



## MATRIX NORMAL BLOCK EM: THE R PACKAGE MNBEM

JAN DE LEEUW AND WEI TAN TSAI

ABSTRACT. The Expectation-Maximization algorithm is derived as a special case of the Majorization Method. We specialize this general derivation to multivariate normal distributions, emphasizing in particular direct sum and direct product structures for the dispersion matrix. A special case important in environmental statistics is missing data imputation in generalized growth curve models for the matrix variate normal distribution. The corresponding algorithms, with [R](#) code, are also given.

### 1. INTRODUCTION

The *majorization method* is a general approach, or family of approaches, to construct optimization methods. Some general publications about majorization are Kiers [1990]; De Leeuw [1994]; Heiser [1995]; Lange et al. [2000]; Hunter and Lange [2004]; De Leeuw and Lange [2009].

Suppose the problem is to minimize  $f : \mathcal{X} \Rightarrow \mathbb{R}$  over  $\mathcal{X} \subseteq \mathbb{R}^n$ . A function  $F : \mathcal{X} \otimes \mathcal{X} \Rightarrow \mathbb{R}$  is a *majorization function* if  $f(x) \leq F(x, y)$  for all  $x, y \in \mathcal{X}$  and  $f(x) = F(x, x)$  for all  $x \in \mathcal{X}$ .

The iterative *majorization algorithm* finds the update of  $x^{(k)}$  by computing

$$\mathcal{X}^{(k)} \triangleq \underset{x \in \mathcal{X}}{\operatorname{argmax}} F(x, x^{(k)}).$$

If  $x^{(k)} \in \mathcal{X}^{(k)}$  we stop. Else we select  $x^{(k+1)} \in \mathcal{X}^{(k)}$ . The *sandwich inequality*

$$f(x^{(k+1)}) \leq F(x^{(k+1)}, x^{(k)}) < F(x^{(k)}, x^{(k)}) = f(x^{(k)})$$

shows that the algorithm either stops, or produces a decreasing sequence of function values. Under compactness and continuity conditions this implies convergence [Zangwill, 1969].

Of course if we are maximizing  $f$ , then we can construct a suitable minorization function and maximize that in each iterative step. To cover both minorization and majorization Lange et al. [2000] propose the name *MM algorithm*, where the first  $M$  stands for either majorization or minorization, and the second  $M$  stands for either maximization or minimization.

Majorization and minorization functions are usually derived from classical inequalities, from Taylor's Theorem, or from convexity considerations. The *Expectation-Maximization* or *EM algorithm* is a family of MM algorithms based on Jensen's Inequality, usually applied in the statistical context of computing maximum likelihood estimates [Dempster et al., 1977; McLachlan and Krishnan, 2008]. The general idea of using MM algorithms in data analysis came about by realizing that the EM algorithm, based on Jensen's Inequality, and the SMACOF method for multi-dimensional scaling [De Leeuw, 1977], based on the Cauchy-Schwartz Inequality, were both examples of a more general approach to algorithm construction.

**1.1. EM as MM.** Suppose that  $g : X \otimes Y \Rightarrow \mathbb{R}^+$ , where  $X \subseteq \mathbb{R}^n$  and  $Y \subseteq \mathbb{R}^m$ . Define  $f : X \rightarrow \mathbb{R}^+$  by

$$f(x) \triangleq \log \int_Y g(x, y) dy.$$

The problem we study in this paper is maximization of  $f$  over  $X$ .

Suppose  $x, \tilde{x} \in X$ . We assume that if  $x \neq \tilde{x}$  then  $g(x, y) \neq g(\tilde{x}, y)$  for all  $y \in Y$ . Now

$$f(x) - f(\tilde{x}) = \log \frac{\int_Y g(x, y) dy}{\int_Y g(\tilde{x}, y) dy} = \log \frac{\int_Y g(\tilde{x}, y) \frac{g(x, y)}{g(\tilde{x}, y)} dy}{\int_Y g(\tilde{x}, y) dy}.$$

Let

$$h(x, y) \triangleq \frac{g(x, y)}{\int_Y g(x, y) dy}.$$

Then  $\int_Y h(x, y) dy = 1$  for all  $x$  and

$$f(x) - f(\tilde{x}) = \log \int_Y h(\tilde{x}, y) \frac{g(x, y)}{g(\tilde{x}, y)} dy.$$

Applying Jensen's Inequality to the right hand side gives

$$f(x) > f(\tilde{x}) + k(x, \tilde{x}) - k(\tilde{x}, \tilde{x}),$$

where we use the abbreviation

$$k(x, \tilde{x}) \triangleq \int_Y h(\tilde{x}, y) \log g(x, y) dy.$$

The function  $F(x, \tilde{x}) = f(\tilde{x}) + k(x, \tilde{x}) - k(\tilde{x}, \tilde{x})$  is the required minorization function.

This leads to the MM algorithm in which

$$\mathcal{X}^{(k)} \triangleq \underset{x \in X}{\operatorname{argmax}} F(x, x^{(k)}) = \underset{x \in X}{\operatorname{argmax}} k(x, x^{(k)}),$$

and  $x^{(k+1)} \in \mathcal{X}^{(k)}$ .

## 2. MULTINORMAL EM

**2.1. General.** In the case of a  $p$ -dimensional multinormal density the parameter space  $X$  is some set of mean vectors  $\mu$  and covariance matrices  $\Sigma$ . We write

$$g(\mu, \Sigma, y) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left\{-\frac{1}{2}(y - \mu)' \Sigma^{-1} (y - \mu)\right\}.$$

Thus

$$f(\mu, \Sigma) = \log \int_Y g(\mu, \Sigma, y) dy,$$

and

$$h(\mu, \Sigma, y) = \frac{g(\mu, \Sigma, y)}{\int_Y g(\mu, \Sigma, y) dy}.$$

The MM algorithm minimizes

$$(1) \quad \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \log |\Sigma| + \int_Y \tilde{h}(\tilde{\mu}, \tilde{\Sigma}, y) (y - \mu)' \Sigma^{-1} (y - \mu) dy.$$

If we let

$$(2a) \quad \tilde{m} \triangleq \int_Y h(\tilde{\mu}, \tilde{\Sigma}, y) y dy,$$

and

$$(2b) \quad \tilde{V} \triangleq \int_Y h(\tilde{\mu}, \tilde{\Sigma}, y) (y - \tilde{m})(y - \tilde{m})' dy,$$

then

$$(3) \quad \ell(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) = \log |\boldsymbol{\Sigma}| + \mathbf{tr} \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{V}} + (\tilde{\mathbf{m}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\tilde{\mathbf{m}} - \boldsymbol{\mu}).$$

**2.2. Marginal Constraints.** Suppose there is a  $z \in \mathbb{R}^p$  such that  $Y = \{z\} \otimes \mathbb{R}^q$ , with  $p + q = m$ . Integration is over the last  $q$  of the  $m$  coordinates of  $y = (z, u)$ , with the first  $p$  coordinates fixed at  $z$ . Thus

$$f(\boldsymbol{\mu}, \boldsymbol{\Sigma}, z) = \log \int_{\mathbb{R}^q} g(\boldsymbol{\mu}, \boldsymbol{\Sigma}, z, u) du.$$

In this case  $h$  is the conditional density of  $u$  given  $z$ . Thus

$$h(\boldsymbol{\mu}, \boldsymbol{\Sigma}, u, z) = \frac{1}{\sqrt{(2\pi)^q |\boldsymbol{\Sigma}_{u|z}|}} \exp\left\{-\frac{1}{2}(u - \boldsymbol{\mu}_{u|z})' \boldsymbol{\Sigma}_{u|z}^{-1} (u - \boldsymbol{\mu}_{u|z})\right\},$$

where

$$(4a) \quad \boldsymbol{\mu}_{u|z} = \boldsymbol{\mu}_u + \boldsymbol{\Sigma}_{uz} \boldsymbol{\Sigma}_{zz}^{-1} (z - \boldsymbol{\mu}_z),$$

$$(4b) \quad \boldsymbol{\Sigma}_{u|z} = \boldsymbol{\Sigma}_{uu} - \boldsymbol{\Sigma}_{uz} \boldsymbol{\Sigma}_{zz}^{-1} \boldsymbol{\Sigma}_{zu}.$$

It follows directly that

$$(5a) \quad \tilde{\mathbf{m}} = \begin{bmatrix} z \\ \tilde{\boldsymbol{\mu}}_u + \tilde{\boldsymbol{\Sigma}}_{uz} \tilde{\boldsymbol{\Sigma}}_{zz}^{-1} (z - \tilde{\boldsymbol{\mu}}_z) \end{bmatrix},$$

$$(5b) \quad \tilde{\mathbf{V}} = \begin{bmatrix} 0 & 0 \\ 0 & \tilde{\boldsymbol{\Sigma}}_{uu} - \tilde{\boldsymbol{\Sigma}}_{uz} \tilde{\boldsymbol{\Sigma}}_{zz}^{-1} \tilde{\boldsymbol{\Sigma}}_{zu} \end{bmatrix}.$$

From the computational point of view we can most easily compute  $\boldsymbol{\mu}_{u|z}$  and  $\boldsymbol{\Sigma}_{u|z}$  by applying the sweep operator, explained for example in Lange [1999, Chapter 7], to the joint dispersion matrix  $\tilde{\boldsymbol{\Sigma}}$ , sweeping the elements corresponding with  $z$ .

**2.3. Direct Sums.** Suppose  $\boldsymbol{\Sigma}$  has direct sum structure, i.e. there are  $n$  positive semi-definite matrices  $\boldsymbol{\Sigma}_i$ , or order  $m_i$ , such that

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & 0 & \cdots & 0 \\ 0 & \boldsymbol{\Sigma}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \boldsymbol{\Sigma}_n \end{bmatrix}.$$

The corresponding partition of  $\mu$  is

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}.$$

Then

$$f(\mu, \Sigma) = \sum_{i=1}^n \log \int_{Y_i} g(\mu_i, \Sigma_i, y) dy,$$

where  $Y_i \subseteq \mathbb{R}^{m_i}$ . Also

$$(6a) \quad \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \sum_{i=1}^n \left\{ \log |\Sigma_i| + \int_{Y_i} h_i(\tilde{\mu}_i, \tilde{\Sigma}_i, y) (y - \mu_i)' \Sigma_i^{-1} (y - \mu_i) dy \right\},$$

where

$$h_i(\tilde{\mu}_i, \tilde{\Sigma}_i, y) = \frac{g(\tilde{\mu}_i, \tilde{\Sigma}_i, y)}{\int_{Y_i} g(\tilde{\mu}_i, \tilde{\Sigma}_i, y) dy}.$$

Thus, with obvious indexing notation,

$$(6b) \quad \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \sum_{i=1}^n \left\{ \log |\Sigma_i| + \mathbf{tr} \Sigma_i^{-1} \tilde{V}_i + (\tilde{m}_i - \mu_i)' \Sigma_i^{-1} (\tilde{m}_i - \mu_i) \right\}.$$

This is basically the same as  $n$  replications of Equation (3).

**2.4. Repeated Independent Trials.** The results from Section 2.3 simplify if all  $\Sigma_i$  are equal to, say, the same positive definite matrix  $\Sigma$  of order  $m$ <sup>1</sup> Define  $\tilde{V} \triangleq \frac{1}{n} \sum_{i=1}^n \tilde{V}_i$ .

Then

$$(7a) \quad \frac{1}{n} \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \log |\Sigma| + \mathbf{tr} \Sigma^{-1} \tilde{V} + \frac{1}{n} \sum_{i=1}^n (\tilde{m}_i - \mu_i)' \Sigma^{-1} (\tilde{m}_i - \mu_i).$$

If we collect all  $\tilde{m}_i$  in the  $n \times m$  matrix  $\tilde{M}$  and all  $\mu_i$  in the  $n \times m$  matrix  $\Xi$ , then

$$(7b) \quad \frac{1}{n} \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \log |\Sigma| + \mathbf{tr} \Sigma^{-1} \tilde{V} + \frac{1}{n} \mathbf{tr} (\tilde{M} - \Xi) \Sigma^{-1} (\tilde{M} - \Xi)'$$

Define  $S(\Xi) \triangleq \frac{1}{n} (\tilde{M} - \Xi)' (\tilde{M} - \Xi)$ . Then

$$(7c) \quad \frac{1}{n} \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \log |\Sigma| + \mathbf{tr} \Sigma^{-1} \tilde{V} + \mathbf{tr} \Sigma^{-1} S(\Xi).$$

<sup>1</sup>Observe we use  $\Sigma$  both for the  $m \times m$  column-covariances, and for  $I \otimes \Sigma = \underbrace{\Sigma \oplus \cdots \oplus \Sigma}_{n \text{ times}}$ . It will be

obvious from the context which  $\Sigma$  we mean.

If the  $\mu_i$  are also equal we define  $\tilde{m} \triangleq \frac{1}{n} \sum_{i=1}^n \tilde{m}_i$ . We then find

$$(7d) \quad \frac{1}{n} \ell(\mu, \Sigma, \tilde{\mu}, \tilde{\Sigma}) = \log |\Sigma| + \mathbf{tr} \Sigma^{-1} \tilde{V} + (\tilde{m} - \mu)' \Sigma^{-1} (\tilde{m} - \mu) + \mathbf{tr} \Sigma^{-1} S,$$

where  $S \triangleq \frac{1}{n} \sum_{i=1}^n (\tilde{m}_i - \tilde{m})(\tilde{m}_i - \tilde{m})'$ .

**2.5. Direct Product (Kronecker) Structure.** Now consider the more general case in which the dispersion matrix is the  $nm \times nm$  matrix  $\Gamma \otimes \Sigma$ , where  $\Gamma$  is of order  $n$  and  $\Sigma$  is of order  $m$ . This defines the matrix variate normal distribution, discussed in detail in Gupta and Nagar [2000]. The repeated independent trials from Section 2.4 are the special case in which  $\Gamma = I$ . In the usual interpretation we have an  $n \times m$  matrix valued random variable  $\underline{Y}$ , and there is a row covariance matrix  $\Gamma$  and a column covariance matrix  $\Sigma$ , which combine to the direct product covariance matrix  $\Gamma \otimes \Sigma$  of all  $nm$  variables  $y_{ij}$ . Note that assuming Kronecker structure reduces the number of parameters in the dispersions from  $\frac{1}{2}nm(nm + 1)$  to  $\frac{1}{2}n(n + 1) + \frac{1}{2}m(m + 1)$ .

In environmental statistics direct product covariance structures are used to approximate the impractically large covariance matrices of space time time or vector-valued time series. In that context, covariance matrices with direct product, or Kronecker, structure are often called separable. See, for example, Matsuda and Yajima [2004], Lu and Zimmerman [2005] or Mitchell et al. [2006]. In the context of growth curve models Kronecker product dispersion matrices have been studied by many authors. Some recent interesting publications are Srivastava et al. [2008a,b]. For repeated measures data, Kronecker product covariance structures have been proposed by Naik and Rao [2001] and Roy and Khattree [2005]. In most of the papers cited, however, emphasis is on estimation and testing, and not on actual computation.

The log-likelihood for the matrix variate normal is

$$(8) \quad \log g(\Xi, \Gamma, \Sigma, Y) = m \log \Gamma + n \log \Sigma + \mathbf{tr} \Gamma^{-1} (Y - \Xi) \Sigma^{-1} (Y - \Xi)',$$

and the majorization function for multinormal EM becomes

$$(9) \quad \ell(\Xi, \Gamma, \Sigma, \tilde{\Xi}, \tilde{\Gamma}, \tilde{\Sigma}) = \\ = m \log |\Gamma| + n \log |\Sigma| + \mathbf{tr} (\Gamma^{-1} \otimes \Sigma^{-1}) \tilde{V} + \mathbf{tr} \Gamma^{-1} (\tilde{M} - \Xi)' \Sigma^{-1} (\tilde{M} - \Xi).$$

Observe that now  $\tilde{V}$  is of order  $nm$ , and it has a non-zero row and column for all missing elements of the  $n \times m$  data matrix.

## 3. MULTINORMAL MAXIMUM LIKELIHOOD USING BLOCK RELAXATION

The block relaxation method for multinormal maximum likelihood estimation [Oberhofer and Kmenta, 1974; De Leeuw, 1994] is designed to minimize

$$\mathcal{D}(\mu, \Sigma) = \log g(\mu, \Sigma, x) = \log |\Sigma| + (x - \mu)' \Sigma^{-1} (x - \mu)$$

over  $\mu \in \mathcal{M} \subseteq \mathbb{R}^m$  and  $\Sigma \in \mathcal{P} \in \mathbb{R}^{m \times m}$ . The algorithm alternates optimization of  $\mu$  and  $\Sigma$ , and computes its updates using the rule

$$(10a) \quad \mu^{(k+1)} \in \underset{\mu \in \mathcal{M}}{\operatorname{argmin}} (x - \mu)' [\Sigma^{(k)}]^{-1} (x - \mu),$$

$$(10b) \quad \Sigma^{(k+1)} \in \underset{\Sigma \in \mathcal{P}}{\operatorname{argmin}} \log |\Sigma| + \operatorname{tr} \Sigma^{-1} (x - \mu^{(k+1)}) (x - \mu^{(k+1)})'.$$

The update rules in (10) can be relaxed to define generalized block methods. It is not necessary to minimize in each of the two subproblems, it is enough to strictly decrease, as long as this decrease is done with continuous maps. So suppose we have two continuous maps  $F : \mathbb{R}^m \otimes \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^m$  and  $G : \mathbb{R}^m \otimes \mathbb{R}^{m \times m} \rightarrow \mathbb{R}^{m \times m}$ . Define  $\mu^{(k+1)} = F(\mu^{(k)}, \Sigma^{(k)})$  and  $\Sigma^{(k+1)} = G(\mu^{(k+1)}, \Sigma^{(k)})$ . Suppose

$$(x - \mu^{(k+1)})' [\Sigma^{(k)}]^{-1} (x - \mu^{(k+1)}) < (x - \mu^{(k)})' [\Sigma^{(k)}]^{-1} (x - \mu^{(k)}),$$

and

$$\begin{aligned} \log |\Sigma^{(k+1)}| + \operatorname{tr} [\Sigma^{(k+1)}]^{-1} (x - \mu^{(k+1)}) (x - \mu^{(k+1)})' &< \\ &< \log |\Sigma^{(k)}| + \operatorname{tr} [\Sigma^{(k)}]^{-1} (x - \mu^{(k+1)}) (x - \mu^{(k+1)})' \end{aligned}$$

This still defines a convergent algorithm.

**3.1. Block EM.** Using our EM results we can extend the block relaxation algorithm to minimizing

$$\mathcal{D}(\mu, \Sigma) = \log \int_Y g(\mu, \Sigma, y) dy$$

One possible sequence of steps is

$$(11a) \quad m^{(k)} = \int_Y h(\mu^{(k)}, \Sigma^{(k)}, y) y dy,$$

$$(11b) \quad V^{(k)} = \int_Y h(\mu^{(k)}, \Sigma^{(k)}, y) (y - m^{(k)}) (y - m^{(k)})' dy,$$

$$(11c) \quad \mu^{(k+1)} \in \underset{\mu \in \mathcal{M}}{\operatorname{argmin}} (m^{(k)} - \mu)' [\Sigma^{(k)}]^{-1} (m^{(k)} - \mu),$$

$$(11d) \quad \Sigma^{(k+1)} \in \underset{\Sigma \in \mathcal{P}}{\operatorname{argmin}} \log |\Sigma| + \operatorname{tr} \Sigma^{-1} \{V^{(k)} + (m^{(k)} - \mu^{(k+1)}) (m^{(k)} - \mu^{(k+1)})'\}.$$

There are many variations possible, however. We can alternate 11c and 11d a number of times, before we go back to 11a and 11b. We can do steps 11a and 11b after each step 11c and after each step 11d. And so on. In some cases it will even be possible to replace 11c and 11d by

$$(\boldsymbol{\mu}^{(k+1)}, \boldsymbol{\Sigma}^{(k+1)}) \in \underset{\boldsymbol{\mu} \in \mathcal{M}, \boldsymbol{\Sigma} \in \mathcal{P}}{\mathbf{argmin}} \log |\boldsymbol{\Sigma}| + \mathbf{tr} \boldsymbol{\Sigma}^{-1} \{V^{(k)} + (m^{(k)} - \boldsymbol{\mu})(m^{(k)} - \boldsymbol{\mu})'\}.$$

For instance if there are no constraints on  $\boldsymbol{\mu}$ , i.e.  $\mathcal{M} = \mathbb{R}^m$ , then  $\boldsymbol{\mu}^{(k+1)} = m^{(k)}$ , and

$$\boldsymbol{\Sigma}^{(k+1)} \in \underset{\boldsymbol{\Sigma} \in \mathcal{P}}{\mathbf{argmin}} \log |\boldsymbol{\Sigma}| + \mathbf{tr} \boldsymbol{\Sigma}^{-1} V^{(k)}.$$

**3.2. Kronecker Weights.** Minimizing (9) requires some extra thought. We will use a block algorithms with five substeps.

$$(12a) \quad m^{(k)} = \mathbf{vec}(M^{(k)}) = \int_Y h(\boldsymbol{\mu}^{(k)}, \Gamma^{(k)}, \boldsymbol{\Sigma}^{(k)}, y) \mathbf{vec}(Y) dy,$$

$$(12b) \quad V^{(k)} = \int_Y h(\boldsymbol{\mu}^{(k)}, \Gamma^{(k)}, \boldsymbol{\Sigma}^{(k)}, y) (\mathbf{vec}(Y) - m^{(k)}) (\mathbf{vec}(Y) - m^{(k)})' dy,$$

$$(12c) \quad \boldsymbol{\Xi}^{(k+1)} \in \underset{\boldsymbol{\Xi} \in \mathcal{M}}{\mathbf{argmin}} \mathbf{tr} \{[\Gamma^{(k)}]^{-1} (M^{(k)} - \boldsymbol{\Xi}) [\boldsymbol{\Sigma}^{(k)}]^{-1} (M^{(k)} - \boldsymbol{\Xi})'\},$$

$$(12d) \quad \Gamma^{(k+1)} \in \underset{\Gamma \in \mathcal{G}}{\mathbf{argmin}} m \log |\Gamma| + \mathbf{tr} (\Gamma^{-1} \otimes [\boldsymbol{\Omega}^{(k)}]^{-1}) V^{(k)} + \mathbf{tr} \Gamma^{-1} \hat{\Gamma}^{(k)},$$

$$(12e) \quad \boldsymbol{\Sigma}^{(k+1)} \in \underset{\boldsymbol{\Sigma} \in \mathcal{S}}{\mathbf{argmin}} n \log |\boldsymbol{\Sigma}| + \mathbf{tr} ([\Gamma^{(k+1)}]^{-1} \otimes \boldsymbol{\Omega}^{-1}) V^{(k)} + \mathbf{tr} \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\Sigma}}^{(k)},$$

where

$$(13a) \quad \hat{\Gamma}^{(k)} = (M^{(k)} - \boldsymbol{\Xi}^{(k+1)}) [\boldsymbol{\Sigma}^{(k)}]^{-1} (M^{(k)} - \boldsymbol{\Xi}^{(k+1)})',$$

$$(13b) \quad \hat{\boldsymbol{\Sigma}}^{(k)} = (M^{(k)} - \boldsymbol{\Xi}^{(k+1)})' [\Gamma^{(k+1)}]^{-1} (M^{(k)} - \boldsymbol{\Xi}^{(k+1)}).$$

The *piece de r sistance* is the term  $\mathbf{tr} (\Gamma^{-1} \otimes \boldsymbol{\Omega}^{-1}) V^{(k)}$ , which involves multiplication of two  $nm \times nm$  matrices. Since we have applications in mind where  $nm$  could be of order  $10^4$ , the matrices could grow unpleasantly large. We can use the Kronecker structure of  $\Gamma^{-1} \otimes \boldsymbol{\Omega}^{-1}$  and the sparseness of  $V^{(k)}$  to get considerable savings in computation and storage.

If we are updating  $\boldsymbol{\Sigma}$  in the block relaxation process we use

$$\mathbf{tr} (\Gamma^{-1} \otimes \boldsymbol{\Sigma}^{-1}) V^{(k)} = \mathbf{tr} \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \sum_{k=1}^n \gamma^{ik} V_{ik}^{(k)},$$



where the  $\gamma^{ik}$  are the elements of  $\Gamma^{-1}$  and the  $V_{ik}^{(k)}$  are the  $m \times m$  submatrices of  $V^{(k)}$  corresponding with row  $i$  and row  $k$  of  $Y$ . Note that element  $j, \ell$  of  $V_{ik}^{(k)}$  is nonzero if and only if both  $y_{ij}$  and  $y_{k\ell}$  are missing.

In the same way if we are updating  $\Gamma$  we use

$$\mathbf{tr} (\Gamma^{-1} \otimes \Sigma^{-1}) V^{(k)} = \mathbf{tr} \Gamma^{-1} \sum_{j=1}^m \sum_{\ell=1}^m \sigma^{j\ell} V_{j\ell}^{(k)},$$

where now the  $V_{j\ell}^{(k)}$  are the  $n \times n$  submatrices of  $V^{(k)}$  corresponding with column  $j$  and column  $\ell$  of  $Y$ .

It is probably a good idea to start the iterations by minimizing the log-likelihood (8) over  $\Sigma$ ,  $\Gamma$ , and  $\Xi$  and over the missing elements of  $Y$ . This avoids the huge Kronecker products altogether and can be assumed to give a good initial estimate of the structural parameters. At some point, we then let the algorithm (12) take over.

## REFERENCES

- J. De Leeuw. Applications of Convex Analysis to Multidimensional Scaling. In J.R. Barra, F. Brodeau, G. Romier, and B. Van Cutsem, editors, *Recent developments in statistics*, pages 133–145, Amsterdam, The Netherlands, 1977. North Holland Publishing Company.
- J. De Leeuw. Block Relaxation Methods in Statistics. In H.H. Bock, W. Lenski, and M.M. Richter, editors, *Information Systems and Data Analysis*, Berlin, 1994. Springer Verlag.
- J. De Leeuw and K. Lange. Sharp Quadratic Majorization in One Dimension. *Computational Statistics and Data Analysis*, 53:2471–2484, 2009.
- A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum Likelihood from Incomplete Data via the EM algorithm (with Discussion). *Journal of the Royal Statistical Society Series B*, 39:1–38, 1977.
- A.K. Gupta and D.K. Nagar. *Matrix Variate Distributions*. Chapman & Hall/CRC, Boca Raton, Florida, 2000.
- W.J. Heiser. Convergent Computing by Iterative Majorization: Theory and Applications in Multidimensional Data Analysis. In W.J. Krzanowski, editor, *Recent Advantages in Descriptive Multivariate Analysis*, pages 157–189. Clarendon Press, Oxford, 1995.
- D.R. Hunter and K. Lange. A Tutorial on MM Algorithms. *American Statistician*, 58(30–37), 2004.
- H. Kiers. Majorization as a Tool for Optimizing a Class of Matrix Functions. *Psychometrika*, 55:417–428, 1990.
- K. Lange. *Numerical Analysis for Statisticians*. Springer, Berlin, Heidelberg, New York, 1999.
- K. Lange, D.R. Hunter, and I. Yang. Optimization Transfer Using Surrogate Objective Functions. *Journal of Computational and Graphical Statistics*, 9:1–20, 2000.
- N. Lu and D.L. Zimmerman. The Likelihood Ratio Test for a Separable Covariance Matrix. *Statistics and Probability Letters*, 73:449–457, 2005.
- Y. Matsuda and Y. Yajima. On Testing for Separable Correlations of Multivariate Time Series. *Journal of Time Series Analysis*, 25:501–528, 2004.
- G.J. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. Wiley, New York, second edition, 2008.

- M.W. Mitchell, M.G. Genton, and M.L. Gumpertz. A Likelihood Ratio Test for Separability of Covariances. *Journal of Multivariate Analysis*, 97:1025–1043, 2006.
- D.N. Naik and S.S. Rao. Analysis of Multivariate Repeated Measures Data with a Kronecker Product Structured Covariance Matrix. *Journal of Applied Statistics*, 28:91–105, 2001.
- W. Oberhofer and J. Kmenta. A General Procedure for Obtaining Maximum Likelihood Estimates in Generalized Regression Models. *Econometrica*, 42:579–590, 1974.
- A. Roy and R. Khattree. On Implementation of a Test for Kronecker Product Covariance Structure for Multivariate Repeated Measures Data. *Statistical Methodology*, 2:297–306, 2005.
- M.S. Srivastava, T. von Rosen, and D. von Rosen. Models with a Kronecker Product Covariance Structure: Estimation and Testing. *Mathematical Methods of Statistics*, 17:357–370, 2008a.
- M.S. Srivastava, T. von Rosen, and D. von Rosen. Estimation in General Multivariate Linear Models with Kronecker Product Covariance Structure. Research Report 2008:1, Center of Biostatistics, Swedish University of Agricultural Sciences, 2008b.
- W. I. Zangwill. *Nonlinear Programming: a Unified Approach*. Prentice-Hall, Englewood-Cliffs, N.J., 1969.

## APPENDIX A. CODE

```

1 #
2 # mnem package
3 # Copyright (C) 2009 Jan de Leeuw <deleeuw@stat.ucla.edu>
4 # UCLA Department of Statistics, Box 951554, Los Angeles, CA 90095-1554
5 #
6 # This program is free software; you can redistribute it and/or modify
7 # it under the terms of the GNU General Public License as published by
8 # the Free Software Foundation; either version 2 of the License, or
9 # (at your option) any later version.
10 #
11 # This program is distributed in the hope that it will be useful,
12 # but WITHOUT ANY WARRANTY; without even the implied warrant(y) of
13 # MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
14 # GNU General Public License for more details.
15 #
16 # You should have received a copy of the GNU General Public License
17 # along with this program; if not, write to the Free Software
18 # Foundation, Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
19 #
20 #####
21 #
22 # version 0.1, 2009-06-06 Initial Release
23 # version 0.2, 2009-06-08 Added some stuff
24 #
25
26 # which(outer(i,j,function(x,y) is.na(x)&is.na(y)))
27
28 imputeK<-function(y,v,w,mu) {
29 n<-nrow(y); m<-ncol(y)
30 vsig<-kronecker(v,w); vy<-as.vector(t(y)); vmu<-as.vector(t(mu))
31 vim<-imputeY(vy,vsig,vmu)
32 return(list(vim=matrix(vim$yimp,n,m,byrow=TRUE),vzv=vim$vim))
33 }
34
35 sigmaHat<-function(vmat,v) {
36 n<-nrow(v); m<-nrow(vmat)/n; mm<-1:m; shat<-matrix(0,m,m)
37 for (i in 1:n) for (k in 1:n)
38 shat<-shat+v[i,k]*vmat[mm+(i-1)*m,mm+(k-1)*m]
39 return(shat)
40 }
41
42 gammaHat<-function(vmat,w) {
43 m<-nrow(w); n<-nrow(vmat)/m; mm<-(0:(n-1))*m; ghat<-matrix(0,n,n)
44 for (j in 1:m) for (l in 1:m)
45 ghat<-ghat+w[j,l]*vmat[mm+j,mm+l]
46 return(ghat)
47 }
48
49 # given a multivariate normal with mean mu and
50 # covariance matrix sig, and a vector of observations y,
51 # impute the missing values in y and the corresponding

```

```

52 # conditional covariance matrix
53
54 imputeY<-function(y,sig,mu) {
55   n<-length(y); nm<-is.na(y)
56   indi<-which(!nm); v<-matrix(0,n,n)
57   cs<-condiStat(sig,mu,y[indi],indi)
58   y[nm]<-cs$cmean; v[nm,nm]<-cs$cdisp
59   return(list(yimp=y,vimp=v))
60 }
61
62 # given a multivariate normal with mean mu and
63 # covariance matrix sig, compute the conditional
64 # mean and variance if we fix the variables indexed
65 # with indi to be equal to z
66
67 condiStat<-function(sig,mu,z,indi) {
68   m<-nrow(sig); jndi<-(1:m)[-indi]
69   mv<-rep(0,m); mv[indi]<-mu[indi]-z; mv[jndi]<-mu[jndi]
70   aa<-sweeper(cbind(sig,mv),indi)
71   return(list(cmean=aa[jndi,m+1],cdisp=aa[jndi,jndi]))
72 }
73
74 # beaton's sweep function
75
76 sweeper<-function(a,indi) {
77   n<-nrow(a); m<-length(indi)
78   for (j in indi) {
79     pv<-a[j,j]
80     if (pv == 0) next()
81     pr<-a[j,-j]; pc<-a[-j,j]
82     a[j,j] <- -1/pv
83     a[j,-j] <- pr/pv
84     a[-j,j] <- pc/pv
85     a[-j,-j] <- a[-j,-j]-outer(pc,pr)/pv
86   }
87   return(a)
88 }
89
90 # this minimizes trace{(y-mu)'(y-mu)} over the
91 # missing elements of y
92
93 imputeOLS<-function(y,mu) return(ifelse(is.na(y),mu,y))
94
95 imputeWLS<-function(y,k,ginv,sinv,mu,eps=1e-6,itmax=100,verbose=FALSE)
96 {
97   n<-nrow(y); m<-ncol(y); itel<-1
98   res<-y-mu; ras<-ginv%*%res%*%sinv; ff<-sum(ras*res)
99   wgt<-outer(diag(ginv),diag(sinv))
100  repeat {
101    thmax<-0
102    for (i in 1:n) for (j in 1:m) {
103      if (!is.na(k[i,j])) next()
104      rij<-ras[i,j]; wij<-wgt[i,j]

```

```

105     th<-rij/wij
106     thmax<-max(thmax,abs(th))
107     z[i,j]<-z[i,j]+th
108     ff<-ff-(rij^2)/wij
109     ras<-ras+th*outer(ginv[,i],sinv[,j])
110   }
111   if (verbose)
112     cat("itel",formatC(itel,format="d",width=4)," maxth",formatC(thmax,
113       format="f",digits=8,width=15)," func",formatC(ff,format="f",digits
114       =8,width=15),"\n")
113   if ((thmax < eps) || (itel == itmax)) break()
114   itel<-itel+1
115   }
116   return(y)
117 }

```

DEPARTMENT OF STATISTICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, CA 90095-1554

*E-mail address*, Jan de Leeuw: [deleeuw@stat.ucla.edu](mailto:deleeuw@stat.ucla.edu)

*URL*, Jan de Leeuw: <http://gifi.stat.ucla.edu>

*E-mail address*, Wei Tan Tsai: [tsai@stat.ucla.edu](mailto:tsai@stat.ucla.edu)

*URL*, Wei Tan Tsai: <http://www.stat.ucla.edu/~tsai/>