



1. Introduction

Longitudinal data is data which is collected from the same individuals sequentially over multiple points in time. We are interested in analysing data of this form to attempt to identify trends and changes in the data over time and to see how the individuals are correlated. This project will introduce a class of models that aims to model longitudinal data whilst taking into account variability between each individual.

2. Building a Linear Mixed Effects Model

One can introduce the structure of a **linear mixed effects model (LMEM)** in two stages [3]. Let Y_{ij} be the response of individual i at time $j = 1, \dots, n_i$, where n_i is the total number of observations for the i^{th} individual. Stage 1 assumes that \mathbf{Y}_i follows a linear regression model, given by

$$\mathbf{Y}_i = \mathbf{Z}_i \boldsymbol{\beta}_i + \boldsymbol{\epsilon}_i, \quad (1)$$

where \mathbf{Z}_i is an $n_i \times q$ dimensional design matrix, $\boldsymbol{\beta}_i$ is a q -dimensional vector of unknown subject specific regression coefficients. We also have $\boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_{n_i})$. In Stage 2, we aim to explain the variability between subjects by including random effects in $\boldsymbol{\beta}_i$ such as

$$\boldsymbol{\beta}_i = \mathbf{M}_i \boldsymbol{\beta} + \mathbf{u}_i, \quad (2)$$

where $\mathbf{M}_i \in \mathbb{R}^{q \times p}$ is a matrix of known covariates and $\boldsymbol{\beta} \in \mathbb{R}^p$ is a vector of unknown fixed regression coefficients. Here \mathbf{u}_i is a vector of subject specific random effects which are assumed to be independent and $\mathbf{u}_i \sim N(\mathbf{0}, \mathbf{Q})$. Substituting (2) into (1) gives

$$\begin{aligned} \mathbf{Y}_i &= \mathbf{Z}_i \mathbf{M}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{u}_i + \boldsymbol{\epsilon}_i \\ &= \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \mathbf{u}_i + \boldsymbol{\epsilon}_i \end{aligned}$$

where we have written $\mathbf{X}_i = \mathbf{Z}_i \mathbf{M}_i$. We have the general framework for the entire model as

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{Z} \mathbf{u} + \boldsymbol{\epsilon} \quad (3)$$

Note that $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n)$ and $\mathbf{u} \sim N(\mathbf{0}, \mathbf{Q})$ where \mathbf{Q} is a block diagonal matrix with the i -th diagonal block being equal to the covariance matrix of \mathbf{u}_i .

3. References

- [1] Verbeke Fitzmaurice Davidian and Molenberghs. *Longitudinal Data Analysis*. Chapman and Hall/CRC, 2009.
- [2] Lynn Roy LaMotte. "A Direct Derivation of the REML Likelihood function". In: *Statistical Papers* (2005).
- [3] Geert Verbeke and Geert Molenberghs. *Linear Mixed Models for Longitudinal Data*. Springer, 2000.
- [4] Geert Verbeke and Geert Molenberghs. *Models for Discrete Longitudinal Data*. Springer.

4. Distribution of a Linear Mixed Effects Model

The distribution of \mathbf{Y} is as follows:

$$\mathbf{Y} \sim N(\mathbf{X} \boldsymbol{\beta}, \boldsymbol{\Sigma}), \quad (4)$$

where $\boldsymbol{\Sigma} = \mathbf{Z} \mathbf{Q} \mathbf{Z}^T + \sigma^2 \mathbf{I}_N$, with $N = \sum_i n_i$, the total number of observations across the whole dataset. This is due to

$$\mathbb{E}(\mathbf{Y}) = \mathbf{X} \boldsymbol{\beta} + \mathbf{Z} \mathbb{E}(\mathbf{u}) + \mathbb{E}(\boldsymbol{\epsilon}) = \mathbf{X} \boldsymbol{\beta}$$

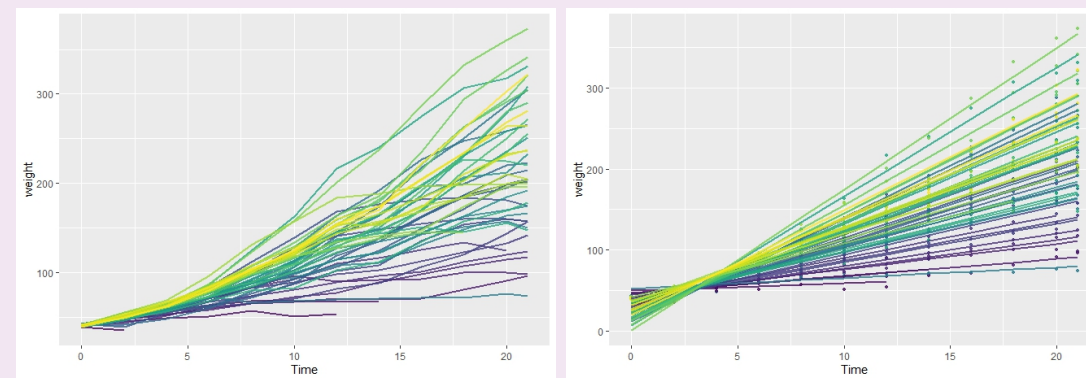
$$\text{Var}(\mathbf{Y}) = \mathbf{Z} \text{Var}(\mathbf{u}) \mathbf{Z}^T + \text{Var}(\boldsymbol{\epsilon}) = \mathbf{Z} \mathbf{Q} \mathbf{Z}^T + \sigma^2 \mathbf{I}_N$$

5. Applying a Linear Mixed Effects Model

We will apply an LMEM with a random intercept and slope to the **ChickWeight** dataset. The data consists of 50 chicks which have been weighed at 12 time points. The chicks were fed one of four different diets, which we denote Diet A, B, C or D. Using R, we fit a model with the equation

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + \beta_2 B_i + \beta_3 C_i + \beta_4 D_i + u_{0i} + u_{1i} t_{ij} + \epsilon_{ij}$$

where B_i, C_i and D_i are diet indicator variables.



(a) ChickWeight plotted

(b) LME fitted to ChickWeight

Figure 1: Plots Showing How the LME Models ChickWeight

The REML estimation (see **Section 7**) for the coefficients is as follows:

$$Y_{ij} = 26.4 + 8.4 t_{ij} + 2.8 B_i + 2.0 C_i + 9.3 D_i + u_{0i} + u_{1i} t_{ij} + \epsilon_{ij}$$

6. Inference for the Random Effects

One can estimate \mathbf{u}_i by taking $\hat{\mathbf{u}}_i(\theta) = \mathbb{E}(\mathbf{u}_i | \mathbf{Y}_i = \mathbf{y}_i)$, its posterior mean. We know that

$$\begin{pmatrix} \mathbf{u} \\ \mathbf{Y} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{X} \boldsymbol{\beta} \end{pmatrix}, \begin{pmatrix} \mathbf{Q} & \text{Cov}(\mathbf{u}, \mathbf{Y}) \\ \text{Cov}(\mathbf{u}, \mathbf{Y})^T & \boldsymbol{\Sigma} \end{pmatrix} \right)$$

where $\text{Cov}(\mathbf{u}, \mathbf{Y}) = \mathbf{Q} \mathbf{Z}^T$. Using a general result for multivariate conditional normal distributions, the posterior mean can be specified as $\hat{\mathbf{u}}_i(\theta) = \mathbf{Q} \mathbf{Z}^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})$ and hence this is our estimate for \mathbf{u}_i .

7. REML Estimator for the Likelihood

One can introduce an estimator for the likelihood which is not dependent on $\boldsymbol{\beta}$ [4]. One does this by multiplying the general equation by the transpose of a matrix \mathbf{A} which is orthogonal to \mathbf{X} . One has

$$\mathbf{A}^T \mathbf{Y} = \mathbf{A}^T \mathbf{X} \boldsymbol{\beta} + \mathbf{A}^T \mathbf{Z} \mathbf{u} + \mathbf{A}^T \boldsymbol{\epsilon} = \mathbf{A}^T \mathbf{Z} \mathbf{u} + \mathbf{A}^T \boldsymbol{\epsilon}$$

$$\mathbf{A}^T \mathbf{Y} \sim N(\mathbf{0}, \mathbf{A}^T \boldsymbol{\Sigma} \mathbf{A}).$$

LaMotte [2] proves that the likelihood of this function can be written in a form which does not depend on \mathbf{A} , with $\hat{\boldsymbol{\beta}}$ being the generalised least squares (GLS) estimate. This is given by:

$$(2\pi)^{-(n-p)} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} |\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X}|^{-\frac{1}{2}} |\mathbf{X}^T \mathbf{X}|^{\frac{1}{2}} \exp \left(-\frac{1}{2} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}) \right). \quad (5)$$

8. Missing Data in Longitudinal Datasets

Longitudinal datasets often include missing data. The **lmer** function in R can still fit an LMEM to a dataset that contains missing values. We randomly set a percentage of the response values equal to NA and calculated the residual sum of squares (RSS) of the model fitted using **lmer**.

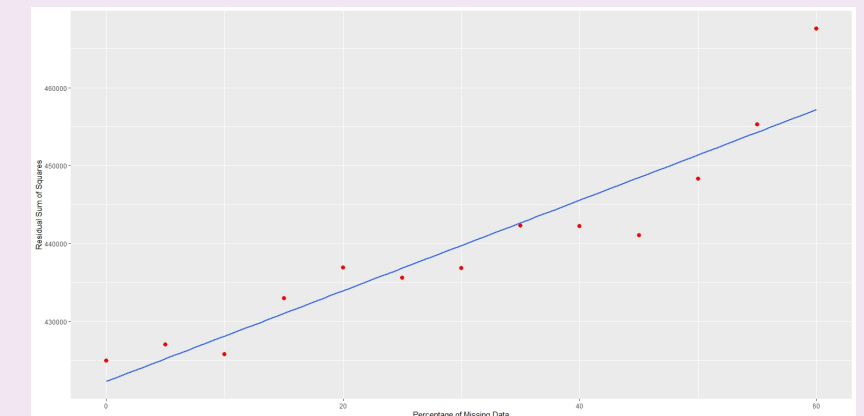


Figure 2: RSS of an LME Fitted to ChickWeight with Missing Data

The higher the percentage of missing data, the more R struggles to fit the model. There are many ways that one can impute missing data from simple methods such as single mean imputation and last observation carried forward to more advanced methods such as multiple imputation.

9. Conclusions and Future Work

We have seen that we can model longitudinal datasets effectively with LMEMs which capture the trends of the data for each individual.

A future step in the project is to create a function on R to fit a selection model following the theory from [1]. This is a hierarchical model which assumes the data follows a multivariate normal distribution and the missingness follows a logistic regression.