

# 1 Cook's Distance for LMEs

Diagnostic methods for fixed effects are generally analogues of methods used in classical linear models. Cook (1977) greatly expanded the study of residuals and influence measures. Cook's Distance, denoted as  $D_{(i)}$ , is a well known diagnostic technique used in classical linear models, used as an overall measure of the combined impact of the  $i$ th case of all estimated regression coefficients. Cook's key observation was the effects of deleting each observation in turn could be calculated with little additional computation. That is to say,  $D_{(i)}$  can be calculated without fitting a new regression coefficient each time an observation is deleted. Consequently deletion diagnostics have become an integral part of assessing linear models.

The focus of this analysis is related to the estimation of point estimates (i.e. regression coefficients). It must be pointed out that the effect on the precision of estimates is separate from the effect on the point estimates. Data points that have a small Cook's distance, for example, can still greatly affect hypothesis tests and confidence intervals, if their influence on the precision of the estimates is large.

As well as individual observations, Cook's distance can be used to analyse the influence of observations in subset  $U$  on a vector of parameter estimates (Cook, 1977).

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

$$\delta_{(U)} = \hat{\beta} - \hat{\beta}_{(U)} \quad (2)$$

Diagnostic methods for variance components are based on 'one-step' methods. Cook (1986) gives a completely general method for assessing the influence of local departures from assumptions in statistical models.

For fixed effects parameter estimates in LME models, the Cook's distance can be extended to measure influence on these fixed effects.

## 1.1 Cook's Distance

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.

Cook's Distance is a well known diagnostic technique used in classical linear models, extended to LME models. For LME models, two formulations exist; a Cook's distance that examines the change in fixed fixed parameter estimates, and another that examines the change in random effects parameter estimates. The outcome of either Cook's distance is a scaled change in either  $\beta$  or  $\theta$ .

Cook's D is a good measure of the influence of an observation and is proportional to the sum of the squared differences between predictions made with all observations in the analysis and predictions made leaving out the observation in question. If the predictions are the same with or without the observation in question, then the observation has no influence on the regression model. If the predictions differ greatly when the observation is not included in the analysis, then the observation is influential.

## 1.2 Cook's Distance

In statistics, Cook's Distance or Cook's D is a commonly used estimate of the influence of a data point when performing least squares regression analysis.[1] In a practical ordinary least squares analysis, Cook's distance can be used in several ways: to indicate data points that are particularly worth checking for validity; to indicate regions of the design space where it would be good to be able to obtain more data points. It is named after the American statistician R. Dennis Cook, who introduced the concept in 1977.

## 1.3 Cook's Distance

Cook's  $D$  statistics (i.e. colloquially Cook's Distance) is a measure of the influence of observations in subset  $U$  on a vector of parameter estimates (Cook, 1977).

$$\delta_{(U)} = \hat{\beta} - \hat{\beta}_{(U)}$$

If  $V$  is known, Cook's D can be calibrated according to a chi-square distribution with degrees of freedom equal to the rank of  $\mathbf{X}$  (?).

For fixed effects parameter estimates in LME models, the Cook's distance can be extended to measure influence on these fixed effects.

$$CD_i(\beta) = \frac{(c_{ii} - r_{ii}) \times t_i^2}{r_{ii} \times p}$$

For random effect estimates, the Cook's distance is

$$CD_i(b) = g'_{(i)}(I_r + \text{var}(\hat{b})D)^{-2}\text{var}(\hat{b})g_{(i)}.$$

Large values for Cook's distance indicate observations for special attention.

As well as individual observations, Cook's distance can be used to analyse the influence of observations in subset  $U$  on a vector of parameter estimates (Cook, 1977).

$$\hat{e}_{i(U)} = y_i - x\hat{\beta}_{(U)} \quad (3)$$

$$\delta_{(U)} = \hat{\beta} - \hat{\beta}_{(U)} \quad (4)$$

For LME models, Cook's distance can be extended to model influence diagnostics by definining.

It is also desirable to measure the influence of the case deletions on the covariance matrix of  $\hat{\beta}$ .

Diagnostic methods for fixed effects are generally analogues of methods used in classical linear models. Diagnostic methods for variance components are based on 'one-step' methods. Cook (1986) gives a completely general method for assessing the influence of local departures from assumptions in statistical models.

## 1.4 Change in the precision of estimates

The effect on the precision of estimates is separate from the effect on the point estimates. Data points that have a small Cook's distance, for example, can still greatly affect hypothesis tests and confidence intervals, if their influence on the precision of the estimates is large.

### 1.4.1 Computational Limitations for Cook's Distance

Application of Cook's Distances are limited by computation tractability.

Application of case-deletion diagnostics offer some interested for Method Comparison Studies

Care must be given when interpreting these plots. For example the position of case 68 on the BSVR indicates that that case 68

Any diagnostic plot may constructed using Overall variability and intermethod bias.

#### 1.4.2 Taxonomy of Cook's Distances for LMEs

Zewotir and Galpin (2005) discusses a taxonomy of Cook's distance when applied to LME models.

- For variance components  $\gamma$ :  $CD(\gamma)_i$ ,
- For fixed effect parameters  $\beta$ :  $CD(\beta)_i$ ,
- For random effect parameters  $\mathbf{u}$ :  $CD(u)_i$ ,
- For linear functions of  $\beta$ :  $CD(\psi)_i$

## 2 Cook's Distance

In classical linear regression, a commonly used measure of influence is Cook's distance. It is used as a measure of influence on the regression coefficients.

For linear mixed effects models, Cook's distance can be extended to model influence diagnostics by defining.

$$C_{\beta i} = \frac{(\hat{\beta} - \hat{\beta}_{[i]})^T (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}) (\hat{\beta} - \hat{\beta}_{[i]})}{p}$$

It is also desirable to measure the influence of the case deletions on the covariance matrix of  $\hat{\beta}$ .

## 3 Cook's Distance - Implementation with R

Cook's Distance is a measure indicating to what extent model parameters are influenced by (a set of) influential data on which the model is based. This function computes the Cook's distance based on the information returned by the `estex()` function.

## 4 Cook's Distance for LMEs

Diagnostic methods for fixed effects are generally analogues of methods used in classical linear models. Diagnostic methods for variance components are based on 'one-step'

methods. Cook (1986) gives a completely general method for assessing the influence of local departures from assumptions in statistical models.

Cook's Distance How to extract/compute leverage and Cook's distances for linear mixed effects models

Does anyone know how to compute (or extract) leverage and Cook's distances for a mer class object (obtained through lme4 package)? I'd like to plot these for a residuals analysis.

You should have a look at the R package influence.ME. It allows you to compute measures of influential data for mixed effects models generated by lme4.

An example model:

```
library(lme4)
model <- lmer(mpg ~ disp + (1 | cyl), mtcars)
```

The function influence is the basis for all further steps:

```
library(influence.ME)
infl <- influence(model, obs = TRUE)
```

Calculate Cook's distance:

```
cooks.distance(infl)
```

Plot Cook's distance:

```
plot(infl, which = "cook")
```

enter image description here

How to extract/compute leverage and Cook's distances for linear mixed effects models

Does anyone know how to compute (or extract) leverage and Cook's distances for a mer class object (obtained through lme4 package)? I'd like to plot these for a residuals analysis.

You should have a look at the R package `influence.ME`. It allows you to compute measures of influential data for mixed effects models generated by `lme4`.

An example model:

```
library(lme4)
model <- lmer(mpg ~ disp + (1 | cyl), mtcars)
```

The function `influence` is the basis for all further steps:

```
library(influence.ME)
infl <- influence(model, obs = TRUE)
Calculate Cook's distance:
```

```
cooks.distance(infl)
Plot Cook's distance:
```

```
plot(infl, which = "cook")
```

## 4.1 Cook's Distance

In classical linear regression, a commonly used measure of influence is Cook's distance. It is used as a measure of influence on the regression coefficients.

For linear mixed effects models, Cook's distance can be extended to model influence diagnostics by defining.

$$C_{\beta i} = \frac{(\hat{\beta} - \hat{\beta}_{[i]})^T (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}) (\hat{\beta} - \hat{\beta}_{[i]})}{p}$$

It is also desirable to measure the influence of the case deletions on the covariance matrix of  $\hat{\beta}$ .

### 4.1.1 Random Effects

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

## 4.2 linear functions

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

Cook (1986) gave a completely general method for assessing influence of local departures from assumptions in statistical models.

For LME models, two formulations exist; a Cook's distance that examines the change in fixed fixed parameter estimates, and another that examines the change in random effects parameter estimates. The outcome of either Cook's distance is a scaled change in either  $\beta$  or  $\theta$ .

## 5 Cook's Distance for LMEs

Diagnostic methods for fixed effects are generally analogues of methods used in classical linear models. Diagnostic methods for variance components are based on 'one-step' methods. Cook (1986) gives a completely general method for assessing the influence of local departures from assumptions in statistical models.

### 5.0.1 Interpretation

Specifically  $D_i$  can be interpreted as the distance one's estimates move within the confidence ellipsoid that represents a region of plausible values for the parameters.[clarification needed] This is shown by an alternative but equivalent representation of Cook's distance in terms of changes to the estimates of the regression parameters between the cases where the particular observation is either included or excluded from the regression analysis.

### 5.1 Interpreting Cook's Distance

A common rule of thumb is that an observation with a value of Cook's D over 1.0 has too much influence. As with all rules of thumb, this rule should be applied judiciously and not thoughtlessly.

### 5.2 Interpreting Cook's Distance

Some texts tell you that points for which Cook's distance is higher than 1 are to be considered as influential. Other texts give you a threshold of  $4/N$  or  $4/(Nk1)$ , where

N is the number of observations and k the number of explanatory variables. In your case the latter formula should yield a threshold around 0.1 .

John Fox (1), in his booklet on regression diagnostics is rather cautious when it comes to giving numerical thresholds. He advises the use of graphics and to examine in closer details the points with "values of D that are substantially larger than the rest". According to Fox, thresholds should just be used to enhance graphical displays.

In your case the observations 7 and 16 could be considered as influential. Well, I would at least have a closer look at them. The observation 29 is not substantially different from a couple of other observations.

(1) Fox, John. (1991). Regression Diagnostics: An Introduction. Sage Publications.

## 6 Exention of Cook's Distance methodology to LME models

Cook's Distance is extended to LME models. For LME models, two formulations exist; a Cook's distance that examines the change in fixed fixed parameter estimates, and another that examines the change in random effects parameter estimates. The outcome of either Cook's distance is a scaled change in either  $\beta$  or  $\theta$ .

Diagnostic methods for variance components are based on 'one-step' methods. Cook (1986) gives a completely general method for assessing the influence of local departures from assumptions in statistical models. For fixed effects parameter estimates in LME models, the Cook's distance can be extended to measure influence on these fixed effects.

$$CD_i(\beta) = \frac{(c_{ii} - r_{ii}) \times t_i^2}{r_{ii} \times p}$$

For random effect estimates, the Cook's distance is

$$CD_i(b) = g'_{(i)}(I_r + \text{var}(\hat{b})D)^{-2}\text{var}(\hat{b})g_{(i)}.$$

Large values for Cook's distance indicate observations for special attention.

Cook's Distance was extended from classical linear models to LME models. For linear mixed effects models, Cook's distance can be extended to model influence diagnostics by definining.



$$CD_{\beta i} = \frac{(\hat{\beta} - \hat{\beta}_{[i]})^T (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X}) (\hat{\beta} - \hat{\beta}_{[i]})}{p}$$

It is also desirable to measure the influence of the case deletions on the covariance matrix of  $\hat{\beta}$ .

## References

- Cook, R. (1977). Detection of influential observations in linear regression. *Technometrics* 19, 15–18.
- Cook, R. (1986). Assessment of local influence. *Journal of the Royal Statistical Society. Series B (Methodological)* 48(2), 133–169.
- Zewotir, T. and J. Galpin (2005). Influence diagnostics for linear mixed models. *Journal of Data Science* 3, 153–177.