

Influence Analysis

The general idea of quantifying the influence of one or more observations relies on computing parameter estimates based on all data points, removing the cases in question from the data, refitting the model, and computing statistics based on the change between full-data and reduced-data estimation.

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.

0.0.1 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.

0.0.2 Well Known Influence Measures

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

0.0.3 Influence Diagnostics: Basic Idea and Statistics

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: Cooks (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

For linear models for uncorrelated data, it is not necessary to refit the model after removing a data point in order to measure the impact of an observation on the model. The change in fixed effect estimates, residuals, residual sums of squares, and the variance-covariance matrix of the fixed effects can be computed based on the fit to the full data alone. By contrast, in mixed models several important complications arise. Data points can affect not only the fixed effects but also the covariance parameter estimates on which the fixed-effects estimates depend.

Furthermore, closed-form expressions for computing the change in important model quantities might not be available. This section provides background material for the various influence diagnostics available with the MIXED procedure. See the section Mixed Models Theory for relevant expressions and definitions. The parameter vector denotes all unknown parameters in the and matrix. The observations whose influence is

being ascertained are represented by the set and referred to simply as "the observations in ." The estimate of a parameter vector, such as β , obtained from all observations except those in the set is denoted $\hat{\beta}_{(-i)}$. In case of a matrix Σ , the notation $\Sigma_{(-i)}$ represents the matrix with the rows in removed; these rows are collected in \mathbf{r}_i . If Σ is symmetric, then notation implies removal of rows and columns. The vector \mathbf{r}_i comprises the responses of the data points being removed, and $\Sigma_{(-i)}$ is the variance-covariance matrix of the remaining observations. When \mathbf{r}_i , lowercase notation emphasizes that single points are removed, such as \mathbf{r}_i .

0.1 LME diagnostic measures

0.1.1 Cook's Distance

- For variance components γ

Diagnostic tool for variance components

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

0.1.2 Variance Ratio

- For fixed effect parameters β .

0.1.3 Cook-Weisberg statistic

- For fixed effect parameters β .

0.1.4 Andrews-Pregibon statistic

- For fixed effect parameters β .

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.

PRESS

Schabenberger (2004) describes the use of the *PRESS* and *DFITS* in determining influence.

The *PRESS* residual is the difference between the observed value and the predicted (marginal) value.

$$\hat{e}_{i(U)} = y_i - x\hat{\beta}_{(U)} \quad (1)$$

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.

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

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

$$PRESS = \sum (y - y^{-Q})^2$$

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

PRESS Residuals and PRESS Statistic

The predicted residual sum of squares (PRESS) statistic is a form of cross-validation used in regression analysis to provide a summary measure of the fit of a model to a sample of observations that were not themselves used to estimate the model. It is calculated as the sums of squares of the prediction residuals for those observations.

A fitted model having been produced, each observation in turn is removed and the model is refitted using the remaining observations. The out-of-sample predicted value is calculated for the omitted observation in each case, and the PRESS statistic is calculated as the sum of the squares of all the resulting prediction errors:[4]

$$\text{PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{i,-i})^2$$

Given this procedure, the PRESS statistic can be calculated for a number of candidate model structures for the same dataset, with the lowest values of PRESS indicating the best structures. Models that are over-parameterised (over-fitted) would tend to give small residuals for observations included in the model-fitting but large residuals for observations that are excluded.

0.1.5 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.

0.2 Measures of Influence

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

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. Cooks Distance (D_i) is an overall measure of the combined impact of the i th case of all estimated regression coefficients. It uses the same structure for measuring the combined impact of the differences in the estimated regression coefficients when the k th case is deleted. $D_{(k)}$ can be calculated without fitting a new regression coefficient each time an observation is deleted.

Cooks D Measure of aggregate impact of each observation on the group of regression coefficients, as well as the group of fitted values. Values larger than $4/n$ are considered highly influential.

DFFITS Measure of how much an observation has effected its fitted value from the regression model. Values larger than $2\sqrt{(k+1)/n}$ in absolute value are considered highly influential.

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.

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

DFBETAS Measure of how much an observation has effected the estimate of a regression coefficient (there is one DFBETA for each regression coefficient, including

the intercept). Values larger than $2/\sqrt{n}$ in absolute value are considered highly influential.

DFBETAS 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.

DFBETAS The measure that measures how much impact each observation has on a particular predictor is DFBETAs The DFBETA for a predictor and for a particular observation is the difference between the regression coefficient calculated for all of the data and the regression coefficient calculated with the observation deleted, scaled by the standard error calculated with the observation deleted.

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.

$$DFBETA_a = \hat{\beta} - \hat{\beta}_{(a)} \quad (3)$$

$$= B(Y - Y_{\bar{a}}) \quad (4)$$

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.

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

DFBETAS

$$DFBETA_a = \hat{\beta} - \hat{\beta}_{(a)} \quad (5)$$

$$= B(Y - Y_{\bar{a}}) \quad (6)$$

Studentized Residuals Residuals divided by their estimated standard errors (like t-statistics). Observations with values larger than 3 in absolute value are considered outliers.

Leverage Values (Hat Diag) Measure of how far an observation is from the others in terms of the levels of the independent variables (not the dependent variable). Observations with values larger than $2(k+1)/n$ are considered to be potentially highly influential, where k is the number of predictors and n is the sample size.

0.2.1 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, β_0 , and 1 β 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.

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.

Bibliography

Belsley, D. A. and E. Kuh. Welsch., re (1980). regression diagnostics: Identifying influential data and sources of collinearity. *Uiley Series in Probability and Mathematical Statistics*.

Cook, R. (1986). Assessment of local influence. *Journal of the Royal Statistical Society. Series B (Methodological)* 48(2), 133–169.

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

Zewotir, T. and J. S. Galpin (2005). Influence diagnostics for linear mixed models. *Journal of Data Science* 3(2), 153–177.