

### 0.0.1 Influence Statistics for LME models

Influence statistics can be coarsely grouped by the aspect of estimation that is their primary target:

- overall measures compare changes in objective functions: (restricted) likelihood distance (Cook and Weisberg 1982, Ch. 5.2)
- influence on parameter estimates: Cook's (Cook 1977, 1979), MDFFITS (Belsley, Kuh, and Welsch 1980, p. 32)
- influence on precision of estimates: CovRatio and CovTrace
- influence on fitted and predicted values: PRESS residual, PRESS statistic (Allen 1974), DFFITS (Belsley, Kuh, and Welsch 1980, p. 15)
- outlier properties: internally and externally studentized residuals, leverage

### 0.0.2 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.

$$PRESS = \sum (y - y^{(k)})^2 \quad (1)$$

- $e_{-Q} = y_Q - x_Q \hat{\beta}^{-Q}$
- $PRESS_{(U)} = y_i - x_i \hat{\beta}_{(U)}$

### 0.0.3 DFBETA

$$DFBETA_a = \hat{\beta} - \hat{\beta}_{(a)} \quad (2)$$

$$= B(Y - Y_{\bar{a}}) \quad (3)$$

## 0.0.4 Influential Observations : DFBeta and DFBetas

Cook's distance refers to how far, on average, predicted y-values will move if the observation in question is dropped from the data set. dfbeta refers to how much a parameter estimate changes if the observation in question is dropped from the data set. Note that with k covariates, there will be k+1 dfbetas (the intercept,  $\beta_0$ , and 1  $\beta$  for each covariate). Cook's distance is presumably more important to you if you are doing predictive modeling, whereas dfbeta is more important in explanatory modeling.

## 0.1 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  $\hat{V}$ .

### 0.1.1 DFBETA

$$DFBETA_a = \hat{\beta} - \hat{\beta}_{(a)} \quad (4)$$

$$= B(Y - Y_{\hat{a}}) \quad (5)$$

### 0.1.2 DFFITS

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

$$DFFITS = \frac{\widehat{y}_i - \widehat{y}_{i(k)}}{s_{(k)}\sqrt{h_{ii}}}$$

## 0.2 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.

### 0.2.1 DFFITS

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

$$DFFITS = \frac{\widehat{y}_i - \widehat{y}_{i(k)}}{s_{(k)}\sqrt{h_{ii}}}$$

### 0.2.2 PRESS

The prediction residual sum of squares (PRESS) is a 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.

$$PRESS = \sum (y - y^{(k)})^2 \quad (6)$$

- $e_{-Q} = y_Q - x_Q\hat{\beta}^{-Q}$
- $PRESS_{(U)} = y_i - x_i\hat{\beta}_{(U)}$

## DFBETA

$$DFBETA_a = \hat{\beta} - \hat{\beta}_{(a)} \quad (7)$$

$$= B(Y - Y_a) \quad (8)$$

### 0.2.3 Residual diagnostics

For classical linear models, residual diagnostics are typically implemented as a plot of the observed residuals and the predicted values. A visual inspection for the presence of trends inform the analyst on the validity of distributional assumptions, and to detect outliers and influential observations.

### Abstract

This chapter is broken into two parts. The first part is a review of diagnostics methods for linear models, intended to acquaint the reader with the subject, and also to provide a basis for material covered in the second part. Particular attention is drawn to graphical methods.

## Influence Analysis

The basic rationale behind measuring influential cases is that when iteratively single units are omitted from the data, models based on these data should not produce substantially different estimates.

### Well Known Influence Measures

*“Regression Diagnostics: Identifying Influential Data and Source of Collinearity (1980)”*

by Belsley, Kuh, & Welsch is a landmark text in the field of residual diagnostics, and provides a foundation for much of the subsequent work.

**Cook’s Distance** Cooks Distance is a measure indicating to what extent model parameters are influenced by (a set of) influential data on which the model is based.

**DFBETAS** DFBETAS (standardized difference of the beta) is a measure that standardizes the absolute difference in parameter estimates between a (mixed effects) regression model based on a full set of data, and a model from which a (potentially influential) subset of data is removed. A value for DFBETAS is calculated for each parameter in the model separately.

### 0.2.4 Quantifying Influence

The basic procedure for quantifying influence is simple as follows:

- Fit the model to the data and obtain estimates of all parameters.
- Remove one or more data points from the analysis and compute updated estimates of model parameters.
- Based on full- and reduced-data estimates, contrast quantities of interest to determine how the absence of the observations changes the analysis.

Cook (1986) introduces powerful tools for local-influence assessment and examining perturbations in the assumptions of a model. In particular the effect of local perturbations of parameters or observations are examined.

## 0.3 Measures 2

### 0.3.1 Cook's Distance

- For variance components  $\gamma$

Diagnostic tool for variance components

$$C_{\theta i} = (\hat{(\theta)}_{[i]} - \hat{(\theta)})^T \text{cov}(\hat{(\theta)})^{-1} (\hat{(\theta)}_{[i]} - \hat{(\theta)})$$

### 0.3.2 Variance Ratio

- For fixed effect parameters  $\beta$ .

### 0.3.3 Variance Ratio

- For fixed effect parameters  $\beta$ .

### 0.3.4 Cook-Weisberg statistic

- For fixed effect parameters  $\beta$ .

### 0.3.5 Andrews-Pregibon statistic

- For fixed effect parameters  $\beta$ .

The Andrews-Pregibon statistic  $AP_i$  is a measure of influence based on the volume of the confidence ellipsoid. The larger this statistic is for observation  $i$ , the stronger the influence that observation will have on the model fit.

## 0.4 Measures 2

### 0.4.1 Cook's Distance

- For variance components  $\gamma$

Diagnostic tool for variance components

$$C_{\theta i} = ((\hat{\theta})_{[i]} - \hat{\theta})^T \text{cov}(\hat{\theta})^{-1} ((\hat{\theta})_{[i]} - \hat{\theta})$$

### 0.4.2 Variance Ratio

- For fixed effect parameters  $\beta$ .

### 0.4.3 Cook-Weisberg statistic

- For fixed effect parameters  $\beta$ .

### 0.4.4 Andrews-Pregibon statistic

- For fixed effect parameters  $\beta$ .

The Andrews-Pregibon statistic  $AP_i$  is a measure of influence based on the volume of the confidence ellipsoid. The larger this statistic is for observation  $i$ , the stronger the influence that observation will have on the model fit.

## 0.5 Measures 2



## Random Effects

A large value for  $CD(u)_i$  indicates that the  $i$ -th observation is influential in predicting random effects.

## linear functions

$CD(\psi)_i$  does not have to be calculated unless  $CD(\beta)_i$  is large.

### 0.5.1 Information Ratio

### 0.5.2 Cook's Distance

- For variance components  $\gamma$

Diagnostic tool for variance components

$$C_{\theta i} = (\hat{(\theta)}_{[i]} - \hat{(\theta)})^T \text{cov}(\hat{(\theta)})^{-1} (\hat{(\theta)}_{[i]} - \hat{(\theta)})$$

## 0.6 Zewotir Measures of Influence in LME Models

Zewotir describes a number of approaches to model diagnostics, investigating each of the following;

- Variance components
- Fixed effects parameters
- Prediction of the response variable and of random effects
- likelihood function

# Bibliography