# Sequential Monte Carlo EM for multivariate probit models[☆]

Giusi Moffa[a], Jack Kuipers[b]

[a]*Institut für funktionelle Genomik, Universität Regensburg,
Josef Engertstraße 9, 93053 Regensburg, Germany*
[b]*Institut für theoretische Physik, Universität Regensburg,
D-93040 Regensburg, Germany*

## Abstract

Multivariate probit models have the appealing feature of capturing some of the dependence structure between the components of multidimensional binary responses. The key for the dependence modelling is the covariance matrix of an underlying latent multivariate Gaussian. Most approaches to maximum likelihood estimation in multivariate probit regression rely on Monte Carlo EM algorithms to avoid computationally intensive evaluations of multivariate normal orthant probabilities. As an alternative to the much used Gibbs sampler a new sequential Monte Carlo (SMC) sampler for truncated multivariate normals is proposed. The algorithm proceeds in two stages where samples are first drawn from truncated multivariate Student $t$ distributions and then further evolved towards a Gaussian. The sampler is then embedded in a Monte Carlo EM algorithm. The sequential nature of SMC methods can be exploited to design a fully sequential version of the EM, where the samples are simply updated from one iteration to the next rather than resampled from scratch. Recycling the samples in this manner significantly reduces the computational cost. An alternative view of the standard conditional maximisation step provides the basis for an iterative procedure to fully perform the maximisation needed in the EM algorithm. The identifiability of multivariate probit models is also thoroughly discussed. In particular, the likelihood invariance can be embedded in the EM algorithm to ensure that constrained and unconstrained maximisation are equivalent. A simple iterative procedure is then derived for either maximisation which takes effectively no computational time. The method is validated by applying it to the widely analysed Six Cities dataset and on a higher dimensional simulated example. Previous approaches to the Six Cities dataset overly restrict the parameter space but, by considering the correct invariance, the maximum likelihood is quite naturally improved when treating the full unrestricted model.

*Keywords:* Maximum likelihood, Multivariate probit, Monte Carlo EM, adaptive sequential Monte Carlo

## 1. Introduction

Multivariate probit models, originally introduced by Ashford and Sowden (1970) for the bivariate case, are particularly useful tools to capture some of the dependence structure of binary, and more generally multinomial, response variables (McCulloch, 1994; McCulloch and Rossi, 1994; Bock and Gibbons, 1996; Chib and Greenberg, 1998; Natarajan et al., 2000; Gueorguieva and Agresti, 2001; Li and Schafer, 2008). Inference for such models is typically computationally involved and often still impracticable in high dimensions. To mitigate these difficulties, Varin and Czado (2010) proposed a pseudo-likelihood approach as a surrogate for a full likelihood analysis. Similar pairwise likelihood approaches were also previously considered by Kuk and Nott (2000) and Renard et al. (2004). A number of Bayesian approaches have also been considered including Chib and Greenberg (1998); McCulloch et al. (2000); Nobile (1998, 2000); Imai and van Dyk (2005) and more recently Talhouk et al. (2012).

Due to the data augmentation or latent variable nature of the problem, the expectation maximisation (EM) algorithm (Dempster et al., 1977) is typically employed for maximising the likelihood as its iterative procedure is usually

---

more attractive than classical numerical optimisation schemes. Each iteration consists of an expectation (E) step and a maximisation (M) step, and both should ideally be easy to implement.

For cases in which the E step is analytically intractable, Wei and Tanner (1990) introduced a Monte Carlo version of the EM algorithm (MCEM). Sampling from the truncated normal distributions involved is often based on Markov chain Monte Carlo (MCMC) methods and the Gibbs sampler in particular (see e.g. Geweke, 1991). As a different option we employ a sequential Monte Carlo (SMC) sampler (Del Moral et al., 2006). A sequence of distributions of interest is then approximated by a collection of weighted random samples, called particles, which are progressively updated by means of sampling and weighting operations. Though originally introduced in dynamical scenarios (Gordon et al., 1993; Kitagawa, 1996; Liu and Chen, 1998; Doucet et al., 2001) as a more general alternative to the well known Kalman filter (Kalman, 1960), SMC algorithms can also be used in static inference (see e.g. Chopin, 2002) where artificial dynamics are introduced. When the target is a truncated multivariate normal, as in our case, an obvious sequence of distributions is obtained by gradually shifting the truncation region to the desired position. Since normal distributions decay very quickly in the tails, we propose to use flatter Student $t$ distributions to drive the SMC particles more efficiently towards the target region, and only then take the appropriate limit to recover the required truncated multivariate normal. The resulting algorithm is compared to the Gibbs sampler (Geweke, 1991; Robert, 1995).

The main difficulty in the M step rests with the computational complexity of standard numerical optimisation over large parameter spaces, for which Meng and Rubin (1993) suggested a conditional maximisation approach. A simple extension of their method allows us to define an iterative procedure to further maximise the likelihood at each M step. Though the likelihood converges, there is no guarantee that the parameters converge to a point (Wu, 1983). Restrictions to the parameter space have then been introduced to treat the identifiability issue where the data does not determine the parameters uniquely (McCulloch and Rossi, 1994; Bock and Gibbons, 1996), raising the problem of constrained maximisation, normally significantly more difficult than unconstrained. When constraints are only introduced to overcome identifiability issues rather than being intrinsic to the problem, they can be regarded as artificial, and similar observations are at the basis of parameter expansion approaches to EM (Liu et al., 1998). In fact we show in our analysis of multivariate probit models that both constrained and unconstrained maximisation can be made identical. Furthermore we describe a simple novel strategy which allows either maximisation to be easily computed.

Building on the fundamental ideas of the SMC methodology it is possible to define a sequential version of the EM, where the particle approximation is simply evolved after the parameter update in the M step, rather than resampled from scratch, so reducing the computational burden of the otherwise expensive E step. Finally we validate our methods by comparison with previous approaches (Chib and Greenberg, 1998; Craig, 2008), and on a simulated higher dimensional example.


## 2. Background and notation

### 2.1. Sequential Monte Carlo samplers

Sequential Monte Carlo samplers (Del Moral et al., 2006) are a class of iterative algorithms to produce weighted sample approximations from a sequence $\{\pi_n\}$ of distributions of interest where the normalising constant $C_n$ need not be known, $\pi_n = \gamma_n/C_n$. For a given probability distribution $\pi$, one obtains a collection of weighted samples $\{W^{(k)}, \mathbf{Z}^{(k)}\}$, also referred to as particle approximation of $\pi$, such that

$$E_\pi(h(\mathbf{Z})) \simeq \sum_{k=1}^{M} W^{(k)} h(\mathbf{Z}^{(k)}),$$

where $M$ is the number of particles and $h$ a function of interest. In a static scenario the main purpose is to obtain such an approximation from the last element of an artificially defined targeted sequence.

In order to control for the degeneracy of the sample, resampling (see Douc et al., 2005, for a review of resampling schemes) is typically performed when the effective sample size (ESS), as defined by Kong et al. (1994) and often approximated as (Doucet et al., 2000):

$$\mathrm{ESS}^{-1} = \sum_{k=1}^{M} (W_n^{(k)})^2, \tag{1}$$

2

falls below a given threshold $ESS^\star = sM$, with $0 < s < 1$ though often $s = 1/2$ is chosen as a trade-off between efficiency and accuracy. The move from the target $\pi_{n-1}$ to the next $\pi_n$ is achieved by means of a transition kernel $K_n$, so that $Z_n^{(k)} \sim K_n(Z_{n-1}^{(k)}, \cdot)$, and updating the normalised weights

$$W_n^{(k)} \propto W_{n-1}^{(k)} \tilde{w}_n^{(k)}, \qquad \tilde{w}_n(\mathbf{Z}_{n-1}^{(k)}, \mathbf{Z}_n^{(k)}) = \frac{\gamma_n(\mathbf{Z}_n^{(k)}) L_{n-1}(\mathbf{Z}_n^{(k)}, \mathbf{Z}_{n-1}^{(k)})}{\gamma_{n-1}(\mathbf{Z}_{n-1}^{(k)}) K_n(\mathbf{Z}_{n-1}^{(k)}, \mathbf{Z}_n^{(k)})}, \qquad k = 1, \dots, M.$$

The quantity $L_{n-1}$ in the formula for the incremental weights $\tilde{w}_n^{(k)}$ is a backward kernel introduced by Del Moral et al. (2006) to address computational issues and should be optimised with respect to the transition kernel $K_n$ in order to minimise the variance of the importance weights. In the same work the authors also discuss a number of choices for $K_n$ suggesting MCMC kernels with $\pi_n$ as an invariant distribution as a convenient choice in many applications. A good approximation for the optimal backward kernel is then (Del Moral et al., 2006)

$$L_{n-1}(z_n, z_{n-1}) = \frac{\pi_n(z_{n-1}) K_n(z_{n-1}, z_n)}{\pi_n(z_n)}.$$

Following standard practice we therefore adopt here in particular a random walk Metropolis Hastings kernel. The samples at a given iteration $n$ are obtained by moving each particle $\mathbf{Z}_{n-1}^{(k)}$ to a new location $\mathbf{Z}_n^{(k)} = Y^k \sim \mathcal{N}(\mathbf{Z}_{n-1}^{(k)}, \Sigma_n^{\text{MH}})$ with probability $\alpha^k = 1 \wedge \rho^k$ and leaving it unchanged otherwise, with $\rho^k = \pi_n(Y^k)/\pi_n(\mathbf{Z}_{n-1}^{(k)})$. The covariance matrix $\Sigma_n^{\text{MH}} = \kappa \widehat{\Sigma}_\pi$ in the random walk proposal is a scaled version of an approximation $\widehat{\Sigma}_\pi$ of the target covariance matrix. In practice we set $\Sigma_n^{\text{MH}} = \kappa \widehat{\Sigma}_{\pi_{n-1}}$ since at iteration $n$, $\pi_{n-1}$ is the best approximation available for $\pi_n$. Unlike MCMC schemes however, in the case of SMC samplers no convergence conditions are required since any discrepancies arising from sampling from the wrong distribution are corrected by means of importance sampling reweighting (see e.g. section 2.2.3 of Del Moral et al., 2007).

As extensively investigated in the MCMC literature (for example the original paper of Gilks et al., 1998; Haario et al., 2001; Atchadé and Rosenthal, 2005, or the review of Andrieu and Thoms, 2008) the scaling factor $\kappa$ can be adaptively tuned by monitoring the average empirical acceptance probability $\hat{\alpha}_n$ at iteration $n$. For some canonical target distributions it has been proved by Roberts et al. (1997) that the asymptotically (with the dimension) optimal acceptance rate is 0.234, while Roberts and Rosenthal (1998) found that it is 0.574 for Metropolis adjusted Langevin algorithms (see also Roberts and Rosenthal, 2001, for a survey of results). There is however no gold standard on how to choose the desired acceptance rate in more realistic situations. Within the SMC framework, where one of the purposes of the MCMC move is to help maintain the sample diversity it seems sensible to fix slightly higher values than those found in theory. Especially high values should nevertheless be avoided when performing local moves as they would only be masking an eventual sample degeneracy, and lead to highly correlated samples (see Chopin, 2002, for a discussion of related issues).

In the case of SMC samplers with a Metropolis Hastings transition kernel, the empirical acceptance rate can be evaluated as

$$\hat{\alpha}_n = \sum_{k=1}^{M} W_n^{(k)} (1 \wedge \pi_n(Y_n^{(k)})/\pi_n(\mathbf{Z}_{n-1}^{(k)})).$$

Hence a stochastic approximation type algorithm can be implemented aiming to keep the above quantity equal (or close) to a prespecified value $\alpha^\star$ (see e.g. section 4.2 of Andrieu and Thoms, 2008) by setting $\Sigma_{n+1}^{\text{MH}} = \kappa_n \widehat{\Sigma}_{\pi_n}$ with scaling factor adapted as

$$\log(\kappa_{n+1}) = \log(\kappa_n) + \xi_n(\hat{\alpha}_n(\log(\kappa_n)) - \alpha^\star), \tag{2}$$

with $\xi_n$ a stepsize and where the logarithm ensures that the scaling factors are positive.

Adaptation of the transition kernel specifically within SMC has also recently been considered by Jasra et al. (2011) and Fearnhead and Taylor (2013).

## 2.2. Monte Carlo EM

An EM algorithm (Dempster et al., 1977) is an iterative procedure for the computation of maximum likelihood or maximum a posteriori estimates in the context of incomplete data problems, where the likelihood is typically

3

intractable. The algorithm relies on the definition of an associated complete data problem for which the object function of the maximisation is tractable and therefore more easily solved. Let $Y$ be a random variable representing the observed data and $\psi$ a vector of unknown parameters. Alternating between an expectation or E step and a maximisation or M step the algorithm provides us with a sequence $\{\psi^m\}$ of parameter estimates such that the observed data likelihood $\mathcal{L}(\psi \mid Y)$ is non decreasing (namely $\mathcal{L}(\psi^{m+1} \mid Y) \geq \mathcal{L}(\psi^m \mid Y)$), and eventually converges to a local maximum. Let $Z$ be a random variable corresponding to the augmented data. Separating the observed data log-likelihood in terms of the complete $(Y, Z)$ and conditional missing data $Z \mid Y, \psi$ distributions and by taking the expectation with respect to the latent variable $Z \mid Y, \psi^m$ conditioned on the observed data $Y$ and the parameter estimate $\psi^m$ at iteration $m$, the log-likelihood can be written as

$$l(\psi \mid Y) = \log(\mathrm{pr}\{Y \mid \psi\}) = Q(\psi, \psi^m) - H(\psi, \psi^m),$$

with

$$Q(\psi, \psi^m) = E_{Z\mid Y, \psi^m}\left[\log(\mathrm{pr}\{Y, Z \mid \psi\})\right], \qquad H(\psi, \psi^m) = E_{Z\mid Y, \psi^m}\left[\log(\mathrm{pr}\{Z \mid Y, \psi\})\right].$$

Jensen's inequality implies that $H(\psi, \psi^m) \leq H(\psi^m, \psi^m)$, so that the likelihood is certainly not decreased at each step if $Q(\psi^{m+1}, \psi^m) \geq Q(\psi^m, \psi^m)$. An iteration of the EM algorithm then comprises the following two steps

**E step.** Evaluate $Q(\psi, \psi^m)$.

**M step.** Maximise $Q(\psi, \psi^m)$ with respect to $\psi$.

When it is not possible to perform the E step analytically a standard solution is given by the MCEM (Wei and Tanner, 1990) where the expectation in the E step is replaced by a Monte Carlo estimate

$$Q(\psi, \psi^m) = E_{Z\mid Y, \psi^m}\left[\log(\mathrm{pr}\{Y, Z \mid \psi\})\right] \simeq \frac{1}{M}\sum_{k=1}^{M}\log(\mathrm{pr}\{Y, Z_k \mid \psi\}), \tag{3}$$

with the samples $Z_k$ drawn from the conditional distribution of the augmented data $Z \mid Y, \psi^m$.

For situations where the maximisation in the M step is not feasible, Dempster et al. (1977) suggested settling for a value that simply increases $Q(\psi, \psi^m)$ at each iteration, and they termed the resulting procedure a *generalised* EM algorithm (see also McLachlan and Krishnan, 2007, section 1.5.5). When the M step cannot be performed analytically, to overcome the difficulties associated with numerical maximisation, Meng and Rubin (1993) suggested replacing the maximisation over the full parameter space by a multi-step conditional maximisation over several subspaces in turn. Ideally we wish to set $\psi^{m+1}$ to the value of $\psi$ which maximises $Q(\psi, \psi^m)$, as required by the actual EM. In Section 3.3.1 we show, for the first time, how such a value can easily be found for the multivariate probit model.

### 2.3. Multivariate probit model

Following the formulation in Chib and Greenberg (1998), denote by $y^j$ a binary vector corresponding to the $j$th observation of a response variable $Y^j$ with $p$ components. Let $x_i^j$ be a size $k_i$ column vector containing the covariates associated to the $i$th component $y_i^j$ of the $j$th observation $Y^j$. The first element of the vector of covariates $x_i^j$ can be set to 1 to account for an intercept. Define the $j$th matrix of covariates

$$X^j \triangleq \mathrm{diag}((x_1^j)^{\mathrm{T}}, \ldots, (x_p^j)^{\mathrm{T}}),$$

as a $p \times k$ block diagonal matrix, with $k = \sum_{i=1}^{p} k_i$. A multivariate probit model with parameters $\beta \in \mathbb{R}^k$ and $\Sigma$, a $p \times p$ covariance matrix, can be specified by setting

$$\mathrm{pr}\{Y^j = y^j \mid X^j, \beta, \Sigma\} = \int_{A_1^j}\cdots\int_{A_p^j}\phi_p(z^j; X^j\beta, \Sigma)\,\mathrm{d}z^j, \qquad A_i^j = \begin{cases} (0, \infty) & \text{if} \quad y_i^j = 1 \\ (-\infty, 0] & \text{if} \quad y_i^j = 0 \end{cases}, \tag{4}$$

where $\phi_p$ is the density function of a multivariate normal random variable with mean vector $\mu = X^j\beta$ and covariance matrix $\Sigma$. The vector of regression coefficient is $\beta = (\beta_1^{\mathrm{T}}, \ldots, \beta_p^{\mathrm{T}})^{\mathrm{T}}$, with each subvector $\beta_i \in \mathbb{R}^{k_i}$ corresponding to the $i$th component of the response variable. Naturally the situation where it is assumed that the same number of covariates

4

are observed for each component of the response variable and the vectors $\boldsymbol{\beta}_i$ are also taken to be all identical can be treated as a special case. When considering particular settings however, care must be taken in reconsidering the model identifiability, as discussed in Section 3.5.

The probit model can also be understood in terms of a continuous latent variable construction, where the binary response $\boldsymbol{Y}$ is obtained by discretization of a multivariate Gaussian variable $\boldsymbol{Z} \sim \mathcal{N}(\boldsymbol{X\beta}, \boldsymbol{\Sigma})$. The observations are then thought of as obtained from an unobserved sample of multivariate Gaussian vectors $\{z^1, \ldots, z^N\}$ as $y_i^j = I_{z>0}(z_i^j)$, where specifically $\boldsymbol{Z}^j \sim \mathcal{N}(\boldsymbol{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})$ and $I$ is the indicator function.

The covariance matrix $\boldsymbol{\Sigma}$ is a crucial parameter for the multivariate probit model since it accounts, though indirectly, for some of the dependence structure among the components of the response variable. The identity matrix corresponds to the assumption of independence and the model reduces to a collection of independent one dimensional probit models, for which the regression coefficients $\boldsymbol{\beta}$ can be easily estimated, component by component, and used as starting point for more elaborate inference strategies. An alternative to the identity for the initial covariance matrix can be obtained (Emrich and Piedmonte, 1991) by pairwise approximations, which are likely however to lead to non positive definite matrices. 'Bending' techniques (Hayes and Hill, 1981; Montana, 2005) are then necessary to ensure the positivity of the eigenvalues.

### 2.3.1. Monte Carlo E step

For the multivariate probit model, letting $\boldsymbol{\psi} = (\boldsymbol{\Sigma}, \boldsymbol{\beta})$ be the parameter vector and $\boldsymbol{Z}^j \sim \mathcal{N}(\boldsymbol{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})$ the latent variables, the complete data log-likelihood function is

$$\log(\text{pr}\{\boldsymbol{Y}, \boldsymbol{Z} \mid \boldsymbol{\psi}\}) = \sum_{j=1}^{N} \log\left[I_{A^j}(z^j)\phi(z^j; \boldsymbol{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})\right]. \tag{5}$$

Using the cyclicity of the trace and ignoring some normalising constants, the corresponding $Q(\boldsymbol{\psi}, \boldsymbol{\psi}^m)$ function can be written as

$$2Q(\boldsymbol{\psi}, \boldsymbol{\psi}^m) = -N\log|\boldsymbol{\Sigma}| - N\text{tr}\left[\boldsymbol{\Sigma}^{-1}\boldsymbol{S}\right], \qquad \boldsymbol{S} = \frac{1}{N}\sum_{j=1}^{N} E_{\boldsymbol{Z}^j|\boldsymbol{Y}^j,\boldsymbol{\psi}^m}\left\{(\boldsymbol{Z}^j - \boldsymbol{X}^j\boldsymbol{\beta})(\boldsymbol{Z}^j - \boldsymbol{X}^j\boldsymbol{\beta})^{\text{T}}\right\}, \tag{6}$$

and for completeness a detailed derivation is provided in Appendix A. The expression in (6) can be transformed into the one provided in Chib and Greenberg (1998) by using the cyclicity of the trace as in (A.5), but the form in (6) is convenient for the maximisation. The second term of (6) is analytically intractable since it involves expectations with respect to high dimensional truncated multivariate Gaussian densities. In a MCEM approach (Wei and Tanner, 1990) the expectations can be approximated as

$$E_{\boldsymbol{Z}^j|\boldsymbol{Y}^j,\boldsymbol{\psi}^m}\left\{(\boldsymbol{Z}^j - \boldsymbol{X}^j\boldsymbol{\beta})(\boldsymbol{Z}^j - \boldsymbol{X}^j\boldsymbol{\beta})^{\text{T}}\right\} \simeq \sum_{k=1}^{M} W^{j(k)}(\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\boldsymbol{\beta})(\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\boldsymbol{\beta})^{\text{T}}, \tag{7}$$

over a weighted sample $\{W^{j(k)}, \boldsymbol{Z}^{j(k)}\}_{k=1}^{M}$ from $\pi(z^j \mid y^j, \boldsymbol{\psi}^m) = \text{TMN}(A^j, \boldsymbol{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})$, a multivariate normal distribution truncated to the domain $A^j$. The weights should be normalised $\sum_{k=1}^{M} W^{j(k)} = 1$ and the samples themselves may be approximate, such as provided by MCMC or importance sampling based algorithms. In our analysis we suggest to use particle approximations provided by the SMC samplers, which are detailed in Section 3.1 for the truncated multivariate normal distribution. The particle approximations so obtained can also be updated in a sequential manner from one EM iteration to the next, without the need to redraw the complete sample from scratch at each E step. The result is a more efficient EM algorithm as presented in Section 3.2 for multivariate probit models.

### 2.3.2. Two-step conditional maximisation

The multivariate normal regression with incomplete data is considered as an example in Meng and Rubin (1993). The parameters $\boldsymbol{\psi}^m$ at step $m$ are split into $\boldsymbol{\Sigma}^m$ and $\boldsymbol{\beta}^m$ leading to a two-step conditional maximisation which can be performed analytically. The solutions can be obtained by setting to zero the derivatives of (6). Using the cyclicity of the trace the maximisation condition becomes

$$2\text{d}Q = -N\text{tr}\left[\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{S}\boldsymbol{\Sigma}^{-1}\right)\text{d}\boldsymbol{\Sigma} + \boldsymbol{\Sigma}^{-1}\text{d}\boldsymbol{S}\right] = 0, \tag{8}$$

5

The function $S$ from (6) only depends on $\boldsymbol{\beta}$ so by fixing $\boldsymbol{\beta}$, the value of $\hat{\boldsymbol{\Sigma}}$ which satisfies equation (8) is simply $S$. Writing the result in terms of the particle approximation in (7), we have

$$\hat{\boldsymbol{\Sigma}}(\boldsymbol{\beta}) = \frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)} (\mathbf{Z}^{j(k)} - \boldsymbol{X}^j \boldsymbol{\beta})(\mathbf{Z}^{j(k)} - \boldsymbol{X}^j \boldsymbol{\beta})^{\mathrm{T}}. \tag{9}$$

For example, by evaluating this at the current regression vector value $\boldsymbol{\beta}^m$ the covariance matrix can be updated to $\boldsymbol{\Sigma}^{m+1} = \hat{\boldsymbol{\Sigma}}(\boldsymbol{\beta}^m)$ for the next step of the EM algorithm. Keeping $\boldsymbol{\Sigma}$ fixed, the cyclicity of the trace allows us to write the conditional maximisation condition from (8) as

$$0 = -N \mathrm{tr}\left[\boldsymbol{\Sigma}^{-1} \mathrm{d}S\right] \simeq 2 \,(\mathrm{d}\boldsymbol{\beta})^{\mathrm{T}} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)} \left(\boldsymbol{X}^j\right)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{Z}^{j(k)} - \boldsymbol{X}^j \boldsymbol{\beta}), \tag{10}$$

where again the Monte Carlo estimate in (7) has been substituted for $S$. The value of $\hat{\boldsymbol{\beta}}$ which satisfies this condition is then

$$\hat{\boldsymbol{\beta}}(\boldsymbol{\Sigma}) = \Big( \sum_{j=1}^{N} (\boldsymbol{X}^j)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{X}^j \Big)^{-1} \sum_{j=1}^{N} (\boldsymbol{X}^j)^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} \sum_{k=1}^{M} \left( W^{j(k)} \mathbf{Z}^{j(k)} \right), \tag{11}$$

so that by using the already updated value $\boldsymbol{\Sigma}^{m+1}$ the regression parameters for the next step can be updated as $\boldsymbol{\beta}^{m+1} = \hat{\boldsymbol{\beta}}\left(\boldsymbol{\Sigma}^{m+1}\right)$ to give the new parameters $\boldsymbol{\psi}^{m+1}$. Though this two-step approach does not maximise $\boldsymbol{\psi}$ at each step, it removes the need for computationally intensive maximisation and (in the large $M$ limit with particle approximations) increases the likelihood at each step to ensure convergence of the (generalised) EM.

## 2.4. Model invariance and identifiability

When the data is 'incomplete', maximisation of the observed data likelihood may not lead to uniquely identified parameters. Imposing constraints is a standard measure to ensure identifiability, but often with the effect of making the M step more involved (e.g Kuk and Chan, 2001, and more specifically for multivariate probit models Bock and Gibbons (1996); Chan and Kuk (1997)). The phenomenon is directly linked to symmetries of the likelihood, where it is invariant under some change of coordinates of the parameters. Both *global* and *local* symmetries can play a role. In the first case the invariance of the likelihood $\mathcal{L}(\boldsymbol{\psi})$ does not depend on the particular value of $\boldsymbol{\psi} \in \Psi$. The parameter space can then be decomposed as $\Psi = \Delta \times \Xi$ into an invariant space $\Delta$ and a reduced parameter space $\Xi$ so that $\boldsymbol{\psi} = (\boldsymbol{\delta}, \boldsymbol{\xi})$ with $\boldsymbol{\delta} \in \Delta$ and $\boldsymbol{\xi} \in \Xi$. Due to the invariance of the likelihood over $\Delta$

$$\mathcal{L}(\boldsymbol{\psi}) = \mathcal{L}(\boldsymbol{\delta}, \boldsymbol{\xi}) = \mathcal{L}(\boldsymbol{\xi}) \Rightarrow \max_{\boldsymbol{\psi}} \mathcal{L}(\boldsymbol{\psi}) = \max_{\boldsymbol{\xi}} \mathcal{L}(\boldsymbol{\xi}),$$

unconstrained maximisation over the whole space $\Psi$ is identical to performing it 'constrained' over the reduced space $\Xi$, with the difference that the parameters maximising the likelihood in the larger space are $\boldsymbol{\psi}^* = \Delta \times \boldsymbol{\xi}^*$. Conversely, if the likelihood depended on some subspace of $\Delta$ then it would be identified during the maximisation process. Therefore the dimension of $\Delta$ is the number of constraints needed to ensure identifiability.

In addition to any global symmetries, the likelihood function could also show a *local* symmetry so that $\hat{\mathcal{L}}(\boldsymbol{\xi})$ is maximised by a higher dimensional manifold rather than a single point (as discussed in Wu, 1983). In principle a local change of variables is possible (for example making the non-zero eigenvalues of the Hessian equal to $-1$ around the maximum) to decompose the space further, but in practice this presumes knowledge of the likelihood function. As above though, maximisation over the subspace or the whole space are exactly equivalent because there will still be (local) dimensions which do not affect the value of the likelihood.

Within the EM algorithm the identifiability issue becomes more subtle since the likelihood is not maximised directly, but by proxy through the function $Q(\boldsymbol{\psi}, \boldsymbol{\psi}^m)$. If this were to share the symmetries of the likelihood, then the simpler unconstrained maximisation would be equivalent to the constrained version, as for the likelihood. If this is not the case, for example due to conditioning on the previous parameter value $\boldsymbol{\psi}^m$, then any changes in $Q$ arising from shifting $\boldsymbol{\psi}$ in the invariant space $\Delta$ of the likelihood must be exactly mimicked by changes in $H$. This spurious dependence can create differences between constrained and unconstrained maximisation. The non decreasing

behaviour of the likelihood remains preserved, since neither maximisation decreases $Q$ nor, because of Jensens's inequality, increases $H$. Hence either choice leads to the EM algorithm finding a maximum of the likelihood (though not necessarily the same one) and explains the conjecture of Bock and Gibbons (1996); Chan and Kuk (1997) and the agreement between constrained and unconstrained maximisation found in Kuk and Chan (2001).

In fact such symmetries can be seen as a natural example of the parameter expansion EM algorithm of Liu et al. (1998) where in general one seeks additional parameters which do not affect the likelihood but which can be incorporated into the EM steps. Here for example the standard EM would be over the constrained space $\Xi$ while the parameter expanded version would include some or all the parameters in $\Delta$. By examining the second differentials of the likelihood, Liu et al. (1998) showed that (at least in a quadratic neighbourhood of the maxima) the parameter expanded EM algorithm converges at least as fast as the standard version, suggesting the more parameters the better. They also gave some examples where the speed up was very significant. In general unconstrained maximisation is also less demanding than constrained maximisation and then more appealing on both fronts.

## 3. Methodology

Most approaches to MCEM for multivariate probit (e.g. Chan and Kuk, 1997; Natarajan et al., 2000) rely on MCMC schemes based on the Gibbs sampler to approximate the expectations in (7). As an alternative in Section 3.1 we propose a SMC sampler for truncated multivariate normal distributions. In Section 3.2 we discuss how to evolve the particle approximation through EM iterations and so avoid to fully draw a new sample at each E step. By taking an alternative view of the conditional maximisation an iterative procedure to complete the maximisation is discussed in Section 3.3. Based on identifiability considerations specifically for multivariate probit models, it is discussed in Section 3.5 how to perform constrained maximisation at almost no computational cost.

### 3.1. SMC sampler for truncated multivariate normal and t distributions

Since the probability of a random walk Metropolis to move towards the tails of a Gaussian distribution decreases exponentially, a SMC method involving normals may be highly inefficient in moving samples towards regions of low probability. To achieve higher rates of acceptance in the tails we suggest starting with a flatter distribution: the multivariate (of dimension $p$) Student $t$ distribution $\mathcal{T}(\nu, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ with degree of freedom $\nu$, a size $p$ vector $\boldsymbol{\mu}$ and a $p \times p$ positive definite matrix $\boldsymbol{\Sigma}$ as location and scale parameters respectively. The probability density function of a variable $\mathbf{Z} \sim \mathcal{T}(\nu, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ can be defined (Nadarajah and Kotz, 2005) as

$$f(z) = \frac{\Gamma(\frac{\nu+p}{2})}{\Gamma(\frac{\nu}{2})(\pi\nu)^{p/2}|\boldsymbol{\Sigma}|^{1/2}} \left[ 1 + \frac{1}{\nu}(z - \boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(z - \boldsymbol{\mu}) \right]^{-\frac{\nu+p}{2}}. \tag{12}$$

Replacing the $\nu$ in the denominator inside the square brackets by $(\nu - 2)$, and correspondingly changing the normalisation factor, would provide the Student distribution with a covariance of $\boldsymbol{\Sigma}$. As it stands, the distribution in (12) actually has a covariance of $\nu\boldsymbol{\Sigma}/(\nu - 2)$ which further increases the acceptance in the tails. Once in the region of low probability we allow the number of degrees of freedom to grow to infinity ($\nu \to \infty$) so the distribution approaches a $p$-variate Gaussian with the same mean and covariance matrix $\boldsymbol{\Sigma}$.

To sample in the region of interest $A$, we define a sequence of target distributions $\{\pi_n\}_0^T$ such that the first target is an unconstrained multivariate Student and the last one is the same distribution truncated to $A$. Quite naturally the intermediate distributions are defined in terms of intermediate target domains $\{A_n\}_0^T$, included in each other $A_{k+1} \subset A_k$, with $A_T \equiv A$ and $A_0 \equiv \mathbb{R}^p$. The local target $\pi_n$ at iteration $n$ of the SMC algorithm is then $\pi_n(z) = \gamma_n(z)/C_n$, with

$$\gamma_n(z) = \left[ 1 + \frac{1}{\nu}(z - \boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(z - \boldsymbol{\mu}) \right]^{-\frac{\nu+p}{2}} I_{A_n}(z),$$

where $C_n$ is a normalising constant which can be estimated (Del Moral et al., 2006) from

$$\widehat{C_n} = C_0 \prod_{i=1}^{n} \widehat{\frac{C_i}{C_{i-1}}}, \qquad \widehat{\frac{C_i}{C_{i-1}}} = \sum_{k=1}^{M} W_{i-1}^{(k)} \tilde{w}_i(\mathbf{Z}_{i-1}^{(k)}, \mathbf{Z}_i^{(k)}),$$
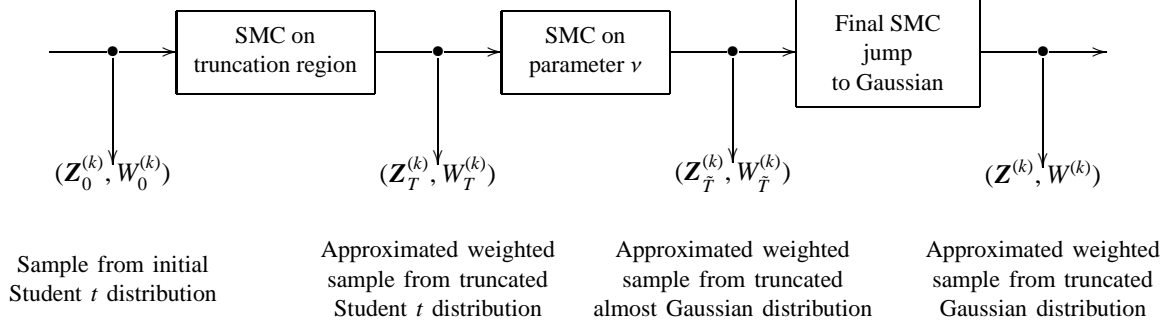
Figure 1: Cascade interpretation of the SMC sampler for truncated multivariate normal via Student $t$.

and $C_0$ follows from (12). It follows that the probability that a random variable $\mathbf{Z} \sim \pi_0$ from the initial distribution falls within region $A_n$ can be approximated by $P(\mathbf{Z} \in A_n) = \frac{C_n}{C_0} \simeq \prod_{i=1}^{n} \widehat{\frac{C_i}{C_{i-1}}}$. This ultimately allows us to obtain the probabilities of the regions in (4) and hence the likelihood for the probit model.

After reaching the required region, we define a new sequence of target distributions $\{\tilde{\pi}_n\}_0^{\tilde{T}}$ which this time starts from the truncated Student $\tilde{\pi}_0 = \pi_T = \gamma_T/C_T$. The following terms of the sequence are defined by increasing the degree of freedom $\nu$ up to a value $\nu_{\tilde{T}}$ large enough so that the truncated Student $\tilde{\pi}_{\tilde{T}} = \tilde{\pi}(\nu_{\tilde{T}})$ cannot be distinguished from the desired truncated multivariate normal within a certain level of accuracy. A final step is then performed to explicitly move to the Gaussian. One could also vary both the truncation region and the degree of freedom concurrently in the sequence of target distributions, but since the main reason for introducing the flatter Student distribution is to aid moving to regions of low probability we chose this two-step approach. A graphical overview of the process of moving from a Student $t$ distribution to a truncated Gaussian is given in Figure 1.

Similar sampling problems have also been studied in the literature dealing with rare event analysis. There smooth sequences of distributions gradually concentrating on the rare set have been suggested, as opposed to simply increasingly truncated distributions (see Johansen et al., 2006, and references therein). The two stage approach based on the Student $t$ distribution is preferred here for its relative conceptual simplicity and the fact that it proved well suited in practice to a linear adaptation framework of the type described in the following section 3.1.1.

### 3.1.1. Adaptive approach to artificial dynamics

Adaptive strategies can be applied not only for tuning the transition kernel $K_n$, as noticed in Section 2.1, but also to define the artificial dynamics leading to the distribution of interest $\pi_T$. The problem of finding the optimal path linking an initial measure $\pi_0$ to the target $\pi_T$ on the space of distributions is not addressed, rather it is assumed that the functional form of the intermediate distributions is given and can be described in terms of a parameter $\theta$. In the examples of the Section 3.1 we have $\theta = A$ for the truncation case and $\theta = \nu$ when moving the truncated Student to a Gaussian. An adaptive strategy to move from $\pi_0$ to $\pi_T$ is one that does not require the sampling points $\{\theta_n\}$ defining the intermediate targets $\{\pi_n = \pi(\theta_n)\}$ to be fixed a priori, but allows us to determine them dynamically on the basis of the local difficulty of the problem.

Adaptation can be achieved by controlling some statistics related to the performance of the algorithm and evolving with the parameter $\theta$. The ESS introduced in (1) is an ideal quantity to monitor. Theoretically we wish to solve

$$\text{ESS}_n(\theta_n) - \text{ESS}_A^\star = 0, \tag{13}$$

where $\text{ESS}_A^\star$ is a value chosen to compromise between efficiency and accuracy, and which can be different (lower) from the resampling threshold $\text{ESS}^\star$. Inspired by the Robbins-Monro recursion (see for example Kushner and Yin, 2003, page 3) for stochastic approximation, and aiming at the dynamical design of a sequence which keeps the ESS on average close to the threshold $\text{ESS}_A^\star$, we define the updating scheme

$$\theta_{n+1} = \left[ \theta_n + \left( \zeta_n \frac{\widetilde{\text{ESS}}_n - \text{ESS}_A^\star}{M} \vee \Delta\theta_{\min} \right) \right] \wedge \theta_T, \tag{14}$$
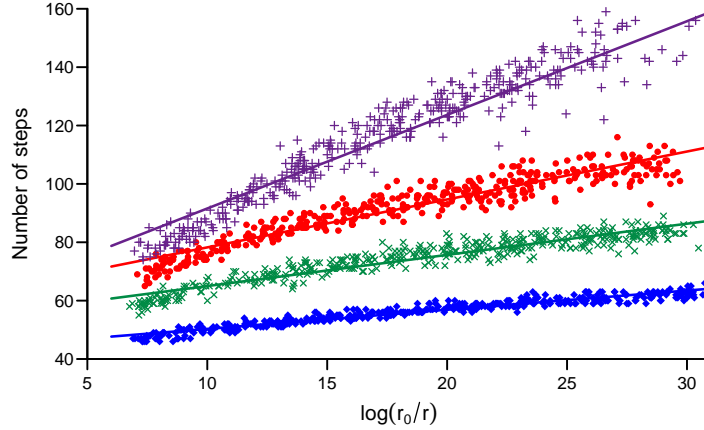
8

Figure 2: The number of steps required for the SMC algorithm to reach a region of probability $r$ for dimensions 2 (diamonds), 4 (crosses), 8 (dots) and 16 (pluses).

where $\widetilde{\mathrm{ESS}}_n$ is the value observed for ESS at iteration $n$ and the division by the number of particles $M$ is only introduced for scaling purposes. Taking the maximum between the correction term and $\Delta\theta_{\min}$ ensures that the resulting sequence approaches the final target monotonically, while taking the minimum with $\theta_T$ ensures that the sequence ends at the desired target $\pi(\theta_T)$. Theoretically the ESS should ideally be equal to the total number of particles $M$ of the SMC sampler. To promote motion and so a quicker progression of the algorithm towards its final target, the threshold $\mathrm{ESS}_A^\star$ can be fixed as a fraction $a \in (0, 1)$ of $M$, namely $\mathrm{ESS}_A^\star = aM$. The fraction $a$ should be slightly smaller than the fraction $s$ defined in 2.1 to control the resampling, say $a = .9s$, to ensure that the resampling threshold $\mathrm{ESS}^\star$ is also crossed while the algorithm runs. The number of iterations needed to reach the target $\pi_T$ is reduced for smaller $a$. Similar adaptive ideas have also been applied to inference for stochastic volatility models by Jasra et al. (2011) and rather recently discussed in Schäfer and Chopin (2013) for sampling in large binary spaces.

The details of the complete SMC sampler are summarised in Algorithm 1 in Appendix B.

### 3.1.2. Scaling behaviour

The advantage of the SMC method, over alternatives which may be more efficient in sampling from truncated multivariate normals in low dimensions, is the scaling behaviour with the dimension $p$. Solving the adaptive equation (13) exactly means that we lose a fixed proportion of the probability mass at each iteration. The number of steps required to reach a target region of low probability $r$, then behaves like $\log(r)$, independently of $p$. This may not be true when using (14) as a numerical adaptive approximation to (13), especially as the number of steps for the adaption to settle grows linearly with $p$, so a weak dependence on the dimension could be expected.

A simulation study with targets of dimensions $2^n$ for $n = 1, \ldots, 4$ was performed. To limit the sources of variability, only one covariance structure was considered for the unconstrained distribution, with unit diagonals and a single non-zero off-diagonal element of $0 \cdot 9$. The SMC algorithm was initialised so that after an initial move the Student $t$ target would be truncated to a region containing one quarter of the probability mass of an independent Gaussian, and we denote by $r_0$ the actual estimated probability. The cutoff for the final target, the same in all directions, was drawn so as to ensure that the log probability of an independent multivariate normal would be uniform on a given interval. The number of steps needed to reach the target are plotted against $\log(r_0/r)$ in Fig. 2, for 400 runs of a SMC sampler with 4000 particles for the different dimensions. A behaviour close to linear can be observed, though the offset increases by a factor of about $1 \cdot 4$ over the range of dimensions and the slope increases roughly linearly with $p$, which is likely due to any inexactness in the adaptation. The theoretical stability of these types of algorithms has recently been investigated in depth by Beskos et al. (2012).

### 3.2. Sequential Monte Carlo EM for multivariate probit models

The SMC sampler of Section 3.1 for truncated multivariate normals has good scaling behaviour for rare events, but depending on the choice of initial parameters $\psi^0 = (\Sigma^0, \beta^0)$ there may be more efficient methods of obtaining the initial sample. If for example, as suggested earlier, the covariance matrix $\Sigma^0$ is set to be the identity matrix $I$, a sample of the corresponding distribution truncated to a region $A = A_1 \times A_2 \times \cdots \times A_p$ with $A_i$ of the type defined in equation (4) can be simply obtained by truncating each of the $p$ components independently. Draws from a univariate truncated normal can for example be quite efficiently obtained via the mixed rejection algorithm of Geweke (1991) or for truncation near the mean by the recent method proposed by Chopin (2011). When a better and much more structured guess is available for the initial covariance matrix $\Sigma^0$, such that the components cannot be truncated independently and the efficient univariate methods cannot be applied, the SMC sampling method proposed above is then of course a valid alternative.

On the other hand once the initial sample has been obtained (from whichever method of choice) sequential Monte Carlo methods provide a natural machinery for efficient parameter updating during the E-step of a Monte Carlo based EM. In fact given a particle approximation from the truncated target distribution corresponding to the initial parameter values a sequential Monte Carlo approach can easily be defined to move between subsequent estimates $\psi^m = (\Sigma^m, \beta^m)$ without the need to perform the complete truncation again. The M step of iteration $m$ provides the newly optimised parameters $(\Sigma^{m+1}, \beta^{m+1})$. While from the previous E step, for each observation $j$, a particle approximation is available from a multivariate normal with mean $X^j \beta^m$ and covariance $\Sigma^m$ truncated to the region $A^j = A_1^j \times A_2^j \times \cdots \times A_p^j$. One wishes to simply move these particles to approximate the truncation to the same region $A^j$ of a multivariate normal with updated mean $X^j \beta^{m+1}$ and covariance $\Sigma^{m+1}$.

Translating the particles could lead to the situation where the mean shift would effectively imply moving the sample to a bigger region, which would prevent us from using a simplified version of the backward kernel $L_n$ (see section 3.3.2.3 of Del Moral et al., 2006). Instead the coordinates of the previous sample can just be rescaled by multiplying them by the diagonal matrix $D^{-1}$. As long as the scaling factors (elements of $D^{-1}$) are all positive, this transformation does not affect the truncation region and the mean vector is likewise scaled to $D^{-1} X^j \beta^m$. By simply choosing $D$ to set this equal to the new required mean vector $X^j \beta^{m+1}$ we have a particle approximation with the correct mean and truncation region, but with a covariance matrix of $D^{-1} \Sigma^m D^{-1}$. A SMC algorithm can then be applied to target a distribution with the new covariance matrix $\Sigma^{m+1}$.

Multiple sub-steps might be needed to update $D^{-1} \Sigma^m D^{-1}$ to $\Sigma^{m+1}$, depending on how different the two corresponding targets are. As the EM algorithm progresses however, these two must tend to approach each other and a single step will start to suffice. For each observation $j$ the local (to the EM iteration) initial and final distributions of an artificial sequence $\{\pi_n\}_0^T$ can be defined via $\pi_0 = \text{TMN}(A^j, X^j \beta^{m+1}, D^{-1} \Sigma^m D^{-1})$ and $\pi_T = \text{TMN}(A^j, X^j \beta^{m+1}, \Sigma^{m+1})$ respectively. The parameter sequence $\{\theta_n\}_0^T$ then defines the intermediate targets moves such that $\theta_0 = D^{-1} \Sigma^m D^{-1}$ and $\theta_T = \Sigma^{m+1}$, and possibly only requires a single step.

The complete SMC EM procedure is outlined in Algorithm 2 of Appendix B.

### 3.3. M-step in the multivariate probit: alternative approach to the conditional maximisation step

The first step of the conditional maximisation in Section 2.3.2 can be interpreted as follows. The expression for the covariance matrix $\hat{\Sigma}(\beta)$ of equation (9) maximises the $Q$ function $Q(\psi, \psi^m)$ over $\Sigma$ for any value of $\beta$. Therefore it can be substituted in the expression of $Q(\psi, \psi^m)$ in (6) providing a function which only depends on $\beta$

$$2\hat{Q}(\beta, \psi^m) = -N \log |\hat{\Sigma}(\beta)| - Np, \tag{15}$$

where the simplification of the second term to a constant derives from the fact that the argument of the trace reduces to the identity matrix. Finding the value $\tilde{\beta}$ which maximises (15) over $\beta$ and setting $\tilde{\Sigma} = \hat{\Sigma}(\tilde{\beta})$ in (9) provides the new parameter $\tilde{\psi} = (\tilde{\beta}, \tilde{\Sigma})$ which maximises the likelihood. Setting to zero the derivative of (15) with respect to $\beta$ leads to the condition $\text{tr}\{\hat{\Sigma}^{-1} d\hat{\Sigma}\} = 0$. Since $\hat{\Sigma} = S$, from the Monte Carlo expression of $S$ obtained when combining (6) and (7) it follows that

$$d\hat{\Sigma} = dS \simeq -\frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)} \left[ (X^j d\beta)(Z^{j(k)} - X^j \beta)^{\mathrm{T}} + (Z^{j(k)} - X^j \beta)(X^j d\beta)^{\mathrm{T}} \right].$$

and again by the cyclicity of the trace the condition $\mathrm{tr}\{\hat{\boldsymbol{\Sigma}}^{-1}\mathrm{d}\hat{\boldsymbol{\Sigma}}\} = 0$ reduces to the condition in (10) with $\hat{\boldsymbol{\Sigma}}$ instead of $\boldsymbol{\Sigma}$

$$2\,(\mathrm{d}\boldsymbol{\beta})^{\mathrm{T}} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)} \left(\boldsymbol{X}^{j}\right)^{\mathrm{T}} \hat{\boldsymbol{\Sigma}}^{-1}(\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^{j}\boldsymbol{\beta}) = 0. \tag{16}$$

Though $\mathrm{d}\hat{\boldsymbol{\Sigma}}$ is linear in the components of $\boldsymbol{\beta}$, the inverse matrix $\hat{\boldsymbol{\Sigma}}^{-1}$ leads to a system of coupled higher order polynomial equations. Solving these is impracticable, but an iterative scheme over a sequence of points $\tilde{\boldsymbol{\beta}}^{n}$ can be followed.

Performing Newton-type iterations would be an option, but when the starting point is not too far from the maximum a simpler approximate maximisation can be employed. Starting from a point $\tilde{\boldsymbol{\beta}}^{n}$ first set $\tilde{\boldsymbol{\Sigma}}^{n+1} = \hat{\boldsymbol{\Sigma}}(\tilde{\boldsymbol{\beta}}^{n})$, then separate $\hat{\boldsymbol{\Sigma}}(\boldsymbol{\beta}) = \tilde{\boldsymbol{\Sigma}}^{n+1} + \Delta\tilde{\boldsymbol{\Sigma}}(\boldsymbol{\beta})$. Recall that $\log|\boldsymbol{G}| = \mathrm{tr}\{\log\boldsymbol{G}\}$ for any matrix $\boldsymbol{G}$ and make the approximation $\log(\boldsymbol{I}+\boldsymbol{G}) \approx \boldsymbol{G}$ for $\boldsymbol{G}$ near $\boldsymbol{0}$ (Petersen and Pedersen, 2012; Higham, 2008) to rewrite the $\hat{Q}$ function as

$$2\hat{Q}(\boldsymbol{\beta}, \boldsymbol{\psi}^{m}) \approx -N\mathrm{tr}\left\{\log\tilde{\boldsymbol{\Sigma}}^{n+1}\right\} - N\mathrm{tr}\left\{(\tilde{\boldsymbol{\Sigma}}^{n+1})^{-1}\Delta\tilde{\boldsymbol{\Sigma}}\right\} - Np.$$

where the only term depending on $\boldsymbol{\beta}$ is $\Delta\tilde{\boldsymbol{\Sigma}}$. Since $\mathrm{d}\Delta\tilde{\boldsymbol{\Sigma}} = \mathrm{d}\hat{\boldsymbol{\Sigma}}$, when differentiating with respect to $\boldsymbol{\beta}$, maximising the previous expression is achieved by solving

$$0 = -N\mathrm{tr}\left\{(\tilde{\boldsymbol{\Sigma}}^{n+1})^{-1}\mathrm{d}\hat{\boldsymbol{\Sigma}}\right\} \simeq 2\,(\mathrm{d}\boldsymbol{\beta})^{\mathrm{T}} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)}(\boldsymbol{X}^{j})^{\mathrm{T}}(\tilde{\boldsymbol{\Sigma}}^{n+1})^{-1}(\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^{j}\boldsymbol{\beta}), \tag{17}$$

which is again just (10) evaluated at a given $\tilde{\boldsymbol{\Sigma}}^{n+1}$. The solution is of the form in (11) when replacing $\boldsymbol{\Sigma}$ with $\tilde{\boldsymbol{\Sigma}}^{n+1}$, so that the next value of $\tilde{\boldsymbol{\beta}}$ can be set as $\tilde{\boldsymbol{\beta}}^{n+1} = \hat{\boldsymbol{\beta}}(\tilde{\boldsymbol{\Sigma}}^{n+1})$. For a given starting point of the full parameter vector $\tilde{\boldsymbol{\psi}}^{n}$ the covariance matrix $\tilde{\boldsymbol{\Sigma}}^{n+1}$ and subsequently the coefficient vector $\tilde{\boldsymbol{\beta}}^{n+1}$ can be updated in this way to provide a point $\tilde{\boldsymbol{\psi}}^{n+1}$ for the next iteration.

### 3.3.1. Complete maximisation

Because of the logarithmic approximation, each value of $\tilde{\boldsymbol{\beta}}^{n+1}$ found in a single step does not yet maximise $\hat{Q}$. It is clear however that the value of $\boldsymbol{\beta}$ can be iteratively adjusted until the maximum is found. The procedure moves along the sequence $\tilde{\boldsymbol{\beta}}^{n}$ using the current value as the starting point for the next iteration, while updating $\tilde{\boldsymbol{\Sigma}}^{n+1}$ at the same time. The maximisation can be completed by iterating until one finds the maximiser $\tilde{\boldsymbol{\psi}} = \lim_{n\to\infty} \tilde{\boldsymbol{\psi}}^{n}$, and numerically stopping the iterations when the Euclidean norm $\|\tilde{\boldsymbol{\beta}}^{n+1} - \tilde{\boldsymbol{\beta}}^{n}\|$ is small. For the EM algorithm one can then set $\boldsymbol{\psi}^{m+1} = \tilde{\boldsymbol{\psi}}$.

In general the surety of convergence or even of not decreasing $\hat{Q}$ is lost with approximations. But choosing $\tilde{\boldsymbol{\beta}}^{0} = \boldsymbol{\beta}^{m}$ (or $\tilde{\boldsymbol{\psi}}^{0} = \boldsymbol{\psi}^{m}$) as a starting point leads to the same expression as that found after the first conditional maximisation over $\boldsymbol{\Sigma}$ in (9) and $\tilde{\boldsymbol{\Sigma}}^{1} = \boldsymbol{\Sigma}^{m+1}$. The logarithmic approximation then gives $\tilde{\boldsymbol{\beta}}^{1} = \boldsymbol{\beta}^{m+1}$, so that neatly, a single iteration using the logarithmic approximation and the two-step conditional maximisation of Meng and Rubin (1993) are equivalent when started at the same point ($\boldsymbol{\psi}^{m}$ for example). That each iteration (without a particle approximation) does not decrease the likelihood follows from the arguments in Meng and Rubin (1993), confirming the convergence of the maximisation. More importantly, completing the maximisation by iterating until convergence can equivalently be achieved by running through the two-step conditional maximisation many times.

### 3.4. From generalised EM to EM

Though the focus of Section 3.3.1 is on multivariate normals, cycling through the conditional maximisations of Meng and Rubin (1993) until convergence can be applied more generally, turning the *generalised* EM of their single round procedure into an EM again. However, as they mention, it may be computationally advantageous to perform an E step between conditional maximisations when these are more demanding, and in such cases the algorithm remains a generalised one.

### 3.5. Invariance and identifiability in the multivariate probit model

The full parameter space $\Psi$ of a multivariate probit model comprises $p(p+1)/2$ entries from the covariance matrix $\boldsymbol{\Sigma}$ and $k$ regression coefficients from $\boldsymbol{\beta}$. Invariance of the likelihood is observed under a rescaling of the coordinates of the latent multivariate normal variable $\boldsymbol{Z}$ by means of a diagonal matrix $\boldsymbol{D}$ with positive entries $(d_1, \ldots, d_p)$ according to the transformation $z^j = \boldsymbol{D}u^j$. The covariance matrix $\boldsymbol{\Sigma}$ gets transformed to $\boldsymbol{\Omega} = \boldsymbol{D}^{-1}\boldsymbol{\Sigma}\boldsymbol{D}^{-1}$ and the vector of regression coefficients $\boldsymbol{\beta}$ to $\boldsymbol{\lambda} = (d_1^{-1}\boldsymbol{\beta}_1^{\mathrm{T}}, \ldots, d_p^{-1}\boldsymbol{\beta}_p^{\mathrm{T}})^{\mathrm{T}}$, but it can easily be checked that the likelihood is left unchanged. Choosing the entries of $\boldsymbol{D}$ to be the square root of the diagonal elements of $\boldsymbol{\Sigma}$ reduces $\boldsymbol{\Omega}$ to correlation form. The invariant space $\Delta$ then has coordinates given by the $p$ diagonal elements of $\boldsymbol{\Sigma}$ (i.e. $\delta_1 = 1/\sqrt{\sigma_{11}}$ etc.) while the reduced space $\Xi$ includes the $p(p-1)/2$ rescaled upper triangular elements of $\boldsymbol{\Omega}$ (i.e. $\omega_{ij} = \delta_i \delta_j \sigma_{ij}$) and the $k$ elements of $\boldsymbol{\lambda} = (\delta_1 \boldsymbol{\beta}_1^{\mathrm{T}}, \ldots, \delta_p \boldsymbol{\beta}_p^{\mathrm{T}})^{\mathrm{T}}$.

The likelihood however is not maximised directly, but through the function

$$Q(\boldsymbol{\psi}, \boldsymbol{\psi}^m) = \sum_{j=1}^{N} \int_{A^j} \mathrm{TMN}(A^j, \boldsymbol{X}^j\boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m)\left[\log\left(\frac{1}{|\boldsymbol{\Sigma}|^{1/2}}\right) - \frac{1}{2}(z^{(j)} - \boldsymbol{X}^j\boldsymbol{\beta})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(z^j - \boldsymbol{X}^j\boldsymbol{\beta})\right] \mathrm{d}z^j. \tag{18}$$

Given a diagonal matrix $\boldsymbol{D}$ the expression in (18) above is only invariant under a change of the integration variables $z^j = \boldsymbol{D}u^j$, if a factor $|\boldsymbol{D}|$ is included inside the log and correspondingly inside the log of $H(\boldsymbol{\psi}, \boldsymbol{\psi}^m)$. Moreover, both $\boldsymbol{\psi}$ and $\boldsymbol{\psi}^m$ need to be scaled by the same matrix so that essentially $\boldsymbol{\delta} = \boldsymbol{\delta}^m$. Therefore $\boldsymbol{\psi}$ and $\boldsymbol{\psi}^m$ are tied together in the $Q$ function in an apparent constraint, though theoretically they have independent invariant spaces for the likelihood. In Chib and Greenberg (1998) maximisation is performed inside the constrained space $\Xi$, while keeping $\delta_i = 1$. Denote by $\boldsymbol{\psi}_{\mathrm{c}}$ the corresponding solution and by $\boldsymbol{\psi}_{\mathrm{u}}$ the one obtained through unconstrained maximisation of $Q$. Clearly $Q(\boldsymbol{\psi}_{\mathrm{u}}, \boldsymbol{\psi}^m) \geq Q(\boldsymbol{\psi}_{\mathrm{c}}, \boldsymbol{\psi}^m)$, but when projecting $\boldsymbol{\psi}_{\mathrm{u}}$ to a point $\boldsymbol{\psi}_{\mathrm{p}}$ in the constrained space $\Xi$ by setting $\delta_i = 1$ then $Q(\boldsymbol{\psi}_{\mathrm{p}}, \boldsymbol{\psi}^m) \leq Q(\boldsymbol{\psi}_{\mathrm{c}}, \boldsymbol{\psi}^m)$. Since the likelihood is invariant under this projection

$$Q(\boldsymbol{\psi}_{\mathrm{u}}, \boldsymbol{\psi}^m) - Q(\boldsymbol{\psi}_{\mathrm{p}}, \boldsymbol{\psi}^m) = H(\boldsymbol{\psi}_{\mathrm{u}}, \boldsymbol{\psi}^m) - H(\boldsymbol{\psi}_{\mathrm{p}}, \boldsymbol{\psi}^m),$$

and without any information on $H(\boldsymbol{\psi}_{\mathrm{u}}, \boldsymbol{\psi}^m) - H(\boldsymbol{\psi}_{\mathrm{c}}, \boldsymbol{\psi}^m)$, for example from the second differential of the likelihood as in parameter expanded EM (Liu et al., 1998), it is impossible to say which maximisation increases the likelihood most and is to be preferred in that respect.

### 3.5.1. Reintroducing the likelihood invariance in the Q function

To remove the above ambiguity, $Q$ can be redefined to respect the invariance of the likelihood, for example by replacing the parameters $(\boldsymbol{\Sigma}, \boldsymbol{\beta})$ in (18) by their projection $(\boldsymbol{\Omega}, \boldsymbol{\lambda})$. Such a replacement effectively enforces invariance of the resulting function $\tilde{Q}$ with respect to a rescaling of $(\boldsymbol{\Sigma}, \boldsymbol{\beta})$, making constrained and unconstrained maximisation identical. However, this is no longer true when a (cyclical) two-step conditional maximisation is performed.

With the replacement in (6), $\tilde{Q}$ becomes

$$\tilde{Q}(\boldsymbol{\psi}, \boldsymbol{\psi}^m) = -\frac{N}{2}\left[\log\frac{|\boldsymbol{\Sigma}|}{|\boldsymbol{D}|^2} + \mathrm{tr}\left\{\boldsymbol{D}\boldsymbol{\Sigma}^{-1}\boldsymbol{D}\tilde{\boldsymbol{S}}\right\}\right], \qquad \tilde{\boldsymbol{S}} = \frac{1}{N}\sum_{j=1}^{N}\sum_{k=1}^{M} W^{j(k)}(\boldsymbol{Z}^{j(k)} - \boldsymbol{D}^{-1}\boldsymbol{X}^j\boldsymbol{\beta})(\boldsymbol{Z}^{j(k)} - \boldsymbol{D}^{-1}\boldsymbol{X}^j\boldsymbol{\beta})^{\mathrm{T}}, \tag{19}$$

in terms of the particle approximation in (7), and where $\boldsymbol{D}$ is a diagonal matrix whose elements are the square roots of the diagonal elements of $\boldsymbol{\Sigma}$ (so that its projection into correlation form is $\boldsymbol{\Omega} = \boldsymbol{D}^{-1}\boldsymbol{\Sigma}\boldsymbol{D}^{-1}$). Though $\tilde{Q}$ may appear to be limited to the constrained space, it depends on the full parameter space when one of $\boldsymbol{\Sigma}$ or $\boldsymbol{\beta}$ are given. Assume that for given $\boldsymbol{\psi}^m$ and $\boldsymbol{\beta}^m = \boldsymbol{\lambda}^m$ we wish to find $\boldsymbol{\Sigma}^{m+1}$. Constrained maximisation enforces $\delta_i = 1$ to find $\boldsymbol{\Omega}_{\mathrm{c}}^{m+1}$. An unconstrained maximisation allows $\delta_i$ to vary, leading to $\boldsymbol{\Sigma}_{\mathrm{u}}^{m+1}$ such that $\tilde{Q}((\boldsymbol{\Sigma}_{\mathrm{u}}^{m+1}, \boldsymbol{\beta}^m), \boldsymbol{\psi}^m) \geq \tilde{Q}((\boldsymbol{\Omega}_{\mathrm{c}}^{m+1}, \boldsymbol{\lambda}^m), \boldsymbol{\psi}^m)$. Because of the invariance, the projection of $(\boldsymbol{\Sigma}_{\mathrm{u}}^{m+1}, \boldsymbol{\beta}^m)$ does not now change $\tilde{Q}$ resulting in a point in the constrained space with a higher value. It can now be unambiguously seen that the unconstrained maximisation is preferable. In fact $\boldsymbol{\beta}$ is only defined up to a scale, which need not be preserved during each conditional maximisation, nor given the stochastic nature of the estimation step.

### 3.5.2. Constrained maximisation

Introducing the invariance of the likelihood into the function $Q$ to obtain the $\tilde{Q}$ in (19) provides a function whose maximisation allows constrained maximisation to be performed over the original $Q$ since they are identical when $\boldsymbol{D}$ is the identity matrix.

First we differentiate (19) to obtain the maximisation conditions

$$2\mathrm{d}\tilde{Q} = -N\mathrm{tr}\left[\boldsymbol{\Sigma}^{-1}\left(1 - \boldsymbol{D}\tilde{\boldsymbol{S}}\boldsymbol{D}\boldsymbol{\Sigma}^{-1}\right)\mathrm{d}\boldsymbol{\Sigma} - 2\boldsymbol{D}^{-1}\mathrm{d}\boldsymbol{D} + \boldsymbol{\Sigma}^{-1}\mathrm{d}\left(\boldsymbol{D}\tilde{\boldsymbol{S}}\boldsymbol{D}\right)\right] = 0, \tag{20}$$

with

$$\mathrm{d}\left(\boldsymbol{D}\tilde{\boldsymbol{S}}\boldsymbol{D}\right) = \frac{1}{N}\sum_{j=1}^{N}\sum_{k=1}^{M}W^{j(k)}\left[(\mathrm{d}\boldsymbol{D}\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\mathrm{d}\boldsymbol{\beta})(\boldsymbol{D}\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\boldsymbol{\beta})^{\mathrm{T}} + (\boldsymbol{D}\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\boldsymbol{\beta})(\mathrm{d}\boldsymbol{D}\boldsymbol{Z}^{j(k)} - \boldsymbol{X}^j\mathrm{d}\boldsymbol{\beta})^{\mathrm{T}}\right], \tag{21}$$

where $\tilde{\boldsymbol{S}}$ now depends on both $\boldsymbol{D}$ and $\boldsymbol{\beta}$.

Performing conditional maximisation by fixing $\boldsymbol{\Sigma}$ (and hence $\boldsymbol{D}$), the value of $\hat{\boldsymbol{\beta}}$ satisfying equations (20) and (21) is

$$\hat{\boldsymbol{\beta}}(\boldsymbol{\Sigma}) = \left(\sum_{j=1}^{N}(\boldsymbol{X}^j)^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{X}^j\right)^{-1}\sum_{j=1}^{N}(\boldsymbol{X}^j)^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{D}\sum_{k=1}^{M}\left(W^{j(k)}\boldsymbol{Z}^{j(k)}\right), \tag{22}$$

which is almost the same as in (11) but with an extra factor $\boldsymbol{D}$ before the sum over $k$.

Maximisation with fixed $\boldsymbol{\beta}$ over $\boldsymbol{\Sigma}$ can in turn be done in two steps. The differential $\mathrm{d}\boldsymbol{\Sigma}$ is split into a diagonal $(2\boldsymbol{D}\mathrm{d}\boldsymbol{D})$ and an off-diagonal part. The condition for the latter to vanish is that $\boldsymbol{\Sigma}^{-1}(1 - \boldsymbol{D}\tilde{\boldsymbol{S}}\boldsymbol{D}\boldsymbol{\Sigma}^{-1})$ be a diagonal matrix, or equivalently that $(\boldsymbol{\Omega}^{-1} - \boldsymbol{\Omega}^{-1}\tilde{\boldsymbol{S}}\boldsymbol{\Omega}^{-1})$ is the diagonal matrix $\boldsymbol{A}$. As long as the diagonal elements of $\tilde{\boldsymbol{S}}$ are not too far from 1, a solution can be found by a simple iterative approach starting from an arbitrary $\boldsymbol{\Omega}_0$ and then solving for the diagonal matrix $\boldsymbol{A}$ the linear equations

$$\boldsymbol{\Omega}_{k+1} = \tilde{\boldsymbol{S}} + \boldsymbol{\Omega}_k\boldsymbol{A}\boldsymbol{\Omega}_k, \tag{23}$$

so that $\boldsymbol{\Omega}_{k+1}$ is in correlation form. Iterations are repeated until numerical convergence provides the required $\boldsymbol{\Omega}$. Since for fixed $\boldsymbol{D}$ the diagonal part of $\mathrm{d}\boldsymbol{\Sigma}$ is identically zero the steps above allow constrained maximisation to be performed for both (19) and (6).

The above procedure leads to a significant speed up with respect to numerical optimisation routines over the off-diagonal elements of $\boldsymbol{\Omega}$. Both methods involve inverting and multiplying $p \times p$ matrices, but the relative complexity of the numerical optimisation would be expected to grow at least as fast as the number of off-diagonal parameters, namely as $p(p-1)/2$ in dimension $p$. In a simple test comparing to the 'nlm' function of the stats package in R the target parameter was set to a noisy version of the identity matrix, which was used as starting point for both algorithms. The method here was over 13 times faster in dimension 4, nearly 40 times faster in dimension 6 and about 100 times for $p = 8$, highlighting the scale of improvement that can be expected, and is consistent with a growth like $p^2$ or better.

If $\boldsymbol{D}$ can vary, for the diagonal elements of $\mathrm{d}\boldsymbol{\Sigma}$ to vanish the matrix

$$\boldsymbol{A} - \boldsymbol{I} + \boldsymbol{\Omega}^{-1}\frac{1}{N}\sum_{j=1}^{N}\sum_{k=1}^{M}W^{j(k)}\boldsymbol{Z}^j(\boldsymbol{Z}^j)^{\mathrm{T}} - \boldsymbol{\Omega}^{-1}\boldsymbol{D}^{-1}\frac{1}{N}\sum_{j=1}^{N}\sum_{k=1}^{M}W^{j(k)}\boldsymbol{X}^j\boldsymbol{\beta}(\boldsymbol{Z}^j)^{\mathrm{T}}, \tag{24}$$

must have zero along the diagonal. The condition translates into a linear equation in the inverse elements of $\boldsymbol{D}$ and so can likewise be solved easily. The solution depends on $\boldsymbol{\Omega}$, which in turn depends (through $\tilde{\boldsymbol{S}}$) on $\boldsymbol{D}$. The unconstrained maximisation of (19) over $\boldsymbol{\Sigma}$ for a given $\boldsymbol{\beta}$ requires then cycling through solving (24) and (23). As such, the difference between constrained and unconstrained maximisation is made transparent.

### 3.6. Identifiability for specific formulations of the multivariate probit model

As pointed out in section 2.4 the identifiability of the parameters of a model is directly related to any invariance of the likelihood. In order to correctly evaluate the identifiability of a given model it is then crucial to account for any constraints explicitly or implicitly imposed on the parameter space. In an attempt to clarify sources of confusion, different formulations of the multivariate probit models are considered in detail. For clarity the different cases are summarised in Table 1.

13

Table 1: Special formulations of the multivariate probit model with a $p$-dimensional response variable $y = (y_1, \ldots, y_p)^\mathrm{T}$. The form of the design matrix with the covariates associated to each observation is provided, where however the observation index $j$ is dropped to simplify the notation. The scaling matrix $D$ is defined as $D = \mathrm{diag}(d_1, \ldots, d_p)$ with $d_i = \sqrt{\sigma_{ii}}$ and $\sigma_{ii}$ the $i$-th diagonal element of $\Sigma$. Likelihood invariance means that $\mathcal{L}(\beta, \Sigma) = \mathcal{L}(\lambda, \Omega)$ holds under the given transformation of the parameters. $E(Y) = \Phi_A(\cdot)$ is shorthand for the model definition in (4).

| | Design matrix | Regression coefficients | Likelihood invariance |
|---|---|---|---|
| General form | block diagonal | size $k = \sum_i k_i$ vector | |
| $E(Y) = \Phi_A(X\beta; \Sigma)$ | $X = \mathrm{diag}((x_1)^\mathrm{T}, \ldots, (x_p)^\mathrm{T})$ | $\beta = (\beta_1^\mathrm{T}, \ldots, \beta_p^\mathrm{T})^\mathrm{T}$ | $\Omega = D^{-1}\Sigma D^{-1}$ |
| | $x_i = (x_{i1}, \ldots, x_{ik_i})^\mathrm{T}$ | $\beta_i = (\beta_{i1}, \ldots \beta_{ik_i})^\mathrm{T}$ | $\lambda = \left(d_1^{-1}\beta_1^\mathrm{T}, \ldots, d_p^{-1}\beta_p^\mathrm{T}\right)^\mathrm{T}$ |
| Shared covariates | size $k$ vector | $p \times k$ matrix | |
| $E(Y) = \Phi_A(\beta X; \Sigma)$ | $X = (x_1, \ldots, x_k)^\mathrm{T}$ | $\beta = (\beta_1, \ldots, \beta_p)^\mathrm{T}$ | $\Omega = D^{-1}\Sigma D^{-1}$ |
| | | $\beta_i = (\beta_{i1}, \ldots \beta_{ik})^\mathrm{T}$ | $\lambda = D^{-1}\beta$ |
| Shared coefficients | $p \times k$ matrix | size $k$ vector | |
| $E(Y) = \Phi_A(X_\mathrm{c}\beta_\mathrm{c}; \Sigma)$ | $X_\mathrm{c} = (x_1, \ldots, x_p)^\mathrm{T}$ | $\beta_\mathrm{c} = (\beta_1, \ldots, \beta_k)^\mathrm{T}$ | $\Omega_\mathrm{c} = d_1^{-2}\Sigma$ |
| | $x_i = (x_{i1}, \ldots, x_{ik})^\mathrm{T}$ | | $\lambda_\mathrm{c} = d_1^{-1}\beta_\mathrm{c}$ |

The most general form of multivariate probit model described in Section 2.3 allows for a different number of covariates for each component of the response variable and consequently different vectors of regression coefficients. The design matrix for each observation is block diagonal, with one block, in the form of a row vector, for each response. Depending on the data at hand the model can be specialised in different ways. The covariates may be shared between the components of the response variables, while keeping the vectors of regression coefficients different. This is the case for example in Talhouk et al. (2012) and Xu and Craig (2010). Indeed this is simply a special case of the most general formulation since sharing the covariates does not change the number of free parameters of the problem, meaning that the dimension of the invariant space remains the same. Having a single vector of covariates however allows the problem to be represented in a slight different form, where the regression coefficients can be packed in a matrix with each row corresponding to a given response component, and the design matrix reduced to a single column vector. The parameter transformation yielding invariance of the likelihood can then be written in a more compact form, as summarised in Table 1. An important practical consequence is that a closed form solution can be derived for the M-step as pursued in Xu and Craig (2010), provided that no constraints are imposed.

Alternatively the regression coefficients, and the number of covariates, may be shared among the responses. The model chosen in Chib and Greenberg (1998) for the Six Cities dataset is of this type, the same number of covariates are observed for each component of the response variables, but they take different values. The model can be represented into a more compact form with a $p \times k$ design matrix $X_\mathrm{c}$, where each row corresponds to one component of the response variable and a $k$-dimensional vector of regression coefficients $\beta_\mathrm{c}$ (see Table 1). The conditional maximisation in Section 3.3.1 can be applied to maximise over the constrained space of $(\Sigma, \beta_\mathrm{c})$. The extreme case of fixing both the covariates and the regression coefficients would lead to a univariate probit model.

Fixing the vector of regression coefficients across response components accounts to imposing constraints at the modelling stage, with the effect of reducing the number of free parameters and hence the dimension of the invariant space. This aspect seems to have been overlooked in the literature where the Six Cities dataset is taken as an example, including Chib and Greenberg (1998). It is simply treated as a special case of the most general form, and a correlation structure is imposed on the covariance matrix, but in fact this is unnecessary for identifiability.

Consider the standard formulation where the design matrix $X$ is block diagonal, with each element of the same length $k$, and the vector of regression coefficients is a vector of length $p \times k$ with $p$ repeated sub-vectors $\beta_\mathrm{c}$. It has been noted in Section 3.5 that the transformation which guarantees invariance is such that each subvector $\beta_i$ is multiplied by a different positive factor $d_i$, violating the desired constraint that they be all equal. In order to avoid that the sub vectors of regression coefficients differ they all need to be multiplied by the same factor. In other words the likelihood is now left unchanged, independently of $X^j$, only when rescaling all the coordinate directions by the same amount, corresponding to a one dimensional invariant space. A reduced space can be defined by fixing the first diagonal element of the covariance matrix to 1, call $(\Omega_\mathrm{c}, \lambda_\mathrm{c})$ the corresponding parameters. An invariant $\tilde{Q}$ is obtained
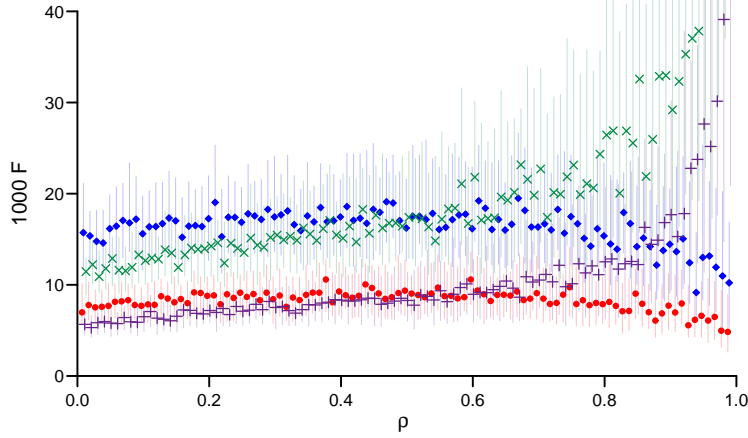
Figure 3: The median error $F$ in estimating the matrix of second moments in 4d as well as the inter-quartile ranges for the SMC sampler (diamonds) and a Gibbs sampler (crosses). Also included are runs with 4 times as many particles for the SMC (dots) and Gibbs (plusses) sampler.

by replacing $X^j, \Sigma$ and $\beta$ in (18) by $X_c^j, \Omega_c$ and $\lambda_c$ respectively and by setting all the elements of $D$ in (19) to be the square root of the first element of $\Sigma$. Constrained and unconstrained maximisation follow from the considerations in Section 3.5 but with the slight changes that only the first element of the matrix $A$ in (23) is non-zero and just the trace of (24) needs to be 0.

Maximising over an overly constrained space leads in general to a lower likelihood than when only imposing the conditions needed to ensure identifiability. Nevertheless, were the correlation form desired for modelling reasons, maximisation can be performed by setting $D$ to be the identity matrix and using $X_c^j$ in the formulae (22) and (23) above.

## 4. Results

### 4.1. Comparison of the SMC and Gibbs samplers for truncated multivariate normals

The SMC method for sampling truncated multivariate normals is now compared to a Gibbs sampler (Geweke, 1991; Robert, 1995) which is a Markov Chain where each component is sampled conditional on all the others. In dimension $p$ each 'pass' of the Gibbs sampler requires drawing $p$ univariate truncated normal variables which can be efficiently achieved by the mixed accept-reject algorithm of Geweke (1991) or for better efficiency near the mean using the tabulated accept-reject of Chopin (2011). The Gibbs sampler starts from the correct truncation region but, as noted in Geweke (1991), it can be rather slow at converging to the correct correlation structure. Convergence however is improved for extreme truncations since these have the effect of reducing the correlation among the components.

The SMC sampler on the other hand starts with the 'correct' correlation structure and moves to the required truncation region. Better performance could then be expected for correlated samples with the Gibbs sampler becoming preferable for more extreme truncation or lower correlation. A simulation study is conducted for the type of truncation regions that occur for the multivariate probit model in four dimensions, $p = 4$. The binary variables are set to $y_i = 1$, corresponding to the quadrant with positive $z$, and the vector of means is chosen as $\mu = (-1, -1, 1, 1)^T$ so that the mean is included in two dimensions and excluded in the other two. All the off-diagonal elements of the correlation matrix $\Sigma$ are set equal to $\rho$. The matrix of second moments is estimated using 10 million samples obtained by rejection sampling. Estimates are then obtained from the SMC and the Gibbs samplers, and a statistic $F$ is defined as the square root of the mean square distance between each estimate and the reference value from the rejection sampling. This process is repeated for a range of $\rho$ values.

Initially both algorithms are run to obtain 10 000 samples; for the Gibbs sampler this is the sample size after discarding the first fifth as burn in. For each value of $\rho$ we run both samplers 100 times and plot the median value of $F$ as well as the inter-quartile range in Fig. 3. The error in the Gibbs sampler increases with $\rho$ and starts to increase rapidly after $\rho \approx 0.75$. The error from the SMC sample on the other hand is fairly constant and actually starts to

decrease for large $\rho$. Both samplers then have similar performance for moderate correlations with the SMC approach notably better for large correlations. In the central range of $\rho \approx 0.5$ the truncation actually reduces the correlation so that the off-diagonal elements of the sample correlation matrix are only around 0.23 on average.

However the Gibbs sampler implementation is faster than the SMC version so that a Gibbs chain of about 50 000 passes can be obtained in the same time as the SMC sampler provides a sample of 10 000. Discarding the first fifth leaves a sample which is four times larger and whose error is roughly halved correspondingly. The growth in the error with $\rho$ still allows the SMC sampler to perform better, but now only for rather large values of $\rho$ above about 0.85, or where the actual sample correlations are above the more moderate value of approximately 0.56. Of course the actual computational time depends upon how efficiently each algorithm is implemented, but the simulation here suggests that the SMC sampler will have an advantage for high correlations.

Finally, SMC samples of size 40 000 are also obtained. The errors can again be observed to halve compared to the samples of 10 000. For both samplers in Fig. 3 the errors observed with the larger sample sizes resemble a scaled version of those from smaller sample sizes.

### 4.2. Application: the Six Cities dataset

To test the validity of our method, we treat the widely analysed data set from the Six Cities longitudinal study on the health effects of air pollution, for which a multivariate probit model was considered for a range of covariance structures by Chib and Greenberg (1998), who conducted both Bayesian and non-Bayesian analysis. Later Song and Lee (2005) proposed a confirmatory factor analysis for the same model, while more recently Craig (2008) used the example as a test case for a new method for geometrically reconstructing multivariate orthant probabilities which leads to an efficient evaluation of probabilities of the type in (4). As opposed to the MCMC procedure of Chib and Greenberg (1998), the SMC method provides estimates of the orthant probabilities as a by-product of the sampling, so that likelihoods are readily available for comparison to the results in Craig (2008). Moreover the SMC sampler produces a sample from the fitted distribution which is useful for further evaluation of intractable expectations of interest.

The Six Cities study was meant to model a probabilistic relation over time between the wheezing status of children, the smoking habit of their mother during the first year of observation and their age. In particular the subset of data considered for analysis refers to the observation of 537 children from Steubenville, Ohio. The wheezing condition $y_i^j$ of each child $j$ at age $i \in \{7, 8, 9, 10\}$ and the smoking habit $h^j$ of their mother are recorded as binary variables, with value 1 indicating the condition (wheezing/smoking) present. Three covariates are assumed for each component $i$, namely the age $x_{i1}^j = i - 9$ of child $j$ centred at 9, the smoking habit $x_{i2}^j = h^j$ and an interaction term $x_{i3}^j = (i - 9)h^j$ between the two. A probit model can then be constructed

$$\mathrm{pr}\{y_i^j = 1\} = \mathrm{pr}(z_i^j > 0) = \Phi\left[(\beta_0 + \beta_1 \cdot x_{i1}^j + \beta_2 \cdot x_{i2}^j + \beta_3 \cdot x_{i3}^j)\sigma_{ii}^{-\frac{1}{2}}\right],$$

where $z_i^j$ is the $i$th component of a multivariate random variable $\mathbf{Z}^j \sim \mathcal{N}(X_c^j \boldsymbol{\beta}, \boldsymbol{\Sigma})$ and $\Phi$ is the cumulative distribution function of a standard normal random variable.

Note that this is an example of a compact model as discussed in Section 3.6 and the invariant space is therefore only 1-dimensional. For identifiability it is then sufficient to fix only one of the diagonal elements of the covariance matrix, which in previous approaches has been overly restricted to be in correlation form instead. For comparison we therefore first perform constrained maximisation with $\boldsymbol{\Sigma}$ in correlation form using the methods in section 3.5.2 and then run the SMC EM algorithm with the correct invariance.

#### 4.2.1. Variance reduction

To reduce the variance associated with the stochastic nature of the Monte Carlo E step, the parameter can be updated according to a stochastic approximation type rule

$$\boldsymbol{\psi}^m = \boldsymbol{\psi}^{m-1} + \zeta_m(\hat{\boldsymbol{\psi}}^m - \boldsymbol{\psi}^{m-1}) \equiv (1 - \zeta_m)\boldsymbol{\psi}^{m-1} + \zeta_m\hat{\boldsymbol{\psi}}^m,$$

where $\hat{\boldsymbol{\psi}}_m$ is the actual estimate obtained from the M-step and $\zeta_m \in (0, 1)$ a stepsize with the purpose of gradually shifting the relative importance from the innovation $(\hat{\boldsymbol{\psi}}^m - \boldsymbol{\psi}^{m-1})$ to the value of the parameter $\boldsymbol{\psi}_{m-1}$ learnt through the previous iterations, and which therefore goes to 0. The scheme is like taking a weighted average of the previous

Table 2: Maximum likelihood estimates for the Six Cities dataset as obtained by using the constrained SMC algorithm with linearly increasing number of particles, after variance reduction and for a single run where the samples are recycled. Included for comparison are the results of Chib and Greenberg (1998) and Craig (2008). The value in brackets next to each estimate is the estimated standard error. The values of the parameters (and their errors) have all been multiplied by 1000. The last lines report the estimated $l(\psi)$ and corrected $\hat{l}(\psi)$ log-likelihoods, and a more accurate value obtained by numerical integrations.

| | Chib and Greenberg (1998) | | Craig (2008) | | linear increase | | variance reduction | | recycled samples | |
|---|---|---|---|---|---|---|---|---|---|---|
| $\beta_0$ | -1118 | (65) | -1122 | (62) | -1122 | (62) | -1123 | (62) | -1122 | (62) |
| $\beta_1$ | -79 | (33) | -78 | (31) | 79 | (31) | -79 | (31) | -78 | (31) |
| $\beta_2$ | 152 | (102) | 159 | (101) | 159 | (101) | 159 | (101) | 158 | (101) |
| $\beta_3$ | 39 | (52) | 37 | (51) | 37 | (51) | 38 | (51) | 37 | (51) |
| $\sigma_{12}$ | 584 | (68) | 585 | (66) | 583 | (66) | 583 | (66) | 581 | (66) |
| $\sigma_{13}$ | 521 | (76) | 524 | (72) | 522 | (71) | 522 | (71) | 523 | (71) |
| $\sigma_{14}$ | 586 | (95) | 579 | (74) | 577 | (74) | 578 | (74) | 579 | (73) |
| $\sigma_{23}$ | 688 | (51) | 687 | (56) | 686 | (56) | 686 | (56) | 682 | (57) |
| $\sigma_{24}$ | 562 | (77) | 559 | (74) | 558 | (74) | 558 | (74) | 558 | (74) |
| $\sigma_{34}$ | 631 | (77) | 631 | (67) | 626 | (67) | 627 | (67) | 625 | (67) |
| $l(\psi)$ | -794·94 | (0·69) | -794·93 | (0·66) | -794·91 | (0·59) | -794·95 | (0·82) | -794·86 | (0·66) |
| $\hat{l}(\psi)$ | -794·70 | | -794·72 | | -794·73 | | -794·61 | | -794·65 | |
| | {-794·749} | | {-794·738} | | {-794·742} | | {-794·740} | | {-794·748} | $(10^{-5})$ |

estimates, so we refer to it as a 'variation reduction' step. This way the monotonicity property of the EM algorithm is not guaranteed, but as long as the parameters remain within a neighbourhood of the maximum likelihood point where it can be approximated quadratically, monotonicity trivially follows from the convexity, so that in many practical cases this matter may not cause any issues.

### 4.2.2. Comparison to alternative approaches with covariance matrix in correlation form

To fit the model, a SMC sampler is implemented with the number of particles increasing linearly from 50 to 2000, over 40 iterations, followed by 10 further steps of variance reduction with 4000 particles. As for the tuning parameters of the algorithm the desired acceptance probability is set to $\alpha^\star = .6$ and the fraction $s$ defining the resampling threshold ESS$^\star$ as $s = .8$. Results for the constrained maximisation are presented in Table 2 along with those of Chib and Greenberg (1998) and Craig (2008). Good agreement can be observed both for the estimates and the standard errors. The latter are the square roots of the diagonal elements of the inverse observed Fisher information matrix, which in the case of missing data can be obtained from Louis' method (Louis, 1982).

Also given in Table 2 are average values of the corresponding log-likelihoods together with the standard deviation estimates over 50 runs. No real differences can be seen, with likelihoods comparable to, but slightly below, the estimate of -794·74 in Craig (2008). Due to the sampling noise the log-likelihood tends to be underestimated. A simple correction, discussed in Appendix C, consisting in adding half the variance over the runs, brings the estimates, $\hat{l}(\psi)$ in Table 2, closer to that in Craig (2008).

In the case of the Six Cities dataset the design matrix $X_c$ only takes two possible values, corresponding respectively to a smoking and non-smoking mother. Numerical calculation of the likelihood requires 32 regions to be evaluated for each of the parameter value $\psi$ estimated by the different methods. Numerical integration in dimension 4 is feasible and gives the results in curly brackets in Table 2, with accuracy $10^{-5}$, confirming that the SMC EM finds better parameter values than Chib and Greenberg (1998). We have also noticed that the estimates from other MCMC methods, such as in Song and Lee (2005), seem to be closer to those in Chib and Greenberg (1998), while ours are closer to the results of the exact method of Craig (2008).

Regardless of the method used for drawing the samples, our approach avoids the computationally expensive constrained maximisation employed by Chib and Greenberg (1998) by replacing it with the simple procedure in section 3.5.2. Despite the advantage with respect to standard numerical optimisation, the cost of either method is not significant compared to the sampling time. However failing to take advantage of the cyclicity of the trace would make the numerical optimisation substantially more involved. The evaluation of $Q$ would require quadratic forms as in

Table 3: Example maximum likelihood estimates for the Six Cities dataset obtained using the unconstrained SMC algorithm for non-invariant $Q$, invariant $\tilde{Q}$ and by fixing $\sigma_{11} = 1$. The standard deviations of the log-likelihood estimates are 0·90, 0·75 and 0·70 respectively, so that the corrected values of the likelihood $\hat{l}(\psi)$ are -792·97, -792·87, -792·825 and the numerical values with $10^{-5}$ accuracy are -792·849, -792·836, -792·834. The values of the parameters have all been multiplied by 1000.

| | $\beta_0$ | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\sigma_{12}$ | $\sigma_{13}$ | $\sigma_{14}$ | $\sigma_{22}$ | $\sigma_{23}$ | $\sigma_{24}$ | $\sigma_{33}$ | $\sigma_{34}$ | $\sigma_{44}$ | $l(\psi)$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $Q$ | -1176 | -84 | 159 | 41 | 647 | 592 | 572 | 1208 | 855 | 619 | 1255 | 715 | 1001 | -793·37 |
| $\tilde{Q}$ | -1235 | -113 | 168 | 47 | 664 | 622 | 612 | 1275 | 921 | 683 | 1383 | 802 | 1146 | -793·15 |
| fixed $\sigma_{11}$ | -1241 | -116 | 169 | 48 | 666 | 626 | 615 | 1279 | 927 | 686 | 1395 | 809 | 1158 | -793·07 |

(A.5) to be recalculated for each value of $\Sigma$ during the numerical routine, with an $O(N)$ increase in cost. This might be the reason behind the suggestion at the beginning of page 355 of Chib and Greenberg (1998) to redraw the latent variables between the two conditional maximisations for efficiency reasons.

Assuming that cyclicity of the trace is indeed exploited, a more significant gain comes instead from completing the maximisation by cycling through the conditional maximisations until convergence, rather then performing a single cycle. Completing the maximisation in fact reduces the number of full EM iterations needed.

A further major benefit of the SMC method is that the particle approximation can be updated after each M step and need not necessarily be resampled at each iteration, as described in section 3.2. The last column of Table 2 shows results obtained when recycling the samples in a SMC EM algorithm with 2000 particles and 40 iterations. Since oscillations before the variance reduction step are around 0·001 between iterations (with 2000 particles), parameter estimates when recycling the sample are essentially equivalent, at a much reduced computational cost. Updating the particle approximation is about 15 times faster than drawing the sample again, and the entire procedure proves to be about 5 times faster than a run with a linear increase of the number of particles drawn from scratch each time, with the latter strategy still enjoying a factor two improvement over keeping the number fixed at 2000. Similar parameter values are obtained when using fewer particles, but obviously with higher variance.

### 4.2.3. Unrestricted model

Strictly speaking the Six Cities model does not require $\Sigma$ to be in correlation form for identifiablity reasons. As discussed in Section 3.6 the invariant space is in fact only one-dimensional. To illustrate an application of the ideas presented in Section 3.5 a more general model which does not impose correlation form is analysed. To respect the invariance, one can either fix the first element of $\Sigma$ (i.e. set $\sigma_{11} = 1$) or run unconstrained maximisation (and project the results). Unconstrained maximisation was performed over 60 iterations with 4000 particles before the variance reduction step, since it may take longer for the EM algorithm to explore a larger space. A fairly robust point is found with the non-invariant $Q$, while the invariant $\tilde{Q}$ seems to lead to a flatter likelihood neighbourhood, with the solution appearing more sensitive to the number of particles during earlier iterations or on imposing the constraint of fixing $\sigma_{11}$ to 1. Results are given in Table 3, and again can be quite closely reproduced by recycling the samples in a sequential manner between parameter updates. Numerical integrations with $10^{-5}$ accuracy produces the values $-792.849, -792.836, -792.834$, when working respectively with $Q$, the invariant $\tilde{Q}$ or fixing $\sigma_{11} = 1$. For the latter two, despite the different parameter estimates, the likelihoods are essentially identical, and interestingly also higher than the one obtained for the non invariant $Q$. This seems practical evidence for an advantage in targeting the likelihood more directly through utilising the invariance.

Since there are only two possible forms for the design matrix $X_c$, a local symmetry arises, along with the global one, when $\beta_0\beta_3 = \beta_1\beta_2$. Moving near this symmetry may allow the EM algorithm to find different final maxima and explain the different parameter values found by using invariance or not. The local symmetry along with the estimation noise may be responsible for the non positive definiteness of the observed Fisher information (resulting from the difference of two positive definite matrices). Fixing a further parameter value, such as another diagonal element of $\Sigma$ to be 1, removes the local symmetry and allows standard errors to be obtained, centred around .10 and ranging from .045 to .16.

## 4.3. Higher dimensional simulated dataset

A higher dimensional example with simulated data is presented to show that the method scales reasonably well. The model chosen has the same formulation as the one used for the Six Cities case, corresponding to the third line in Table 1. The response variable is 8-dimensional with 7 covariates (including the intercept) associated to each component, resulting in a $8 \times 7$ design matrix. The entries different from the intercept are drawn from a uniform distribution on the interval $(-.5, .5)$. The parameters are set to

$$
\beta_c = \begin{pmatrix} 1.00 \\ 0.30 \\ -0.30 \\ 0.20 \\ -0.20 \\ 0.10 \\ -0.10 \end{pmatrix}, \Sigma = \begin{pmatrix} 1.00 & 0.10 & 0.10 & 0.10 & 0.10 & 0.20 & 0.20 & 0.40 \\ 0.10 & 1.20 & 0.10 & 0.10 & 0.10 & 0.20 & 0.30 & 0.40 \\ 0.10 & 0.10 & 1.20 & 0.10 & 0.20 & 0.20 & 0.30 & 0.40 \\ 0.10 & 0.10 & 0.10 & 1.10 & 0.20 & 0.20 & 0.30 & 0.50 \\ 0.10 & 0.10 & 0.20 & 0.20 & 1.10 & 0.20 & 0.30 & 0.50 \\ 0.20 & 0.20 & 0.20 & 0.20 & 0.20 & 0.90 & 0.40 & 0.60 \\ 0.20 & 0.30 & 0.30 & 0.30 & 0.30 & 0.40 & 0.90 & 0.60 \\ 0.40 & 0.40 & 0.40 & 0.50 & 0.50 & 0.60 & 0.60 & 0.80 \end{pmatrix}
$$

and 1000 observations are generated from the resulting model. Inference is then performed using both our SMC EM method and a Gibbs based MCEM approach. Since, as noted in Section 4.2.2, the cost of the M step with respect to the E step is relatively low in both cases, we retain our implementation of the M step rather than using numerical optimisation routines, despite marginally penalising the SMC EM algorithm in the comparison. The number of particles M is set to 4000 for both the SMC and the Gibbs sampler, and 40 iterations of the EM are performed. The square root of the mean squared distance of the estimated parameters from the real ones are found to be .079 and .080. However the distance between them is about .007, but the two clouds from different runs are barely distinguishable, since they have variations around .01 within them. Local symmetries are excluded in the simulated example so the standard errors could be estimated, ranging between about .045 and .172, and centred at .074 for the SMC, with similar values for the parameter values estimated via Gibbs, and in agreement with the actual distance to the real values. The SMC sampler has the advantage over Gibbs of automatically providing estimates of the likelihoods. Over 50 runs the average log-likelihoods were found to be around $-3280.3$ and $-3280.4$, with a standard deviation of .4. Although the likelihood for the parameters estimated via the SMC EM method happens to be marginally better in this case, the noise is too high for an accurate comparison and standard numerical integration is not an option due to the high dimensionality of the problem. A computational advantage still remains, since with the same number of particles the run time for the SMC EM is about two and a half times shorter than for the Gibbs based EM.

## 5. Conclusions

A new method based on sequential Monte Carlo samplers is introduced for the maximum likelihood estimation of multivariate probit models. In particular an adaptive sequential Monte Carlo algorithm is proposed to sample from high dimensional truncated normals. The proposal builds upon the property that a Student $t$ distribution approaches a Gaussian as the degrees of freedom go to infinity. When comparing to a Gibbs sampler the quality of the sample produced by the SMC sampler seems to be better for high correlation.

The typical iterative procedure of the EM algorithm appears like the ideal setting for SMC methods, which provide a natural machinery to evolve the particle approximation from one iteration to the next when updating the parameters in the M step. Performing the truncation to the current region starting from scratch can so be avoided. This way the computational cost is greatly reduced, with no particular loss in the performance, as seen for the example of the Six Cities dataset where similar parameter estimates are obtained when restarting at each iteration and with the fast sequential updating scheme.

Since by construction the sequential Monte Carlo sampler also provides samples from truncated Student distributions, it is clear that the method can be easily extended to a scenario where a Student $t$ distribution is assumed for the underlying latent variable of the probit model, rather than a normal distribution. Extensions to models with multinomial response variables are of course also possible.

Furthermore some of the confusion that has arisen around the maximisation step is clarified and the first complete EM algorithm for multivariate probit models is presented. Previously, methods typically proposed in the literature

have inevitably resulted in a generalised EM, while here the full maximisation is both easy to implement and efficient, with almost no computational cost. By examining the identifiability of such models we show that there is in fact a simple way to perform constrained maximisation, a process which is normally more computationally demanding. More importantly, we demonstrate how to tweak the EM algorithm so that it more directly targets increasing the likelihood. This is achieved by mimicking the invariance of the likelihood in the function $Q$ at the basis of the maximisation process, a strategy that should be of interest for other models.

An interesting alternative to EM for point estimation in the context of latent variable models, when neither the E step nor the M step are analytically tractable, is provided by a set of methods combining multiple imputation and simulated annealing ideas, as in Doucet et al. (2002); Gaetan and Yao (2003); Johansen et al. (2008). Sampling is then performed not only in the E step to impute the latent variables, but also in the M step to draw parameter values which are expected to converge to the maxima of the object function of interest. A desirable property of algorithms based on a stochastic version of the M step with respect to its deterministic counterpart is that they have a chance to escape local maxima. Obtaining multiple copies of the latent variables is essentially equivalent to drawing a sample in the E step of a standard Monte Carlo EM algorithm, therefore the same sampler can be applied. In the case of multivariate probit models then the SMC sampler of Section 3.1 would also be an option for the multiple imputations. Drawing the parameters to mimic the M step on the other hand may be non trivial, especially in higher dimensions. The difficulties lie in particular with ensuring that the identification constraints on the covariance matrix are met, as already noted by Chib and Greenberg (1998), and further discussed for example by McCulloch et al. (2000); Nobile (2000), in relation to multinomial probit models. More recently a parameter expanded method to simulate correlation matrices has been suggested by Liu and Daniels (2006). However in the context of multivariate probit models we show that performing the M step is actually pretty straightforward.

### Acknowledgements

### Appendix A. The $Q$ function for the probit model

For the multivariate probit model, substituting (5) into the expectation in (3) gives

$$Q(\psi, \psi^m) = \mathbb{E}_{\mathbf{Z}|\mathbf{Y},\psi^m}\left[l(\psi|\mathbf{Y},\mathbf{Z})\right] = \int_{z|\mathbf{y},\psi^m} \sum_{j=1}^{N} \log\left[I_{A^j}(z^j)\phi(z^j; \mathbf{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})\right] \cdot \prod_{l=1}^{N} \pi(z^l|\mathbf{y}^l, \psi^m)\, \mathrm{d}z^1 \cdots z^N$$

(A.1)

where $\phi(z^j; \mathbf{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})$ is the density of a multivariate normal distribution and $\pi(z^l|\mathbf{y}^l, \psi^m)$ is the density of a truncated multivariate normal constrained to the domain $A^l$. Denote it by $\mathrm{TMN}(z^l; A^l, \mathbf{X}^l\boldsymbol{\beta^m}, \boldsymbol{\Sigma^m})$ in the following. After inverting the order of integration and summation, and accounting for the fact that the integrals with respect to all the variables $z^l$ for $l \neq j$ can be independently evaluated and are normalised

$$\int_{\substack{z^l|\mathbf{y}^j,\psi^m \\ l \neq j}} \prod_{\substack{l=1 \\ l \neq j}} \mathrm{TMN}(z^l; A^l, \mathbf{X}^l\boldsymbol{\beta^m}, \boldsymbol{\Sigma^m}) \prod_{\substack{l=1 \\ l \neq j}} \mathrm{d}z^l = 1,$$

the integral in (A.1) then simplifies to

$$Q(\psi, \psi^m) = \sum_{j=1}^{N} \int_{z^j|\mathbf{y}^j,\psi^m} \log\left[I_{A^j}(z^j)\phi(z^j; \mathbf{X}^j\boldsymbol{\beta}, \boldsymbol{\Sigma})\right] \mathrm{TMN}(z^j, A^j, \mathbf{X}^j\boldsymbol{\beta^m}, \boldsymbol{\Sigma^m})\, \mathrm{d}z^j,$$

(A.2)

with $I_{A^j}(z^j) \equiv 1$ on the domain of integration since $\mathrm{TMN}(A^j, X^j \boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m)$ is only different from zero for $z^j \in A^j$. Substituting into (A.2) the expression for the density of a multivariate Gaussian density, and neglecting the proportionality constant term which is irrelevant for the maximisation,

$$\phi(z^j; X^j \boldsymbol{\beta}, \boldsymbol{\Sigma}) \propto |\boldsymbol{\Sigma}|^{(-1/2)} \exp\left(-\frac{1}{2}(z^{(j)} - X^j \boldsymbol{\beta})' \boldsymbol{\Sigma}^{-1}(z^j - X^j \boldsymbol{\beta})\right),$$

the $Q$ function becomes

$$Q(\psi, \psi^m) = -\frac{1}{2} \sum_{j=1}^{N} \int_{z^j|y^j, \psi^m} \left[\log |\boldsymbol{\Sigma}| + (z^j - X^j \boldsymbol{\beta})' \boldsymbol{\Sigma}_L^{-1}(z^j - X^j \boldsymbol{\beta})\right] \cdot \mathrm{TMN}(z^j, A^j, X^j \boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m) \, \mathrm{d}z^j. \tag{A.3}$$

The addends in the square brackets of (A.3) lead to two terms, the first of which can be simplified as

$$-\frac{1}{2} \log |\boldsymbol{\Sigma}| \sum_{j=1}^{N} \int_{z^j y^j, \psi^m} \mathrm{TMN}(z^j, A^j, X^j \boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m) \, \mathrm{d}z^j = -\frac{N}{2} \log |\boldsymbol{\Sigma}|. \tag{A.4}$$

By the cyclicity property of the trace of a matrix

$$(z^j - X^j \boldsymbol{\beta})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}(z^j - X^j \boldsymbol{\beta}) = \mathrm{tr}\{\boldsymbol{\Sigma}^{-1}(z^j - X^j \boldsymbol{\beta})(z^j - X^j \boldsymbol{\beta})^{\mathrm{T}}\}, \tag{A.5}$$

hence the second term of (A.3) can be written as

$$-\frac{1}{2} \sum_{j=1}^{N} \int_{z^j|y^j, \psi^m} \mathrm{tr}\{\boldsymbol{\Sigma}^{-1}(z^j - X^j \boldsymbol{\beta})(z^j - X^j \boldsymbol{\beta})^{\mathrm{T}}\} \cdot \mathrm{TMN}(A^j, X^j \boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m) \, \mathrm{d}z^j$$

$$= -\frac{1}{2} \mathrm{tr}\left\{\boldsymbol{\Sigma}^{-1} \sum_{j=1}^{N} \int_{z^j|y^j, \psi^m} (z^j - X^j \boldsymbol{\beta})(z^j - X^j \boldsymbol{\beta})^{\mathrm{T}} \cdot \mathrm{TMN}(A^j, X^j \boldsymbol{\beta}^m, \boldsymbol{\Sigma}^m) \, \mathrm{d}z^j\right\}$$

$$\equiv -\frac{1}{2} \mathrm{tr}\left\{\boldsymbol{\Sigma}^{-1} \sum_{j=1}^{N} \mathbb{E}_{\boldsymbol{Z}^j|\boldsymbol{Y}^j, \psi^m}\left[(\boldsymbol{Z}^j - X^j \boldsymbol{\beta})(\boldsymbol{Z}^j - X^j \boldsymbol{\beta})^{\mathrm{T}}\right]\right\}. \tag{A.6}$$

By combining equations (A.4) and (A.6) we obtain the final expression for the $Q$ function as in equation (6)

$$Q(\psi, \psi^m) = -\frac{N}{2}\left[\log |\boldsymbol{\Sigma}| + \mathrm{tr}\left\{\boldsymbol{\Sigma}^{-1} \frac{1}{N} \sum_{j=1}^{N} \mathbb{E}_{\boldsymbol{Z}^j|\boldsymbol{Y}^j, \psi^m}\left\{(\boldsymbol{Z}^j - X^j \boldsymbol{\beta})(\boldsymbol{Z}^j - X^j \boldsymbol{\beta})^{\mathrm{T}}\right\}\right\}\right].$$

## Appendix B. Algorithms

---

**Algorithm 1:** SMC SAMPLER
Key steps of a SMC sampler with a Random Walk Metropolis transition kernel and normalising constant estimation

---

Initialisation:     obtain a weighted particle approximation from the initial distribution

$$(W_0^{(k)}, \boldsymbol{Z}_0^{(k)}) \sim \pi_0(\theta_0) \tag{1}$$

set parameters $\theta_1, \boldsymbol{\Sigma}_1^{\mathrm{MH}}, \kappa_1, \zeta_1, \xi_1$ for a first move, $n = 1$
SMC core:     Repeat the following loop until $\theta_n \equiv \theta_T$ ($\pi_n \equiv \pi_T$)

Loop:        evaluate incremental weights

$$\tilde{w}_n(\mathbf{Z}_{n-1}^{(k)}, \mathbf{Z}_{n-1}^{(k)}) = \frac{\gamma_n(\mathbf{Z}_{n-1}^{(k)})}{\gamma_{n-1}(\mathbf{Z}_{n-1}^{(k)})} \tag{2}$$

update normalised weights

$$W_n^{(k)} \propto W_{n-1}^{(k)} \tilde{w}_n^{(k)} \tag{3}$$

update normalising constant estimate

$$\widehat{C}_n = \widehat{C}_{n-1} \sum_{k=1}^{M} W_{n-1}^{(k)} \tilde{w}_n^{(k)} \tag{4}$$

evaluate $\text{ESS}_n$ as a measure of the degree of degeneracy

$$\text{ESS}_n = \frac{1}{\sum\limits_{k=1}^{M} (W_n^{(k)})^2} \tag{5}$$

if $\text{ESS} < \text{ESS}^*$ resample

$$(W_n^{(k)}, \mathbf{Z}_{n-1}^{(k)}) \rightarrow \left(\frac{1}{M}, \tilde{\mathbf{Z}}_{n-1}^{(k)}\right) \sim \pi_n \tag{6}$$

MCMC step: $\forall k \in \{1, 2, \ldots, M\}$
      sample $\mathbf{Y}^k \sim \mathcal{N}(\mathbf{Z}_{n-1}^{(k)}, \mathbf{\Sigma}_n^{\text{MH}})$
      set $\mathbf{Z}_n^{(k)} = \mathbf{Y}^k$ with probability

$$\alpha^k = 1 \wedge \rho^k \tag{7}$$

where

$$\rho^k = \frac{\pi_n(\mathbf{Y}^k)}{\pi_n(\mathbf{Z}_{n-1}^{(k)})} \equiv \frac{\gamma_n(\mathbf{Y}^k)}{\gamma_n(\mathbf{Z}_{n-1}^{(k)})} \tag{8}$$

adapt scaling factor

$$\log(\kappa_{n+1}) = \log(\kappa_n) + \xi_n(\hat{\alpha}_n(\log(\kappa_n)) - \alpha^\star) \tag{9}$$

set new proposal covariance matrix $\mathbf{\Sigma}_{n+1}^{\text{MH}} = \kappa_n \widehat{\mathbf{\Sigma}}_{\pi_n}$
update the parameter identifying the next target

$$\theta_{n+1} = \left[\theta_n + \left(\zeta_n \frac{\text{ESS}_n - \text{ESS}_A^\star}{M} \vee \Delta\theta_{\min}\right)\right] \wedge \theta_T, \tag{10}$$

current particle approximation

$$(W_n^{(k)}, \mathbf{Z}_n^{(k)}) \text{ `` $\sim$ '' } \pi_n \tag{11}$$

go to next iteration

$$n = n + 1 \tag{12}$$

End of loop:    particle approximation available

$$(W_T^{(k)}, \mathbf{Z}_T^{(k)}) \text{ `` $\sim$ '' } \pi_T \tag{13}$$

---

*Further implementation details..* When targeting the multivariate Student $t$ distribution truncated to a domain $A$ as described in Section 3.1 the parameter $\theta = A$ can be defined as a vector of components $\theta = (\pm a_1, \ldots, \pm a_p)^{\text{T}}$ with the signs and direction of trunctation given by the observations. Similarly a vector $\theta_n = (\pm a_{1n}, \ldots, \pm a_{pn})^{\text{T}}$ defines the target region $A_n$ at iteration $n$. In practice the algorithm cycles through the dimensions one at the time, so we can focus on one particular component for a more detailed description. Drawing from an untruncated distribution effectively means to fix $a_{i0} = -\infty$. The first truncation points $a_{i1}$ can then be chosen for example by ensuring that a certain proportion of the probability mass of a multivariate normal distribution with independent components is preserved after

the truncation (to make sure that a non-negligible number of particles is kept). After the initialization the algorithm proceeds by updating each component according to equation (10) in Algorithm 1. The initial covariance matrix for the random walk Metropolis $\Sigma_1^{\text{MH}}$ is set equal to the covariance matrix target of the multivariate normal distribution (untruncated). Further tuning parameters in our runs were set as $\kappa_1 = 1, \zeta_1 = 2, \xi_1 = 7, \Delta\theta_{\min} = .02$, however they will depend on the particular scale of the problem, but they are not automatically tuned in our implementation. The resampling and adaptive thresholds $\text{ESS}^*$ and $\text{ESS}_A^\star$ are both set to $.8M$.

---

**Algorithm 2:** PROBIT SMC EM
Key steps of the SMC EM algorithm for multivariate probit models

---

Initialisation:     Set parameters

$$\psi^0 = (\Sigma^0, \beta^0), \quad m = 0 \tag{1}$$

obtain a sample (possibly weighted) from the initial distribution

$$(W_0^{(k)}, \mathbf{Z}_0^{(k)}) \sim \pi_0(\psi^0), \tag{2}$$

EM core:     Repeat the following loop until $\|\psi^m - \psi^{m-1}\|$ converges up to noise

Loop:     $m = m + 1, \tilde{\beta}^0 = \beta^0, \tilde{\Sigma}^0 = \Sigma^0, n = 0$

M-step:     Cycle through conditional maximisation until $\|\tilde{\beta}^n - \tilde{\beta}^{n-1}\| < \epsilon$

$n = n + 1$

update covariance matrix

$$\tilde{\Sigma}^n = \frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{M} W^{j(k)} (\mathbf{Z}^{j(k)} - X^j\tilde{\beta}^{n-1})(\mathbf{Z}^{j(k)} - X^j\tilde{\beta}^{n-1})^{\text{T}} \tag{3}$$

update regression coefficients

$$\tilde{\beta}^n = \Big( \sum_{j=1}^{N} (X^j)^{\text{T}} (\tilde{\Sigma}^n)^{-1} X^j \Big)^{-1} \sum_{j=1}^{N} (X^j)^{\text{T}} (\tilde{\Sigma}^n)^{-1} \sum_{k=1}^{M} \Big( W^{j(k)} \mathbf{Z}^{j(k)} \Big) \tag{4}$$

update parameter $\psi_m$ before E-step

$$\Sigma^m = \tilde{\Sigma}^n, \beta^m = \tilde{\beta}^n \tag{5}$$

E-step:     Implement a SMC sampler to move samples from $\pi_{m-1}(\psi^{m-1})$ to target $\pi_m(\psi^m)$
$\forall j \in \{1, 2, \ldots, N\}$
Rescale sample: $\forall k \in \{1, 2, \ldots, M\}$

$$\tilde{Z}_{m-1}^{(k)} = D^{-1} Z_{m-1}^{(k)} \tag{6}$$

with scaling $D$ such that

$$X^j\beta^m = D^{-1} X^j\beta^{m-1} \tag{7}$$

set covariance matrix

$$\tilde{\Sigma}^{m-1} = D^{-1} \Sigma^{m-1} D^{-1} \tag{8}$$

current particle approximation

$$(W_{m-1}^{(k)}, \tilde{\mathbf{Z}}_{m-1}^{(k)}) \text{ `` } \sim \text{ '' } \text{TMN}(A^j, X^j\beta^m, \tilde{\Sigma}^{m-1}) \tag{9}$$

build a SMC sampler to move from

$$\pi_0 = \text{TMN}(A^j, X^j\beta^m, \tilde{\Sigma}^{m-1}), \quad \theta_0 = \tilde{\Sigma}^{m-1} \tag{10}$$

to

$$\pi_T = \text{TMN}(A^j, X^j\beta^m, \Sigma^m), \quad \theta_T = \Sigma^m \tag{11}$$

## Appendix C. Log-likelihood correction

The log-likelihood is the sum of the log-probabilities of the regions corresponding to each observation, for which the SMC sampler provides noisy estimates. Assume the value returned is $p_j(1 + \xi_j)$ for each observation $j$, where $\xi_j$ are random relative noise variables with zero mean so that any bias is kept in the value $p_j$. The total log-likelihood is then

$$l = \sum_j \log\left[p_j(1 + \xi_j)\right] = \sum_j \log(p_j) + \log(1 + \xi_j) = \sum_j \log(p_j) + \xi_j - \frac{\xi_j^2}{2} + \ldots \tag{C.1}$$

with the logarithms expanded up to second order. Consider the sums over the random noises and their squares. From the central limit theorem, these should tend towards normal distributions with variances related to the sums of the second and fourth moments of the $\xi_j$. Assuming the relative errors $\xi_j \ll 1$, fourth and higher moments will decay quickly compared to the second, leading to the approximation

$$\sum_j \log(1 + \xi_j) \approx \zeta - \frac{\sigma^2}{2}, \tag{C.2}$$

where $\sigma^2$ corresponds to the variance between different runs of the log-likelihood estimate and $\zeta$ is a normally distributed random variable with the same variance $\zeta \sim N(0, \sigma^2)$. The log-likelihood should then be corrected to

$$\hat{l} = \sum_j \log(p_j) \approx \sum_j \log\left[p_j(1 + \xi_j)\right] + \frac{\sigma^2}{2}. \tag{C.3}$$

## References

Andrieu, C., Thoms, J., 2008. A tutorial on adaptive MCMC. Statistics and Computing 18.

Ashford, J.R., Sowden, R.R., 1970. Multi-variate probit analysis. Biometrics 26, 535–546.

Atchadé, Y.F., Rosenthal, J.S., 2005. On adaptive Markov chain Monte Carlo algorithms. Bernoulli 11, 815–828.

Beskos, A., Crisan, D., Jasra, A., 2012. On the stability of sequential Monte Carlo methods in high dimensions. Preprint, arXiv:1103.3965v2.

Bock, R.D., Gibbons, R.D., 1996. High-dimensional multivariate probit model. Biometrics 52, 1183–1194.

Chan, J.S.K., Kuk, A.Y.C., 1997. Maximum likelihood estimation for probit-linear mixed models with correlated random effects. Biometrics 53, 86–97.

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

Chopin, N., 2002. A sequential particle filter method for static models. Biometrika 89, 539–551.

Chopin, N., 2011. Fast simulation of truncated Gaussian distributions. Statistics and Computing 21, 275–288.

Craig, P., 2008. A new reconstruction of multivariate normal orthant probabilities. Journal of the Royal Statistical Society Series B 70, 227–243.

Del Moral, P., Doucet, A., Jasra, A., 2006. Sequential Monte Carlo samplers. Journal of the Royal Statistical Society Series B 68, 411–436.

Del Moral, P., Doucet, A., Jasra, A., 2007. Sequential Monte Carlo for Bayesian computations, in: Bayesian Statistics 8. Oxford University Press, pp. 1–34.

Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society Series B 39, 1–38.

Douc, R., Cappe, O., Moulines, E., 2005. Comparison of resampling schemes for particle filtering. Proceedings of the 4th International Symposium on Image and Signal Processing and Analysis , 64–69.

Doucet, A., de Freitas, N., Gordon, N., 2001. Sequential Monte Carlo methods in practice. Statistics for engineering and information science, Springer.

Doucet, A., Godsill, S., Andrieu, C., 2000. On sequential Monce Carlo sampling methods for Bayesian filtering. Statistics and Computing 10.

Doucet, A., Godsill, S.J., Robert, C.P., 2002. Marginal maximum a posteriori estimation using Markov chain Monte Carlo. Statistics and Computing 12, 77–84.

Emrich, L.J., Piedmonte, M.R., 1991. A method of generating high-dimensional multivariate binary variables. The American Statistician 45, 302–304.

Fearnhead, P., Taylor, B.M., 2013. An adaptive sequential Monte Carlo sampler. Bayesian Analysis 8, 1–28.

Gaetan, C., Yao, J.F., 2003. A multiple-imputation Metropolis version of the EM algorithm. Biometrika 90, 643–654.

Geweke, J., 1991. Efficient simulation from the multivariate normal and student-$t$ distributions subject to linear constraints. Computing Science and Statistics: Proceedings of the Twenty-Third Symposium on the Interface , 571–578.

Gilks, W.R., Roberts, G.O., Sahu, S.K., 1998. Adaptive Markov chain Monte Carlo through regeneration. Journal of the American Statistical Association 93, 1045 – 1054.

Gordon, N.J., Salmond, D.J., Smith, A.F.M., 1993. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. IEE-F 140, 107–113.

Gueorguieva, R.V., Agresti, A., 2001. A correlated probit model for joint modeling of clustered binary continuous responses. Journal of the American Statistical Association 96, 1102–1112.

Haario, H., Saksman, E., Tamminen, J., 2001. An adaptive Metropolis algorithm. Bernoulli 7, 223–242.

Hayes, J.F., Hill, W.G., 1981. Modification of estimates of parameters in the construction of genetic selection indices ('bending'). Biometrics 37, 483–493.

Higham, N.J., 2008. Functions of Matrices: Theory and Computation. Society for Industrial and Applied Mathematics.

Imai, K., van Dyk, D.A., 2005. A Bayesian analysis of the multinomial probit model using marginal data augmentation. Journal of Econometrics 124, 311–334.

Jasra, A., Stephens, D.A., Doucet, A., Tsagaris, T., 2011. Inference for Lévy driven stochastic volatility models via adaptive sequential Monte Carlo. Scandinavian Journal of Statistics 38, 1–22.

Johansen, A., Doucet, A., Davy, M., 2008. Particle methods for maximum likelihood estimation in latent variable models. Statistics and Computing 18, 47–57.

Johansen, A.M., Del Moral, P., Doucet, A., 2006. Sequential Monte Carlo samplers for rare events, in: Proceedings of 6th International Workshop on Rare Event Simulation, Bamberg, Germany. pp. 256–267.

Kalman, R.E., 1960. A new approach to linear filtering and prediction problems. Transactions of the ASME–Journal of Basic Engineering 82, 35–45.

Kitagawa, G., 1996. Monte Carlo filter and smoother for non-Gaussian non-linear state space model. Journal of Computational and Graphical Statistics 5, 1–25.

Kong, A., Liu, J.S., Wong, W.H., 1994. Sequential imputations and Bayesian missing data problems. Journal of the American Statistical Association 89.

Kuk, A.Y.C., Chan, J.S.K., 2001. Three ways of implementing the EM algorithm when parameters are not identifiable. Biometrical Journal 43, 207–218.

Kuk, A.Y.C., Nott, D.J., 2000. A pairwise likelihood approach to analyzing correlated binary data. Statistics and Probability Letters 47, 329–335.

Kushner, H.J., Yin, G.G., 2003. Stochastic approximation and recursive algorithms and applications. Applications of mathematics, Springer. second edition.

Li, Y., Schafer, D.W., 2008. Likelihood analysis of the multivariate ordinal probit regression model for repeated ordinal responses. Computational Statistical and Data Analysis 52, 3474–3492.

Liu, C., Rubin, D.B., Wu, Y.N., 1998. Parameter expansion to accelerate EM: The PX-EM algorithm. Biometrika 85, 755–770.

Liu, J.S., Chen, R., 1998. Sequential Monte Carlo methods for dynamic systems. Journal of the American Statistical Association 93, 1032–44.

Liu, X., Daniels, M.J., 2006. A new algorithm for simulating a correlation matrix based on parameter expansion and reparameterization. Journal of Computational and Graphical Statistics 15, 897–914.

Louis, T.A., 1982. Finding the observed information matrix when using the EM algorithm. Journal of the Royal Statistical Society Series B 44, 226–233.

McCulloch, C.E., 1994. Maximum likelihood variance components estimation for binary data. Journal of the American Statistical Association 89, 330–335.

McCulloch, R., Rossi, P.E., 1994. An exact likelihood analysis of the multinomial probit model. Journal of Econometrics 64, 207–240.

McCulloch, R.E., Polson, N.G., Rossi, P.E., 2000. A Bayesian analysis of the multinomial probit model with fully identified parameters. Journal of Econometrics 99, 173–193.

McLachlan, G.J., Krishnan, T., 2007. The EM algorithm and Extensions. Wiley Series in Probability and Statistics, Wiley. second edition.

Meng, X.L., Rubin, D.B., 1993. Maximum likelihood estimation via the ECM algorithm: A general framework. Biometrika 80, 267–278.

Montana, G., 2005. Hapsim: a simulation tool for generating haplotype data with pre-specified allele frequencies and ld coefficients. Bioinformatics 21, 4309–4311.

Nadarajah, S., Kotz, S., 2005. Sampling distributions associated with the multivariate $t$ distribution. Statistica Neerlandica 59, 214–234.

Natarajan, R., McCulloch, C.E., Kiefer, N.M., 2000. A Monte Carlo EM method for estimating multinomial probit models. Computational Statistics & Data Analysis 34, 33–50.

Nobile, A., 1998. A hybrid markov chain for the bayesian analysis of the multinomial probit model. Statistics and Computing 8, 229–242.

Nobile, A., 2000. Comment: Bayesian multinomial probit models with a normalization constraint. Journal of Econometrics 99, 335–345.

Petersen, K.B., Pedersen, M.S., 2012. The matrix cookbook.

Renard, D., Molenberghs, G., Geys, H., 2004. A pairwise likelihood approach to estimation in multilevel probit models. Computational Statistics & Data Analysis 44, 649–667.

Robert, C.P., 1995. Simulation of truncated normal variables. Statistics and Computing 5, 121–125.

Roberts, G.O., Gelman, A., Gilks, W.R., 1997. Weak convergence and optimal scaling of random walk Metropolis algorithms. The Annals of Applied Probability 7, 110–120.

Roberts, G.O., Rosenthal, J.S., 1998. Optimal scaling of discrete approximations to Langevin diffusions. Journal of the Royal Statistical Society Series B 60, 255–268.

Roberts, G.O., Rosenthal, J.S., 2001. Optimal scaling for various Metropolis-Hastings algorithms. Statistical Sciences 16, 351–367.

Schäfer, C., Chopin, N., 2013. Sequential Monte Carlo on large binary sampling spaces. Statistics and Computing 23, 163–184.

Song, X.Y., Lee, S.Y., 2005. A multivariate Probit latent variable model for analyzing dichotomous responses. Statistica Sinica 15, 645–664.

Talhouk, A., Doucet, A., Murphy, K., 2012. Efficient Bayesian inference for multivariate probit models with sparse inverse correlation matrices. Journal of Econometrics 21, 739–757.

Varin, C., Czado, C., 2010. A mixed autoregressive probit model for ordinal longitudinal data. Biostatistics 11, 127–138.

Wei, G.C.G., Tanner, M.A., 1990. A Monte Carlo implementation of the EM algorithm and the poor man's data augmentation algorithms. Journal

of the American Statistical Association 85, 699–704.

Wu, C.G.J., 1983. On the convergence properties of the EM algorithm. The Annals of Statistics 1, 95–103.

Xu, H., Craig, B.A., 2010. Likelihood analysis of multivariate probit models using a parameter expanded MCEM algorithm. Technometrics 52, 340–348.