a study designed to assess the association of the severity of the normal tissue reactions after radiotherapy in women with endometrial or cervical cancer and their genetic characteristics. Previous analyses of the data can be found in Grigorova [18]. In this manuscript we focus on modeling the severity of two types of reactions. The variables of main interest are the severity of skin reactions and the severity of urogenital reactions. The variables take values from the following levels: absent (1), weak (2), moderate or severe (3) reactions. We examine how the genotype of polymorphism XRCC3 codon 241 (C>T) is related to the severity of the reactions.

In the analysis we include 121 individuals in the study who have a complete set of observations. The variable XRCC3_241 takes values 0 for genotype {C,C} (45 observations) and 1 for genotype {C,T} or {T,T} (76 observations). We merged genotypes {C,T} and {T,T} in one level of the categorical variable because we have too few women with genotype {T,T}. This way we assess the influence of the T allele on the severity of early adverse reactions. A summary of the data is presented in Table 2 and Table 3.

Table 2. Contingency table of skin and urogenital reactions

Skin reactions	Urogenital reactions		
	absent	weak	moderate or severe
absent	24	17	7
weak	12	16	14
moderate or severe	7	12	12

Table 3. Contingency table of skin reactions versus genotype and contingency table of urogenital reactions versus genotype

XRCC3 241	Skin reactions		
	absent	weak	moderate or severe
{C,C}	11	19	15
{C,T} or {T,T}	37	23	16
XRCC3 241	Urogenital reactions		
	absent	weak	moderate or severe
{C,C}	16	17	12
{C,T} or {T,T}	27	28	21

We fit the following correlated probit model to the data:

$$y_{ij} = \beta_{j0} + \beta_{j1}XRCC3_241_i + b_{ij} + \epsilon_{ij},$$

Table 4. Table of estimates and standard errors for the model fitted to the radiotherapy data

	β_{10}	β_{11}	β_{20}	β_{21}	$\delta_{1,2}$	$\delta_{2,2}$	σ_{11}	σ_{12}	σ_{22}
estimates	0.880	-0.772	0.554	0.025	1.417	1.500	1.237	0.801	1.366
standard errors	0.263	0.264	0.242	0.268	0.187	0.181	0.235	0.307	0.256
z-score	3.342	-2.928	2.292	0.093	7.578	8.275	5.255	2.610	5.327

$$\text{Reactions}_{ij} = \begin{cases} \text{absent, } y_{ij} \leq \alpha_{j,1}, (\alpha_{j,1} = 0), \\ \text{weak, } 0 < y_{ij} \leq \alpha_{j,2}, \\ \text{moderate or severe, } y_{ij} > \alpha_{j,2}, \end{cases}$$

where $(\epsilon_{i1}, \epsilon_{i2})' \sim N(\mathbf{0}, \mathbf{I}_2)$, $j = 1$ for skin reactions, $j = 2$ for urogenital reactions and

$$\Sigma = \text{Var} \begin{pmatrix} b_{i1} \\ b_{i2} \end{pmatrix} = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}.$$

The estimates of the parameters, their standard errors and z-scores are presented in Table 4. Z-scores are computed before rounding off the estimates and their standard errors, and then rounded to the third decimal point. The results show that all of the parameters in the model are statistically significantly different from zero, except the regression parameter β_{21} for the genotype in the sub-model for the urogenital reactions.

The parameters of most interest are the regression coefficients β_{11} and β_{21} . The estimate for β_{11} is negative. The two-sided z-test statistic for this parameter is $z = -2.928$ with p -value = 0.0034 and thus we conclude that genotype {C,C} of polymorphism XRCC3 codon 241 increases the risk of adverse skin reactions. Further study including additional covariates may reveal an association between the severity of the side effects and other particular subject characteristics.

We note that the variances of the random intercepts are significantly greater than 0. Because under the null hypothesis, the variances take values at the boundary of the parameter space, the actual distributions of the squared z-scores are mixtures of chi-square distributions with 0 and 1 degrees of freedom. Using the normal distribution rather than the mixture for calculating the p-value of the test is a conservative approach. The p-values for both variances are < 0.0001 . The covariance of the random effects is statistically significantly different than zero (p -value = 0.009 for the two-sided test). The estimate for the covariance of the random effects is 0.801 and the conclusion is that the association between skin and urogenital reactions is significantly positive.

5. Discussion. In this paper we considered a correlated probit model for the analysis of multiple ordinal outcomes. We proposed an extension of the EM algorithm of Chan and Kuk [11] and the ECM algorithm of Grigorova and Gueorguieva [19] for obtaining maximum likelihood estimates. The algorithm is implemented in the free software environment for statistical computing and graphics R [33]. We studied its performance via simulations. We illustrated the approach on a data set previously analyzed in Grigorova [18]. Our approach has advantages over alternative estimation methods in that it can handle a large dimension of the multivariate outcome, it can be easily extended to any combination of binary, ordinal and continuous outcomes, and it provides asymptotically unbiased estimates. It is also easily implemented in the free open-source software environment R.

There are several possible directions in which the algorithm implementation can be improved. There is a possible extension of the algorithm, called parameter expanded EM algorithm (PX-EM algorithm, [25]), that can accelerate the speed of convergence of the algorithm. Rather than restricting some parameters (e.g., the variance of the error terms) in order to achieve parameter identifiability up front, this extension allows estimation of all or some parameters free of restrictions. At the last iteration of the PX-EM algorithm fully identifiable functions of the parameters are calculated (e.g., the ratios of the regression parameters and the squared root of the variance of the errors estimate). An example of an implementation of this algorithm can be found in Gueorguieva and Agresti [22].

We used a bootstrap method for standard error estimation, which 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 Louis's approximation method [27].

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

The observed positive association of the early skin and urogenital reactions is a novel finding. Although statistically significant, it is not clear what the clinical significance is. One possible explanation is the clinical assessment of the side effects. Some urogenital reactions reported by the patient could be related to the skin reactions in the irradiated area. For example dysuria is one of the urogenital radiation side effects and is a condition when the patient feels pain when urinating. But when a patient has skin side effects around the genital area urination also could cause itching and pain, so we could wrongly assess this adverse

event as urogenital instead of skin reaction.

Our analysis revealed that genotype {C,C} of polymorphism XRCC3 241 increases the risk of adverse skin normal tissue reactions. Our result is close to the findings of a Danish post-mastectomy study that has found a significant relationship between allele C at codon 241 and increased risk of subcutaneous fibrosis in a subgroup analysis of 41 patients. Such association was also found for teleangiectasia [5]. In other studies, however, carried out specifically to confirm the results from those 41 breast cancer patients, no replication of the initial results on association between this SNP and the risk of late radiation effects was obtained [12, 3]. Further extended studies are needed.

Acknowledgements. This work was supported by the European Social Fund through the Human Resource Development Operational Programme under contract BG051PO001-3.3.06-0052 (2012/2014).

The research is partially supported by appropriated state fund for research allocated to Sofia University (contract 111/2013), Bulgaria.

REFERENCES

- [1] AITCHISON J., S. D. SILVEY. The generalization of probit analysis to the case of multiple responses. *Biometrika*, **44** (1957), No 1–2, 131–140.
- [2] AKA P., R. MATEUCA, J. P. BUCHET, H. THIERENS, M. KIRSCH-VOLDERS. Are genetic polymorphisms in OGG1, XRCC1 and XRCC3 genes predictive for the DNA strand break repair phenotype and genotoxicity in workers exposed to low dose ionising radiations? *Mutat. Res.*, **556** (2004), No 1–2, 169–181.
- [3] ANDREASSEN C. N. Can risk of radiotherapy-induced normal tissue complications be predicted from genetic profiles? *Acta Oncol.*, **44** (2005), 801–815.
- [4] ANDREASSEN C. N., J. ALSNER, J. OVERGAARD. Does variability in normal tissue reactions after radiotherapy have a genetic basis—here and how to look for it? *Radiother. Oncol.*, **64** (2002), 131–140.
- [5] ANDREASSEN C. N., J. ALSNER, J. OVERGAARD. Prediction of normal tissue radiosensitivity from polymorphisms in candidate genes. *Radiother. Oncol.*, **69** (2003), 127–135.
- [6] ASHFORD J. R., R. R. SOWDEN. Multi-variate probit analysis. *Biometrics*, **26** (1970), No 3, 535–546.

- [7] BLISS C. I. The method of probits. *Science*, **79** (1934), 38–39.
- [8] BLISS, C. I. The method of probits — a correction. *Science*, **79** (1934), 409–410.
- [9] BRESLOW N. E., D. G. CLAYTON. Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, **88** (1993), 9–25.
- [10] CATALANO, P. J. Bivariate modelling of clustered continuous and ordered categorical outcomes. *Statistics in Medicine*, **16** (1997), No 8, 883–900.
- [11] CHAN, J. S. K., A. Y. C. KUK Maximum likelihood estimation for probit-linear mixed models with correlated random effects. *Biometrics*, **53** (1997), 86–97.
- [12] DE RUYCK K., M. VAN EIJKEREN, K. CLAES, R. MORTIER, A. DE PAEPE, A. VRAL, L. DE RIDDER, H. THIERENS. Radiation-induced damage to normal tissues after radiotherapy in patients treated for gynecologic tumors: association with single nucleotide polymorphisms in XRCC1, XRCC3, and OGG1 genes and in vitro chromosomal radiosensitivity in lymphocytes. *Int. J. Radiat. Oncol. Biol. Phys.*, **62** (2005), 1140–1149.
- [13] DEMPSTER A. P., 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 2, 1–22.
- [14] EFRON B., R. J. TIBSHIRANI. An Introduction to the Bootstrap. 1 ed., Vol. 57 of Monographs on Statistics & Applied Probability, Chapman & Hall, New York, 1994.
- [15] FAHRMEIR L., G. TUTZ. Multivariate Statistical Modelling Based on Generalized Linear Models, second ed., Springer-Verlag, New York, 2001.
- [16] GADDUM J. H. Methods of biological assay depending on a quantal response. Reports on biological standards, III, 1933.
- [17] GIBBONS R. D., D. HEDEKER. Application of random-effects probit regression models. *Journal of Consulting and Clinical Psychology*, **62** (1994), No 2, 285–296.
- [18] GRIGOROVA D. Assessing the effect of genetic factors and other factors on the normal tissue reactions after radiotherapy in patients with cancer. In: Proceedings of the 16th EYSM, Bucharest, Romania, 2009, 108–112. Short paper.

- [19] GRIGOROVA D., R. GUEORGUEVA. Implementation of the EM algorithm for maximum likelihood estimation of a random effects model for one longitudinal ordinal outcome. *Pliska Stud. Math. Bulgar.* **22** (2013), 41–56.
- [20] GRILLI L., C. RAMPICHINI. Alternative specifications of multivariate multi-level probit ordinal response models. *Journal of Educational and Behavioral Statistics*, **28** (2003), 31–44.
- [21] GUEORGUEVA R. V. Correlated probit model. Encyclopedia of Biopharmaceutical Statistics, 2006, ch. 59, 355–362.
- [22] GUEORGUEVA R. V., 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.
- [23] GUEORGUEVA R. V., G. SANACORA. Joint analysis of repeatedly observed continuous and ordinal measures of disease severity. *Statistics in Medicine*, **25** (2006), 1307–1322.
- [24] KAWAKATSU H., A. G. LARGEY. EM algorithms for ordered probit models with endogenous regressors. *Econometrics Journal*, **12** (2009), 164–186.
- [25] LIU C., D. RUBIN, Y. WU. Parameter expansion to accelerate EM: The PX-EM algorithm. *Biometrika*, **85** (1998), No 4, 755–770.
- [26] LIU Q., D. A. PIERCE. A note on Gauss-Hermite quadrature. *Biometrika*, **81** (1994), No 3, 624–629.
- [27] LOUIS T. A. Finding the observed information matrix when using the EM algorithm. *Journal of the Royal Statistical Society. Series B*, **44** (1982), No 2, 226–233.
- [28] MANJUNATH B. G., S. WILHELM. Moments calculation for the double truncated multivariate normal density. <http://ssrn.com/abstract=1472153>, September 11, 2009
- [29] MCLACHLAN G. J., T. KRISHNAN. The EM Algorithm and Extensions. Wiley Series in Probability and Statistics, 2 ed., Wiley-Interscience, 2008.
- [30] MENG X.-L., D. B. RUBIN. Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika*, **80** (1993), No 2, 267–278.
- [31] OCHI Y., R. L. PRENTICE. Likelihood inference in a correlated probit regression model. *Biometrika*, **71** (1984), No 3, 531–543.
- [32] POPANDA O., J. U. MARQUARDT, J. CHANG-CLAUDE, P. SCHMEZER. Genetic variation in normal tissue toxicity induced by ionizing radiation. *Mutat. Res.*, **667** (2009), 58–69.

- [33] THE R DEVELOPMENT CORE TEAM. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, ISBN 3-900051-07-0, Vienna, Austria, 2013.
- [34] RUUD P. A. Extensions of estimation methods using the EM algorithm. *Journal of Econometrics*, **49** (1991), No 3, 305–341.
- [35] WILHELM S. Gibbs sampler for the truncated multivariate normal distribution. Electronic, April 6 2012. <http://cran.r-project.org/web/packages/tmvtnorm/vignettes/GibbsSampler.pdf>.
- [36] WILHELM S., B. G. MANJUNATH. Tmvtnorm: Truncated Multivariate Normal and Student t Distribution, R package version 1.4-7, 2012.
- [37] WOLFINGER R., M O'CONNELL. Generalized linear mixed models: A pseudo-likelihood approach. *Journal of Statistical Computation and Simulation*, **48** (1993), 233–243.
- [38] WOOD R. D., M. MITCHELL, J. SGOUROS, T. LINDAHL. Human DNA repair genes. *Science*, **291** (2001), 1284–1289.

Denitsa Grigorova
Department of Probability,
Operational Research and Statistics
Faculty of Mathematics and Informatics
“St. Kl. Ohridski” University of Sofia
5, J. Bourchier Blvd, P.O. Box 48
1164 Sofia, Bulgaria
e-mail: dgrigorova@fmi.uni-sofia.bg

Ralitza Gueorguieva
School of Public Health
Department of Biostatistics
Yale University
New Haven, Connecticut, USA
e-mail: ralitza.gueorguieva@yale.edu

Elitsa Encheva
Department of Radiation Oncology
Saint Marina University Hospital
Medical University Varna
9002 Varna, Bulgaria
e-mail: dr.encheva@gmail.com

Received November 11, 2013
Final Accepted December 12, 2013