

STAT 420: Methods of Applied Statistics

Model Diagnostics — Outliers

Ruoqing Zhu, Ph.D. <rqzhu@illinois.edu>

Course website: <https://sites.google.com/site/teazrq/teaching/STAT420>

Department of Statistics
University of Illinois at Urbana-Champaign
March 25, 2017

Detecting Outliers

- Previously we learned how to check the normality and constant variance assumptions.
- Sometimes there are observations that are distant from others, and we want to detect and remove them from the data.
- Usually these are observations that are “far away” from the regression line, and have large influence on the estimated $\hat{\beta}$.
- We will introduce some statistical methods for detecting them.

An important notation change

- In previous lecture notes, we use p to denote the number of **predictors**, and $p + 1$ is the number of parameters in a linear model if an intercept term is used.
- From now on, we don't want to discuss the two cases (with or without intercept) separately. So, we will use p as **the total number of parameters**, regardless of whether an intercept term is used or not.
- This is mainly for simplification of the notation.

Cook's Distance as a Measure of Influence

- Cook's distance, D_i is a measure of the distance between predicted values when all of the observations are used versus when a given observation i is omitted.
- Let $\hat{\beta}_{(-i)}$ denote the OLS estimate when the i^{th} observations is deleted.
- This new parameter estimate should produce a new set of predicted values:

$$\hat{\mathbf{y}}_{(-i)} = \mathbf{X}\hat{\beta}_{(-i)}$$

- This should be different from the original predicted values:

$$\hat{\mathbf{y}}_{(-i)} - \hat{\mathbf{y}} = \mathbf{X}(\hat{\beta}_{(-i)} - \hat{\beta})$$

- Hence, the squared distance is

$$(\hat{\mathbf{y}}_{(-i)} - \hat{\mathbf{y}})^T (\hat{\mathbf{y}}_{(-i)} - \hat{\mathbf{y}}) = (\hat{\beta}_{(-i)} - \hat{\beta})^T \mathbf{X}^T \mathbf{X} (\hat{\beta}_{(-i)} - \hat{\beta})$$

Cook's Distance as a Measure of Influence

- It is interesting that this squared distance can be simplified into something that does not require refitting of the model:
- The Cook's distance (1977), D_i is

$$\begin{aligned} D_i &= \frac{(\hat{\beta}_{(-i)} - \hat{\beta})^T \mathbf{X}^T \mathbf{X} (\hat{\beta}_{(-i)} - \hat{\beta})}{p \hat{\sigma}^2} \\ &= \frac{e_i^2}{p \hat{\sigma}^2} \left[\frac{h_{ii}}{(1 - h_{ii})^2} \right] \end{aligned}$$

where all components in the second line are from the **original model**:

- e_i is the residual of the i th subject
 - h_{ii} is the i th diagonal element in the “hat matrix” \mathbf{H}
 - $\hat{\sigma}^2$ is the estimated variance of residual
- The second equality is due to Beckman and Trussel (1974).

Cook's distance: remove subject

```
1 > fit = lm(Work ~ FPC, data = ATT)
2 # Calculate Cook's distance for subject 58
3 # Refit by removing the subject
4 > ATTi = ATT[-58,]
5 > fiti = lm(Work ~ FPC, data = ATTi)
6 # Calculate Di using the first approach
7 > X = cbind(1, ATT$FPC)
8 > diff = X %*% (fit$coefficients - fiti$coefficients)
9 > sigma2 = sum(fit$residuals^2)/(nrow(ATT)-2)
10 > t(diff) %*% diff / 2 / sigma2
11           [,1]
12 [1,] 0.05399123
```

- Example: Calculate the Cook's distance D_i statistic for all subjects in the ATT data, and identify the subject with the largest D_i
- The subject ID is:

A : 58 B : 24 C : 16

Cook's distance using \mathbf{H}

- The easier and default way is to use the “hat matrix”, because it does not involve refitting the model (computational cost is high).

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- h_{ii} is also called leverage. Large values of h_{ii} are due to extreme values in \mathbf{X} .
- Recall the interpretation of \mathbf{H} , and by matrix operations, we have

$$\hat{y}_i = \sum_{j=1}^n h_{ij} y_j.$$

Hence a point with high leverage has the potential to greatly influence the fit.

- A “rule of thumb” is that leverages of more than $2\bar{h}_{ii}$ (some suggests $3\bar{h}_{ii}$) should be looked at more closely, where $\bar{h}_{ii} = n^{-1} \sum_i h_{ii}$

- Some properties about \mathbf{H} should be noted: $0 \leq h_{ii} \leq 1$.
- This is due to the fact that \mathbf{H} is idempotent, i.e., $\mathbf{H}^2 = \mathbf{H}$. Hence
$$h_{ii} = \sum_{j=1}^n h_{ij}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$$
- Hence $h_{ii} \geq 0$ because $\sum_{j=1}^n h_{ij}^2 \geq 0$.
- $h_{ii} \leq 1$ can be shown by contradiction: if $h_{ii} > 1$, then $\sum_{j \neq i} h_{ij}^2$ has to be negative.

Cook's distance: use H

```
1 > H = X %*% solve(t(X) %*% X) %*% t(X)
2 > h58 = H[58, 58]
3 > (h58 / (1 - h58)^2) * fit$residuals[58]^2 / 2 / sigma2
4      58
5 0.05399123
```

- Example: Calculate the Cook's distance D_i statistic for all subjects in the ATT data using the hat matrix approach.
- The subject that has the smallest D_i :

A : 74 B : 57 C : 79

Build-in approaches

```
1 # leverage of all subjects
2 > hat(X, intercept =FALSE)
3 # Di of all subjects
4 > cooks.distance( fit )
5
6 > hat(X, intercept =FALSE)[58]
7 [1] 0.0113612
8 > cooks.distance( fit )[58]
9          58
10 0.05399123
11 > which.max(cooks.distance( fit ))
12 24
13 24
```

- In practice, any observation with $D_i > 0.5$ could be influential. > 1 is quite likely to be influential. Or, if it stands out from the other D_i values, it is almost certainly influential.
- Cook's distance D_i is related to the F test statistic with degrees of freedom $(p, n - p)$.
- In practice, one compares D_i with $F_{0.5, p, n-p}$, and investigate any observation above that threshold.

- DFFITS is another practically useful approach.

$$\text{DFFITS}_i = \frac{\hat{y}_i - \hat{y}_{(-i)}}{\sqrt{\hat{\sigma}_{(-i)}^2 h_{ii}}}$$

where $\sigma_{(-i)}^2$ is the MSE estimated by removing subject i .

- “Rule of thumb”: if DFFITS_i is greater than $2\sqrt{\frac{p+1}{n-p-1}}$, its worth investigating.

- **Warning:** an observation is influential does NOT mean that it should be removed immediately!!
- Iteratively reweighting the observations may reduce their influence (`rlm` in `MASS` package).
- Sometimes we may consider giving up the squared error (ℓ_2 norm) in OLS and use alternative approaches to deal with outliers
 - Use least absolute residuals instead of squared loss (`VGAM` package)
 - Use quantile regression (`quantreg` package)
- Sometimes we may also consider giving up the linear functional form and use nonparametric regressions such as random forests (`randomForest` package), generalized additive model (`gam` package), kernel method (`loess` function), etc.