# Graphical Models for Ordinal Data

**4 authors**, including:

Jian Guo
Harvard University
**8** PUBLICATIONS **524** CITATIONS

George Michailidis
University of Michigan
**397** PUBLICATIONS **7,993** CITATIONS

**Some of the authors of this publication are also working on these related projects:**

Change Point Detection in Networks View project

Control of Network Systems View project

# Graphical Models for Ordinal Data

**Jian Guo**[*], **Elizaveta Levina**[†], **George Michailidis**[†], and **Ji Zhu**[†]

[*]Department of Biostatistics, Harvard University

[†]Department of Statistics, University of Michigan, Ann Arbor

## Abstract

A graphical model for ordinal variables is considered, where it is assumed that the data are generated by discretizing the marginal distributions of a latent multivariate Gaussian distribution. The relationships between these ordinal variables are then described by the underlying Gaussian graphical model and can be inferred by estimating the corresponding concentration matrix. Direct estimation of the model is computationally expensive, but an approximate EM-like algorithm is developed to provide an accurate estimate of the parameters at a fraction of the computational cost. Numerical evidence based on simulation studies shows the strong performance of the algorithm, which is also illustrated on data sets on movie ratings and an educational survey.

### Keywords

Graphical model; lasso; ordinal variable; probit model

## 1 Introduction

Graphical models have been successful in identifying directed and undirected structures from high dimensional data. In a graphical model, the nodes of the network correspond to random variables and the edges represent their corresponding associations (Lauritzen, 1996). Two canonical classes of graphical models are the Gaussian one, where the dependence structure is fully specified by the inverse covariance matrix and the Markov one, where the dependence structure is captured by the interaction effects in an exponential family model. In the latter model, each interaction effect can be interpreted as the conditional log-odds-ratio of the two associated variables given all other variables. In both models, a zero element in the inverse covariance matrix or a zero interaction effect determines a conditionally independent relationship between the corresponding nodes in the network.

Estimation of such models from high dimensional data under a sparsity assumption has attracted a lot of interest in the statistics and machine learning literature, including regularized likelihood and regression methods, for example, see Yuan and Lin (2007); Banerjee et al. (2008); Friedman et al. (2008); Rothman et al. (2008); Fan et al. (2009); Meinshausen and Buhlmann (2006); Rocha et al. (2008); Peng et al. (2009) and references therein. For a Markov network, direct estimation of a regularized likelihood is infeasible due to the intractable partition function in the likelihood. Instead, existing methods in the literature employ variants of approximation estimation methods. Examples include the surrogate likelihood methods (Banerjee et al., 2008; Kolar and Xing, 2008) and the pseudo-

likelihood methods (Höefling and Tibshirani, 2009; Ravikumar et al., 2010; Guo et al., 2010).

In many applications involving categorical data, an ordering of the categories can be safely assumed. For example, in marketing studies consumers rate their preferences for a wide range of products. Similarly, computer recommender systems utilize customer ratings to make purchase recommendations to new customers; this constitutes a key aspect of the business model behind Netflix, Amazon, Tripadvisor, etc (Koren et al., 2009).

Ordinal variables are also an integral part of survey data, where respondents rate items or express level of agreement/disagreement on issues/topics under consideration. Such responses correspond to Likert items and a popular model to analyze such data is the polychotomous Rasch model (von Davier and Carstensen, 2010) that obtains interval level estimates on a continuum, an idea that we explore in this work as well. Ordinal response variables in regression analysis give rise to variants of the classical linear model, including the proportional odds model (Walker and Duncan, 1967; McCullagh, 1980), the partial proportional odds model (Peterson, 1990), the probit model (Bliss, 1935; Albert and Chib, 1993; Chib and Greenberg, 1998), etc. A comprehensive review of ordinal regression is given in McCullagh and Nelder (1989) and O'Connell (2005).

In this paper, we introduce a graphical model for ordinal variables. It is based on the assumption that the ordinal scales are generated by discretizing the marginal distributions of a latent multivariate Gaussian distribution and the dependence relationships of these ordinal variables are induced by the underlying Gaussian graphical model. In this context, an EM-like algorithm is appropriate for estimating the underlying latent network, which presents a number of technical challenges that have to be addressed for successfully pursuing this strategy.

Our work is related to Albert and Chib (1993), Chib and Greenberg (1998) and Stern et al. (2009) in the sense that they are all built on the probit model and/or the EM algorithmic framework. Albert and Chib (1993) proposed an MCMC algorithm for the probit-model-based univariate ordinal regression problem, where an ordinal response is fitted on a number of covariates, while Chib and Greenberg (1998) can be considered an extension to the multivariate case. Stern et al. (2009) aims in building an online recommender system via collaborative filtering and applied the discretization/thresholding idea in the probit model to the ordinal matrix factorization problem. Our model, on the other hand, has a completely different motivation from these works. Our objective is explore associations between a set of ordinal variables, rather than prediction and/or regression problems. Nevertheless, the EM framework employed is related to that in Chib and Greenberg (1998), but due to the different goal, the form of the likelihood function of the proposed model is different from that of the ordinal regression problem. Further, as seen in Section 2, we do not use any MCMC or Gibbs sampling scheme.

The remainder of the paper is organized as follows. Section 2 presents the probit graphical model and discusses algorithmic and model selection issues. Section 3 evaluates the performance of the proposed method on several synthetic examples and Section 4 applies

the model to two data examples, one on movie ratings and the other on a national educational longitudinal survey study.

## 2 Methodology

### 2.1 The probit graphical model

Suppose we have $p$ ordinal random variables $X_1, \ldots, X_p$, where $X_j \in \{1, 2, \ldots, K_j\}$ for some integer $K_j$, which is the number of the ordinal levels in variable $j$. In the proposed probit graphical model, we assume that there exist $p$ latent random variables $Z_1, \ldots, Z_p$ from a joint Gaussian distribution with mean zero and covariance matrix $\Sigma = (\sigma_{j,j'})_{p \times p}$. Without loss of generality, we further assume that $Z_j$'s have unit variances ($\sigma_{j,j} = 1$ for $j = 1, \ldots, p$), i.e., the $Z_j$'s marginally follow standard Gaussian distributions. Each observed variable $X_j$ is discretized from its latent counterpart $Z_j$. Specifically, for the $j$-th variable ($j = 1, \ldots, p$), we assume that $(-\infty, +\infty)$ is split into $K_j$ disjointed intervals by a set of thresholds $-\infty = \theta_0^{(j)} \langle \theta_1^{(j)} \langle \ldots \langle theta_{K_j-1}^{(j)} \langle \theta_{K_j}^{(j)} = +\infty$, such that $X_j = k$ if and only if $Z_j$ falls in the interval $\left[\theta_{k-1}^{(j)}, \theta_k^{(j)}\right)$. Thus,

$$\Pr\left(X_j = k\right) = \Pr\left(\theta_{k-1}^{(j)} \leq Z_j < \theta_k^{(j)}\right) = \Phi\left(\theta_k^{(j)}\right) - \Phi\left(\theta_{k-1}^{(j)}\right), \quad (1)$$

where $\Phi(\cdot)$ denotes the cumulative density function of the standard normal distribution.

Let $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1} = \left(\omega_{\mathbf{j},\mathbf{j}'}\right)_{\mathbf{p} \times \mathbf{p}}$, $\boldsymbol{\Theta} = \left\{\theta_{\mathbf{k}}^{(\mathbf{j})} : \mathbf{j} = \mathbf{1}, \ldots, \mathbf{p}; \mathbf{k} = \mathbf{1}, \ldots, \mathbf{K_j}\right\}$, $X = (X_1, \ldots, X_p)$, $\mathbf{Z} = (Z_1, \ldots, Z_p)$. Let $C(\mathbf{X}, \Theta)$ be the hyper-cube defined by $\left[\theta_{X_1-1}^{(1)}, \theta_{X_1}^{(1)}\right) \times \ldots \times \left[\theta_{X_p-1}^{(p)}, \theta_{X_p}^{(p)}\right)$. Then we can write the joint density function of $(\mathbf{X}, \mathbf{Z})$ as:

$$\mathrm{f}_{\boldsymbol{X},\boldsymbol{Z}}\left(\boldsymbol{x}, \boldsymbol{z}; \Omega, \Theta\right) = \mathrm{f}\left(\boldsymbol{z}; \Omega\right) \prod_{j=1}^{p} \mathrm{f}_{\Theta}\left(x_j | z_j; \Theta\right) = \frac{\det\left(\Omega\right)}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2}\boldsymbol{z}\Omega\boldsymbol{z}^{\mathrm{T}}\right) \mathrm{I}\left(\boldsymbol{z} \in C\left(\boldsymbol{x}, \Theta\right)\right) \quad (2)$$

where $\mathrm{I}(\cdot)$ is the indicator function. Thus, the marginal probability density function of the observed $\boldsymbol{X}$ is given by

$$\mathrm{f}_{\boldsymbol{X}}\left(\boldsymbol{x}; \Omega, \Theta\right) = \int_{z \in \mathbb{R}^p} \mathrm{f}_{\boldsymbol{X},\boldsymbol{Z}}\left(\boldsymbol{x}, \boldsymbol{z}; \Omega, \Theta\right) dz \quad (3)$$

We refer to (1)–(3) as the *probit graphical model*, which is motivated by the probit regression model (Bliss, 1935; Albert and Chib, 1993; Chib and Greenberg, 1998) and the polychotomous Rasch model (von Davier and Carstensen, 2010).

To fit the probit graphical model, we propose maximizing an $l_1$-regularized log-likelihood of the observed data. Let $x_{i,j}$ and $z_{i,j}$ be the $i$-th realizations of the observed variable $X_j$ and the latent variable $Z_j$, respectively, with $\boldsymbol{x}_i = (x_{i,1}, \ldots, x_{i,p})$ and $\boldsymbol{z}_i = (z_{i,1}, \ldots, z_{i,p})$. The criterion is given by

$$\sum_{i=1}^{n} log \, \mathrm{f}_{\boldsymbol{X}} \, (\boldsymbol{x}_i; \Omega, \Theta) - \lambda \sum_{j \neq j'} |\omega_{j,j'}|. \quad (4)$$

The tuning parameter $\lambda$ in (4) controls the degree of sparsity in the underlying network. When $\lambda$ is large enough, some $\omega_{j,j'}$'s can be shrunken to zero, resulting in the removal of the corresponding links in the underlying network. Numerically, it is difficult to maximize criterion (4) directly, because of the integral in (3). Next, we introduce an EM-type algorithm to maximize (4) in an iterative manner.

### 2.2 An algorithm for fitting the probit graphical model

Criterion (4) depends on the parameters $\Theta$ and $\Omega$ and the latent variable $\boldsymbol{Z}$. The former has a closed-form estimator. Specifically, for each $j = 1, \dots, p$, we set

$$\hat{\theta}_k^{(j)} = \begin{cases} -\infty, & \text{if } k=0; \\ \Phi^{-1} \left( n^{-1} \sum_{i=1}^{n} \mathrm{I} \left( x_{i,j} < k \right) \right), & \text{if } k=1, \dots, K_j - 1; \\ +\infty, & \text{if } k=K_j. \end{cases} \quad (5)$$

where $\Phi$ is the cumulative distribution function of the standard normal. One can show that $\hat{\Theta}$ consistently estimates $\Theta$. The estimation of $\Omega$, on the other hand, is nontrivial due to the multiple integrals in (3). To address this problem, we apply the EM algorithm to optimizing (4), where the latent variables $z_{i,j}$'s ($i = 1, \dots, n; j = 1, \dots, p$) are treated as "missing data" and are imputed in the E-step, and the parameter $\Omega$ is estimated in the M-step.

**E-step**. Suppose $\hat{\Omega}$ is the updated estimate of $\Omega$ from the M-step. Then the E-step computes the conditional expectation of the joint log-likelihood given the estimates $\hat{\Theta}$ and $\hat{\Omega}$, which is usually called the *Q*-function in the literature:

$$\mathrm{Q} \left( \Theta, \Omega \right) = \sum_{i=1}^{n} \mathrm{E}_{\boldsymbol{Z}} \left[ log \, \mathrm{f}_{\boldsymbol{X}, \boldsymbol{Z}} \left( \boldsymbol{x}_i, \boldsymbol{Z}; \hat{\Theta}, \hat{\Omega} \right) \right] = \frac{n}{2} \left[ log \det \left( \Omega \right) - \mathrm{trace} \left( \boldsymbol{S} \Omega \right) - p \, log \left( 2\pi \right) \right] \quad (6)$$

Here $\boldsymbol{S}$ is a $p \times p$ matrix whose $(j, j')$-th element is $s_{j,j'} = n^{-1} \sum_{i=1}^{n} \mathrm{E} \left( z_{i,j} z_{i,j'} \mid \boldsymbol{x}_i; \hat{\Theta}, \hat{\Omega} \right) \, \left( 1 \leq j, j' \leq p \right)$. The distribution of $\boldsymbol{z}_i$ conditional on $\boldsymbol{x}_i$ is equal to that of $\boldsymbol{z}_i$ conditional on $\boldsymbol{z}_i \in C(\boldsymbol{x}_i, \Theta)$, which follows a truncated multivariate Gaussian distribution on the hyper-cube $C(\boldsymbol{x}_i, \Theta)$. Therefore, $\mathrm{E} \left( z_{i,j} z_{i,j'} \mid \boldsymbol{x}_i; \hat{\Theta}, \hat{\Omega} \right)$ is the second moment of a truncated multivariate Gaussian distribution and it can be directly estimated using the algorithms proposed by Tallis (1961), Lee (1979), Leppard and Tallis (1989) and Manjunath and Wilhelm (2012). Nevertheless, the computational cost of these direct estimation algorithms is extremely high and thus not suitable for even moderate size problems. An alternative approach is based on the Markov-chain-Monte-Carlo (MCMC) method. Specifically, we randomly generate a sequence of samples from the conditional distribution $\mathrm{f}_{\boldsymbol{Z}|\boldsymbol{X}} \left( \boldsymbol{z}_i \mid \boldsymbol{x}_i; \hat{\Theta}, \hat{\Omega} \right)$ using a Gibbs sampler from a multivariate truncated normal distribution (Kotecha and Djuric, 1999) and then $\mathrm{E} \left( z_{i,j} z_{i,j'} \mid \boldsymbol{x}_i; \hat{\Theta}, \hat{\Omega} \right)$ is estimated by the

empirical conditional second moment from these samples. Although the MCMC approach is faster than the direct estimation method, it is still not efficient for large scale problems. To address this computational issue, we develop an efficient approximate estimation algorithm, discussed in Section 2.3.

**M-step.** The M-step updates $\Omega$ by maximizing the $l_1$-regularized $Q$-function (up to a constant and a factor):

$$\tilde{\Omega}=\arg \max_{\Omega} \quad log \det (\Omega) - \text{trace}\,(\boldsymbol{S}\Omega) - \lambda \sum_{j \neq j'} |\omega_{j,j'}|. \quad (7)$$

The optimization problem (7) can be solved efficiently by existing algorithms such as the graphical lasso (Friedman et al., 2008) and SPICE (Rothman et al., 2008). However, the estimated covariance matrix, $\widetilde{\boldsymbol{\Sigma}}=\widetilde{\boldsymbol{\Omega}}^{-1}$, does not necessarily have unit diagonal elements postulated by the probit graphical model. Therefore, we post-process $\widetilde{\boldsymbol{\Sigma}}$ by scaling it to a unit-diagonal matrix $\hat{\boldsymbol{\Sigma}}$ and update $\hat{\boldsymbol{\Omega}}=\hat{\boldsymbol{\Sigma}}^{-1}$, which will be used in the E-step of the next iteration.

### 2.3 Approximating the conditional expectation

Note that when $j = j'$, the corresponding conditional expectation is the conditional second moment $\text{E}\left(z_{i,j}^2 \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)$; when $j \neq j'$, we use a mean field theory approach (Peterson and Anderson, 1987) to approximate it as

$E\left(z_{i,j}z_{i,j'} \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right) \approx E\left(z_{i,j} \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right) E\left(z_{i,j'} \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)$. Note that the approximation decouples the "interaction" between the two variables $z_{i,j}$ and $z_{i,j}'$. Therefore, one would expect that the approximation performs well when $z_j$ and $z_j'$ are close to independence given all other random variables, which often holds when $\Omega$ or the corresponding graph is sparse.

With this approximation, it is sufficient to estimate the first moment $E\left(z_{i,j} \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)$ and the second moment $E\left(z_{ij}^2 \mid \boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)$. In general, the latent variable $z_{i,j}$ not only depends on $x_{i,j}$, but also on all other observed variables $\boldsymbol{x}_{i,-j} = (x_{i,1}, \ldots, x_{i,j-1}, x_{i,j+1}, \ldots, x_{i,p})$. We can write the first and second conditional moments as

$$E\left(z_{i,j}|\boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)=E\left[ E\left(z_{i,j}|\boldsymbol{z}_{i,-j},x_{i,j};\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right) |\boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right], \quad (8)$$

$$E\left(z_{i,j}^2|\boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right)=E\left[ E\left(z_{i,j}^2|\boldsymbol{z}_{i,-j},x_{i,j};\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right) |\boldsymbol{x}_i;\hat{\boldsymbol{\Theta}},\hat{\boldsymbol{\Omega}}\right], \quad (9)$$

where $\boldsymbol{z}_{i,-j} = (z_{i,1}, \ldots, z_{i,j-1}, z_{i,j+1}, \ldots, z_{i,p})$. The inner expectations in (8) and (9) are relatively straightforward to compute: given the parameter estimate $\hat{\boldsymbol{\Omega}}$, $z_{i,1}, \ldots, z_{i,p}$ jointly follow a multivariate Gaussian distribution with mean zero and covariance matrix $\hat{\boldsymbol{\Sigma}}=\hat{\boldsymbol{\Omega}}^{-1}$. A property of the Gaussian distribution is that the conditional distribution of $z_{i,j}$ given $\boldsymbol{z}_{i,-j}$ is also Gaussian, with mean $\widetilde{\boldsymbol{\mu}}_{i,j}=\hat{\boldsymbol{\Sigma}}_{j,-j}\hat{\boldsymbol{\Sigma}}_{-j,-j}^{-1}\boldsymbol{z}_{i,-j}^{\top}$ and variance $\widetilde{\boldsymbol{\sigma}}_{i,j}^2=1 - \hat{\boldsymbol{\Sigma}}_{j,-j}\hat{\boldsymbol{\Sigma}}_{-j,-j}^{-1}\hat{\boldsymbol{\Sigma}}_{j,-j}$. Moreover, given the observed data $x_{i,j}$, conditioning $z_{i,j}$ on

$z_{i,-j}$, $x_{i,j}$ in (8) is equivalent to conditioning on $z_{i,-j}, \theta^{(j)}_{x_{i,j}-1} \leq z_{i,j} \left\langle \theta^{(j)}_{x_{i,j}} \right.$, which follows a truncated Gaussian distribution on the interval $\theta^{(j)}_{x_{i,j}-1}, \theta^{(j)}_{x_{i,j}}$. The following lemma gives the closed-form expressions for the first and second moments of the truncated Gaussian distribution.

**Lemma 1** Suppose that a random variable Y follows the Gaussian distribution with mean $\mu_0$ and variance $\sigma_0^2$. For any constants $t_1 < t_2$, let $\xi_1 = (t_1 - \mu_0)/\sigma_0$ and $\xi_2 = (t_2 - \mu_0)/\sigma_0$. Then the first and second moments of Y truncated to the interval $(t_1, t_2)$ are given by

$$E\left(Y | t_1 < Y < t_2\right) = \mu_0 + \frac{\phi(\xi_1) - \phi(\xi_2)}{\Phi(\xi_2) - \Phi(\xi_1)}\sigma_0 \quad (10)$$

$$E\left(Y^2 | t_1 < Y < t_2\right) = \mu_0^2 + \sigma_0^2 + 2\frac{\phi(\xi_1) - \phi(\xi_2)}{\Phi(\xi_2) - \Phi(\xi_1)}\mu_0\sigma_0 + \frac{\xi_1\phi(\xi_1) - \xi_2\phi(\xi_2)}{\Phi(\xi_2) - \Phi(\xi_1)}\sigma_0^2 \quad (11)$$

where $\phi(\cdot)$ is the probability density function of the standard normal. For more properties of the truncated Gaussian distribution, see Johnson et al. (1994).

Letting $\delta_{i,j,k} = \left(\theta^{(j)}_k - \tilde{\mu}_{i,j}\right)/\tilde{\sigma}_{i,j}$ and applying Lemma 1 to the conditional expectations in (8) and (9), we obtain

$$E\left(z_{i,j} | z_{i,-j}, x_{i,j}; \hat{\Theta}, \hat{\Omega}\right) = \tilde{\mu}_{i,j} + a_{i,j}\tilde{\sigma}_{i,j}, \quad (12)$$

$$E\left(z_{i,j}^2 | z_{i,-j}, x_{i,j}; \hat{\Theta}, \hat{\Omega}\right) = \tilde{\mu}_{i,j}^2 + \tilde{\sigma}_{i,j}^2 + 2a_{i,j}\tilde{\mu}_{i,j}\tilde{\sigma}_{i,j} + b_{i,j}\tilde{\sigma}_{i,j}^2. \quad (13)$$

where

$$a_{i,j} = \frac{\phi\left(\delta_{i,j,x_{i,j}-1}\right) - \phi\left(\delta_{i,j,x_{i,j}}\right)}{\Phi\left(\delta_{i,j,x_{i,j}}\right) - \Phi\left(\delta_{i,j,x_{i,j}-1}\right)}, \quad b_{i,j} = \frac{\delta_{i,j,x_{i,j}-1}\phi\left(\delta_{i,j,x_{i,j}-1}\right) - \delta_{i,j,x_{i,j}}\phi\left(\delta_{i,j,x_{i,j}}\right)}{\Phi\left(\delta_{i,j,x_{i,j}}\right) - \Phi\left(\delta_{i,j,x_{i,j}-1}\right)}.$$

Next, we plug equations (12) and (13) into (8) and (9), respectively. Since $\tilde{\mu}_{i,j}$, $a_{i,j}$ and $b_{i,j}$ depend on the latent variables $z_{i,-j}$'s, the outer expectations in (8) and (9) depend on $E\left(\tilde{\mu}_{i,j} | x_i; \hat{\Theta}, \hat{\Omega}\right)$, $E\left(a_{i,j} | x_i; \hat{\Theta}, \hat{\Omega}\right)$, $E\left(b_{i,j} | x_i; \hat{\Theta}, \hat{\Omega}\right)$ and $E\left(a_{i,j}\tilde{\mu}_{i,j} | x_i; \hat{\Theta}, \hat{\Omega}\right)$. Note that $\tilde{\mu}_{i,j}$ is a linear function of $z_{i,-j}$ and $\tilde{\sigma}_{i,j}$ is a constant irrelevant to the latent data. For each $i = 1, \ldots, n$ and $j = 1, \ldots, p$, the conditional expectation of $\tilde{\mu}_{i,j}$ is

$$E\left(\tilde{\mu}_{i,j} | x_i; \hat{\theta}, \hat{\Omega}\right) = \hat{\Sigma}_{j,-j}\hat{\Sigma}^{-1}_{-j,-j}E\left(z^{\mathrm{T}}_{i,-j} | x_i; \hat{\theta}, \hat{\Omega}\right). \quad (14)$$

However, $a_{i,j}$ and $b_{i,j}$ are nonlinear functions of $\tilde{\mu}_{i,j}$ and thus of $z_{i,-j}$. Using the first order delta method, we approximate their conditional expectations by

$$E\left(a_{i,j}|x_i;\hat{\theta},\hat{\Omega}\right) \approx \frac{\phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right) - \phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right)}{\Phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right) - \Phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right)} \quad (15)$$

$$E\left(b_{i,j}|x_i;\hat{\theta},\hat{\Omega}\right) \approx \frac{\tilde{\delta}_{i,j,x_{i,j}-1}\phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right) - \tilde{\delta}_{i,j,x_{i,j}}\phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right)}{\Phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right) - \Phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right)} \quad (16)$$

where $\tilde{\sigma}_{i,j,x_{i,j}} = \left[\theta_k^{(j)} - E\left(\tilde{\mu}_{i,j} \mid \boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right)\right]/\tilde{\sigma}_{i,j}$. Finally, we approximate $E\left(a_{i,j}\tilde{\mu}_{i,j} \mid \boldsymbol{x}_i;\hat{\theta},\hat{\Omega}\right) \approx E\left(a_{i,j} \mid \boldsymbol{x}_i;\hat{\theta},\hat{\Omega}\right)E\left(\tilde{\mu}_{i,j} \mid \boldsymbol{x}_i;\hat{\theta},\hat{\Omega}\right)$. Therefore (8) and (9) can be approximated by

$$E\left(z_{i,j}|x_i;\hat{\Theta},\hat{\Omega}\right) \approx \hat{\Sigma}_{j,-j}\hat{\Sigma}_{-j,-j}^{-1}E\left(z_{i,-j}^{\mathrm{T}}|x_i;\hat{\Theta},\hat{\Omega} + \frac{\phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right) - \phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right)}{\Phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right) - \Phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right)}\tilde{\sigma}_{i,j}\right) \quad (17)$$

$$\begin{aligned} E\left(z_{i,j}^2|x_i;\hat{\Theta},\hat{\Omega}\right) &\approx \hat{\Sigma}_{j,-j}\hat{\Sigma}_{-j,-j}^{-1}E\left(z_{i,-j}^{\mathrm{T}}z_{i,-j}|x_i;\hat{\Theta},\hat{\Omega}\right)\hat{\Sigma}_{-j,-j}^{-1}\hat{\Sigma}_{j,-j}^{\mathrm{T}} + \tilde{\sigma}_{i,j}^2 \\ &+ 2\frac{\phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right) - \phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right)}{\Phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right) - \Phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right)}\left[\hat{\Sigma}_{j,-j}\hat{\Sigma}_{-j,-j}^{-1}E\left(z_{i,-j}^{\mathrm{T}}|x_i;\hat{\Theta},\hat{\Omega}\right)\right]\tilde{\sigma}_{i,j} \\ &+ \frac{\tilde{\delta}_{i,j,x_{i,j}-1}^{(j)}\phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right) - \tilde{\delta}_{i,j,x_{i,j}}\phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right)}{\Phi\left(\tilde{\delta}_{i,j,x_{i,j}}\right) - \Phi\left(\tilde{\delta}_{i,j,x_{i,j}-1}\right)}\tilde{\sigma}_{i,j}^2 \end{aligned} \quad (18)$$

Equations (17) and (18) establish the recursive relationships among the elements in $E\left(\boldsymbol{z}_i|\boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right)$ and $E\left(\boldsymbol{z_i}^{\mathsf{T}}\boldsymbol{z_i}|\boldsymbol{x_i};\hat{\Theta},\hat{\Omega}\right)$, respectively, giving a natural iterative procedure for estimating these quantities. Algorithm 1 summarizes the main steps of the proposed combined estimation procedure outlined in Sections 2.2 and 2.3.

---

**Algorithm 1** The EM Algorithm for estimating $\boldsymbol{\Omega}$

---

1: Initialize $E\left(z_{i,j} \mid \boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right) \approx E\left(z_{i,j} \mid x_{i,j};\hat{\Theta}\right)$, $E\left(z_{i,j}^2 \mid \boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right) \approx E\left(z_{i,j}^2 \mid x_{i,j};\hat{\Theta}\right)$ and $E\left(z_{i,j}z_{i,j'} \mid \boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right) \approx E\left(z_{i,j} \mid x_{i,j};\hat{\Theta}\right)E\left(z_{i,j'} \mid x_{i,j'};\hat{\Theta}\right)$ for $i = 1, ..., n$ and $j,j' = 1, ..., p$;

2: Initialize $s_{j,j'}$ for $1 \leq j,j' \leq p$ using the Line 1 above, and then estimate $\hat{\boldsymbol{\Omega}}$ by maximizing criterion (7); {Start outer loop}

3: **repeat**

4:     E-step: estimate $S$ in (6); {Start inner loop}

5:     **repeat**

6:         **for** $i = 1$ to n **do**

7:             **if** $j = j'$ **then**

8:                 Update $E\left(z_{i,j}^2 \mid \boldsymbol{x}_i;\hat{\Theta},\hat{\Omega}\right)$ using RHS of equation (18) for $j = 1, ..., p$;

9:             **else**

---

10:

Update $\text{E}\left(z_{i,j} \mid \boldsymbol{x}_i; \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\Omega}}\right)$ using RHS of equation (17) for $j = 1, \dots, p$ and then set
$$\text{E}\left(z_{i,j} z_{i,j'} \mid \boldsymbol{x}_i; \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\Omega}}\right) = \text{E}\left(z_{i,j} \mid \boldsymbol{x}_i; \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\Omega}}\right) \text{E}\left(z_{i,j'} \mid \boldsymbol{x}_i; \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\Omega}}\right)$$ for $1 \leq j \leq j' \leq p$;

11: 	**end if**

12: 	**end for**

13:

Update $s_{j,j'} = 1/n \sum_{i=1}^{n} \text{E}\left(z_{i,j} z_{i,j'} \mid \boldsymbol{x}_i; \hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\Omega}}\right)$ for $1 \leq j,j' \leq p$;

14: 	**until** The inner loop converges;

15:

M-step: update $\hat{\boldsymbol{\Omega}}$ by maximizing criterion (7);

16: 	**until** The outer loop converges.

In Algorithm 1, Lines 1–2 initialize the conditional expectation $E(z_{i,j} \mid \boldsymbol{x}_i)$ and the parameter estimate $\hat{\boldsymbol{\Omega}}$. Lines 3–16 establish the outer loop which iteratively computes the E-step and the M-step. In the E-step, Lines 5–14 consist of the inner loop which recursively estimates the first and second moments of $z_{i,j}$ conditional on $\boldsymbol{x}_i$. The complexity of the inner loop is $O(np^2)$, which is the same as that of the graphical lasso algorithm in the M-step. Therefore, the overall complexity of Algorithm 1 is $O(Mnp^2)$, where $M$ is the number of EM steps required for convergence. In our numerical studies, we found $M$ is often smaller than 50. For a more concrete idea about the computational cost, we note that on a linux server with four 1G Dual-Core AMD Opteron Processors and 4GB RAM, it takes about 2 minutes for the proposed algorithm to complete the fitting on a simulated dataset in Section 3 with $n = 200$ observations and $p = 50$ variables.

## 2.4 Model selection

In the probit graphical model, the tuning parameter $\lambda$ controls the sparsity of the resulting estimator and it can be selected using cross-validation. Specifically, we randomly split the observed data $\boldsymbol{X}$ into $D$ subsets of similar sizes and denote the index set of the observations in the $d$-th subset by $\mathscr{T}_d \, (d = 1, \dots, D)$. For any pre-specified $\lambda$, we denote by $\hat{\Omega}_\lambda^{[-d]}$ the maximizer of the criterion (4) estimated by Algorithm 1 using all observations except those in $\mathscr{T}_d$. We also denote by $\hat{\boldsymbol{\Theta}}^{[-d]}$ and $\mathbf{S}^{[d]} = \left(s_{j,j'}^{[d]}\right)_{p \times p}$ the analogs of $\hat{\boldsymbol{\Theta}}$ $\boldsymbol{S}$ in Section 2.2, but computed from the data in $\mathscr{T}_d^c$ and $\mathscr{T}_d$, respectively. In particular, an element of $\boldsymbol{S}^{[d]}$ is defined as $s_{j,j'}^{[d]} = |\mathscr{T}_d|^{-1} \sum_{i \in \mathscr{T}_d} \text{E}\left(z_i, z_{i,j'} | \boldsymbol{x_i}; \hat{\boldsymbol{\Theta}}^{[-d]}, \hat{\boldsymbol{\Omega}}_\lambda^{[-d]}\right)$, for $1 \leq j, j' \leq p$, where $|\mathscr{T}_d|$ is the cardinality of $\mathscr{T}_d$. Given $\hat{\boldsymbol{\Theta}}^{[-d]}$ and $\hat{\boldsymbol{\Theta}}_\lambda^{[-d]}$, $\boldsymbol{S}^{[d]}$ can be estimated by the algorithm introduced in Section 2.3, i.e., the inner loop of Algorithm 1. Thus, the optimal tuning parameter can be selected by maximizing the following criterion:

$$\max_\lambda \sum_{d=1}^{D} log \det\left(\hat{\Omega}_\lambda^{[-d]}\right) - \text{trace}\left(S^{[d]} \hat{\Omega}_\lambda^{[-d]}\right) - p \, log\,(2\pi). \quad (19)$$

We note that we have also considered the AIC and BIC type criteria for choosing the tuning parameter λ. We found that AIC performs the worst among the three due to estimating many zero parameters as non-zeros (Lian, 2011); BIC and cross-validation tend to have similar performances in estimating zero parameters as zeros, but BIC also tends to estimate the non-zero parameters as zeros. Therefore, we choose to use cross-validation. Due to space limitation, the results are not included.

## 3 Numerical Examples

In this section, we use two sets of simulated experiments to illustrate the performance of the probit graphical model. The first set aims at comparing the computational cost of the three methods that estimate the $Q$-function in the E-step; namely the direct computation, the MCMC sampling and the approximation algorithm. The second set compares the performance of the probit graphical model using the approximation algorithm to that of the Gaussian graphical model.

### 3.1 Computational cost and performance

Note that the computational costs of the direct estimation and the MCMC sampling are extremely high when $p$ is even of moderate size. Therefore, in this experiment, we simulate a low-dimensional data set with $p = 5$ variables and $n = 10$ observations. Specifically, we define the underlying inverse covariance matrix $\Omega$ as a tri-diagonal matrix with 1's on the main diagonal and 0.5 on the first sub-diagonal. The corresponding covariance matrix is then scaled so that all the variances are equal to 1. Then, for $i = 1, \ldots, n$, we generate the latent data $z_i = (z_{i,1}, \ldots, z_{i,p})$ from $N(\mathbf{0}, \Sigma)$ and discretize them as follows: for each $j = 1, \ldots, p$, set

$$\theta_k^{(j)\theta} = \begin{cases} -\infty, & \text{if } k=0; \\ \Phi^{-1}(0.2) & \text{if } k=1; \\ \Phi^{-1}(0.4) & \text{if } k=2; \\ +\infty, & \text{if } k=3. \end{cases} \quad (20)$$

and $x_{i,j} = \sum_{k=0}^{2} I\left(z_{i,j} \geq \theta_k^{(j)}\right) (i = 1, \ldots, n; j, \ldots, p)$, i.e., the value of $x_{i,j}$ is $k$ if it locates in interval $\left[\theta_{k-1}^{(j)}, \theta_k^{(j)}\right)$.

The probit graphical model is applied using four estimation methods in the E-step, namely the direct computation, a standard Gibbs sampling, the Gibbs sampler proposed by Pakman and Paninski (2012) and the approximation algorithm proposed in this manuscript. The procedure is repeated for 20 times, and the computational costs are shown in Table 1. We can see that the median CPU time of the approximation algorithm is only about 1/1,000 of that of the Gibbs sampling and about 1/80,000 of that of the direct computation. To further compare the estimation accuracy of these methods, we use the Frobenius and entropy loss metrics that are defined next:

$$FL = \frac{\sum_{1 \le j < j' \le p} \left( \omega_{j,j'} - \hat{\omega}_{j,j'} \right)^2}{\sum_{1 \le j < j' \le p} \omega_{j,j'}^2} \quad (21)$$

$$EL = \text{trace}\left(\Omega^{-1}\hat{\Omega}\right) - log\left[ \det\left(\Omega^{-1}\hat{\Omega}\right) \right] - p \quad (22)$$

where $\hat{\Omega}$ denotes the estimated network.

The performance of the three estimation methods is depicted in Figure 1. It can be seen that the direct computation and Gibbs sampling methods are fairly similar in performance (The result using the R package "tmg" is almost identical to that of the standard Gibbs sampling and not shown); this is expected since they can all be considered "exact" approaches. In terms of the Frobenius and entropy losses, the approximation algorithm lags slightly behind its competitors when the tuning parameter $\lambda$ is relatively small, whereas for larger $\lambda$ it outperforms them. This is due to the fact that in this simulation study, the true $\Omega$ is very sparse and the mean field approximation also happens to implicitly enforce a conditional independence structure on the $S$ matrix. These findings suggest that the proposed approximation algorithm achieves its orders of magnitude computational savings over the competitors with minimal degradation in performance.

## 3.2 Experiments with different types of graphs

In this section, we evaluate the performance of the proposed method by simulation studies. These examples simulate four types of network structures: a scale-free graph, a hub graph, a nearest-neighbor graph and a block graph. Each network consists of $p = 50$ nodes. The details of these networks are described as follows:

**Example 1** **Scale-free graph.** A scale-free graph has a power-law degree distribution and can be simulated by the Barabasi-Albert algorithm (Barabasi and Albert, 1999). A realization of a scale-free network is depicted in Figure 2 (A).

**Example 2** **Hub graph.** A hub graph consists of a few high-degree nodes (hubs) and a large amount of low-degree nodes. In this example, we follow the simulation setting in Peng et al. (2009) and generate a hub graph by inserting a few hub nodes into a very sparse graph. Specifically, the graph consists of three hubs with degrees around eight, and the other 47 nodes with degrees at most three. An example of the hub graph is shown in Figure 2 (B).

**Example 3** **Nearest-neighbor graph.** To generate nearest neighbor graphs, we slightly modify the data generating mechanism described in Li and Gui (2006). Specifically, we generate $p$ points randomly on a unit square, calculate all $p(p-1)/2$ pairwise distances, and find the $m$ nearest neighbors of each point in terms of these distances. The nearest neighbor network is obtained by

linking any two points that are *m*-nearest neighbors of each other. The integer *m* controls the degree of sparsity of the network and the value $m = 5$ was chosen in the simulation study. Figure 2 (C) exhibits one realization of the nearest-neighbor network.

**Example 4** **Block graph.** In this setting, we generate a graph using a random adjacency matrix generated from the stochastic block model. Specifically, for nodes 1–20 the probability of being linked is 0.2, for nodes 21–30 the probability of being linked is 0.5, whereas for all other pairs of nodes the probability of having a link is 0.02. Figure 2 (D) illustrates such a random graph.

The ordinal data are generated as follows. First, we generate the inverse covariance matrix $\Omega$ of the latent multivariate Gaussian distribution. Specifically, each off-diagonal element $\omega_{j,j'}$ is drawn uniformly from $[-1, -0.5] \bigcup [0.5, 1]$ if nodes $j$ and $j'$ are linked by an edge, otherwise $\omega_{j,j'} = 0$. Further, the diagonal elements were all set to be 2 to ensure positive definiteness, and the corresponding covariance matrix is scaled so that all the variances are equal to 1. Second, we generate the latent data $z_i = (z_{i,1}, \ldots, z_{i,p})$ as an i.i.d. sample from $N(\mathbf{0}, \Sigma)$. Finally, the continuous latent data $z_i$'s are discretized into ordinal scale with three levels by thresholding. Specifically, for each $j = 1, \ldots, p$, we set

$$\theta_k^{(j)} = \begin{cases} -\infty, & \text{if } k=0; \\ \Phi^{-1}(0.1) & \text{if } k=1; \\ \Phi^{-1}(0.2) & \text{if } k=2; \\ +\infty, & \text{if } k=3. \end{cases} \quad (23)$$

and set $x_{i,j} = \sum_{k=0}^{2} I\left(z_{i,j} \geq \theta_k^{(j)}\right)$ $(i = 1, \ldots, n; j = 1, \ldots, p)$. For each example, we considered different sample sizes, with $n$=50, 100, 200 and 500.

We compare the proposed probit graphical model with two other methods. One consists of direct application of the graphical lasso to the ordinal data $X$, ignoring their discrete nature. The second uses the graphical lasso on the latent continuous data $Z$. We refer to the first one as the naive method and the second one as an oracle method because it represents an ideal situation where $Z$ is exactly recovered. Of course, the latter never occurs with real data, but serves as a benchmark for comparison purposes. The receiver operating characteristic curve (ROC) was used to evaluate the accuracy of network structure estimation. The ROC curve plots the sensitivity (the proportion of correctly detected links) against the false positive rate (the proportion of mis-identified zeros) over a range of values of the tuning parameter $\lambda$. The sensitivity and the false positive rate are defined as follows:

$$\text{Sensitivity} = \frac{\sum\limits_{1 \leq j < j' \leq p} \mathscr{I}\left(\omega_{j,j'} \neq 0, \quad \hat{\omega}_{j,j'} \neq 0\right)}{\sum\limits_{1 \leq j < j' \leq p} \mathscr{I}\left(\omega_{i,j'} \neq 0\right)} \quad (24)$$

$$\text{False Positive Rate} = \frac{\sum\limits_{1 \le j < j' \le p} \mathscr{I}\left(\omega_{j,j'} = 0, \quad \hat{\omega}_{j,j'} \ne 0\right)}{\sum\limits_{1 \le j < j' \le p} \mathscr{I}\left(\omega_{j,j'} = 0\right)} \quad (25)$$

where $\mathscr{I}(\cdot)$ is an indicator function whose value is one if the statement in the parenthesis is true, and is zero if it is false. In addition, the Frobenius loss and the entropy loss defined in (21) were used to evaluate the performance of parameter estimation.

Figure 3 shows the ROC curves for all simulated examples. The curves are averaged over 50 replications. The oracle method provides a benchmark curve for each setting (blue dotted line in each panel). We can see that when the sample size is relatively small ($n$=50, 100 or 200), the probit model (dark solid line) dominates the naive method (red dashed line). When the sample size gets larger, the two methods exhibit similar performance.

Table 2 summarizes the parameter estimation measured by the Frobenius loss and the entropy loss. The results were again averaged over 50 repetitions and the tuning parameter λ was selected using the cross-validation introduced in Section 2.4. The oracle method evidently performs the best, as it should. Comparing the two methods based on the observed data $X$, we can see that the Frobenius losses from the probit model are consistently lower than those from the naive method. The advantage is more significant when the sample size is moderate ($n$=100 or 200). In terms of the entropy loss, we can see that the probit model outperforms the naive method for relatively large sample sizes, such as $n$=200 and 500.

## 4 Data Examples

### 4.1 Application to movie rating records

In this section, we apply the probit graphical model to Movielens, a data set containing rating scores for 1682 movies by 943 users. The rating scores have five levels, where 1 corresponds to strong dissatisfaction and 5 to strong satisfaction. More than 90% of the entries are missing in the full data matrix; for this reason, we consider a subset of the data containing 193 users and 32 movies, with 15% missing values. The missing values were imputed by the median of the observed movie ratings.

The estimated network for these 32 movies is shown in Figure 4. We can see that the estimated network consists of a large connected community as well as a few isolated nodes. The large community mainly consists of mass marketed commercial movies, dominated by science fiction and action films. These movies are characterized by high production budgets, state of the art visual effects, and famous directors and actors. Examples in this data subset include the Star Wars franchise ("Star Wars" (1977), "The Empire Strikes Back" (1980) and "Return of the Jedi" (1983), directed/produced by Lucas), the Terminator series (1984, 1991) directed by Cameron, the Indiana Jones franchise ("Raiders of Lost Ark" (1981), "The Last Crusade" (1989), directed by Spielberg), the Alien series, etc. As expected, movies within the same series are most strongly associated. Further, "Raiders of the Lost Ark" (1981) and "Back to the Future" (1985) form two hub nodes each having 16 connections to other movies and their common feature is that they were directed/produced by Spielberg.

On the other hand, isolated nodes tend to represent "artsier" movies, such as crime films and comedies whose popularity relies more on the plot and the cast than on big budgets and special effects, many with cult status among their followers. Examples include "Pulp Fiction" (1994) (one of the most popular Tarantino movies), "Fargo" (1996) (a quintessential Coen brothers movie), "When Harry Met Sally" (1989) and "Princess Bride" (1987). These films have no significant connections in the network, either with each other or with the commercial movies in the large community. This is likely due to two reasons: (1) we restricted the dataset to movies rated by a substantial fraction of the users, so while there probably are connections from "Fargo" to other Coen brothers movies, the other ones did not appear in this set; and (2) there is a greater heterogeneity of genres in this set than among the large group of science-fiction and action films. In other words, liking "When Harry Met Sally" (a romantic comedy) does not make one more likely to enjoy "Silence of the Lambs" (a thriller/horror movie), whereas liking "Terminator" suggests you are more likely to enjoy "The Alien". A more complete analysis of this dataset is an interesting topic for future work and requires a more sophisticated way of dealing with missing data, which is not the focus of the current manuscript.

## 4.2 National education longitudinal survey study

The data for the second example come from the National Educational Longitudinal Study of 1988 (NELS:88), whose objective was to assess student attitudes towards a number of questions about their school, education, and activities. The data used were obtained from the study's website http://nces.ed.gov/surveys/nels88/ and correspond to a sample of 12144 students of eighth-graders. We selected 218 questions with ordinal and/or binary responses that focused on diverse issues, including school, work and home experiences, educational and occupational aspirations, access to educational resources and other support, as well as student background and school characteristics. Ordinal responses were chosen from the following options: "OFTEN", "SOMETIMES", "RARELY", and "NEVER", while binary ones corresponded to a "YES/NO" answer. Figure 5 depicts the histogram of the frequency of options in 218 survey questions.

The estimated network of the selected 218 survey questions is shown in Figure 6. It is apparent that the estimated network exhibits a strong clustering structure. For example, the set of the following nodes "F1S33A", "F1S33B", "F1S33C", "F1S33D", and "F1S33E" forms a cluster, separated from the remaining nodes. These five questions are a part of a sequence of similar questions, focusing on vocational coursework. Specifically, the question inquires whether "In your most recent or current VOCATIONAL course, how much emphasis did/does your teacher place on the following objectives?" and the specific objectives are listed in Table 3. It can be seen that questions "F1S33A"–"F1S33E" reflect different aspects of knowledge and analytical ability that a student should acquire from a vocational course, and therefore it is reasonable that they form a tight cluster. Similar clustering patterns can be observed in other parts of the graph, for example, serial "F1S7", serial "F1S8", serial "F1S12", serial "F1S25", etc.

Next, we focus on broad patterns revealed by the model, as depicted in Figure 6. The upper right corner captures relationships between serial questions broadly related to coursework

(F1S22–F1S25) in various disciplines (mathematics, science, English, computer education), whereas in the lower left corner there are questions related to overall attitude and study patterns regarding mathematics and science classes (F1S26–F1S32). It is interesting to observe that the model does not discover any relationships between these two question clusters. In the center of the plot we find questions related to various life aspects and being successful/accomplishing them (F1S46) which is negatively associated with a cluster of questions related to working hard in school for good grades (F1S11). In the center, we also find a cluster of serial questions related to different ways of interacting with friends (F1S44) which is negatively correlated to questions related to students awards (F1S8). In the upper left corner we see the serial cluster on grades performance (F1S39) which is also negatively correlated with some of the questions related to amount of coursework in various subjects (F1S22 and F1S24). Finally, in the bottom right corner we encounter questions related to school attendance and attitude towards it (F1S10, F1S12).

Overall, the model reveals interesting and informative patterns, much more so than its Gaussian counterpart shown in Figure 7.

Next, we examined pairs of questions exhibiting the largest positive partial correlations (based on the theory of Gaussian graphical models, the partial correlation of variables $j$ and $j$′ is defined as $\rho_{j,j'} = -\omega_{j,j'}/\sqrt{\omega_{j,j}\omega_{j',j'}}$). The results are shown in Table 4. Among the top five ones, four pairs correspond to serial questions. The only exception is pair "F1S44D—F1S43", although it inquires about extra reading, outside school. Analogously, Table 5 lists the pair of questions exhibiting the strongest negative partial correlations. Note that question pairs "F1S8F—F1S8A", "F1S15B—F1S15A", "F1S16B—F1S16D" are composed of two opposite questions. It is interesting to observe that the model identifies the pair `F1S10B—F1S12B", which can be interpreted that although students may skip class often they do not feel good about their action. A similar negative partial correlation is present in pair "F1S10A—F1S12A" that addresses a "coming to school late" issue. Overall, the proposed model identifies strong clustering patterns in the questions being asked in this survey, which primarily correspond to series of related in intent and purpose questions, thus indirectly validating its usefulness.

## 5 Summary and Discussion

Ordinal data occur often in practice and are usually treated as continuous for most analyses, including estimating dependencies between the variables under consideration by fitting a graphical model. Our proposed model, explicitly takes into account the ordinal nature of the data in the graphical modeling step. While direct computation for the proposed model is expensive, the approximations employed allow us to efficiently fit high-dimensional models. On those datasets that the model can be fitted directly, our numerical results show that the approximations we make result in a minimal loss of accuracy. We leave the theoretical properties of both the exact estimator and its approximate version as a topic for future work.

The method proposed in this paper can also be extended to fit the multivariate ordinal regression model, where multiple ordinal responses are fitted on a number of covariates. Specifically, suppose $W_{j1}, \ldots, W_{jm_j}$ are the covariates associated with the $j$th response.

Following the notation in Section 2.1, let $X_j$ denote the $j$th response, which is an ordinal variable, and $Z_j$ the corresponding latent continuous variable. We may assume $Z_j = \alpha_{j0} + \alpha_{j1}W_{j1} + \ldots + \alpha_{jm_j}W_{jm_j} + \epsilon_j$, where $\alpha_{j0}$ is the intercept, and $\alpha_{j1}, \ldots, \alpha_{jm_j}$ are regression coefficients. In addition, we assume that $\epsilon_1, \ldots, \epsilon_p$ jointly follow a Gaussian distribution with mean zero and covariance matrix $\Sigma = \Omega^{-1}$. To estimate the regression coefficients, we may modify the M-step in Section 2.2 to estimate $\Omega$ and $\alpha_{j\ell}$'s simultaneously. Rothman et al. (2010) discussed a similar problem as the modified M-step, and the algorithm there can be directly applied.

## Acknowledgments

## References

Albert J, Chib S. Bayesian analysis of binary and polychotomous response data. Journal of the American Statistical Asscociation. 1993; 88:669–679.

Banerjee O, El Ghaoui L, d'Aspremont A. Model selection through sparse maximum likelihood estimation. Journal of Machine Learning Research. 2008; 9:485–516.

Barabasi A-L, Albert R. Emergence of scaling in random networks. Science. 1999; 286:509–512. [PubMed: 10521342]

Bliss C. The calculation of the dosage-mortality curve. Annals of Applied Biology. 1935; 22:134–167.

Chib S, Greenberg E. Analysis of multivariate probit models. Biometrika. 1998; 85:347–361.

Fan J, Feng Y, Wu Y. Network exploration via the adaptive LASSO and SCAD penalties. Annals of Applied Statistics. 2009; 3:521–541. [PubMed: 21643444]

Friedman J, Hastie T, Tibshirani R. Sparse inverse covariance estimation with the graphical lasso. Biostatistics. 2008; 9:432–441. [PubMed: 18079126]

Guo, J.; Levina, E.; Michailidis, G.; Zhu, J. Tech. rep. Department of Statistics, University of Michigan; Ann Arbor: 2010. Joint structure estimation for categorical Markov networks.

Höefling H, Tibshirani R. Estimation of sparse binary pairwise Markov networks using pseudo-likelihoods. Journal of Machine Learning Research. 2009; 10:883–906. [PubMed: 21857799]

Johnson, N.; Kotz, S.; Balakrishnan, N. Continuous Univariate Distributions. 2nd ed. Vol. 1. John Wiley & Sons; New Jersey: 1994.

Kolar, M.; Xing, E. Eprint arXiv:0811.1239. 2008. Improved estimation of high-dimensional Ising models.

Koren Y, Bell R, Volinsky C. Matrix factorization techniques for recommender systems. IEEE Computer. 2009; 42:30–37.

Kotecha J, Djuric P. Gibbs sampling approach for generation of truncated multivariate Gaussian random variables. IEEE Computer Society. 1999; 3:1757–1760.

Lauritzen, S. Graphical Models. Oxford University Press; Oxford: 1996.

Lee L-F. On the first and second moments of the truncated multi-normal distribution and a simple estimator. Economics Letters. 1979; 3:165–169.

Leppard P, Tallis G. Evaluation of the mean and covariance of the truncated multinormal. Applied Statistics. 1989; 38:543–553.

Li H, Gui J. Gradient directed regularization for sparse Gaussian concentration graphs, with applications to inference of genetic networks. Biostatistics. 2006; 7:302–317. [PubMed: 16326758]

Lian H. Shrinkage tuning parameter selection in precision matrices estimation. Journal of Statistical Planning and Inference. 2011; 14:2839–2848.

Manjunath, B.; Wilhelm, S. ArXiv e-prints. 2012. Moments calculation for the doubly truncated multivariate normal density.

McCullagh P. Regression models for ordinal data. Journal of the Royal Statistical Society, Series B. 1980; 42:109–142.

McCullagh, P.; Nelder, J. Generalized Linear Models. 2nd ed. Chapman and Hall/CRC; London, UK: 1989.

Meinshausen N, Buhlmann P. High-dimensional graphs with the lasso. Annals of Statistics. 2006; 34:1436–1462.

O'Connell, A. Logistic Regression Models for Ordinal Response Variables. 1st ed. Sage Publications, Inc; 2005.

Pakman, A.; Paninski, L. ArXiv e-prints. 2012. Exact Hamiltonian monte carlo for truncated multivariate Gaussians.

Peng J, Wang P, Zhou N, Zhu J. Partial correlation estimation by joint sparse regression model. Journal of the American Statistical Asscociation. 2009; 104:735–746.

Peterson B. Partial proportional odds models for ordinal response variables. Applied Statistics. 1990; 39:205–217.

Peterson C, Anderson J. A mean field theory learning algorithm for neural networks. Complex systems. 1987; 1:995–1019.

Ravikumar P, Wainwright M, Lafferty J. High-dimensional Ising model selection using $\ell_1$-regularized logistic regression. Annals of Statistics. 2010; 38:1287–1319.

Rocha, G.; Zhao, P.; Yu, B. Tech. rep. Department of Statistics, University of California; Berkeley: 2008. A path following algorithm for Sparse Pseudo-Likelihood Inverse Covariance Estimation (SPLICE).

Rothman A, Bickel P, Levina E, Zhu J. Sparse permutation invariant covariance estimation. Electronic Journal of Statistics. 2008; 2:494–515.

Rothman A, Levina L, Zhu J. Sparse multivariate regression with covariance estimation. Journal of Computational and Graphical Statistics. 2010; 19:947–962. [PubMed: 24963268]

Stern, D.; Herbrich, R.; Graepel, T. Matchbox: large scale online Bayesian recommendations. Proceedings of Wourd Wide Web 2009; Madrid, Spain. 2009. p. 111-120.

Tallis G. The moment generating function of the truncated multinormal distribution. Journal of the Royal Statistical Society, Series B. 1961; 23:223–229.

von Davier, M.; Carstensen, C. Multivariate and Mixture Distribution Rasch Models: Extensions and Applications. 1st ed. Springer; New York: 2010.

Walker S, Duncan D. Estimation of the probability of an event as a function of several independent variables. Biometrika. 1967; 54:167–179. [PubMed: 6049533]

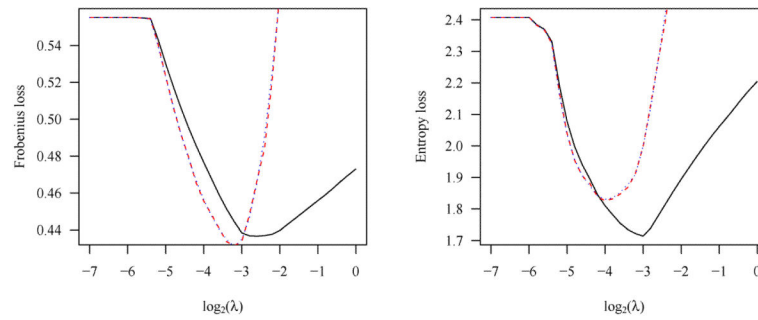Yuan M, Lin Y. Model selection and estimation in the Gaussian graphical model. Biometrika. 2007; 94:19–35.

**Figure 1.**
Comparison of Frobenius loss and Entropy loss over different values of the tuning parameter. The direct computation, the MCMC sampling and the approximation algorithm are respectively represented by blue dotted, red dashed and black solid lines.
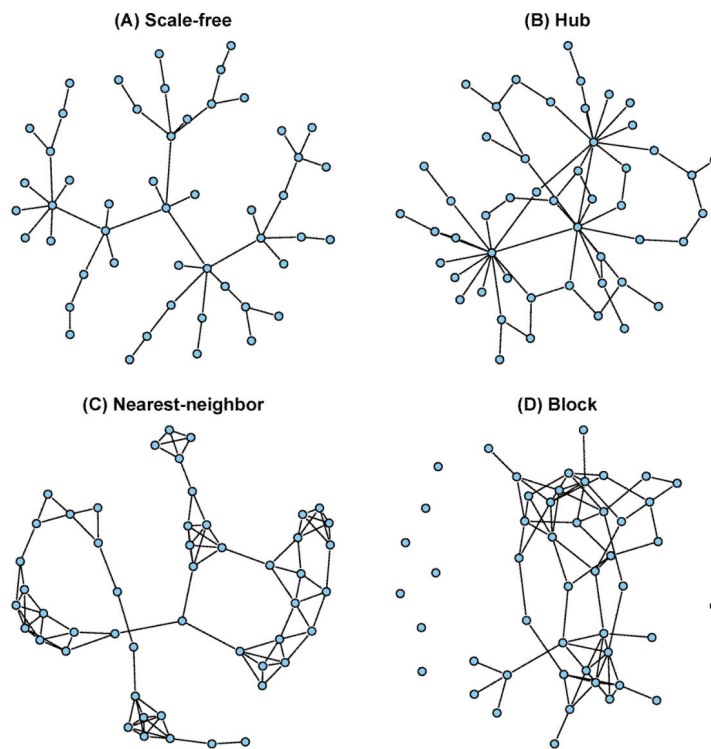
**Figure 2.**
Illustration of the networks used in four simulated examples: scale-free graph, hub graph, nearest-neighbor graph and block graph.
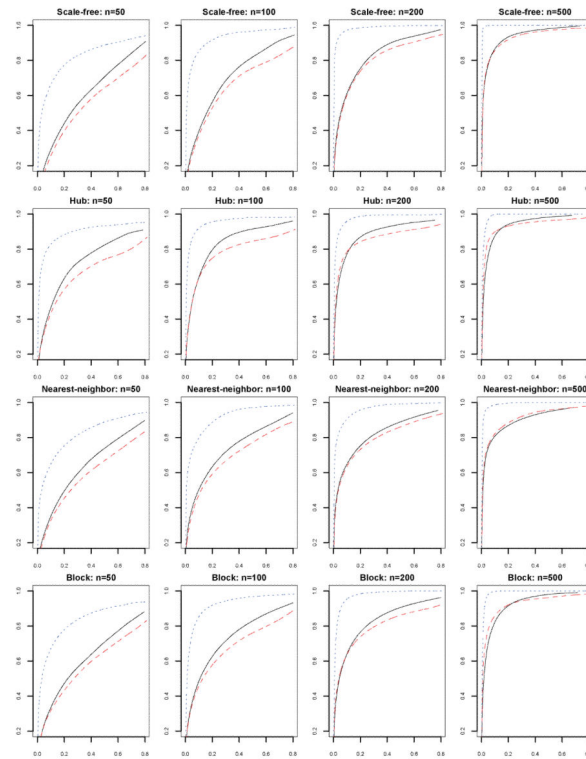
**Figure 3.**
The ROC curves estimated by the probit graphical model (solid dark line), the oracle method (dotted blue line) and the naive method (dashed red line). The oracle method and the naive method simply apply the graphical lasso algorithm to the latent continuous data $Z$ and the observed discrete data $X$, respectively.

**Figure 4.**
The network estimated by the probit graphical model. The nodes represent the movies labeled by their titles. The area of a node is proportional to its degree and the width of a link is proportional to the magnitude of the corresponding partial correlations.
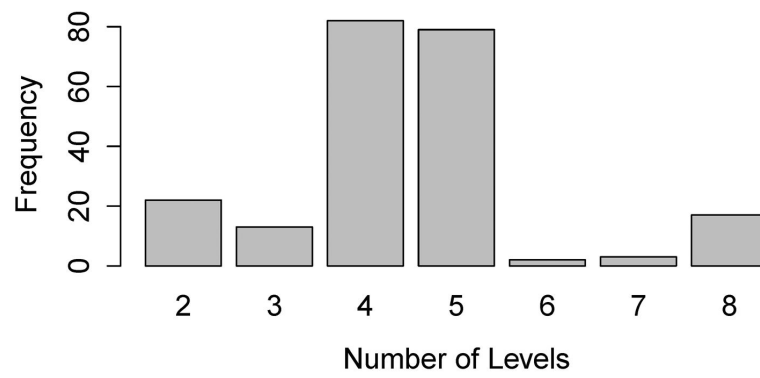
**Figure 5.**
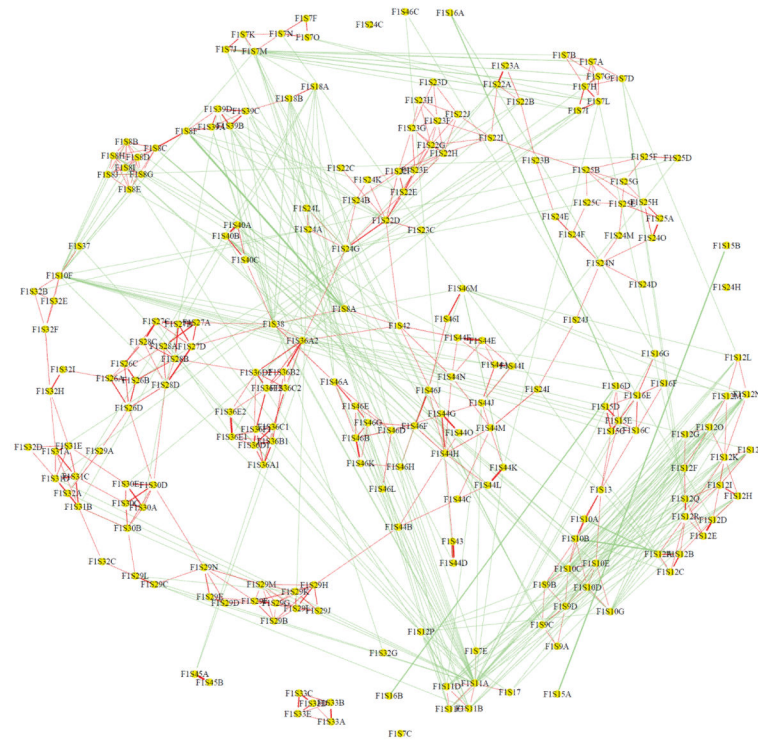Histogram of the number of options in 218 survey questions.

**Figure 6.**
Layout of the network estimated by the proposed probit graphical model. The nodes represent the survey questions labeled by their code. The area of a node is proportional to its degree and the width of a link is proportional to the magnitude of the corresponding partial correlations. The red lines represent positive associations, while the light green lines negative ones.

**Figure 7.**
Layout of the estimated network by the graphical lasso algorithm. The nodes represent the survey questions labeled by their code. The area of a node is proportional to its degree and the width of a link is proportional to the magnitude of the corresponding partial correlations. The red lines represent positive associations, while the light green lines negative ones.
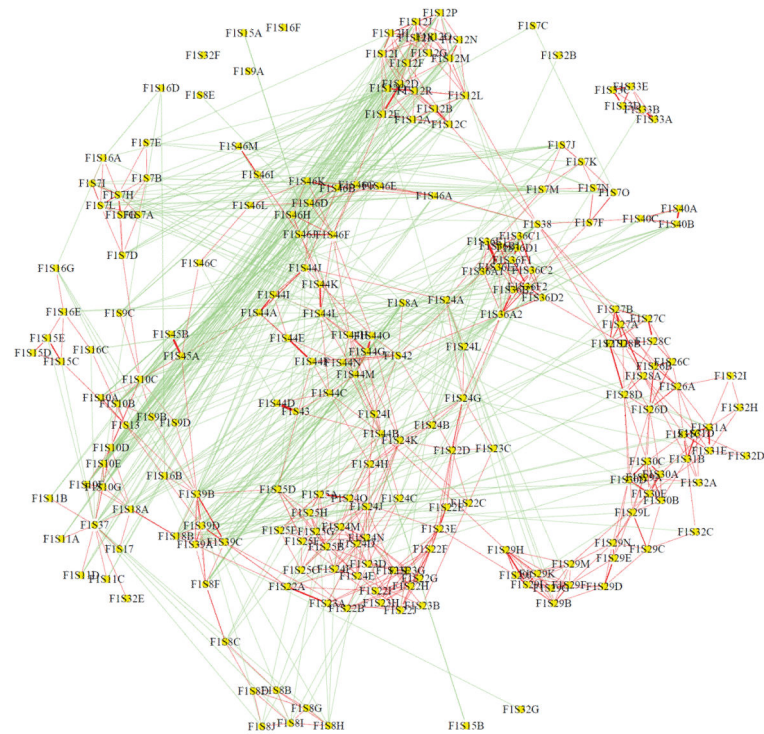
**Table 1**

The numbers are the mean CPU times for different tuning parameter values and 20 replications, with median absolute deviations in parentheses (in second). Direct: direct computation. Gibbs sampler: the regular Gibbs sampler; TMG: the Gibbs sampler proposed by Pakman and Paninski (2012) via the R package "tmg"; Proposed Approximation: the approximation approach proposed in our manuscript

| Method | CPU time in seconds |
|---|---|
| Direct | 3310.21 (199.95) |
| Gibbs sampler | 46.17 (1.51) |
| TMG | 303.94 (11.05) |
| Proposed Approximation | 0.04 (0.03) |

**Table 2**

The Frobenius loss and the entropy loss estimated by the probit graphical model, the oracle method and the naive method. The oracle method and the naive method simply apply the graphical lasso algorithm to the latent continuous data $Z$ and the observed discrete data $X$, respectively. The results are averaged over 50 repetitions and the corresponding standard deviations are recorded in the parentheses.

| Example | $n$ | Frobenius Loss | | | Entropy Loss | | |
|---|---|---|---|---|---|---|---|
| | | Gaussian | Oracle | Probit | Gaussian | Oracle | Probit |
| Scale-free | 50 | 2.3 (0.12) | 0.7 (0.05) | 2.2 (0.13) | 12.0 (0.73) | 3.1 (0.29) | 23.1 (1.83) |
| | 100 | 2.2 (0.13) | 0.4 (0.08) | 1.7 (0.09) | 9.4 (0.68) | 1.9 (0.29) | 10.1 (0.45) |
| | 200 | 1.7 (0.12) | 0.3 (0.02) | 1.2 (0.04) | 6.4 (0.33) | 1.1 (0.10) | 5.4 (0.26) |
| | 500 | 0.9 (0.05) | 0.1 (0.01) | 0.7 (0.04) | 3.3 (0.19) | 0.5 (0.05) | 2.7 (0.19) |
| Hub | 50 | 1.2 (0.06) | 0.3 (0.02) | 1.1 (0.04) | 21.2 (1.32) | 5.8 (0.70) | 29.4 (1.76) |
| | 100 | 1.1 (0.10) | 0.1 (0.01) | 0.8 (0.03) | 15.9 (1.03) | 3.2 (0.27) | 15.1 (0.64) |
| | 200 | 0.8 (0.05) | 0.1 (0.01) | 0.6 (0.01) | 11.9 (0.39) | 1.8 (0.23) | 10.4 (0.33) |
| | 500 | 0.6 (0.02) | 0.0 (0.00) | 0.5 (0.01) | 9.1 (0.16) | 0.7 (0.06) | 7.5 (0.16) |
| Nearest-neighbor | 50 | 1.4 (0.04) | 0.6 (0.02) | 1.3 (0.06) | 16.5 (0.80) | 5.6 (0.30) | 25.6 (2.04) |
| | 100 | 1.3 (0.08) | 0.4 (0.02) | 1.0 (0.02) | 12.1 (0.52) | 3.5 (0.36) | 12.4 (0.76) |
| | 200 | 1.0 (0.04) | 0.2 (0.01) | 0.7 (0.03) | 8.6 (0.32) | 2.0 (0.11) | 7.5 (0.17) |
| | 500 | 0.6 (0.03) | 0.1 (0.01) | 0.5 (0.02) | 5.5 (0.12) | 0.8 (0.02) | 4.5 (0.19) |
| Random-block | 50 | 1.8 (0.05) | 0.7 (0.05) | 1.7 (0.04) | 14.8 (1.04) | 4.7 (0.46) | 23.5 (1.76) |
| | 100 | 1.6 (0.16) | 0.4 (0.02) | 1.3 (0.03) | 10.7 (1.10) | 2.9 (0.27) | 11.3 (0.46) |
| | 200 | 1.3 (0.05) | 0.2 (0.03) | 0.9 (0.05) | 7.2 (0.19) | 1.6 (0.11) | 6.3 (0.32) |
| | 500 | 0.7 (0.03) | 0.1 (0.01) | 0.6 (0.03) | 4.1 (0.15) | 0.7 (0.06) | 3.5 (0.13) |

**Table 3**

Objectives in survey questions "In your most recent or current vocational course, how much emphasis did/ does your teacher place on the following objectives?".

| | |
|---|---|
| F1S33A | Teaching you skills you can use immediately |
| F1S33B | Teaching you facts, rules, and steps |
| F1S33C | Helping you understand how scientific ideas and mathematics are used in work |
| F1S33D | Thinking about what a problem means and the ways it might be solved |
| F1S33E | Helping you to understand mathematical and scientific ideas by helping you to manipulate physical objects (tools, machines, lab equipment) |

**Table 4**

List of pairs of questions with strongest positive partial correlations.

| Connection | Partial Correlation | Description |
| --- | --- | --- |
| F1S44D—F1S43 | 0.617981 | F1S44D: How often do you spend time on reading for pleasure? |
| | | F1S43: How much additional reading do you do each week on your own outside of school - not in connection with schoolwork? |
| F1S45A—F1S45B | 0.443995 | F1S45A: During the school year, how many hours a day do you on weekdays? |
| | | F1S45B: During the school year, how many hours a day do you on weekends? |
| F1S36E1—F1S36E2 | 0.416786 | F1S36E1: How much time do you spend on History homework in school each week? |
| | | F1S36E2: How much time do you spend on History homework out of school each week? |
| F1S44E—F1S44F | 0.398257 | F1S44E: How often do you spend time on going to the park, gym, beach, or pool outside of school? |
| | | F1S44F: How often do you spend time on playing ball or other sports with friends outside of school? |
| F1S12D—F1S12E | 0.388861 | F1S12D: How often do you feel it is "OK" for you to cheat on tests? |
| | | F1S12E: How often do you feel it is "OK" for you to copy someone else's homework? |

**Table 5**

The list of pairs of questions with strongest negative partial correlations.

| Connection | Partial Correlation | Description |
|---|---|---|
| F1S8F—F1S8A | −0.376025 | F1S8F: Did you win any special recognition for good grades or honor roll? |
| | | F1S8A: Haven't you won any awards or received recognition? |
| F1S10B—F1S12B | −0.281428 | F1S10B: How many times did you cut or skipped classes? |
| | | F1S12B: How often do you feel it is "OK" for you to cut a couple of classes? |
| F1S15B—F1S15A | −0.259550 | F1S15B: If does anyone from school called my home on your last absence from school. |
| | | F1S15A: The school did not do anything on your last absence from school. |
| F1S10A—F1S12A | −0.216677 | F1S10A: How many times were you late for school in the first half of the current school year? |
| | | F1S12A: How often do you feel it is "OK" for you to be late for school? |
| F1S16B—F1S16D | −0.214770 | F1S16B: When you came back to school after your last absence, other students helped you catch up on the work you missed. |
| | | F1S16D: When you came back to school after your last absence, you didn't need to catch up on work. |