

Efficient Algorithms for Robust Estimation in Linear Mixed-Effects Models Using the Multivariate t -Distribution

José C. Pinheiro, Chuanhai Liu and

Bell Laboratories

Lucent Technologies

Murray Hill, NJ 07974

Yingnian Wu

Department of Statistics

University of Michigan

Ann Arbor, MI 48109

Abstract

Linear mixed-effects models are frequently used to analyze repeated measures data, because they model flexibly the within-subject correlation often present in this type of data. The most popular linear mixed-effects model for a continuous response assumes normal distributions for the random effects and the within-subject errors, making it sensitive to outliers. Such outliers are more problematic for mixed-effects models than for fixed-effects models, because they may occur in the random effects, in the within-subject errors, or in both, making them harder to be detected in practice. Motivated by a real dataset from an orthodontic study, we propose a robust hierarchical linear mixed-effects model in which the random effects and the within-subject errors have multivariate t -distributions, with known or unknown degrees-of-freedom, which are allowed to vary with subject. By using a gamma-normal hierarchical structure, our model allows the identification and classification of both types of outliers, comparing favorably to other multivariate t models for robust estimation in mixed-effects models previously described in the literature, which use only the marginal distribution of the responses. Allowing for unknown degrees-of-freedom, which may vary with subject and are estimated from the data, our model provides a balance between robustness and efficiency, leading to reliable results for valid inference. We describe and compare efficient EM-type algorithms, including ECM, ECME, and PX-EM, for maximum likelihood estimation in the robust multivariate t model. We compare the performance of the Gaussian and the multivariate t models under different patterns of outliers. Simulation results indicate that the multivariate t substantially outperforms the Gaussian model when outliers are present in the data, even in moderate amounts.

Key words: EM; ECM; ECME; PX-EM; Random effects; Repeated measures; Longitudinal data; Outliers.

1 Introduction

Linear mixed-effects models (Hartley and Rao, 1967) have become a popular tool for analyzing repeated measures data which arise in many areas as diverse as agriculture, biology, economics, and geophysics. The increasing popularity of these models is explained by the flexibility they offer in modeling the within-subject correlation often present in repeated measures data, by the handling of both balanced and unbalanced data, and by the availability of reliable and efficient software for fitting them (Wolfinger, Tobias and Sall, 1991; MathSoft, 1997). The most commonly used linear mixed-effects model for a continuous response was proposed by Laird and Ware (1982) and is expressed as

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \mathbf{e}_i, \quad i = 1, \dots, m, \quad (1)$$

where i is the subject index, \mathbf{y}_i is an n_i -dimensional vector of observed responses, \mathbf{X}_i and \mathbf{Z}_i are known $n_i \times p$ and $n_i \times q$ design matrices corresponding to the p -dimensional fixed effects vector $\boldsymbol{\beta}$ and the q -dimensional random effects vector respectively, and \mathbf{e}_i is an n_i -dimensional vector of within-subject errors. The \mathbf{b}_i are assumed to be independent with distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ and the \mathbf{e}_i are assumed to be independent with distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}_i)$, independent of the \mathbf{b}_i . The $\boldsymbol{\Psi}$ covariance matrix may be unstructured or structured – e.g. diagonal (Jennrich and Schluchter, 1986). The $\boldsymbol{\Lambda}_i$ matrices are typically assumed to depend on i only through their dimensions, being parametrized by a fixed, generally small, set of parameters $\boldsymbol{\rho}$ – e.g. an AR(1) covariance structure (Box, Jenkins and Reinsel, 1994). The most popular estimation methods for the parameters in model (1) are maximum likelihood and restricted maximum likelihood (Lindstrom and Bates, 1988). Confidence intervals and hypothesis tests for the parameters are generally based on asymptotic results (Miller, 1977).

Though model (1) offers great flexibility for modeling the within-subject correlation frequently present in repeated measures data, it suffers from the same lack of robustness against outlying observations as other statistical models based on the Gaussian distribution. An interesting feature of mixed-effects models is that outliers may occur either at the level of the within-subject error \mathbf{e}_i , called *e*–outliers, or at the level of the random effects \mathbf{b}_i , called *b*–outliers. In the first case, some unusual within-subject values are observed, whereas in the second case some unusual subjects are observed. Depending on the percentage of *e*–outliers and the number of observations per subject, it may not be possible to distinguish between the two cases.

A vast statistical literature exists on robust modeling methods, with some authors concentrating more on methods for outlier identification (Barnett and Lewis, 1994) and others on methods for outlier accommodation (Huber, 1981; Hampel, Ronchetti, Rousseeuw and Stahel, 1986). We follow here the robust statistical modeling approach described in Lange, Little and Taylor (1989) and consider a version of model (1) in which the multivariate normal distributions for the \mathbf{b}_i and the \mathbf{e}_i are replaced by multivariate *t*-distributions, with known or unknown degrees-of-freedom, which are allowed to vary with subject. This approach can be regarded as outlier-accommodating, though it also provides useful information for outlier identification.

A multivariate t linear mixed-effects model has been described by Welsh and Richardson (1997), but using only the marginal distribution of the response vectors, without reference to the hierarchical structure of the model. In particular, they do not derive, or discuss, the distributions of the random effects and the error terms under the multivariate t model, which help understanding the robustness of the model. In their description of estimation procedures, the degrees-of-freedom are assumed fixed and computational algorithms are not addressed.

A similar approach to the multivariate t model, but restricted to the distribution of the \mathbf{b}_i , has been considered by Wakefield, Smith, Racine-Poon and Gelfand (1994) and Racine-Poon (1992), within a Bayesian framework. Pendergast and Broffitt (1986) also have mentioned the multivariate t -distribution in connection with M-estimation for growth curve models. Robust estimation in mixed-effects models with variance components only (i.e. without covariance among random effects) using bounded influence estimators has been considered by Richardson and Welsh (1995) and Richardson (1997).

In Section 2, we describe growth curve data in which both \mathbf{b} - and \mathbf{e} -outliers seem to be present. The multivariate t version of model (1) is described in Section 3. In Section 4 we describe efficient EM-type algorithms for maximum likelihood estimation in the multivariate t linear mixed-effects model. We compare the robust maximum likelihood estimators obtained under the multivariate t -distribution to the Gaussian maximum likelihood estimators corresponding to model (1) in Section 5. Our conclusions and suggestions for further research are presented in Section 6.

2 An example: orthodontic distance growth in boys and girls

Our data come from an orthodontic study of 16 boys and 11 girls between the ages of 8 and 14 years and were originally reported in Potthoff and Roy (1964). The response variable is the distance (in millimeters) between the pituitary and the pterygomaxillary fissure, which was measured at 8, 10, 12, and 14 years for each boy and girl. Figure 1 presents a Trellis display (Becker, Cleveland and Shyu, 1996) of the data, along with individual least-squares fits of the simple linear regression model.

Figure 1 about here

Figure 1 reveals that the estimated slope for subject M13 is considerably larger than the remaining estimated slopes and that the responses for subject M09 are more variable around the fitted line. Overall, the responses for the boys vary more around the least squares lines, than do those for the girls. These features are more evident in the residuals plots by gender, displayed in Figure 2 and in the normal plots of the individual coefficients estimates, displayed in Figure 3. These plots suggest that two of the observations on subject M09 are \mathbf{e} -outliers and that subject M13 is a \mathbf{b} -outlier. Subject M10 is also identified in Figure 3 because he is indicated as a possible \mathbf{b} -outlier later in Section 5.1.

Figures 2 and 3 about here

Because both intercept and slope seem to vary with subject and the within-subject variation is larger among boys than girls, the following linear mixed-effects model can be used to describe the orthodontic distance growth with age.

$$y_{ij} = \beta_0 + \delta_0 I_i(F) + (\beta_1 + \delta_1 I_i(F)) t_j + b_{0i} + b_{1i} t_j + e_{ij}, i = 1, \dots, 27 \text{ and } j = 1, \dots, 4, \quad (2)$$

where y_{ij} denotes the orthodontic distance for the i th subject at age t_j , β_0 and β_1 denote respectively the intercept and the slope fixed effects for boys, δ_0 and δ_1 denote respectively the difference in intercept and slope fixed effects between girls and boys, $I_i(F)$ denotes an indicator variable for females, $\mathbf{b}_i = (b_{0i}, b_{1i})$ is the random effects vector for the i th subject, and e_{ij} is the within-subject error.

In Section 5.1, we compare the maximum likelihood estimates (MLEs) under the Gaussian version of the linear mixed-effects model (2) to the MLEs obtained under the multivariate t model described in Section 3.

3 A multivariate t linear mixed-effects model

The Gaussian linear mixed-effects model (1) can alternatively be written as:

$$\begin{bmatrix} \mathbf{y}_i \\ \mathbf{b}_i \end{bmatrix} \stackrel{\text{ind}}{\sim} \mathcal{N}_{n_i+q} \left(\begin{bmatrix} \mathbf{X}_i \boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \boldsymbol{\Lambda}_i & \mathbf{Z}_i \boldsymbol{\Psi} \\ \boldsymbol{\Psi} \mathbf{Z}_i' & \boldsymbol{\Psi} \end{bmatrix} \right), \quad i = 1, \dots, m, \quad (3)$$

with $\boldsymbol{\Lambda}_i = \boldsymbol{\Lambda}_i(\boldsymbol{\rho})$. For robust estimation of $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and $\boldsymbol{\rho}$, we proceed as in Lange et al. (1989) and replace the multivariate normal distribution in (3) with the multivariate t -distribution:

$$\begin{bmatrix} \mathbf{y}_i \\ \mathbf{b}_i \end{bmatrix} \stackrel{\text{ind}}{\sim} t_{n_i+q} \left(\begin{bmatrix} \mathbf{X}_i \boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \boldsymbol{\Lambda}_i & \mathbf{Z}_i \boldsymbol{\Psi} \\ \boldsymbol{\Psi} \mathbf{Z}_i' & \boldsymbol{\Psi} \end{bmatrix}, \nu_i \right), \quad i = 1, \dots, m, \quad (4)$$

where ν_i represents the multivariate t -distribution degrees-of-freedom (d.f.) for the i th subject. It follows from (4) that the \mathbf{y}_i are independent and marginally distributed as

$$\mathbf{y}_i \stackrel{\text{ind}}{\sim} t_{n_i}(\mathbf{X}_i \boldsymbol{\beta}, \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \boldsymbol{\Lambda}_i, \nu_i), \quad (5)$$

which provides yet another characterization of the multivariate t linear mixed-effects model. If $\boldsymbol{\Psi}$ is assumed to be diagonal and $\nu_i = \nu$ are fixed for all subjects, (5) reduces to the model considered in Welsh and Richardson (1997).

The multivariate t model (4) can also be expressed as the marginal distribution of $[\mathbf{y}_i', \mathbf{b}_i']'$ in the following hierarchical models:

$$\begin{aligned} \begin{bmatrix} \mathbf{y}_i \\ \mathbf{b}_i \end{bmatrix} \Big| \tau_i &\stackrel{\text{ind}}{\sim} \mathcal{N}_{n_i+q} \left(\begin{bmatrix} \mathbf{X}_i \boldsymbol{\beta} \\ \mathbf{0} \end{bmatrix}, \frac{1}{\tau_i} \begin{bmatrix} \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \boldsymbol{\Lambda}_i & \mathbf{Z}_i \boldsymbol{\Psi} \\ \boldsymbol{\Psi} \mathbf{Z}_i' & \boldsymbol{\Psi} \end{bmatrix} \right) \quad \text{and} \\ \tau_i &\stackrel{\text{ind}}{\sim} \Gamma\left(\frac{\nu_i}{2}, \frac{\nu_i}{2}\right), \quad i = 1, \dots, m, \end{aligned} \quad (6)$$

or

$$\begin{aligned} \mathbf{y}_i | \mathbf{b}_i, \tau_i &\stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i, \frac{1}{\tau_i} \boldsymbol{\Lambda}_i), \quad \mathbf{b}_i | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{0}, \frac{1}{\tau_i} \boldsymbol{\Psi}), \quad \text{and} \\ \tau_i &\stackrel{\text{ind}}{\sim} \Gamma\left(\frac{\nu_i}{2}, \frac{\nu_i}{2}\right), \quad i = 1, \dots, m. \end{aligned} \quad (7)$$

As described in the sequel, this gamma-normal hierarchical representation of the multivariate t model leads not only to natural EM implementations for maximum likelihood estimation of the unknown parameters, but also to diagnostic statistics that are useful for identification and classification of outliers.

It follows from (6) and (7) that the multivariate t model can be written as

$$\begin{aligned} \mathbf{y}_i &= \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \mathbf{e}_i, \quad i = 1, \dots, m \\ \mathbf{b}_i &\stackrel{\text{ind}}{\sim} t_q(\mathbf{0}, \boldsymbol{\Psi}, \nu_i) \quad \mathbf{e}_i \stackrel{\text{ind}}{\sim} t_{n_i}(\mathbf{0}, \boldsymbol{\Lambda}_i, \nu_i) \end{aligned} \quad (8)$$

with $\mathbf{b}_i | \tau_i$ independent of $\mathbf{e}_i | \tau_i$, implying that \mathbf{b}_i and \mathbf{e}_i are uncorrelated, but not independent, when $\nu_i < \infty$. The multivariate t model assumes that the random effects and the within-subject errors have multivariate t distributions and, therefore, can accommodate both \mathbf{b} -outliers and \mathbf{e} -outliers.

From standard properties of the multivariate t -distribution (Johnson and Kotz, 1972), it follows that, for $\nu_i > 2$,

$$\text{var}(\mathbf{b}_i) = \frac{\nu_i}{\nu_i - 2} \boldsymbol{\Psi} \quad \text{and} \quad \text{var}(\mathbf{e}_i) = \frac{\nu_i}{\nu_i - 2} \boldsymbol{\Lambda}_i, \quad i = 1, \dots, m.$$

Therefore, the interpretation of $\boldsymbol{\Psi}$ and $\boldsymbol{\Lambda}_i$ is different in the Gaussian model (1) and in the multivariate t model (4). Note, in particular, that $\text{var}(\mathbf{b}_i)$ is allowed to change with i in the multivariate t model, while it is independent of i in the Gaussian model. Provided $\nu_i > 1$ in (4), both models have $E(\mathbf{y}_i) = \mathbf{X}_i \boldsymbol{\beta}$, so that the fixed effects have the same interpretation: they represent the population average of the individual parameters.

Generally, some constraints are needed on the ν_i when these are to be estimated from the data. Common constraints are $\nu_i = \nu$ for all $i = 1, \dots, m$, or, more generally,

$$\nu_i = \lambda_{h(i)}, \quad i = 1, \dots, m, \quad (9)$$

where $h(i) \in \{1, \dots, l\}$ denotes the group to which the i th subject belongs and $\lambda_1, \dots, \lambda_l$ are l distinct positive scalar parameters, which can be treated as known, or unknown. We shall focus here on the t linear mixed-effects model (4) with the constraints (9).

Integrating out the \mathbf{b}_i in (7), we can express the distribution of \mathbf{y}_i as the marginal distribution of the following hierarchical model.

$$\mathbf{y}_i | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}\left(\mathbf{X}_i \boldsymbol{\beta}, \frac{1}{\tau_i} (\boldsymbol{\Lambda}_i + \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i')\right) \quad \text{and} \quad \tau_i \stackrel{\text{ind}}{\sim} \Gamma\left(\frac{\nu_i}{2}, \frac{\nu_i}{2}\right), \quad i = 1, \dots, m. \quad (10)$$

A useful consequence of (10) is that

$$\tau_i | \mathbf{y}_i \stackrel{\text{ind}}{\sim} \Gamma\left(\frac{\nu_i + n_i}{2}, \frac{\nu_i + \delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho})}{2}\right),$$

where

$$\delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho}) = (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})' (\mathbf{Z}_i\boldsymbol{\Psi}\mathbf{Z}_i' + \boldsymbol{\Lambda}_i)^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}). \quad (11)$$

Note that, in particular,

$$E(\tau_i | \mathbf{y}_i) = \frac{\nu_i + n_i}{\nu_i + \delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho})}.$$

4 Efficient EM algorithms for maximum likelihood estimation

In this section, we consider the maximum likelihood (ML) estimation of the parameters in the multivariate t linear mixed-effects model (4). We describe three EM-type algorithms for ML estimation with known and unknown degrees-of-freedom, based on two types of missing data structures. The first two algorithms use the hierarchical model (7) with both the \mathbf{b}_i and the τ_i treated as missing. The third algorithm is based on the hierarchical model (10) which, by integrating out the \mathbf{b}_i , has just the τ_i as missing data. The first two algorithms are computationally simpler, with closed-form expressions for the estimates of $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and $\boldsymbol{\rho}$ in the maximization step, but require additional assumptions about the structure of the $\boldsymbol{\Lambda}_i$ matrices. The last algorithm has a more computationally intensive maximization step, but allows more generality in the model specification and only requires minor modifications to existing software for fitting the Gaussian linear mixed-effects model (1). It should be noted that all three algorithms lead to the same MLEs (up to numerical round-off error) under the same structure of the $\boldsymbol{\Lambda}_i$ matrices.

Letting $\boldsymbol{\psi}$ denote a minimal set of parameters to determine $\boldsymbol{\Psi}$ (e.g. the upper triangular elements of $\boldsymbol{\Psi}$ in the unstructured case), we define the population parameters vector $\boldsymbol{\theta} = [\boldsymbol{\beta}', \boldsymbol{\psi}', \boldsymbol{\rho}', \boldsymbol{\lambda}']'$. Compared to the Gaussian linear mixed-effects model (1), the multivariate t model (4) allows each subject to have its own scale τ_i , which is unobserved and needs to be *imputed* from the data. The different individual scales result in different *weights* for estimating the population parameters $\boldsymbol{\theta}$. For example, conditional on $\boldsymbol{\Psi}$, $\boldsymbol{\rho}$, and the τ_i , the ML estimate of $\boldsymbol{\beta}$ minimizes $\sum_{i=1}^m \tau_i \delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho})$, with δ_i^2 as defined as (11). Because $E(\tau_i | \mathbf{y}_i)$ decreases with δ_i^2 , subjects with larger *residual sum of squares* δ_i^2 will have less weight in the determination of the ML estimates. The influence of δ_i^2 on the τ_i scales is controlled by the individual degrees-of-freedom ν_i – the smaller ν_i the larger the influence of δ_i^2 on τ_i .

4.1 The EM algorithm

The EM algorithm (Dempster, Laird and Rubin, 1977) is a popular iterative algorithm for ML estimation in models with incomplete data. More specifically, let \mathbf{y}_{obs} denote the observed data and \mathbf{y}_{mis} denote the missing data. The complete data $\mathbf{y}_{\text{com}} = (\mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}})$ is \mathbf{y}_{obs} augmented with \mathbf{y}_{mis} . We denote by $f(\mathbf{y}_{\text{com}} | \boldsymbol{\theta})$ the complete-data likelihood function of a parameter vector $\boldsymbol{\theta} \in \Theta$, by $L(\boldsymbol{\theta}) = f(\mathbf{y}_{\text{obs}} | \boldsymbol{\theta})$ the log-likelihood function and by $Q(\boldsymbol{\theta} | \boldsymbol{\theta}')$ the expected complete-data

log-likelihood

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}') = \text{E} \{ \ln [f(\mathbf{y}_{\text{com}}|\boldsymbol{\theta})] | \mathbf{y}_{\text{obs}}, \boldsymbol{\theta}' \}.$$

Each iteration of the EM algorithm consists of two steps, the *Expectation* step and the *Maximization* step:

E-step: Compute $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ as a function of $\boldsymbol{\theta}$;

M-step: Find $\boldsymbol{\theta}^{(t+1)}$ such that $Q(\boldsymbol{\theta}^{(t+1)}|\boldsymbol{\theta}^{(t)}) = \max_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$.

Each iteration of the EM algorithm increases the likelihood function $L(\boldsymbol{\theta})$ and, under mild conditions, the EM algorithm converges to a local or global maximum of $L(\boldsymbol{\theta})$ (Dempster et al., 1977; Wu, 1983).

When the M-step in the EM algorithm is difficult to implement, it is often useful to replace it with a sequence of constrained maximization (CM) steps, each of which maximizes $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ over $\boldsymbol{\theta}$ with some function of $\boldsymbol{\theta}$ held fixed. The sequence of CM-steps is such that the overall maximization is over the full parameter space. This leads to a simple extension of the EM algorithm, called the ECM algorithm (Meng and Rubin, 1993). A further extension of the EM algorithm is the ECME algorithm (Liu and Rubin, 1994). This algorithm replaces each CM-step of ECM with a CM-step that maximizes either the constrained Q function, as in ECM, or the correspondingly constrained L function. Liu and Rubin (1994) showed that ECME typically shares with EM the simplicity and stability, but has a faster rate of convergence, especially for the t distribution with unknown degrees of freedom.

4.2 EM algorithms with \mathbf{b}_i and τ_i as missing data

First consider the hierarchical multivariate t model (7) with both the \mathbf{b}_i and the τ_i as missing data. For simplicity, assume that

$$\boldsymbol{\Lambda}_i = \sigma_i^2 \mathbf{R}_i, \quad \sigma_i^2 = \sigma_{g(i)}^2, \quad i = 1, \dots, m, \quad (12)$$

with $g(i) \in \{1, \dots, k\}$ representing the group to which the i th subject belongs. The \mathbf{R}_i are known matrices, usually equal to the identity. We denote by $\boldsymbol{\sigma}^2$ the unique elements in $\{\sigma_1^2, \dots, \sigma_m^2\}$.

The within-subject covariance structure (12) allows for variance heterogeneity among different groups of subjects, but does not include serial correlation structures such as in ARMA models (Box et al., 1994).

4.2.1 ML estimation with known degrees-of-freedom using ECM

Let $\mathbf{y} = [\mathbf{y}'_1, \dots, \mathbf{y}'_m]'$, $\mathbf{b} = [\mathbf{b}'_1, \dots, \mathbf{b}'_m]'$, and $\boldsymbol{\tau} = [\tau_1, \dots, \tau_m]$. Under the constraints (12), the log-likelihood for the complete data in the multivariate t linear mixed-effects model (4) is

$$L(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\sigma}^2 | \mathbf{y}, \mathbf{b}, \boldsymbol{\tau}) = L_1(\boldsymbol{\beta}, \boldsymbol{\sigma}^2 | \mathbf{y}, \mathbf{b}, \boldsymbol{\tau}) + L_2(\boldsymbol{\Psi} | \mathbf{b}, \boldsymbol{\tau}) + \text{constant},$$

where

$$\begin{aligned}
L_1(\boldsymbol{\beta}, \boldsymbol{\sigma}^2 | \mathbf{y}, \mathbf{b}, \boldsymbol{\tau}) &= \sum_{i=1}^m \left[-\frac{n_i}{2} \ln \sigma_i^2 - \frac{\tau_i}{2\sigma_i^2} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i)' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i) \right] \\
&= -\sum_{i=1}^m \frac{n_i}{2} \ln \sigma_i^2 - \sum_{i=1}^m \frac{\tau_i}{2\sigma_i^2} \text{trace} [\mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{Z}_i \mathbf{b}_i) (\mathbf{y}_i - \mathbf{Z}_i \mathbf{b}_i)'] \\
&\quad + \sum_{i=1}^m \frac{\tau_i}{\sigma_i^2} \boldsymbol{\beta}' \mathbf{X}_i' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{Z}_i \mathbf{b}_i) - \sum_{i=1}^m \frac{\tau_i}{2\sigma_i^2} \boldsymbol{\beta}' \mathbf{X}_i' \mathbf{R}_i^{-1} \mathbf{X}_i \boldsymbol{\beta}
\end{aligned}$$

and

$$L_2(\boldsymbol{\Psi} | \mathbf{b}, \boldsymbol{\tau}) = -\frac{m}{2} \ln |\boldsymbol{\Psi}| - \frac{1}{2} \text{trace} \left(\boldsymbol{\Psi}^{-1} \sum_{i=1}^m \tau_i \mathbf{b}_i \mathbf{b}_i' \right).$$

Letting

$$\hat{\tau}_i = E(\tau_i | \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}, \mathbf{y}), \quad \hat{\mathbf{b}}_i = E(\mathbf{b}_i | \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}, \mathbf{y}), \quad \text{and} \quad \hat{\boldsymbol{\Omega}}_i = \tau_i \text{cov}(\mathbf{b}_i | \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}, \mathbf{y}),$$

we obtain

$$\begin{aligned}
\hat{\boldsymbol{\Omega}}_i &= \hat{\boldsymbol{\Psi}} - \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' (\mathbf{Z}_i \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' + \hat{\sigma}_i^2 \mathbf{R}_i)^{-1} \mathbf{Z}_i \hat{\boldsymbol{\Psi}} = \left(\hat{\boldsymbol{\Psi}}^{-1} + \frac{1}{\hat{\sigma}_i^2} \mathbf{Z}_i' \mathbf{R}_i^{-1} \mathbf{Z}_i \right)^{-1}, \\
\hat{\mathbf{b}}_i &= \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' (\mathbf{Z}_i \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' + \hat{\sigma}_i^2 \mathbf{R}_i)^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}) \\
&= \left(\hat{\boldsymbol{\Psi}}^{-1} + \frac{1}{\hat{\sigma}_i^2} \mathbf{Z}_i' \mathbf{R}_i^{-1} \mathbf{Z}_i \right)^{-1} \frac{1}{\hat{\sigma}_i^2} \mathbf{Z}_i' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}), \quad \text{and} \\
\hat{\tau}_i &= \frac{\nu_i + n_i}{\nu_i + \delta_i^2(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Psi}}, \hat{\boldsymbol{\sigma}}^2)}.
\end{aligned} \tag{13}$$

From standard multivariate analysis results (Fang and Zhang, 1990, p. 4) we have

$$\begin{aligned}
&[(\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})', \mathbf{b}_i'] \begin{bmatrix} \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \sigma_i^2 \mathbf{R}_i & \mathbf{Z}_i \boldsymbol{\Psi} \\ \boldsymbol{\Psi} \mathbf{Z}_i' & \boldsymbol{\Psi} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} \\ \mathbf{b}_i \end{bmatrix} \\
&= \mathbf{b}_i' \boldsymbol{\Psi}^{-1} \mathbf{b}_i + \frac{1}{\sigma_i^2} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i)' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i).
\end{aligned}$$

Replacing $\boldsymbol{\theta}$ and \mathbf{b}_i with their current estimates, we obtain the following useful decomposition:

$$\begin{aligned}
\delta_i^2(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Psi}}, \hat{\boldsymbol{\sigma}}^2) &= (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}})' \left(\mathbf{Z}_i \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' + \hat{\sigma}_i^2 \mathbf{R}_i \right)^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}}) \\
&= \hat{\mathbf{b}}_i' \hat{\boldsymbol{\Psi}}^{-1} \hat{\mathbf{b}}_i + \frac{1}{\hat{\sigma}_i^2} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i)' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i) = \hat{\delta}_{\mathbf{b}_i}^2 + \hat{\delta}_{\mathbf{e}_i}^2.
\end{aligned} \tag{14}$$

Equation (14) provides a simple way to compute $\delta_i^2(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\Psi}}, \hat{\boldsymbol{\sigma}}^2)$ as well as the weights $\hat{\tau}_i$. It also gives some insight on how the estimated random effects $\hat{\mathbf{b}}_i$ and the estimated residuals $\hat{\mathbf{e}}_i = \mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i$ affect the individual weights $\hat{\tau}_i$.

Using simple algebra we get

$$\begin{aligned} & \mathbb{E} \left[L_1(\boldsymbol{\beta}, \boldsymbol{\sigma}^2 | \mathbf{y}, \mathbf{b}, \tau) | \mathbf{y}, \hat{\boldsymbol{\theta}} \right] \\ &= - \sum_{i=1}^m \frac{n_i}{2} \ln \sigma_i^2 - \sum_{i=1}^m \frac{1}{2\sigma_i^2} \text{trace} \left[\mathbf{R}_i^{-1} \left(\hat{\tau}_i (\mathbf{y}_i - \mathbf{Z}_i \hat{\mathbf{b}}_i) (\mathbf{y}_i - \mathbf{Z}_i \hat{\mathbf{b}}_i)' + \mathbf{Z}_i \hat{\boldsymbol{\Omega}}_i \mathbf{Z}_i' \right) \right] \\ & \quad + \sum_{i=1}^m \frac{\hat{\tau}_i}{\sigma_i^2} \boldsymbol{\beta}' \mathbf{X}_i' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{Z}_i \hat{\mathbf{b}}_i) - \sum_{i=1}^m \frac{\hat{\tau}_i}{2\sigma_i^2} \boldsymbol{\beta}' \mathbf{X}_i' \mathbf{R}_i^{-1} \mathbf{X}_i \boldsymbol{\beta} \end{aligned}$$

and

$$\mathbb{E} \left[L_2(\boldsymbol{\Psi} | \mathbf{b}, \tau) | \mathbf{y}, \hat{\boldsymbol{\theta}} \right] = -\frac{m}{2} \ln |\boldsymbol{\Psi}| - \frac{1}{2} \text{trace} \left[\boldsymbol{\Psi}^{-1} \sum_{i=1}^m \left(\hat{\tau}_i \hat{\mathbf{b}}_i \hat{\mathbf{b}}_i' + \hat{\boldsymbol{\Omega}}_i \right) \right].$$

We then have the following ECM algorithm:

E-step: Given $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$, compute $\hat{\mathbf{b}}_i$, $\hat{\tau}_i$, and $\hat{\boldsymbol{\Omega}}_i$ for $i = 1, \dots, m$, using (13).

CM-step 1: Fix $\sigma_i^2 = \hat{\sigma}_i^2$ for $i = 1, \dots, m$ and update $\hat{\boldsymbol{\beta}}$ by maximizing $\mathbb{E} \left[L_1(\boldsymbol{\beta}, \hat{\boldsymbol{\sigma}}^2 | \mathbf{y}, \mathbf{b}, \tau) | \mathbf{y}, \hat{\boldsymbol{\theta}} \right]$ over $\boldsymbol{\beta}$, which leads to

$$\hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^m \frac{\hat{\tau}_i}{\hat{\sigma}_i^2} \mathbf{X}_i' \mathbf{R}_i^{-1} \mathbf{X}_i \right)^{-1} \sum_{i=1}^m \frac{\hat{\tau}_i}{\hat{\sigma}_i^2} \mathbf{X}_i' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{Z}_i \hat{\mathbf{b}}_i).$$

CM-step 2: Fix $\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}$ and update $\hat{\sigma}_i^2$ for $i = 1, \dots, m$ by maximizing $\mathbb{E} \left[L_1(\hat{\boldsymbol{\beta}}, \boldsymbol{\sigma}^2 | \mathbf{y}, \mathbf{b}, \tau) | \mathbf{y}, \hat{\boldsymbol{\theta}} \right]$ over σ_i^2 , which gives, for $j = 1, \dots, k$

$$\hat{\sigma}_j^2 = \sum_{i:g(i)=j} \left[\hat{\tau}_i (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i)' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i) + \text{trace}(\hat{\boldsymbol{\Omega}}_i \mathbf{Z}_i' \mathbf{R}_i^{-1} \mathbf{Z}_i) \right] / \sum_{i:g(i)=j} n_i.$$

CM-step 3: Update $\hat{\boldsymbol{\Psi}}$ by maximizing $\mathbb{E} \left[L_2(\boldsymbol{\Psi} | \mathbf{b}, \tau) | \mathbf{y}, \hat{\boldsymbol{\theta}} \right]$ over $\boldsymbol{\Psi}$, that is,

$$\hat{\boldsymbol{\Psi}} = \frac{1}{m} \sum_{i=1}^m \left(\hat{\tau}_i \hat{\mathbf{b}}_i \hat{\mathbf{b}}_i' + \hat{\boldsymbol{\Omega}}_i \right).$$

4.2.2 ML estimation with unknown degrees-of-freedom using ECME

When some, or all, of the degrees-of-freedom ν_1, \dots, ν_m are unknown, we can use the ECME algorithm that has the same E and CM steps as the ECM algorithm described in 4.2.1 for updating the estimates of $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and $\boldsymbol{\sigma}^2$ and an additional CML step that maximizes the constrained likelihood over the degrees-of-freedom with $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and $\boldsymbol{\sigma}^2$ fixed at their current estimates. When λ_j is unknown, the constrained likelihood is computed using

$$\mathbf{y}_i \stackrel{\text{ind}}{\sim} \mathbf{t}_{n_i} \left(\mathbf{X}_i \hat{\boldsymbol{\beta}}, \mathbf{Z}_i \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' + \hat{\sigma}_i^2 \mathbf{R}_i, \lambda_j \right), \quad \text{for } i \in \{i : g(i) = j\}.$$

More specifically, we have

CML-step: Update each unknown λ_j ($j = 1, \dots, l$) by maximizing

$$L_3(\lambda|\mathbf{y}, \boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\sigma}^2) = \sum_{i:h(i)=j} \left\{ \ln \left[\Gamma \left(\frac{\lambda + n_i}{2} \right) \right] - \ln \left[\Gamma \left(\frac{\lambda}{2} \right) \right] + \frac{\lambda}{2} \ln(\lambda) - \frac{\lambda + n_i}{2} \ln \left[\lambda + \delta_i^2 \left(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\Psi}}, \widehat{\boldsymbol{\sigma}}^2 \right) \right] \right\}$$

over λ . This requires only a one-dimensional search and can be obtained, for example, using the Newton-Raphson method (Thisted, 1988, §4.2.2).

4.2.3 Accelerating EM via parameter expansion

Liu, Rubin and Wu (1998) proposed the method of *Parameter Expansion* (PX) to accelerate EM-type algorithms and showed that the PX-EM algorithm shares the simplicity and stability of ordinary EM, but has a faster rate of convergence. The intuitive idea behind PX-EM is to use a *covariance adjustment* to correct the analysis of the M step, capitalizing on extra information captured in the imputed complete data. Technically, PX-EM expands the complete-data model $f(\mathbf{y}_{\text{com}} | \boldsymbol{\theta})$ to a larger model, $f_{\text{x}}(\mathbf{y}_{\text{com}} | \boldsymbol{\Theta})$, with $\boldsymbol{\Theta} = (\boldsymbol{\theta}_{\star}, \alpha)$, where $\boldsymbol{\theta}_{\star}$ plays the same role in $f_{\text{x}}(\mathbf{y}_{\text{com}} | \boldsymbol{\Theta})$ that $\boldsymbol{\theta}$ plays in $f(\mathbf{y}_{\text{com}} | \boldsymbol{\theta})$, and α is an auxiliary parameter which value is fixed at α_0 in the original model. Formally, two conditions must be satisfied. First, the observed-data model is preserved in the sense that, for all $\boldsymbol{\Theta}$, there is a common many-to-one reduction function R , such that $\mathbf{y}_{\text{obs}} | \boldsymbol{\Theta} \sim f\{\mathbf{y}_{\text{obs}} | \boldsymbol{\theta} = R(\boldsymbol{\Theta})\}$. Secondly, the complete-data model is preserved at the null value of α , α_0 , in the sense that, for all $\boldsymbol{\theta}$, $f_{\text{x}}\{\mathbf{y}_{\text{com}} | \boldsymbol{\Theta} = (\boldsymbol{\theta}_{\star}, \alpha_0)\} = f(\mathbf{y}_{\text{com}} | \boldsymbol{\theta} = \boldsymbol{\theta}_{\star})$. These conditions imply that if $\boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2$ then $\boldsymbol{\Theta}_1 \neq \boldsymbol{\Theta}_2$, and that, for all $\boldsymbol{\theta}$, there exists at least one $\boldsymbol{\Theta}$ such that $\mathbf{y}_{\text{obs}} | \boldsymbol{\Theta} \sim f\{\mathbf{y}_{\text{obs}} | \boldsymbol{\theta} = R(\boldsymbol{\Theta})\}$.

The PX-EM algorithm uses $f_{\text{x}}(\mathbf{y}_{\text{com}} | \boldsymbol{\Theta})$ to generate an EM algorithm, by iteratively maximizing the expected log-likelihood of $f_{\text{x}}(\mathbf{y}_{\text{com}} | \boldsymbol{\Theta})$. Specifically, let $\boldsymbol{\Theta}^{(t)} = (\boldsymbol{\theta}^{(t)}, \alpha_0)$ be the estimate of $\boldsymbol{\Theta}$ with $\alpha^{(t)} = \alpha_0$ from the t^{th} iteration. Then, at the $(t+1)^{\text{th}}$ iteration:

PX-E step: Compute $Q_{\text{x}}(\boldsymbol{\Theta} | \boldsymbol{\Theta}^{(t)}) = E_{\mathbf{y}_{\text{com}}} \{\log f_{\text{x}}(\mathbf{y}_{\text{com}} | \boldsymbol{\Theta}) | \mathbf{y}_{\text{obs}}, \boldsymbol{\Theta}^{(t)}\}$.

PX-M step: Find $\boldsymbol{\Theta}^{(t+1)} = \arg \max_{\boldsymbol{\Theta}} Q_{\text{x}}(\boldsymbol{\Theta} | \boldsymbol{\Theta}^{(t)})$; then apply the reduction function $R(\boldsymbol{\theta})$ to obtain $\boldsymbol{\theta}^{(t+1)} = R(\boldsymbol{\Theta}^{(t+1)})$.

The PX-EM algorithm can be used in the context of the multivariate t model to accelerate the EM algorithms described in Sections 4.2.1 and 4.2.2, by adjusting the M step using parameter expansions based on the imputed weights $\widehat{\tau}_i$ and the imputed random effects $\widehat{\mathbf{b}}_i$.

The imputed values of τ_i are only used in the ECM and ECME algorithms of Sections 4.2.1 and 4.2.2 to update the estimates of $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and σ_j^2 . The goodness-of-fit of the model $\tau_i \stackrel{\text{ind}}{\sim} \Gamma(\nu_i/2, \nu_i/2)$ to these values is ignored by the EM algorithms. We make use of this information to adjust the current estimates, by expanding the parameter space to include the scale parameter γ such that

$$\frac{\tau_i}{\gamma} \stackrel{\text{ind}}{\sim} \Gamma\left(\frac{\nu_i}{2}, \frac{\nu_i}{2}\right), \quad i = 1, \dots, m.$$

With the current estimate of γ fixed at $\gamma_0 = 1$, routine algebraic operations lead to the following CM-step for updating γ :

$$\hat{\gamma} = \frac{\sum_{i=1}^m \nu_i \hat{\tau}_i}{\sum_{i=1}^m \nu_i}.$$

Because

$$\mathbf{y}_i | \boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\sigma}^2, \gamma \stackrel{\text{ind}}{\sim} \mathbf{t}_{n_i} \left(\mathbf{X}_i \boldsymbol{\beta}, \frac{1}{\gamma} (\mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i' + \sigma_i^2 \mathbf{R}_i), \nu_i \right), \quad i = 1, \dots, m,$$

the application of the reduction function in the PX-EM algorithm leads to adjustments in the estimates of $\boldsymbol{\Psi}$ and $\boldsymbol{\sigma}^2$, which correspond to replacing their CM-steps in the previous ECM and ECME algorithms with

CM-step 2.X1:

$$\hat{\sigma}_j^2 = \frac{\sum_{i=1}^m \nu_i}{\sum_{i=1}^m \nu_i \hat{\tau}_i} \frac{\sum_{i:g(i)=j} \left[\hat{\tau}_i (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i)' \mathbf{R}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}} - \mathbf{Z}_i \hat{\mathbf{b}}_i) + \text{trace}(\hat{\boldsymbol{\Omega}}_i \mathbf{Z}_i' \mathbf{R}_i^{-1} \mathbf{Z}_i) \right]}{\sum_{i:g(i)=j} n_i}$$

for $j = 1, \dots, k$.

CM-step 3.X1:

$$\hat{\boldsymbol{\Psi}} = \frac{\sum_{i=1}^m \nu_i}{\sum_{i=1}^m \nu_i \hat{\tau}_i} \frac{\sum_{i=1}^m \left(\hat{\tau}_i \hat{\mathbf{b}}_i \hat{\mathbf{b}}_i' + \hat{\boldsymbol{\Omega}}_i \right)}{m}.$$

The PX-EM algorithm can also be used to adjust the current parameter estimates by making use of the information on the covariance matrices between \mathbf{y}_i and \mathbf{b}_i , given τ_i , that is, $\mathbf{Z}_i \boldsymbol{\Psi} / \tau_i$, for all $i = 1, \dots, m$. To do this, we expand the parameter space to include a $q \times q$ matrix $\boldsymbol{\zeta}$ in such a way that the complete-data model for \mathbf{y}_i becomes

$$\mathbf{y}_i | \mathbf{b}_i, \tau_i \stackrel{\text{ind}}{\sim} \mathbf{N}(\mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\zeta} \mathbf{b}_i, \frac{\sigma_i^2}{\tau_i} \mathbf{R}_i), \quad i = 1, \dots, m.$$

The covariance matrix between \mathbf{y}_i and \mathbf{b}_i given τ_i is then $\mathbf{Z}_i \boldsymbol{\zeta} \boldsymbol{\Psi} / \tau_i$.

Letting the current estimate of $\boldsymbol{\zeta}$ be $\boldsymbol{\zeta}_0 = \mathbf{I}_q$, the $q \times q$ identity matrix, and the other parameters be fixed at their current estimates, a CM-step for updating $\boldsymbol{\zeta}$ (together with $\boldsymbol{\beta}$) is obtained as follows.

CM-step 1.X1:

$$\begin{bmatrix} \hat{\boldsymbol{\beta}} \\ \text{vec}(\hat{\boldsymbol{\zeta}}) \end{bmatrix} = \left[\sum_{i=1}^m \frac{1}{\hat{\sigma}_i^2} \begin{pmatrix} \hat{\tau}_i \mathbf{X}_i' \mathbf{R}_i^{-1} \mathbf{X}_i & \hat{\tau}_i \mathbf{X}_i' \mathbf{R}_i^{-1} (\hat{\mathbf{b}}_i' \otimes \mathbf{Z}_i) \\ \hat{\tau}_i (\hat{\mathbf{b}}_i \otimes \mathbf{Z}_i') \mathbf{R}_i^{-1} \mathbf{X}_i & (\hat{\tau}_i \hat{\mathbf{b}}_i \hat{\mathbf{b}}_i' + \hat{\boldsymbol{\Omega}}_i) \otimes (\mathbf{Z}_i' \mathbf{R}_i^{-1} \mathbf{Z}_i) \end{pmatrix} \right]^{-1} \sum_{i=1}^m \frac{\hat{\tau}_i}{\hat{\sigma}_i^2} \begin{pmatrix} \mathbf{X}_i' \\ \hat{\mathbf{b}}_i \otimes \mathbf{Z}_i' \end{pmatrix} \mathbf{R}_i^{-1} \mathbf{y}_i,$$

where $\text{vec}(\widehat{\boldsymbol{\zeta}}) = (\widehat{\zeta}_{1,1}, \dots, \widehat{\zeta}_{q,1}, \dots, \widehat{\zeta}_{1,q}, \dots, \widehat{\zeta}_{q,q})'$ and \otimes stands for the Kronecker, or direct product, operator.

The application of the reduction function in PX-EM replaces the current estimate of $\boldsymbol{\Psi}$, $\widehat{\boldsymbol{\Psi}}$, with $\widehat{\boldsymbol{\zeta}}\widehat{\boldsymbol{\Psi}}\widehat{\boldsymbol{\zeta}}'$.

4.3 ML estimation integrating out the \mathbf{b}_i

The EM algorithms described in Section 4.2 provides closed form expressions for updating the estimates of $\boldsymbol{\theta}$, but require that the within-subject covariance matrices $\boldsymbol{\Lambda}_i$ be constrained to the form given in (12). A more flexible formulation, with no constraints on the $\boldsymbol{\Lambda}_i$, can be used when the \mathbf{b}_i are integrated out of the complete data likelihood, so that only the τ_i are treated as missing data, at the expense of a more computationally intensive CM-step. We describe here an ECME algorithm for this missing data scheme.

The log-likelihood of the complete data $[\mathbf{y}', \boldsymbol{\tau}']'$ in the multivariate t model (4) is

$$L(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho} | \mathbf{y}, \boldsymbol{\tau}) = L_1(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho} | \mathbf{y}, \boldsymbol{\tau}) + \text{constant},$$

where

$$L_1(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho} | \mathbf{y}, \boldsymbol{\tau}) = -\frac{1}{2} \sum_{i=1}^m [n_i \log |\mathbf{V}_i| + \tau_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})],$$

with $\mathbf{V}_i = \boldsymbol{\Lambda}_i + \mathbf{Z}_i \boldsymbol{\Psi} \mathbf{Z}_i'$. Letting $\widehat{\tau}_i$ be defined as in (13), it follows that

$$E \left[L_1(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho} | \mathbf{y}, \widehat{\boldsymbol{\theta}}) \right] = -\frac{1}{2} \sum_{i=1}^m [n_i \log |\mathbf{V}_i| + \widehat{\tau}_i (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})' \mathbf{V}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta})]$$

and, therefore, the following ECME algorithm can be used to obtain the MLEs of $\boldsymbol{\theta}$.

E-step: Given $\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}$, compute $\widehat{\tau}_i = (\nu_i + n_i) / \left[\nu_i + \delta_i^2(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\Psi}}, \widehat{\boldsymbol{\rho}}) \right]$, with $\delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho})$ as defined in (11).

CM-step: For fixed $\widehat{\boldsymbol{\tau}}$, update $\widehat{\boldsymbol{\beta}}$, $\widehat{\boldsymbol{\Psi}}$, and $\widehat{\boldsymbol{\rho}}$ by maximizing the function $E \left[L_1(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho} | \mathbf{y}, \widehat{\boldsymbol{\theta}}) \right]$ over $\boldsymbol{\beta}$, $\boldsymbol{\Psi}$, and $\boldsymbol{\rho}$.

The CM-step in this ECME algorithm is equivalent to maximum likelihood estimation in the Gaussian linear mixed-effects model $\mathbf{y}_i^* = \mathbf{X}_i^* \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{b}_i + \mathbf{e}_i$, $i = 1, \dots, m$, where $\mathbf{y}_i^* = \sqrt{\tau_i} \mathbf{y}_i$ and $\mathbf{X}_i^* = \sqrt{\tau_i} \mathbf{X}_i$. Reliable and efficient implementations of Newton-Raphson algorithms for obtaining the MLEs in the general Gaussian linear mixed-effects model (1) are available in commercial products such as SAS (PROC MIXED) and S-PLUS (lme function). These programs can be used to implement the ECME algorithm described here at low additional cost.

The decomposition of $\delta_i^2(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\Psi}}, \widehat{\boldsymbol{\rho}})$ given in (14) remains valid for general $\boldsymbol{\Lambda}_i$. That is,

$$\widehat{\delta}_i^2 = \delta_i^2(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\Psi}}, \widehat{\boldsymbol{\rho}}) = \widehat{\mathbf{b}}_i' \widehat{\boldsymbol{\Psi}}^{-1} \widehat{\mathbf{b}}_i + (\mathbf{y}_i - \mathbf{X}_i \widehat{\boldsymbol{\beta}} - \mathbf{Z}_i \widehat{\mathbf{b}}_i)' \widehat{\boldsymbol{\Lambda}}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \widehat{\boldsymbol{\beta}} - \mathbf{Z}_i \widehat{\mathbf{b}}_i) = \widehat{\delta}_{\mathbf{b}_i}^2 + \widehat{\delta}_{\mathbf{e}_i}^2,$$

where $\hat{\mathbf{b}}_i = \text{E}(\mathbf{b}_i | \mathbf{y}_i, \hat{\boldsymbol{\theta}}) = \hat{\boldsymbol{\Psi}} \mathbf{Z}_i' \hat{\mathbf{V}}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \hat{\boldsymbol{\beta}})$.

When the degrees-of-freedom ν_i are unknown, an additional CML-step, identical to the one described in Section 4.2.2, can be used to estimate the λ_j , $j = 1, \dots, l$. We have used the `lme` function to implement the ECME algorithm described here in S-PLUS. This implementation allows the degrees-of-freedom to be fixed in advance, or estimated from the data.

4.4 Inference based on the maximum likelihood estimates

One is generally interested in using MLEs to obtain confidence intervals and test hypotheses about the parameters. Because the distribution of the MLEs cannot be explicitly derived, approximate inference methods must be employed. The most common method uses the asymptotic normal approximation to the distribution of the MLEs (Miller, 1977; Lange et al., 1989). Other methods include the bootstrap (Efron and Tibshirani, 1993) and likelihood profiling (Bates and Watts, 1988). These last two methods usually give more accurate approximations, but are computationally intensive for the multivariate t model (4). This paper considers only confidence intervals and tests based on the normal approximation, concentrating on methods for the fixed effects $\boldsymbol{\beta}$.

Asymptotic confidence intervals and tests based on the MLEs can be obtained using either the observed or the expected Fisher information matrix. For the multivariate t model, these can be derived using the results in Appendix B of Lange et al. (1989). Let \mathbf{J} denote the expected Fisher information matrix for the marginal log-likelihood L of the multivariate t model and $\boldsymbol{\omega}$ denote the set of parameters excluding the fixed effects, so that $\boldsymbol{\theta} = (\boldsymbol{\beta}', \boldsymbol{\omega}')'$. It can be shown that

$$\mathbf{J}_{\boldsymbol{\beta}\boldsymbol{\beta}} = \text{E} \frac{\partial^2 L}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = \sum_{i=1}^m \frac{\nu_i + n_i}{\nu_i + n_i + 2} \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{X}_i \quad \text{and} \quad \mathbf{J}_{\boldsymbol{\beta}\boldsymbol{\omega}} = \text{E} \frac{\partial^2 L}{\partial \boldsymbol{\beta} \partial \boldsymbol{\omega}'} = \mathbf{0}.$$

It follows that the expected Fisher information matrix is block diagonal and, in particular, $[\mathbf{J}^{-1}]_{\boldsymbol{\beta}\boldsymbol{\beta}} = \mathbf{J}_{\boldsymbol{\beta}\boldsymbol{\beta}}^{-1}$. Asymptotic confidence intervals and hypothesis tests for the fixed effects are obtained assuming that the MLE $\hat{\boldsymbol{\beta}}$ has approximately a $\mathcal{N}_p(\boldsymbol{\beta}, \mathbf{J}_{\boldsymbol{\beta}\boldsymbol{\beta}}^{-1})$ distribution. In practice, $\mathbf{J}_{\boldsymbol{\beta}\boldsymbol{\beta}}$ is usually unknown and has to be replaced by its MLE $\hat{\mathbf{J}}_{\boldsymbol{\beta}\boldsymbol{\beta}}$.

4.5 Choosing starting values for the parameters

As with most iterative optimization procedures, initial values for the parameters in the multivariate t model must be provided to any of the EM-type algorithms described previously. A simple and generally successful algorithm for deriving initial estimates for the fixed effects $\boldsymbol{\beta}$ and the variance-covariance components $\boldsymbol{\Psi}$ and $\boldsymbol{\rho}$ is to fit separate regression models to each subject in the sample and to form “method of moments” estimates of the population parameters by averaging out the individual estimates. That is, letting $\hat{\boldsymbol{\beta}}_i$ and $\hat{\boldsymbol{\rho}}_i$ denote the individual parameter estimates obtained by fitting a linear regression to the data of the i th subject, $i = 1, \dots, m$, the initial values for the

EM-type algorithms are calculated as

$$\hat{\beta}_0 = \sum_{i=1}^m \hat{\beta}_i/m \quad \hat{\Psi}_0 = \sum_{i=1}^m \left(\hat{\beta}_i - \hat{\beta}_0 \right) \left(\hat{\beta}_i - \hat{\beta}_0 \right)' / (m-1) \quad \hat{\rho}_0 = \sum_{i=1}^m \hat{\rho}_i/m. \quad (15)$$

If the parameters in $\mathbf{\Lambda}_i$ vary according to which group $g(i) \in \{1, \dots, k\}$ subject i belongs (e.g. model (12)), separate initial estimates are obtained averaging over the separate groups

$$\hat{\rho}_j = \sum_{i:g(i)=j} \hat{\rho}_i/m_j,$$

where m_j denotes the number of subject in group j .

If the degrees-of-freedom λ_j for the multivariate t distributions are assumed unknown, initial values for them also need to be provided. It is generally enough to use a relative large initial value for the λ_j , say $\hat{\lambda}_0 = 40$, which corresponds to an initial assumption of near-normality for the random effects and within-subject errors.

The EM-type algorithms described in the previous sections tend to be robust to the choice of starting values for the parameters but, depending on characteristics of the data and of the model being used, it is possible that convergence to local optima occurs. Therefore, it is recommended that different starting values be used with the algorithms to assess the stability of the resulting estimates.

5 Comparing the Gaussian and the multivariate t MLEs

In this section we compare the MLEs under the Gaussian model (1) to the MLEs obtained under the multivariate t model (4). Firstly, we compare the Gaussian MLEs and the multivariate t MLEs for the orthodontic growth example of Section 2. The performance of the two sets of estimators are then compared under different outlier patterns, using results of a simulation.

5.1 The orthodontic growth example revisited

The distributional assumptions for the Gaussian version of the orthodontic growth model (2) are: $\mathbf{b}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{\Psi})$ and $e_{ij} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma_{g(i)}^2)$, with the \mathbf{b}_i independent of the e_{ij} . $g(i) = I_i(F) + 1$ denotes the gender group for the i th subject. The corresponding MLEs are given below.

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\delta}_0 \\ \hat{\beta}_1 \\ \hat{\delta}_1 \end{bmatrix} = \begin{bmatrix} 16.34 \\ 1.03 \\ 0.78 \\ -0.31 \end{bmatrix}, \quad \hat{\Psi} = \begin{bmatrix} 3.20 & -0.11 \\ -0.11 & 0.02 \end{bmatrix}, \quad \begin{bmatrix} \hat{\sigma}_1^2 \\ \hat{\sigma}_2^2 \end{bmatrix} = \begin{bmatrix} 2.63 \\ 0.45 \end{bmatrix}. \quad (16)$$

The corresponding approximate standard errors for the MLEs of the fixed effects, given by the square-roots of the diagonal elements of $(\sum_{i=1}^m \mathbf{X}_i' \hat{\mathbf{V}}_i^{-1} \mathbf{X}_i)^{-1}$, are

$$\hat{\sigma}(\hat{\beta}_0) = 1.111 \quad \hat{\sigma}(\hat{\delta}_0) = 0.097 \quad \hat{\sigma}(\hat{\beta}_1) = 1.334 \quad \hat{\sigma}(\hat{\delta}_1) = 0.115$$

The multivariate t version of model (2) has the following distributional assumptions:

$$\mathbf{b}_i | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{0}, \tau_i^{-1} \mathbf{\Psi}) \quad e_{ij} | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \tau_i^{-1} \sigma_{g(i)}^2) \quad \tau_i \stackrel{\text{ind}}{\sim} \Gamma(\lambda_{g(i)}/2, \lambda_{g(i)}/2)$$

which imply that $\mathbf{b}_i \stackrel{\text{ind}}{\sim} t(\mathbf{0}, \mathbf{\Psi}, \lambda_{g(i)})$, $e_{ij} \stackrel{\text{ind}}{\sim} t(0, \sigma_{g(i)}^2, \lambda_{g(i)})$.

As mentioned in Section 3, the parameters $\mathbf{\Psi}$, σ_1^2 , and σ_2^2 in the Gaussian model do not have the same interpretation as in the multivariate t model. To make the MLEs comparable, we consider the parameters $\text{var}(\mathbf{b}_i) = \mathbf{\Psi}_{g(i)}(t) = [\lambda_{g(i)} / (\lambda_{g(i)} - 2)] \mathbf{\Psi}$ and $\text{var}(e_{ij}) = \sigma_{g(i)}^2(t) = [\lambda_{g(i)} / (\lambda_{g(i)} - 2)] \sigma_{g(i)}^2$. The fixed effects $\beta_0, \delta_0, \beta_1, \delta_1$ have the same interpretation in both models: they represent the population average of the individual parameters and establish the growth patterns for an “average girl” and an “average boy” in the population. The MLEs for the multivariate t model are shown below.

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\delta}_0 \\ \hat{\beta}_1 \\ \hat{\delta}_1 \end{bmatrix} &= \begin{bmatrix} 16.83 \\ 0.54 \\ 0.73 \\ -0.25 \end{bmatrix}, \quad \hat{\mathbf{\Psi}}_1(t) = \begin{bmatrix} 4.79 & -0.16 \\ -0.16 & 0.03 \end{bmatrix}, \quad \hat{\mathbf{\Psi}}_2(t) = \begin{bmatrix} 3.13 & -0.11 \\ -0.11 & 0.02 \end{bmatrix} \\ \begin{bmatrix} \hat{\sigma}_1^2(t) \\ \hat{\sigma}_2^2(t) \end{bmatrix} &= \begin{bmatrix} 2.43 \\ 0.45 \end{bmatrix}, \quad \begin{bmatrix} \hat{\lambda}_1 \\ \hat{\lambda}_2 \end{bmatrix} = \begin{bmatrix} 5.78 \\ 6 \times 10^6 \end{bmatrix}. \end{aligned} \quad (17)$$

The corresponding approximate standard errors for the MLEs of the fixed effects, given by the square-roots of the diagonal elements of the $\hat{\mathbf{J}}_{\beta\beta}^{-1}$ matrix defined in Section 4.4, are

$$\hat{\sigma}(\hat{\beta}_0) = 0.895 \quad \hat{\sigma}(\hat{\delta}_0) = 0.078 \quad \hat{\sigma}(\hat{\beta}_1) = 1.158 \quad \hat{\sigma}(\hat{\delta}_1) = 0.099$$

These are consistently smaller than the corresponding estimated standard errors in the Gaussian model.

The multivariate t MLEs for the orthodontic growth model with unknown degrees-of-freedom were obtained using the three EM algorithms described in Section 4: the ECME algorithm of Section 4.2.2, its PX-EM version presented in Section 4.2.3, and the ECME algorithm of Section 4.3. Stand-alone implementations of the first two algorithms, written in C, were used to obtain the corresponding MLEs, while a modified version of the `lme` function in S-PLUS was used for the third algorithm, denoted by S-PLUS-ECME. Table 1 presents the *number of EM iterations* and the *user time* (on an SGI Challenge XL workstation running Iris 5.3) used to obtain the multivariate t MLEs in the orthodontic growth example, for each algorithm implementation. A relative tolerance of 10^{-7} for the parameter estimates was used as the convergence criterion for the three algorithms.

Table 1 about here

Because the implementations use languages with very different characteristics (compiled C and interpreted S-PLUS), the user times in Table 1 are not directly comparable, but give a sense of the actual performance of the algorithms in a practical setting.

Comparing these estimates to the Gaussian MLEs in (16), we see that the estimates of the incremental parameters (δ_0 and δ_1) fixed effects and the boys' random effects covariance matrix $\Psi_1(t)$ are considerable different. The multivariate t MLEs of δ_0 and δ_1 are respectively 50% smaller and 20% larger than the corresponding Gaussian MLEs. The boys' random effects variances multivariate t MLEs are 50% larger than the Gaussian MLEs. The MLEs of the girls' parameters are essentially unchanged. Using

$$t(\boldsymbol{\mu}, \boldsymbol{\sigma}^2, \nu) \xrightarrow{\nu \rightarrow \infty} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2), \quad (18)$$

it is clear that the estimated degrees-of-freedom $\hat{\lambda}_2$ indicate that a Gaussian model is adequate for the girls' orthodontic growth. The multivariate t linear mixed-effects model (4) can be easily modified into a hybrid model in which some of the subjects have a multivariate t -distribution, while others follow a Gaussian distribution (by setting $\tau_i = 1$ for these subjects).

To better understand the differences between the MLEs under the Gaussian and the multivariate t models, we consider the approximate distributions of the fixed effects estimators (corresponding to the asymptotic distributions evaluated at the MLEs) for each model, presented in Figure 4.

Figure 4 about here

The incremental parameters δ_0 and δ_1 have estimates closer to zero in the multivariate t model and the slope for the girls β_1 appears to be overestimated under the Gaussian model. The estimated variability for the MLEs is smaller in the multivariate t fit (the 95% confidence intervals are between 12 and 16% smaller than in the Gaussian model), suggesting that the parameters are estimated with greater precision.

Because of (18), the Gaussian linear mixed-effects model (1) can be viewed as a particular case of the multivariate t model (4). In the orthodontic growth example, the maximum log-likelihood for the Gaussian model is -203.021 and for the multivariate t model the maximum log-likelihood is -184.555 , corresponding to likelihood ratio statistic of 36.932 (p-value of 10^{-8}). This indicates that the multivariate t model fits the data substantially better than the Gaussian model.

The estimated average distances δ_i^2 , $\delta_{b_i}^2$, and $\delta_{e_i}^2$, defined in (14), provide useful diagnostic statistics for identifying subjects with outlying observations. Note that, under the Gaussian model (1), $E(\delta_{b_i}) = E(\mathbf{b}_i' \boldsymbol{\Psi}^{-1} \mathbf{b}_i) = q$, $E(\delta_{e_i}) = E[(\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i)' \boldsymbol{\Lambda}_i^{-1} (\mathbf{y}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{Z}_i \mathbf{b}_i)] = n_i$, and $E[\delta_i^2(\boldsymbol{\beta}, \boldsymbol{\Psi}, \boldsymbol{\rho})] = n_i$. Therefore, $\hat{\delta}_i^2/n_i$, $\hat{\delta}_{b_i}^2/q$, and $\hat{\delta}_{e_i}^2/n_i$ are expected to be close to 1 under the Gaussian model, and can be used as diagnostics statistics for identifying subjects with outliers (under this Gaussian model). Figure 5 presents these diagnostic statistics for the boys (because of the large value of $\hat{\lambda}_2$, the girls' estimated weights $\hat{\tau}_i$ are all essentially equal to 1). Subjects M09 and M13 present large values of $\hat{\delta}_i^2$ and $\hat{\delta}_{e_i}^2$, suggesting outlying observations at the within-subject level. This is consistent with the preliminary plot of the data, included in Figure 1, which suggests that both subjects have unusual growth patterns. The $\hat{\delta}_{b_i}^2$ plot gives some indication that subject M10

is possibly a **b**-outlier, which can not be concluded from Figure 3. Inspection of Figure 1 reveals that this subject has an unusually high orthodontic distance at the time of the first measurement.

Figure 5 about here

5.1.1 Influence of a single outlier

The robustness of the multivariate t MLEs with respect to the Gaussian MLEs can also be assessed through the influence of a single outlying observation (corresponding to a single **e**-outlier) on the estimated parameters. To simplify, we consider only the model for the girls, which can be represented as

$$y_{ij} = \beta_0 + \beta_1 t_j + b_{0i} + b_{1i} t_j + e_{ij}, \quad i = 1, \dots, 11 \quad j = 1, \dots, 4$$

with $\mathbf{b}_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{\Psi})$ and $e_{ij} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \sigma^2)$ in the Gaussian model and $\mathbf{b}_i | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\mathbf{0}, \tau_i^{-1} \mathbf{\Psi})$, $e_{ij} | \tau_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \tau_i^{-1} \sigma_{g(i)}^2)$, and $\tau_i \stackrel{\text{ind}}{\sim} \Gamma(\lambda/2, \lambda/2)$ in the multivariate t model.

We consider the influence of a change of Δ units in a single measurement on the estimated parameters. That is, we replace a single data point y_{ij} by the contaminated value $y_{ij}(\Delta) = y_{ij} + \Delta$, re-estimate the parameters, and record the relative change in the estimates $(\hat{\theta}(\Delta) - \hat{\theta})/\hat{\theta}$, where $\hat{\theta}$ denotes the original estimate and $\hat{\theta}(\Delta)$ the estimate for the contaminated data. In this example, we contaminated a *typical* value, the fourth observation (age = 14 years) on subject F01, and varied Δ between -20mm and 20mm by increments of 2mm. The Gaussian and the multivariate t fits were identical for the uncontaminated data in this case. Because $\mathbf{\Psi}$ and σ^2 have different interpretations under the Gaussian model (1) and the multivariate t model (4), and even within the multivariate t model for different degrees-of-freedom, we concentrate here on the estimation of the fixed effects β , which have the same interpretation under both models and for different degrees-of-freedom within the multivariate t model. We study the influence of the single outlier $y_{ij}(\Delta)$ on the estimation of $\hat{\beta}$ and of its estimated covariance matrix $\mathbf{V}_{\hat{\beta}}$.

Figure 6 presents the percent change curves for $\hat{\beta}$ and the upper-triangular elements of $\mathbf{V}_{\hat{\beta}}$ for different values of Δ .

Figure 6 about here

The influence of the single outlier is unbounded in the case of the Gaussian model, but clearly bounded in the multivariate t model. In the Gaussian model, the outlying observation has considerable more impact on the estimates of $\mathbf{V}_{\hat{\beta}}$ (changes between -2000% and 1800%), than on the fixed effects $\hat{\beta}$ (changes up to $\pm 60\%$). This has a direct impact on inferences drawn from the fit: confidence intervals increase unboundedly and test statistics go to zero. In the multivariate t fit, the influence of the single outlier for the fixed effects estimates remains bounded between -10% and 6% and for the estimates of $\mathbf{V}_{\hat{\beta}}$ it remains between -107% and 86%.

For closer contamination values ($|\Delta| \leq 2$), the multivariate t fit and the Gaussian fit are essentially identical and therefore have the same influence curves. This occurs because the contaminated observation is not distant enough from the *typical* data to be identified as an outlier, resulting in $\hat{\lambda} = \infty$. Therefore, the two estimation methods will have about the same efficiency for no or close contamination cases.

5.2 Comparing the MLEs under different outlier patterns

To compare the performance of the maximum likelihood estimators under the Gaussian model (1) and the multivariate t model (4), we conducted a simulation study involving different patterns of \mathbf{b} - and \mathbf{e} -outliers.

The linear mixed-effects model used to simulate the data is

$$\mathbf{y}_i = \mathbf{X}(\boldsymbol{\beta} + \mathbf{b}_i) + \mathbf{e}_i, \quad i = 1, \dots, 27, \quad \mathbf{X} = \begin{bmatrix} 1 & 8 \\ 1 & 10 \\ 1 & 12 \\ 1 & 14 \end{bmatrix}, \quad (19)$$

with the following mixture of normals models being used to *contaminate* the distributions of the \mathbf{b}_i and the \mathbf{e}_i .

$$\begin{aligned} \mathbf{b}_i &\stackrel{\text{ind}}{\sim} (1 - p_{\mathbf{b}}) \cdot \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi}) + p_{\mathbf{b}}f \cdot \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi}), \\ e_{ij} &\stackrel{\text{ind}}{\sim} (1 - p_{\mathbf{e}}) \cdot \mathcal{N}(0, \sigma^2) + p_{\mathbf{e}}f \cdot \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, 27, \quad j = 1, \dots, 4, \end{aligned} \quad (20)$$

where $p_{\mathbf{b}}$ and $p_{\mathbf{e}}$ denote, respectively, the expected percentage of \mathbf{b} - and \mathbf{e} -outliers in the data and f denotes the contamination factor. This model is a simplified version of the orthodontic growth model (2), with no gender differences. The parameters in the uncontaminated distributions are similar to the MLEs (17). It follows from (20) that $\text{var}(\mathbf{b}_i) = [1 + (f^2 - 1)p_{\mathbf{b}}] \boldsymbol{\Psi}$ and $\text{var}(e_{ij}) = [1 + (f^2 - 1)p_{\mathbf{e}}] \sigma^2$.

All thirty-two combinations of $p_{\mathbf{b}}$, $p_{\mathbf{e}} = 0, 0.05, 0.1, 0.25$, and $f = 2, 4$ were used in the simulation study. The $f = 2$ case corresponds to a close contamination pattern, while $f = 4$ illustrates a more distant contamination pattern. A total of 500 Monte Carlo replications were obtained for each $(p_{\mathbf{b}}, p_{\mathbf{e}}, f)$ combination.

An S-PLUS implementation of the ECME algorithm of Section 4.3 was used to obtain the multivariate t MLEs at each replication. For the Gaussian MLEs, the `lme` function in S-PLUS (MathSoft, 1997) was used. To enhance the comparability of the results, the same data set was used to obtain the multivariate t estimates and the Gaussian estimates, at each replication. The degrees-of-freedom for the multivariate t -distribution were assumed unknown, being estimated in the ECME algorithm.

As mentioned in Section 5.1, $\boldsymbol{\Psi}$ and σ^2 have different interpretations under the Gaussian model (1) and the multivariate t model (4) and their corresponding MLEs under the two models cannot be directly compared. As before, we concentrate on the estimation of the fixed effects $\boldsymbol{\beta}$, which have the same interpretation under both models.

For the simulation model (19), under estimation method E , the approximate covariance matrix $V_{\hat{\beta}}$ of the fixed effects estimates has the form

$$V_{\hat{\beta}} = \left[\hat{\sigma}_E^2 (\mathbf{X}'\mathbf{X})^{-1} + \hat{\Psi}_E \right] / m.$$

Under the Gaussian model, $\hat{\sigma}_G^2 = \hat{\sigma}^2$ and $\hat{\Psi}_G = \hat{\Psi}$, while under the multivariate t model $\hat{\sigma}_T^2 = m [\sum_{i=1}^m (\hat{\nu}_i + n_i) / (\hat{\nu}_i + n_i + 2)]^{-1} \hat{\sigma}^2$ and $\hat{\Psi}_T = m [\sum_{i=1}^m (\hat{\nu}_i + n_i) / (\hat{\nu}_i + n_i + 2)]^{-1} \hat{\Psi}$, with $\hat{\theta}$ denoting the MLE of θ under the appropriate estimation method. For the purpose of the simulation study, robustness is determined by how close the estimated values are to the parameters of the uncontaminated distribution. The asymptotic covariance matrix for the MLE of β based on the uncontaminated data only is $[\sigma^2 (\mathbf{X}'\mathbf{X})^{-1} + \Psi] / m$. Therefore, we define σ^2 as the target value for $\hat{\sigma}_G^2$ and $\hat{\sigma}_T^2$, σ^2 and Ψ as the target value for $\hat{\Psi}_G$ and $\hat{\Psi}_T$. These estimators can then be used to compare the performance of the two estimation methods with respect to the variance-covariance components Ψ and σ^2 .

The following parameters, with respective *target* values, are used in the comparison of the two estimation methods:

$$\beta_0 = 17, \quad \beta_1 = 0.8, \quad \Psi_{11} = 4, \quad \Psi_{22} = 0.0225, \quad \Psi_{12} = 0, \quad \text{and} \quad \sigma^2 = 1.$$

For the Gaussian model, the MLEs are considered and for the multivariate t model the MLEs of the fixed effects and the modified estimators $\hat{\Psi}_T$ and $\hat{\sigma}_T^2$ of the variance-covariance components are considered.

Let θ denote a parameter of interest, with target value $\theta_0 \neq 0$, estimated by $\hat{\theta}$. The efficiency of the Gaussian estimator $\hat{\theta}_G$ relative to the multivariate t estimator $\hat{\theta}_T$ is defined as the ratio of the respective mean square errors, $E(\hat{\theta}_G - \theta_0)^2 / E(\hat{\theta}_T - \theta_0)^2$. Expectations are taken with respect to the simulation distribution, that is, $E(\hat{\theta} - \theta_0)^2 = \sum_{i=1}^{500} (\hat{\theta}_i - \theta_0)^2 / 500$.

Figures 7 and 8 present the relative efficiency of the multivariate t estimators with respect to the Gaussian estimators. There are substantial gains in efficiency for all parameters under the more distant contamination patterns ($f = 4$) and moderate gains under the close contamination patterns ($f = 2$). The efficiency gains are bigger for the variance-covariance components than for the fixed effects. The two methods have about the same efficiency under the no-contamination case. For the close contamination patterns (Figure 7), the efficiency increases with the percentage of \mathbf{b} - and \mathbf{e} -outliers (except for the Ψ_{11} parameter, for which there is a slight efficiency decrease when the percentage of \mathbf{e} -outliers increases from 10% to 25%). In the case of distant contamination (Figure 8), the efficiency shows a non-monotone behavior with respect to the percentage of \mathbf{b} - and \mathbf{e} -outliers. This pattern suggests that the multivariate t model is more robust than the Gaussian model especially for moderate percentages (5-10%) of outliers.

Figures 7 and 8 about here

The simulation results for the mean square error (not shown here) indicate that outliers affect the variance-covariance components estimates more than they affect the fixed effects estimates. The precision of the estimator of σ^2 seems to be affected only by the percentage of \mathbf{e} -outliers, while the fixed effects and random effects variance-covariance components estimators are affected by both types of outliers.

The MLEs of the fixed effects are nearly unbiased (relative bias $\leq 0.7\%$) for both estimation methods under all contamination patterns. The bias for the variance-covariance components follows the same basic pattern as the mean square error: it increases with the percentage of \mathbf{e} -outliers, is insensitive to the percentage of \mathbf{b} -outliers for σ^2 , and increases in absolute value with the percentage of both types of outliers for the random effects variance-covariance components.

The coverage probabilities of the approximate 95% confidence intervals for the fixed effects, not included here, are generally close to the nominal level for both estimation methods, with the smallest coverage probability 90.4% and the largest 97%. The coverage probabilities tend to increase with the percentage of outliers, because the fixed effects estimators remain unbiased and the confidence intervals get larger. The average length of the 95% confidence intervals is about the same under Gaussian and multivariate t estimation for the close contamination patterns, but 10% to 25% larger in the Gaussian model for the more distant contamination patterns.

6 Conclusion

This article describes a robust version of the linear mixed-effects model of Laird and Ware (1982) in which the Gaussian distributions for the random effects and the within-subject errors are replaced by multivariate t -distributions. Analysis of examples and simulation results indicate that the multivariate t linear mixed-effects model substantially outperforms the Gaussian model when outliers are present in the data, even in moderate amounts. Gains in efficiency for the multivariate t MLEs relative to the Gaussian MLEs, under outlier contamination, are observed for all parameters, being particularly high in the estimation of variance-covariance components, ranging from 20%–30% in the case of close contamination (two standard deviations) to 200%–400% in the case of distant contamination (four standard deviations). This has a direct impact on confidence intervals and test statistics obtained from the fit, which determine all inferences drawn from the estimated model. The influence function is bounded for the multivariate t model and unbounded for the Gaussian model. The multivariate t model also provides diagnostics tools for graphically identifying subjects with outlying observations.

We describe EM-type algorithms for efficient maximum likelihood estimation under two missing data structures: with both the random effects and the individual weights treated as missing and with only the individual weights treated as missing. The former leads to algorithms with closed form expressions for both the E- and the M-step, but imposes some restrictions on the correlation structure of the within-subject errors. The algorithm corresponding to the latter missing data structure, which allows general correlation structures for the within subject errors, involves a more

computationally intensive M-step, but can be implemented using existing, reliable software.

The robust estimation approach described in this article can also be extended to nonlinear mixed-effects models (Lindstrom and Bates, 1990). The computations become considerably more complex, but algorithms based on linear approximations to the marginal distribution of the \mathbf{y}_i can, in principle, be used in conjunction with the methods described here.

References

- Barnett, V. and Lewis, T. (1994). *Outliers in statistical data (Third edition)*, John Wiley & Sons.
- Bates, D. M. and Watts, D. G. (1988). *Nonlinear Regression Analysis and Its Applications*, Wiley, New York.
- Becker, R. A., Cleveland, W. S. and Shyu, M.-J. (1996). The visual design and control of trellis graphics displays, *J. of Computational and Graphical Statistics* **5**(2): 123–156.
- Box, G. E. P., Jenkins, G. M. and Reinsel, G. C. (1994). *Time Series Analysis: Forecasting and Control*, 3rd edn, Holden-Day, San Francisco.
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (c/r: P22-37), *Journal of the Royal Statistical Society, Ser. B* **39**: 1–22.
- Efron, B. and Tibshirani, R. J. (1993). *An Introduction to the Bootstrap*, Chapman & Hall, New York.
- Fang, K. T. and Zhang, Y. Y. (1990). *Generalized Multivariate Analysis*, Science Press, Beijing and Springer-Verlag, Berlin.
- Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J. and Stahel, W. A. (1986). *Robust Statistics: The Approach Based on Influence Functions*, John Wiley & Sons, New York.
- Hartley, H. O. and Rao, J. N. K. (1967). Maximum likelihood estimation for the mixed analysis of variance model, *Biometrika* **54**: 93–108.
- Huber, P. J. (1981). *Robust Statistics*, Wiley, New York.
- Jennrich, R. I. and Schluchter, M. D. (1986). Unbalanced repeated measures models with structural covariance matrices, *Biometrics* **42**(4): 805–820.
- Johnson, N. L. and Kotz, S. (1972). *Distributions in Statistics: Continuous Multivariate Distributions*, Wiley, New York.
- Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data, *Biometrics* **38**: 963–974.

- Lange, K. L., Little, R. J. A. and Taylor, J. M. G. (1989). Robust statistical modeling using the t -distribution, *Journal of the American Statistical Association* **84**: 881–896.
- Lindstrom, M. J. and Bates, D. M. (1988). Newton-Raphson and EM algorithms for linear mixed-effects models for repeated-measures data, *Journal of the American Statistical Association* **83**: 1014–1022.
- Lindstrom, M. J. and Bates, D. M. (1990). Nonlinear mixed effects models for repeated measures data, *Biometrics* **46**: 673–687.
- Liu, C. and Rubin, D. B. (1994). The ECME algorithm: A simple extension of EM and ECM with faster monotone convergence, *Biometrika* **81**: 633–648.
- Liu, C., Rubin, D. B. and Wu, Y. N. (1998). Parameter expansion for EM acceleration – the PX-EM algorithm, *Biometrika* . To appear.
- Maddala, G. S. and Rao, C. R. (1997). *Handbook of Statistics*, Vol. 15, Elsevier Science B. V., Amsterdam.
- MathSoft (1997). *S-PLUS 4 Guide to Statistics*, Data Analysis Products, MathSoft Inc., Seattle, WA.
- Meng, X. L. and Rubin, D. B. (1993). Maximum likelihood estimation via the ECM algorithm: a general framework, *Biometrika* **80**: 267–278.
- Miller, J. J. (1977). Asymptotic properties of maximum likelihood estimates in the mixed model of the analysis of variance, *Ann. of Statistics* **5**: 746–762.
- Pendergast, J. F. and Broffitt, J. D. (1986). Robust estimation in growth curve models, *Communications in Statistics: Theory and Methods* **14**: 1919–1939.
- Potthoff, R. F. and Roy, S. N. (1964). A generalized multivariate analysis of variance model useful especially for growth curve problems, *Biometrika* **51**: 313–326.
- Racine-Poon, A. (1992). Saga: Samples assisted graphical analysis (disc: P401-404), *Bayesian Statistics 4. Proceedings of the Fourth Valencia International Meeting*, pp. 389–401.
- Richardson, A. M. (1997). Bounded influence estimation in the mixed linear model, *Journal of the American Statistical Association* **92**(437): 154–161.
- Richardson, A. M. and Welsh, A. H. (1995). Robust estimation in the mixed linear model, *Biometrics* **51**: 1429–1439.
- Thisted, R. A. (1988). *Elements of Statistical Computing*, Springer-Verlag, London.

- Wakefield, J. C., Smith, A. F. M., Racine-Poon, A. and Gelfand, A. E. (1994). Bayesian analysis of linear and nonlinear population models using the Gibbs sampler, *Applied Statistics* . Accepted for publication.
- Welsh, A. H. and Richardson, A. M. (1997). *Approaches to the Robust Estimation of Mixed Models*, Vol. 15 of Maddala and Rao (1997), chapter 13, pp. 343–384.
- Wolfinger, R., Tobias, R. and Sall, J. (1991). Mixed models: A future direction, *SUGI* **16**: 1380–1388.
- Wu, C. F. J. (1983). On the convergence properties of the EM algorithm, *The Annals of Statistics* **11**: 95–103.

Table 1: Number of iterations and user time to obtain the multivariate t maximum likelihood estimates in the orthodontic growth model.

Algorithm	Iterations	Time (sec)
ECME for missing \mathbf{b}_i and τ_i	268	3.01
PX-EM	134	2.51
S-PLUS-ECME for missing τ_i	16	78.16

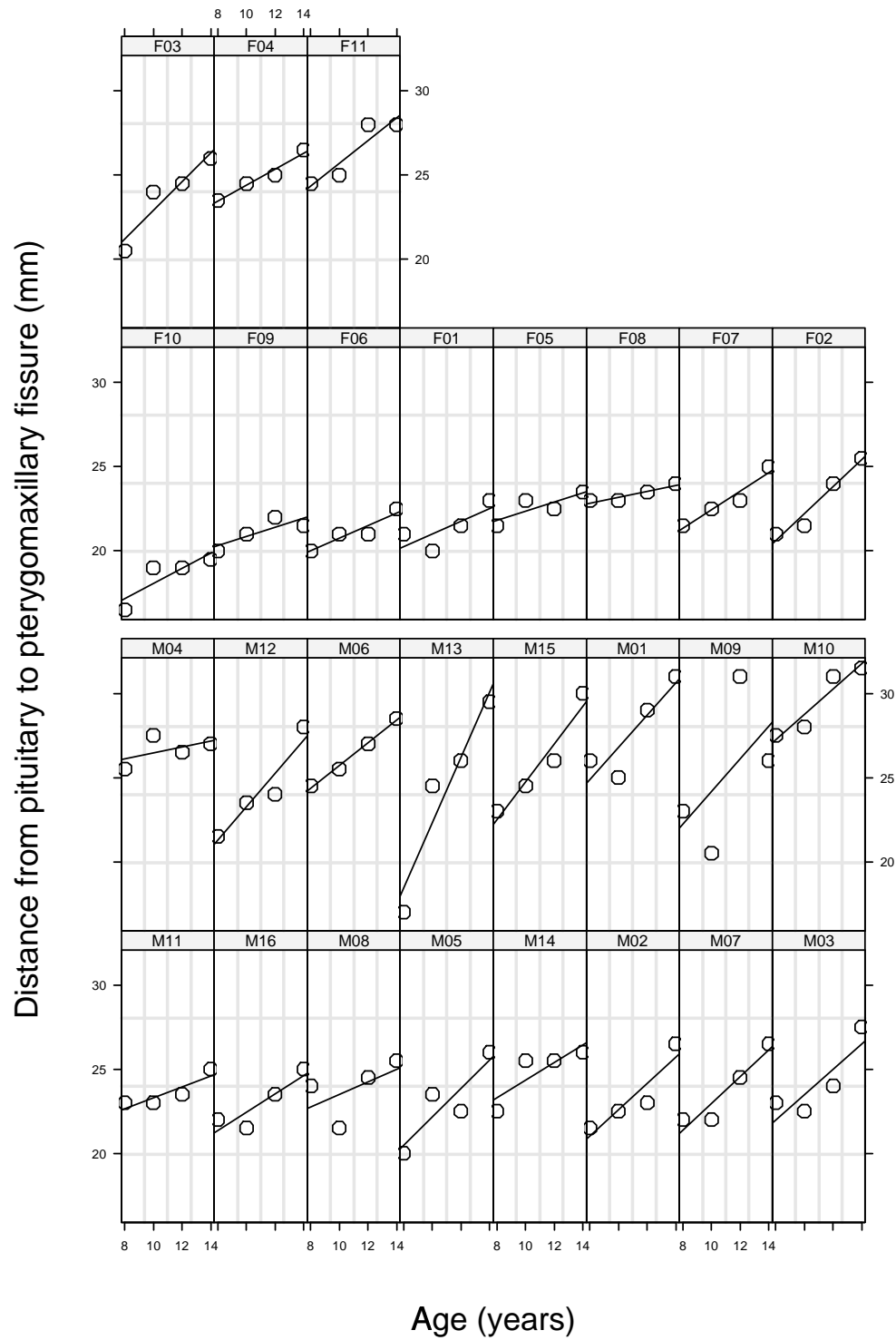


Figure 1: Orthodontic growth patterns in 16 boys(M) and 11 girls(F) between 8 and 14 years of age. Lines represent the individual least squares fits of the simple linear regression model.

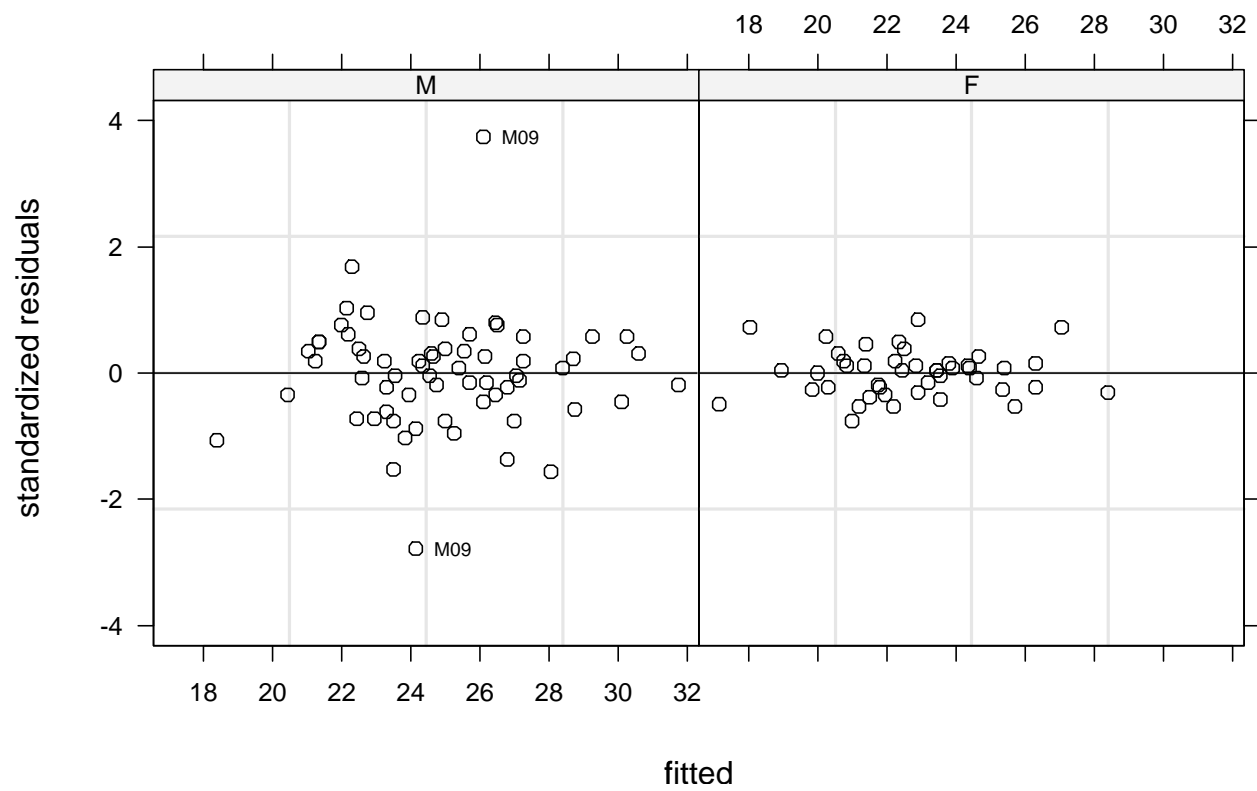


Figure 2: Residuals versus fitted values plots by gender, corresponding to individual least squares fits of the orthodontic growth data.

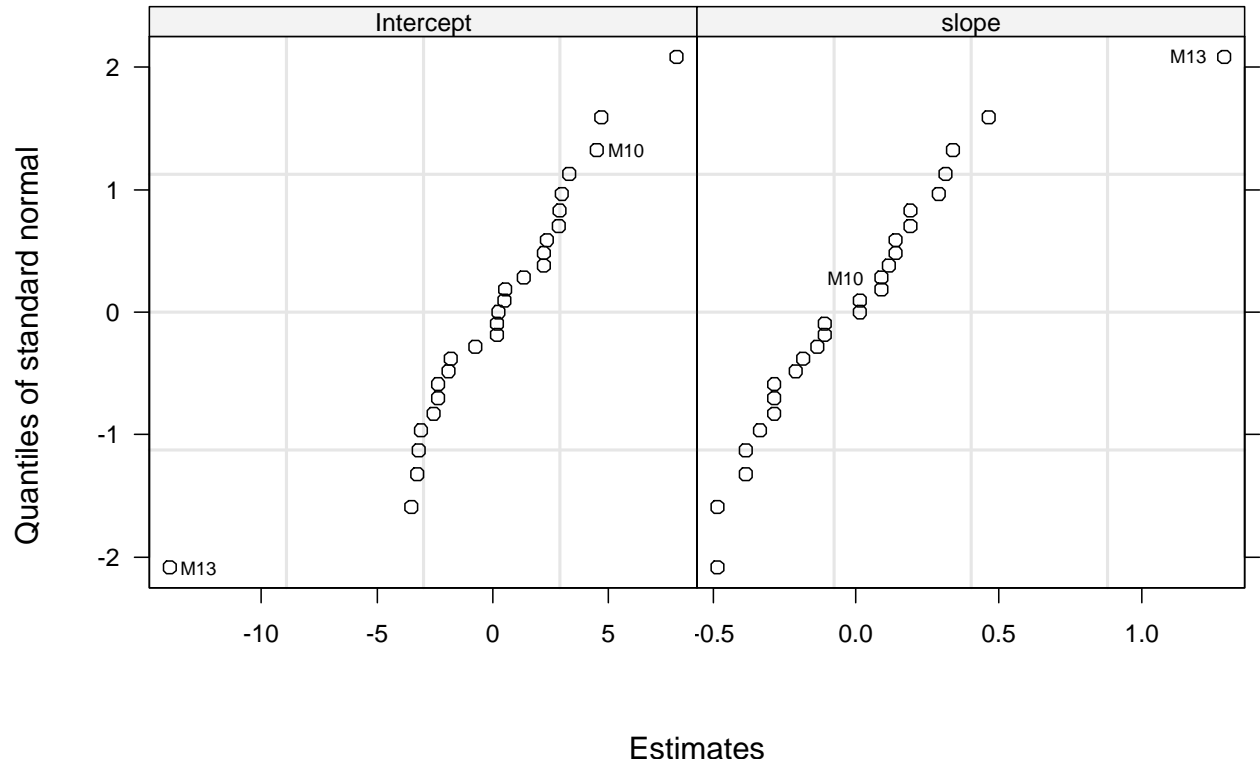


Figure 3: Normal plots of estimated coefficients corresponding to individual least squares fits of the orthodontic growth data.

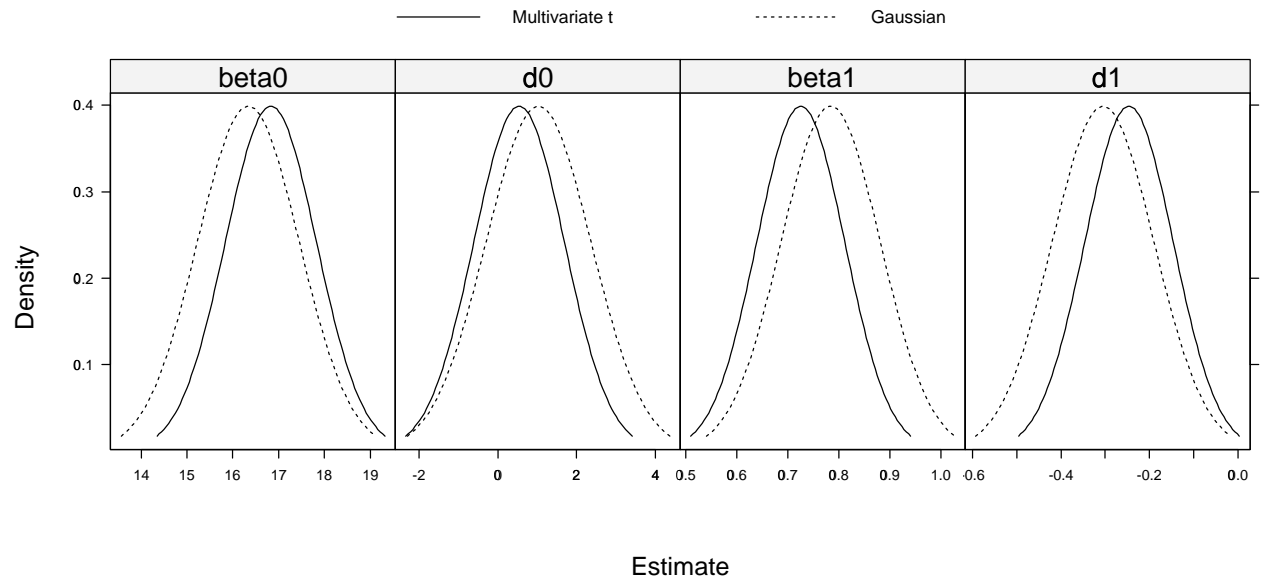


Figure 4: Approximate densities of the fixed effects MLEs in the orthodontic growth model (2) under Gaussian and multivariate t estimation.

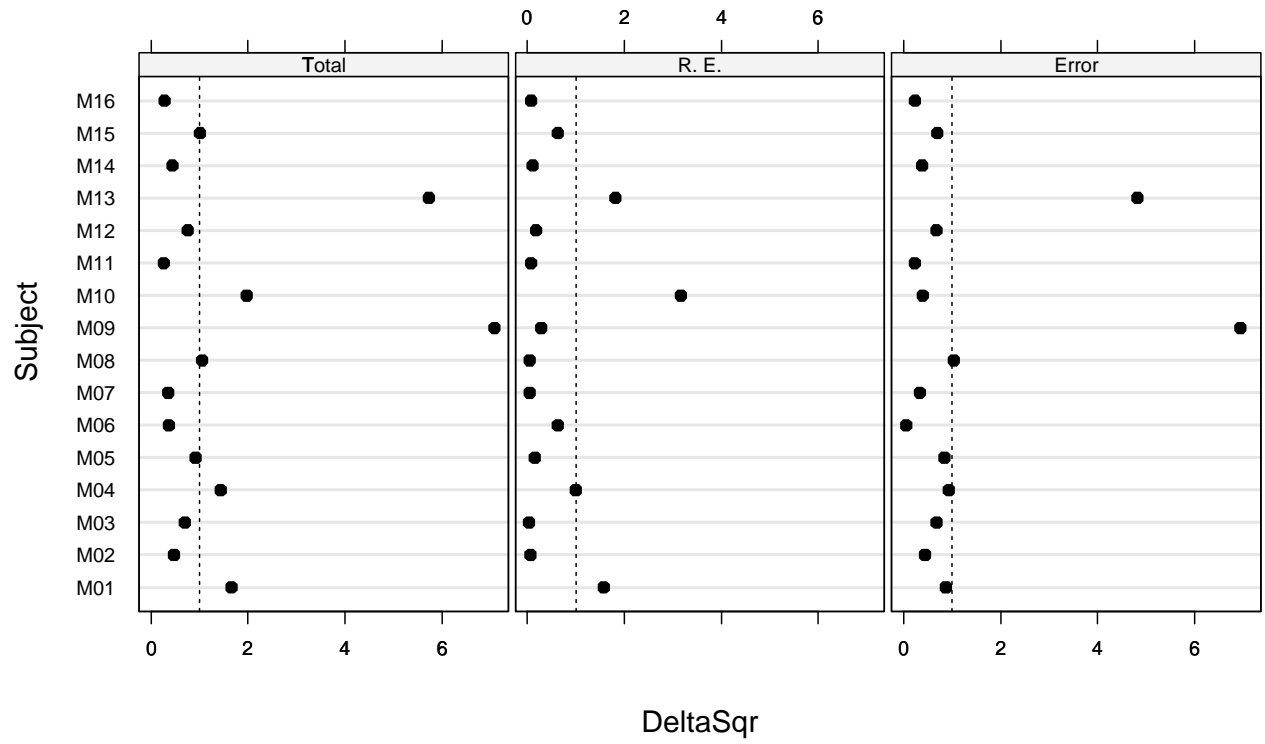


Figure 5: Estimated δ_i^2 (Total), δ_{b_i} (R.E.), and δ_{e_i} (Error) for boys in the multivariate t fit of the orthodontic distance data.

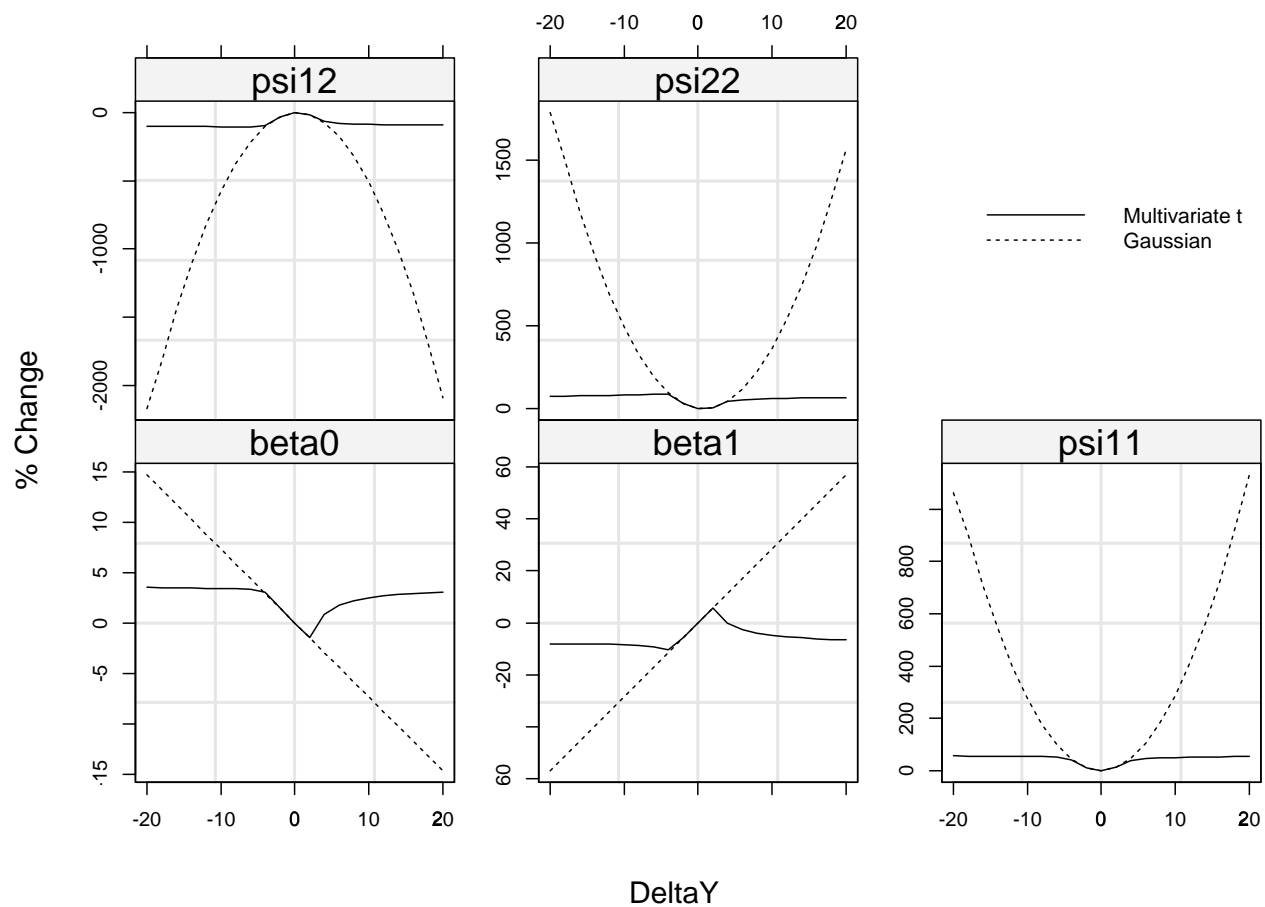


Figure 6: Percent change in maximum likelihood estimates under the Gaussian and multivariate t models for different contaminations Δ of a single observation.

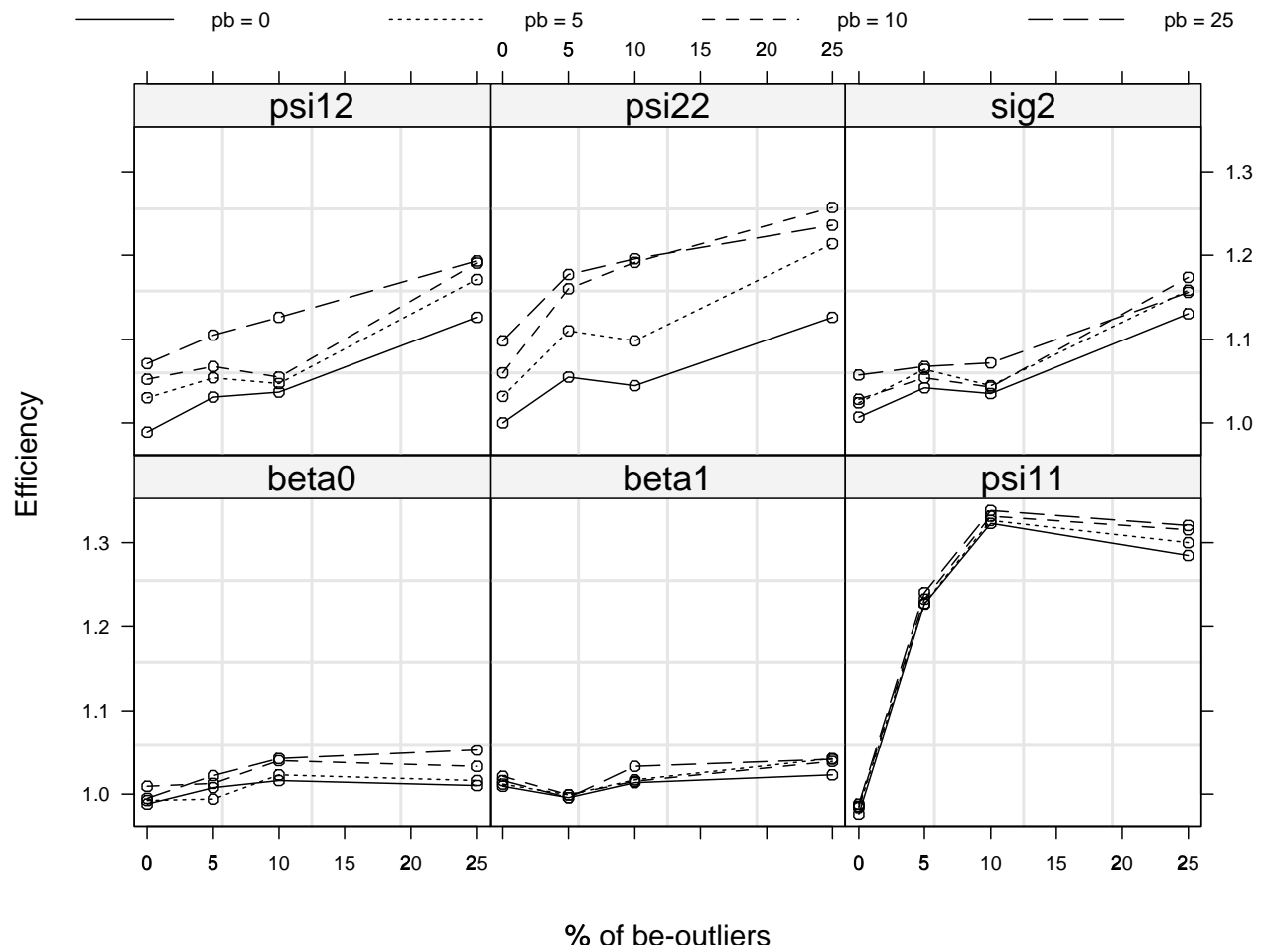


Figure 7: Relative efficiencies of the multivariate t MLEs with respect to the Gaussian MLEs under close outlier contamination patterns.

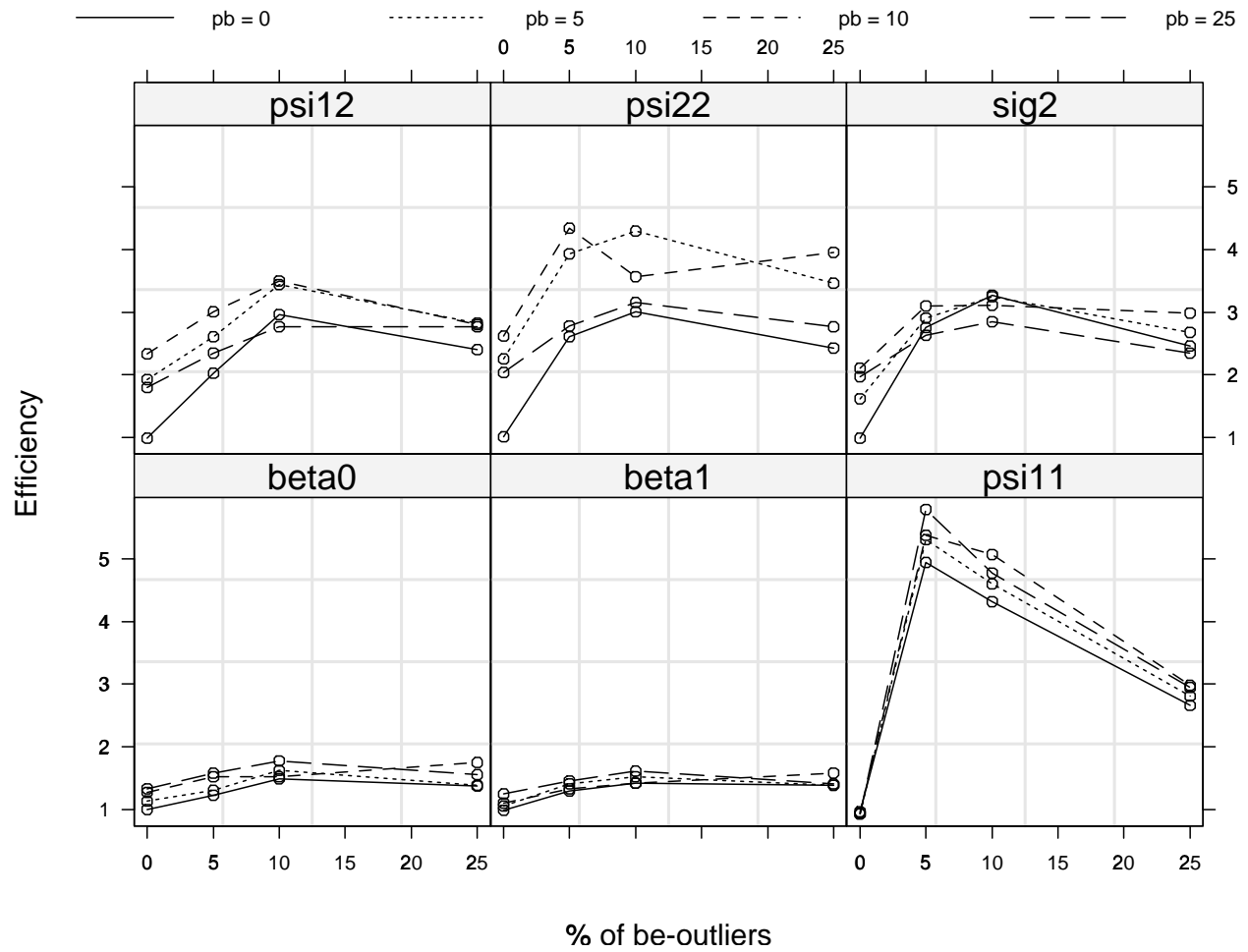


Figure 8: Relative efficiencies of the multivariate t MLEs with respect to the Gaussian MLEs under distant outlier contamination patterns.