# Chapter 1

# Model Diagnostics

## 1.1 Introduction

In classical linear models model diagnostics have been become a required part of any statistical analysis, and the methods are commonly available in statistical packages and standard textbooks on applied regression. However it has been noted by several papers that model diagnostics do not often accompany LME model analyses. Model diagnostic techniques determine whether or not the distributional assumptions are satisfied, and to assess the influence of unusual observations.

### 1.1.1 Model Data Agreement

Schabenberger (2004) describes the examination of model-data agreement as comprising several elements; residual analysis, goodness of fit, collinearity diagnostics and influence analysis.

$$r_{mi} = x_i^T \hat{\beta} \tag{1.1}$$

### 1.1.2 Marginal Residuals

$$\hat{\beta} \;=\; (X^T R^{-1} X)^{-1} X^T R^{-1} Y$$

$$= BY$$

## 1.2 Standardized and studentized residuals

To alleviate the problem caused by inconstant variance, the residuals are scaled (i.e. divided) by their standard deviations. This results in a 'standardized residual'. Because true standard deviations are frequently unknown, one can instead divide a residual by the estimated standard deviation to obtain the 'studentized residual.

### 1.2.1 Standardization

A random variable is said to be standardized if the difference from its mean is scaled by its standard deviation. The residuals above have mean zero but their variance is unknown, it depends on the true values of $\theta$. Standardization is thus not possible in practice.

### 1.2.2 Studentization

Instead, you can compute studentized residuals by dividing a residual by an estimate of its standard deviation.

### 1.2.3 Internal and External Studentization

If that estimate is independent of the $i-$th observation, the process is termed 'external studentization'. This is usually accomplished by excluding the $i-$th observation when computing the estimate of its standard error. If the observation contributes to the standard error computation, the residual is said to be internally studentized.

Externally studentized residual require iterative influence analysis or a profiled residuals variance.

### 1.2.4 Computation

The computation of internally studentized residuals relies on the diagonal entries of $\boldsymbol{V}(\hat{\theta})$ - $\boldsymbol{Q}(\hat{\theta})$, where $\boldsymbol{Q}(\hat{\theta})$ is computed as

$$\boldsymbol{Q}(\hat{\theta}) = \boldsymbol{X}(\boldsymbol{X}'\boldsymbol{Q}(\hat{\theta})^{-1}\boldsymbol{X})\boldsymbol{X}^{-1}$$

### 1.2.5 Pearson Residual

Another possible scaled residual is the 'Pearson residual', whereby a residual is divided by the standard deviation of the dependent variable. The Pearson residual can be used when the variability of $\hat{\beta}$ is disregarded in the underlying assumptions.

## 1.3 Covariance Parameters

The unknown variance elements are referred to as the covariance parameters and collected in the vector $\theta$.

### 1.3.1 Methods and Measures

The key to making deletion diagnostics useable is the development of efficient computational formulas, allowing one to obtain the case deletion diagnostics by making use of basic building blocks, computed only once for the full model.

Zewotir and Galpin (2005) lists several established methods of analyzing influence in LME models. These methods include

- Cook's distance for LME models,

- likelihood distance,

- the variance (information) ration,

- the Cook-Weisberg statistic,

- the Andrews-Prebigon statistic.

## 1.4 Iterative and non-iterative influence analysis

Schabenberger (2004) highlights some of the issue regarding implementing mixed model diagnostics.

A measure of total influence requires updates of all model parameters.

however, this doesnt increase the procedures execution time by the same degree.

### 1.4.1 Iterative Influence Analysis

For linear models, the implementation of influence analysis is straightforward. However, for LME models, the process is more complex. Update formulas for the fixed effects are available only when the covariance parameters are assumed to be known. A measure of total influence requires updates of all model parameters. This can only be achieved in general is by omitting observations, then refitting the model.

Schabenberger (2004) describes the choice between iterative influence analysis and non-iterative influence analysis.

## 1.5 Measures of Influence

The impact of an observation on a regression fitting can be determined by the difference between the estimated regression coefficient of a model with all observations and the estimated coefficient when the particular observation is deleted. The measure DFBETA is the studentized value of this difference.

Influence arises at two stages of the LME model. Firstly when $V$ is estimated by $\hat{V}$, and subsequent estimations of the fixed and random regression coefficients $\beta$ and $u$, given $\hat{V}$.

### 1.5.1 DFFITS

DFFITS is a statistical measured designed to a show how influential an observation is in a statistical model. It is closely related to the studentized residual.

$$DFFITS = \frac{\widehat{y_i} - \widehat{y_{i(k)}}}{s_{(k)}\sqrt{h_{ii}}}$$

### 1.5.2 PRESS

The prediction residual sum of squares (PRESS) is an value associated with this calculation. When fitting linear models, PRESS can be used as a criterion for model selection, with smaller values indicating better model fits.

$$PRESS = \sum (y - y^{(k)})^2 \tag{1.2}$$

- $e_{-Q} = y_Q - x_Q\hat{\beta}^{-Q}$

- $PRESS_{(U)} = y_i - x\hat{\beta}_{(U)}$

### 1.5.3 DFBETA

$$\begin{aligned} DFBETA_a &= \hat{\beta} - \hat{\beta}_{(a)} \tag{1.3} \\ &= B(Y - Y_{\bar{a}} \tag{1.4} \end{aligned}$$

# Chapter 2

# Application to Method Comparison Studies

## 2.1 Application to MCS

Let $\hat{\beta}$ denote the least square estimate of $\beta$ based upon the full set of observations, and let $\hat{\beta}^{(k)}$ denoted the estimate with the $k^{th}$ case excluded.

## 2.2 Grubbs' Data

For the Grubbs data the $\hat{\beta}$ estimated are $\hat{\beta}_0$ and $\hat{\beta}_1$ respectively. Leaving the fourth case out, i.e. $k = 4$ the corresponding estimates are $\hat{\beta}_0^{-4}$ and $\hat{\beta}_1^{-4}$

$$Y^{-Q} = \hat{\beta}^{-Q} X^{-Q} \tag{2.1}$$

When considering the regression of case-wise differences and averages, we write $D^{-Q} = \hat{\beta}^{-Q} A^{-Q}$

|    | F      | C      | D     | A      |
|----|--------|--------|-------|--------|
| 1  | 793.80 | 794.60 | -0.80 | 794.20 |
| 2  | 793.10 | 793.90 | -0.80 | 793.50 |
| 3  | 792.40 | 793.20 | -0.80 | 792.80 |
| 4  | 794.00 | 794.00 | 0.00  | 794.00 |
| 5  | 791.40 | 792.20 | -0.80 | 791.80 |
| 6  | 792.40 | 793.10 | -0.70 | 792.75 |
| 7  | 791.70 | 792.40 | -0.70 | 792.05 |
| 8  | 792.30 | 792.80 | -0.50 | 792.55 |
| 9  | 789.60 | 790.20 | -0.60 | 789.90 |
| 10 | 794.40 | 795.00 | -0.60 | 794.70 |
| 11 | 790.90 | 791.60 | -0.70 | 791.25 |
| 12 | 793.50 | 793.80 | -0.30 | 793.65 |

$$Y^{(k)} = \hat{\beta}^{(k)} X^{(k)} \tag{2.2}$$

Consider two sets of measurements , in this case F and C , with the vectors of case-wise averages $A$ and case-wise differences $D$ respectively. A regression model of differences on averages can be fitted with the view to exploring some characteristics of the data.

When considering the regression of case-wise differences and averages, we write

$$D^{-Q} = \hat{\beta}^{-Q} A^{-Q} \tag{2.3}$$

Let $\hat{\beta}$ denote the least square estimate of $\beta$ based upon the full set of observations, and let $\hat{\beta}^{(k)}$ denoted the estimate with the $k^{th}$ case excluded.

For the Grubbs data the $\hat{\beta}$ estimated are $\hat{\beta}_0$ and $\hat{\beta}_1$ respectively. Leaving the fourth case out, i.e. $k = 4$ the corresponding estimates are $\hat{\beta}_0^{-4}$ and $\hat{\beta}_1^{-4}$

$$Y^{(k)} = \hat{\beta}^{(k)} X^{(k)} \tag{2.4}$$

Consider two sets of measurements , in this case F and C , with the vectors of case-wise averages $A$ and case-wise differences $D$ respectively. A regression model of differences on averages can be fitted with the view to exploring some characteristics of the data.

```
Call: lm(formula = D ~ A)
```

```
Coefficients: (Intercept)            A
-37.51896        0.04656
```

When considering the regression of case-wise differences and averages, we write

$$D^{-Q} = \hat{\beta}^{-Q} A^{-Q} \qquad (2.5)$$

## 2.2.1 Influence measures using R

R provides the following influence measures of each observation.

|     | dfb.1_ | dfb.A | dffit | cov.r | cook.d | hat  |
| --- | ------ | ----- | ----- | ----- | ------ | ---- |
| 1   | 0.42   | -0.42 | -0.56 | 1.13  | 0.15   | 0.18 |
| 2   | 0.17   | -0.17 | -0.34 | 1.14  | 0.06   | 0.11 |
| 3   | 0.01   | -0.01 | -0.24 | 1.17  | 0.03   | 0.08 |
| 4   | -1.08  | 1.08  | 1.57  | 0.24  | 0.56   | 0.16 |
| 5   | -0.14  | 0.14  | -0.24 | 1.30  | 0.03   | 0.13 |
| 6   | -0.00  | 0.00  | -0.11 | 1.31  | 0.01   | 0.08 |
| 7   | -0.04  | 0.04  | -0.08 | 1.37  | 0.00   | 0.11 |
| 8   | 0.02   | -0.02 | 0.15  | 1.28  | 0.01   | 0.09 |
| 9   | 0.69   | -0.68 | 0.75  | 2.08  | 0.29   | 0.48 |
| 10  | 0.18   | -0.18 | -0.22 | 1.63  | 0.03   | 0.27 |
| 11  | -0.03  | 0.03  | -0.04 | 1.53  | 0.00   | 0.19 |
| 12  | -0.25  | 0.25  | 0.44  | 1.05  | 0.09   | 0.12 |

# Bibliography

Schabenberger, O. (2004). Mixed model influence diagnostics. 18929.

Zewotir, T. and J. Galpin (2005). Influence diagnostics for linear mixed models. *Journal of Data Science 3*, 153–177.