

Lecture 13

Model Diagnostics

STAT 8020 Statistical Methods II
September 18, 2019

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Whitney Huang
Clemson University

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

1 Leverage

2 Studentized & Jackknife Residuals

3 DFFITS

4 Non-Constant Variance & Transformation

Recall in MLR that $\hat{Y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where \mathbf{H} is the hat-matrix

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Recall in MLR that $\hat{Y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where \mathbf{H} is the hat-matrix

- The leverage score for the i_{th} observation is defined as:

$$h_i = \mathbf{H}_{ii}$$

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Recall in MLR that $\hat{Y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where \mathbf{H} is the hat-matrix

- The leverage score for the i_{th} observation is defined as:

$$h_i = \mathbf{H}_{ii}$$

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Recall in MLR that $\hat{Y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where \mathbf{H} is the hat-matrix

- The leverage score for the i_{th} observation is defined as:

$$h_i = \mathbf{H}_{ii}$$

- Can show that $\text{Var}(e_i) = \sigma^2(1 - h_i)$, where $e_i = Y_i - \hat{Y}_i$ is the residual for the i_{th} observation

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Recall in MLR that $\hat{Y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{H} \mathbf{Y}$ where \mathbf{H} is the hat-matrix

- The leverage score for the i_{th} observation is defined as:

$$h_i = \mathbf{H}_{ii}$$

- Can show that $\text{Var}(e_i) = \sigma^2(1 - h_i)$, where $e_i = Y_i - \hat{Y}_i$ is the residual for the i_{th} observation

Recall in MLR that $\hat{Y} = X(X^T X)^{-1} X^T Y = HY$ where H is the hat-matrix

- The leverage score for the i_{th} observation is defined as:

$$h_i = H_{ii}$$

- Can show that $\text{Var}(e_i) = \sigma^2(1 - h_i)$, where $e_i = Y_i - \hat{Y}_i$ is the residual for the i_{th} observation
- $\frac{1}{n} \leq h_i \leq 1$, $1 \leq i \leq n$ and $\bar{h}_i = \frac{p}{n} \Rightarrow$ a “rule of thumb” is that leverages of more than $\frac{2p}{n}$ should be looked at more closely

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

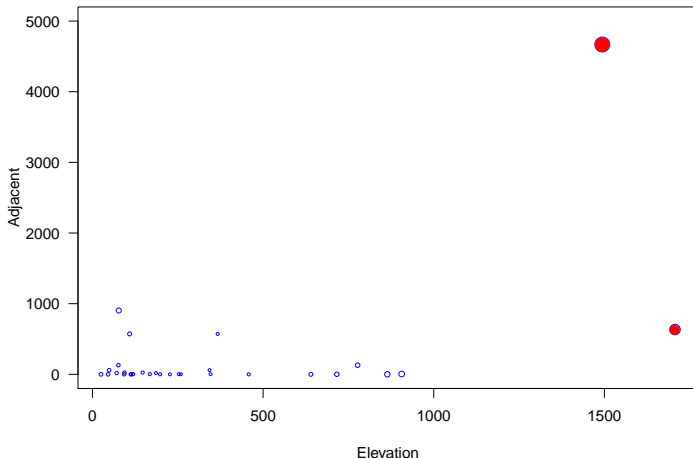
Leverage Scores of Species ~ Elev + Adj

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation



As we have seen $\text{Var}(e_i) = \sigma^2(1 - h_i)$, this suggests the use of

$$r_i = \frac{e_i}{\hat{\sigma}\sqrt{(1-h_i)}}$$

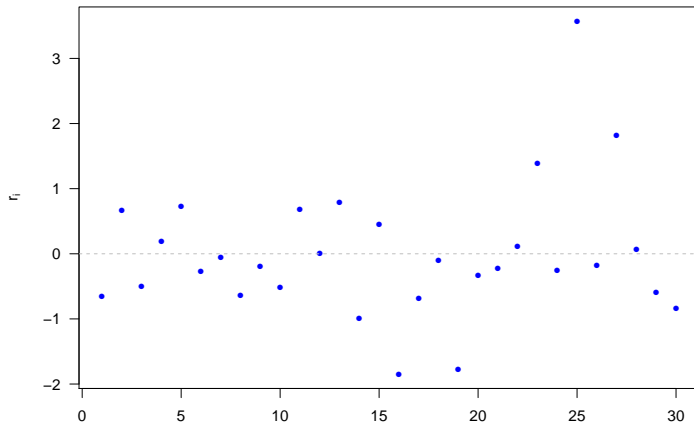
- r_i 's are called **studentized residuals**. r_i 's are sometimes preferred in residual plots as they have been standardized to have equal variance.
- If the model assumptions are correct then $\text{Var}(r_i) = 1$ and $\text{Corr}(e_i, e_j)$ tends to be small

Studentized Residuals of $\text{Species} \sim \text{Elev} + \text{Adj}$

Model Diagnostics

CLEMSON
UNIVERSITY

Studentized Residuals



Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$
- The observation i is an outlier if $\hat{Y}_{i(i)} - Y_i$ is “large”

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$
- The observation i is an outlier if $\hat{Y}_{i(i)} - Y_i$ is “large”

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$
- The observation i is an outlier if $\hat{Y}_{i(i)} - Y_i$ is “large”
- Can show $\text{Var}(\hat{Y}_{i(i)} - Y_i) = \sigma^2 \left(1 + \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right)$

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$
- The observation i is an outlier if $\hat{Y}_{i(i)} - Y_i$ is “large”
- Can show $\text{Var}(\hat{Y}_{i(i)} - Y_i) = \sigma^2 \left(1 + \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right)$

- For a given model, exclude the observation i and recompute $\hat{\beta}_{(i)}$, $\hat{\sigma}_{(i)}$ to obtain $\hat{Y}_{i(i)}$
- The observation i is an outlier if $\hat{Y}_{i(i)} - Y_i$ is “large”
- Can show $\text{Var}(\hat{Y}_{i(i)} - Y_i) = \sigma^2 \left(1 + \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right)$
- Define the **jackknife residuals** as

$$t_i = \frac{\hat{Y}_{i(i)} - Y_i}{\sqrt{\hat{\sigma}^2 \left(1 + \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right)}}$$

which are distributed as a t_{n-p} if the model is correct and $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

Jackknife Residuals of Species ~ Elev + Adj

Model Diagnostics

CLEMSON
UNIVERSITY

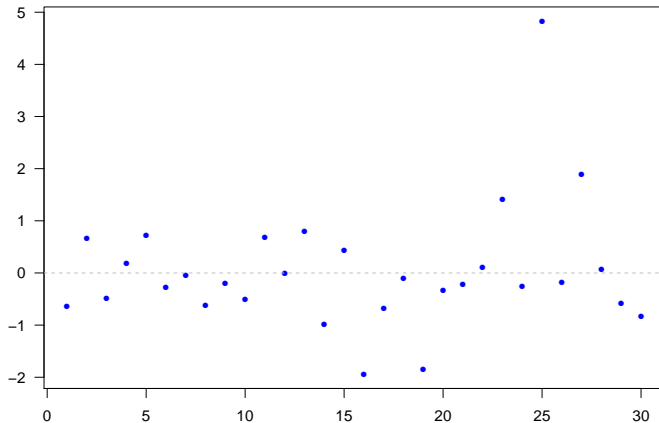
Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Jackknife Residuals

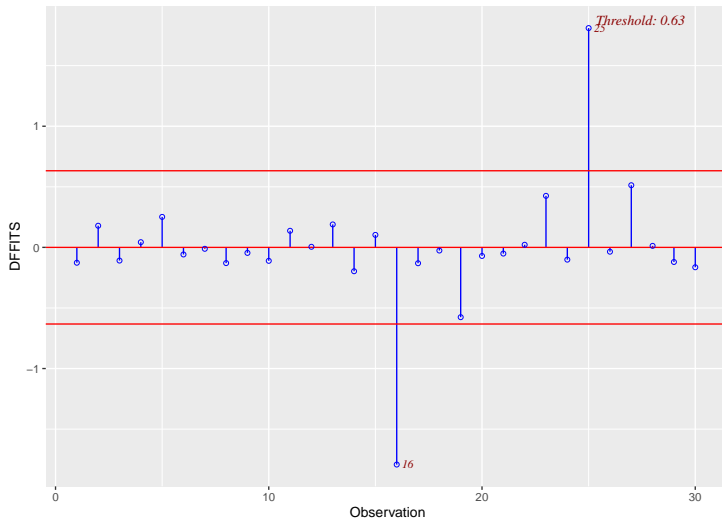


DFFITS

- Difference between the fitted values \hat{Y}_i and the predicted values $\hat{Y}_{i(i)}$
- $$DFFITS_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{MSE_{(i)} h_i}}$$
- Concern if absolute value greater than 1 for small data sets, or greater than $2\sqrt{p/n}$ for large data sets

DFFITS of Species ~ Elev + Adj

Influence Diagnostics for Species



Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Residual Plot of Species ~ Elev + Adj

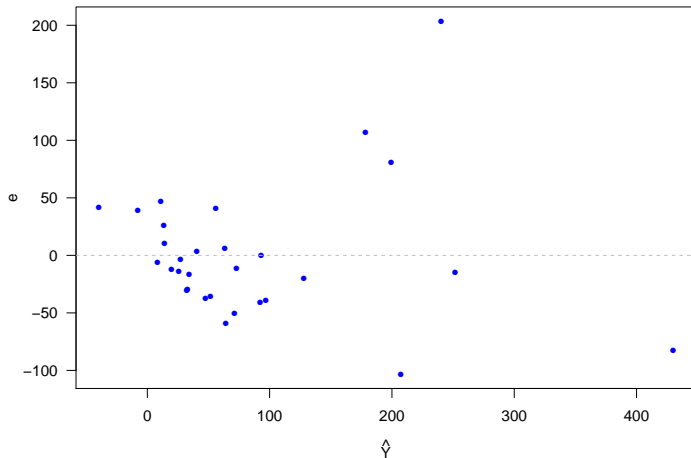
Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Residuals



Residual Plot After Square Root Transformation

Leverage

Studentized &
Jackknife Residuals

DFFITS

Non-Constant
Variance &
Transformation

Residuals

