

Regression Diagnostics

Merlise Clyde

September 2, 2019

Outline

- ▶ Leverage

Outline

- ▶ Leverage
- ▶ Standardized Residuals

Outline

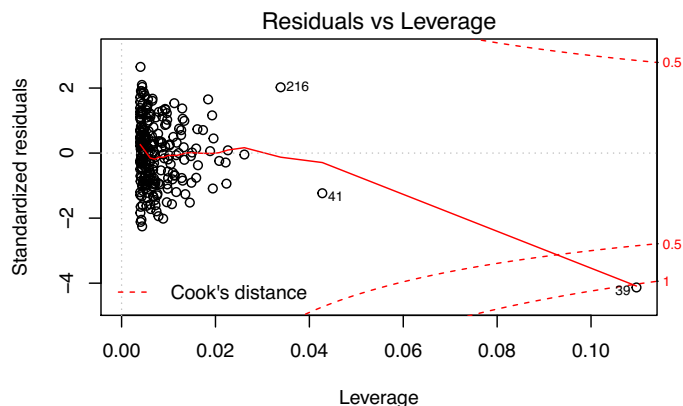
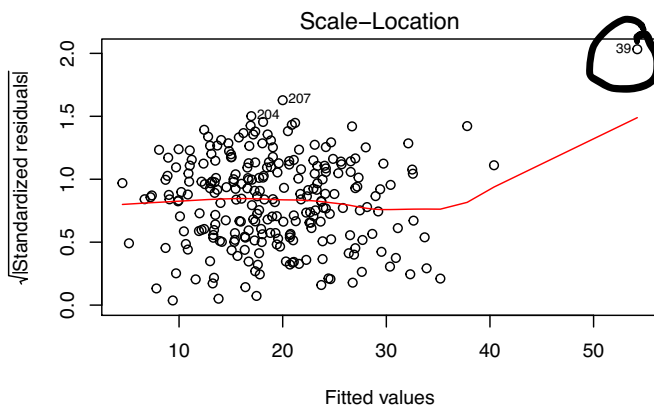
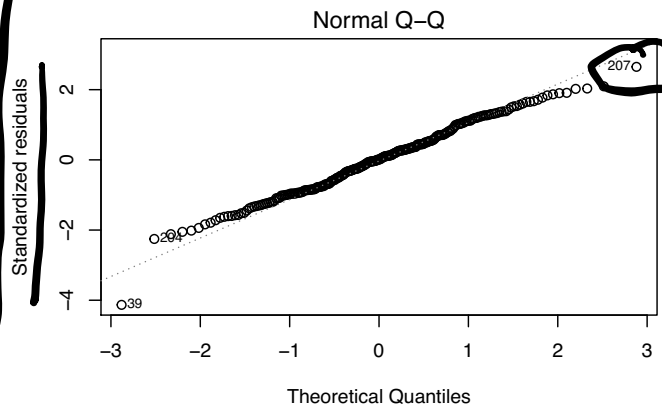
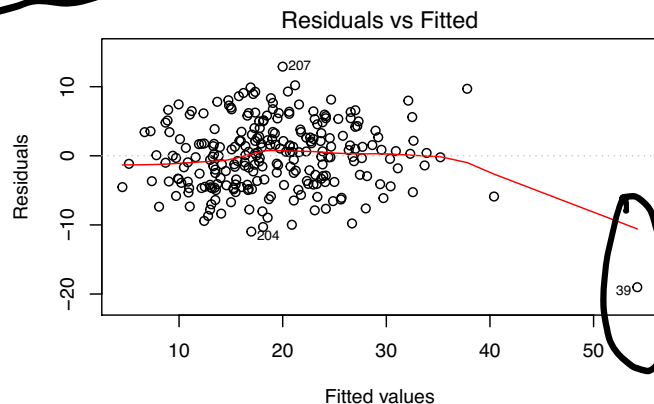
- ▶ Leverage
- ▶ Standardized Residuals
- ▶ Outlier Test

Outline

- ▶ Leverage
- ▶ Standardized Residuals
- ▶ Outlier Test
- ▶ Cook's Distance

Residual Plots

```
bodyfat.lm = lm(Bodyfat ~ Abdomen, data=bodyfat)
par(mfrow=c(2,2))
plot(bodyfat.lm, ask=F)
```



Features of Plots

- ▶ Residuals versus fitted values

Features of Plots

- ▶ Residuals versus fitted values
- ▶ Normal Quantile: check normality of residuals or look for heavier tails than normal where observed quantiles are larger than expected under normality

Features of Plots

- ▶ Residuals versus fitted values
- ▶ Normal Quantile: check normality of residuals or look for heavier tails than normal where observed quantiles are larger than expected under normality
- ▶ Scale-Location plot:
Detect if the spread of the residuals is constant over the range of fitted values. (Constant variance with mean)

Features of Plots

- ▶ Residuals versus fitted values
- ▶ Normal Quantile: check normality of residuals or look for heavier tails than normal where observed quantiles are larger than expected under normality
- ▶ Scale-Location plot:
Detect if the spread of the residuals is constant over the range of fitted values. (Constant variance with mean)
- ▶ standardized residuals versus leverage with contours of Cook's distance: shows influential points where points greater than 1 or $4/n$ are considered influential

Features of Plots

- ▶ Residuals versus fitted values
- ▶ Normal Quantile: check normality of residuals or look for heavier tails than normal where observed quantiles are larger than expected under normality
- ▶ Scale-Location plot:
Detect if the spread of the residuals is constant over the range of fitted values. (Constant variance with mean)
- ▶ standardized residuals versus leverage with contours of Cook's distance: shows influential points where points greater than 1 or $4/n$ are considered influential
- ▶ Case 39 appears to be influential and have a large standardized residual!

Hat Matrix

► predictions

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY$$

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

$$X\hat{\beta} = \hat{Y} \quad \text{compares to } Y$$

Hat Matrix

- predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- Hat Matrix or Projection Matrix

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix
 - ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix
 - ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$
 - ▶ symmetric

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix

- ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$

- ▶ symmetric

- ▶ leverage values are the diagonal elements h_{ii} : $0 \leq h_{ii} \leq 1$

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix
 - ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$
 - ▶ symmetric
 - ▶ leverage values are the diagonal elements h_{ii} : $0 \leq h_{ii} \leq 1$
- ▶ Predictions

$$\hat{Y}_i = h_{ii}Y_i + \sum_{j \neq i} h_{ij}Y_j$$

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

$$() = [] ()$$

$$\hat{Y}_i = h_{ii}Y_i + \sum_{j \neq i} h_{ij}Y_j$$

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix
 - ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$
 - ▶ symmetric
 - ▶ leverage values are the diagonal elements h_{ii} : $0 \leq h_{ii} \leq 1$
- ▶ Predictions

$$\hat{Y}_i = h_{ii} Y_i + \sum_{i \neq j} h_{ij} Y_j$$

- ▶ leverage values near 1 imply $\hat{Y}_i = Y_i$
 \hat{Y}_i Y_i

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix
 - ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$
 - ▶ symmetric
 - ▶ leverage values are the diagonal elements h_{ii} : $0 \leq h_{ii} \leq 1$
- ▶ Predictions

$$\hat{Y}_i = h_{ii}Y_i + \sum_{i \neq j} h_{ij}Y_j$$

- ▶ leverage values near 1 imply $\hat{Y}_i = Y_i$
- ▶ potentially influential

Hat Matrix

- ▶ predictions

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

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

- ▶ Hat Matrix or Projection Matrix

- ▶ idempotent $\mathbf{H}\mathbf{H} = \mathbf{H}$

- ▶ symmetric

- ▶ leverage values are the diagonal elements h_{ii} : $0 \leq h_{ii} \leq 1$

- ▶ Predictions

$$\hat{Y}_i = h_{ii}Y_i + \sum_{i \neq j} h_{ij}Y_j$$

- ▶ leverage values near 1 imply $\hat{Y}_i = Y_i$

- ▶ potentially influential

- ▶ Leverage: measure of how far x_i is from center of data

$$h_{ii} = 1/n + \frac{(\mathbf{x}_i - \bar{\mathbf{x}})^T \left((\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^T)^T (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^T) \right)^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})}{1}$$

$$\frac{1}{n} (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^T)^T (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^T) = \hat{\mathbf{Z}}^T \hat{\mathbf{Z}} = \mathbf{I}$$
$$\hat{Y}_i = h_{i1}Y_1 + h_{i2}Y_2$$

Residual Analysis

► residuals

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$$

$$\text{Var}(\mathbf{e}) = \sigma^2(\mathbf{I} - \mathbf{H})$$

$$\text{var}(e_i) = \sigma^2(1 - h_{ii})$$

Residual Analysis

- residuals

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$$

$$\text{Var}(\mathbf{e}) = \sigma^2(\mathbf{I} - \mathbf{H})$$

$$\text{var}(e_i) = \sigma^2(1 - h_{ii})$$

- Standardized residuals:

$$r_i = e_i / \sqrt{\widehat{\text{var}(e_i)}} = \underline{\underline{e_i / \{\hat{\sigma}^2(1 - h_{ii})\}}}$$

Residual Analysis

- ▶ residuals

$$\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$$

$$\text{Var}(\mathbf{e}) = \sigma^2(\mathbf{I} - \mathbf{H})$$

$$\text{var}(e_i) = \sigma^2(1 - h_{ii})$$

- ▶ Standardized residuals:

$$r_i = e_i / \sqrt{\widehat{\text{var}(e_i)}} = e_i / \{\hat{\sigma}^2(1 - h_{ii})\}$$

- ▶ if leverage is near 1 then residual is near 0 and variance is near 0 and r_i is approximately 0 (may not be helpful)

Predicted Residual

- Estimates without Case (i):

$$\hat{\beta}_{(i)} = (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)}$$

$$= \hat{\beta} - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i e_i}{1 - h_{ii}}$$

$$\hat{\beta}_{D^i}$$

$$\hat{\beta}_{(i)} = \beta - \dots$$

$$(x_1, y_1), \dots, (x_n, y_n) = D$$

$$D^i = D \ominus (x_i, y_i). \text{ Leave-one-out cross validation}$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}_{D^i}(x_i))^2$$

$$\hat{f}_{D^i} = \text{model fit with } D^i$$

*Standardized predicted residual is

$$\frac{e_{(i)}}{\sqrt{\text{var}(e_{(i)})}} = \frac{e_i / (1 - h_{ii})}{\hat{\sigma} / \sqrt{1 - h_{ii}}} = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

these are the same as standardized residual!

Predicted Residual

- Estimates without Case (i):

$$\hat{\beta}_{(i)} = (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)}$$

$$= \hat{\beta} - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i e_i}{1 - h_{ii}},$$

- Predicted residual

$$e_{(i)} = y_i - \mathbf{x}_i^T \hat{\beta}_{(i)} = \frac{e_i}{1 - h_{ii}}$$

$$\hat{\beta} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{\sum} \mathbf{X}^T \mathbf{Y}$$

\mathbf{X} is $n \times p$
 $(p \times n) (n \times 1)$

*Standardized predicted residual is

$$\frac{e_{(i)}}{\sqrt{\text{var}(e_{(i)})}} = \frac{e_i / (1 - h_{ii})}{\hat{\sigma} / \sqrt{1 - h_{ii}}} = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

these are the same as standardized residual!

p^2

$\sum x_i e_i$

$n p^3$

$n^2 p^2$

p^3

$\frac{e_{(i)}}{e_i}$
 \uparrow

Predicted Residual

- Estimates without Case (i):

$$\begin{aligned}\hat{\beta}_{(i)} &= (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{X}_{(i)}^T \mathbf{Y}_{(i)} \\ &= \hat{\beta} - \frac{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i e_i}{1 - h_{ii}}\end{aligned}$$

- Predicted residual

$$\underline{e_{(i)} = y_i - \mathbf{x}_i^T \hat{\beta}_{(i)} = \frac{e_i}{1 - h_{ii}}}$$

$h_{ii} \approx 1$

- variance

$$\underline{\text{var}(e_{(i)}) = \frac{\sigma^2}{1 - h_{ii}}}$$

*Standardized predicted residual is

$$\frac{e_{(i)}}{\sqrt{\text{var}(e_{(i)})}} = \frac{e_i / (1 - h_{ii})}{\hat{\sigma} / \sqrt{1 - h_{ii}}} = \frac{e_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

these are the same as standardized residual!

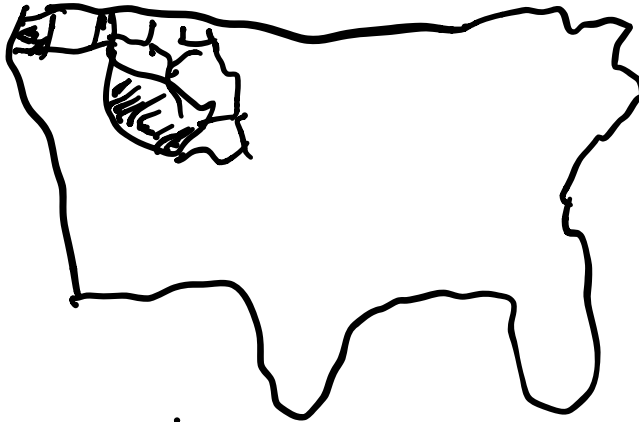
Standardized Residuals with External Estimate of σ

- ▶ Both the standardized residual and standardized predicted residual use all of the data in estimating σ

Standardized Residuals with External Estimate of σ

- ▶ Both the standardized residual and standardized predicted residual use all of the data in estimating σ
- ▶ if case i is an outlier, should also exclude it from estimating σ^2

rural
counties
have
lowest
cancer rates



that
rural
counties
have highest
cancer rate

Standardized Residuals with External Estimate of σ

- ▶ Both the standardized residual and standardized predicted residual use all of the data in estimating σ
- ▶ if case i is an outlier, should also exclude it from estimating σ^2
- ▶ Estimate $\hat{\sigma}_{(i)}^2$ using data with case i deleted

$$SSE_{(i)} = SSE - \frac{e_i^2}{1 - h_{ii}}$$

$$\hat{\sigma}_{(i)}^2 = MSE_{(i)} = \frac{SSE_{(i)}}{\underbrace{n - p - 1}}$$

Standardized Residuals with External Estimate of σ

- ▶ Both the standardized residual and standardized predicted residual use all of the data in estimating σ
- ▶ if case i is an outlier, should also exclude it from estimating σ^2
- ▶ Estimate $\hat{\sigma}_{(i)}^2$ using data with case i deleted

x is standard normal

y is χ^2

$z = \frac{x}{\sqrt{y}}$ is t -dist

$$SSE_{(i)} = SSE - \frac{e_i^2}{1 - h_{ii}}$$

$$\hat{\sigma}_{(i)}^2 = MSE_{(i)} = \frac{SSE_{(i)}}{n - p - 1}$$

t -statistic
 t -dist.

- ▶ Externally Standardized residuals

leave i-out
error

$$t_i = \frac{e_{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2 / (1 - h_{ii})}} = \frac{y_i - \mathbf{x}_i^T \hat{\beta}_{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2 / (1 - h_{ii})}} = r_i \left(\frac{n - p - 1}{n - p - r_i^2} \right)^{1/2}$$

Distribution of Externally Standardized Residuals

$$t_i = \frac{e_{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2/(1 - h_{ii})}} = \frac{y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2/(1 - h_{ii})}} \sim \text{St}(n - p - 1)$$

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus
 - ▶ $H_a: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta} + \alpha_i$ (different mean)

(

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus
 - ▶ $H_a: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta} + \alpha_i$ (different mean)
- ▶ Show that t-test for testing $H_0: \alpha_i = 0$ is equal to t_i

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus
 - ▶ $H_a: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta} + \alpha_i$ (different mean)
- ▶ Show that t-test for testing $H_0: \alpha_i = 0$ is equal to t_i
- ▶ if p-value is small declare the i th case to be an outlier: $E[Y_i]$ not given by $\mathbf{X}\boldsymbol{\beta}$ but $\mathbf{X}\boldsymbol{\beta} + \delta_i \alpha_i$

Outlier Test

- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus
 - ▶ $H_a: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta} + \alpha_i$ (different mean)
- ▶ Show that t-test for testing $H_0: \alpha_i = 0$ is equal to t_i
- ▶ if p-value is small declare the i th case to be an outlier: $E[Y_i]$ not given by $\mathbf{X}\boldsymbol{\beta}$ but $\mathbf{X}\boldsymbol{\beta} + \delta_i \alpha_i$
- ▶ Can extend to include multiple δ_i and δ_j to test that case i and j are both outliers

Outlier Test

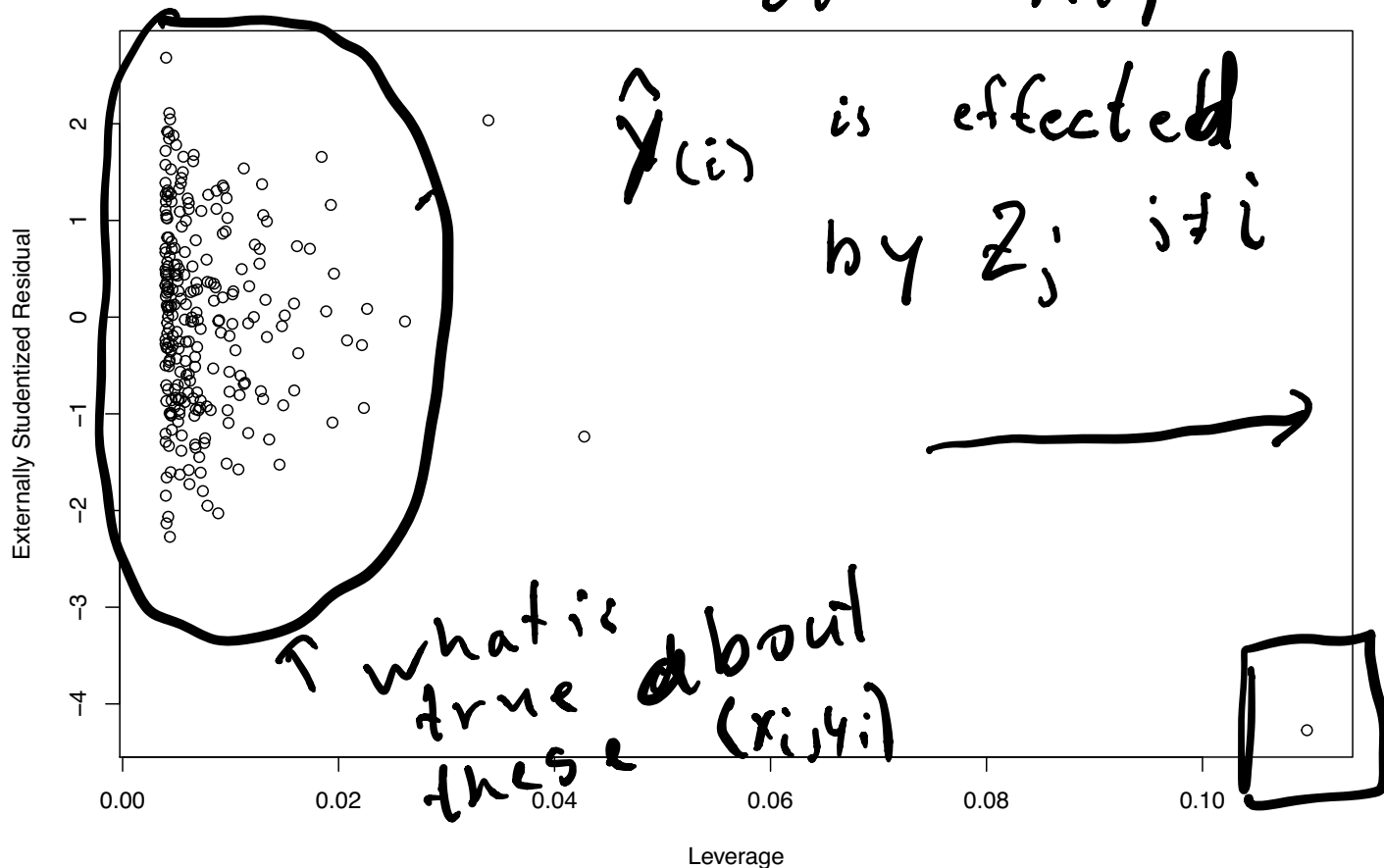
- ▶ Regression $E[Y_i] = \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- ▶ Hypotheses:
 - ▶ $H_0: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta}$ versus
 - ▶ $H_a: \mu_i = \mathbf{x}_i^T \boldsymbol{\beta} + \alpha_i$ (different mean)
- ▶ Show that t-test for testing $H_0: \alpha_i = 0$ is equal to t_i
- ▶ if p-value is small declare the i th case to be an outlier: $E[Y_i]$ not given by $\mathbf{X}\boldsymbol{\beta}$ but $\mathbf{X}\boldsymbol{\beta} + \delta_i \alpha_i$
- ▶ Can extend to include multiple δ_i and δ_j to test that case i and j are both outliers
- ▶ Extreme case $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta} + \mathbf{I}_n \boldsymbol{\alpha}$ all points have their own mean!

$$\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \boldsymbol{\mu} = \mathbf{X} \boldsymbol{\beta} + \mathbf{I}_n \boldsymbol{\alpha} \quad \text{dom.} \quad \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \mathbf{0}$$

R Code

```
plot(rstudent(bodyfat.lm) ~ hatvalues(bodyfat.lm),  
     ylab="Externally Studentized Residual",  
     xlab="Leverage")
```

$$z_i = (x_i, y_i)$$



P-Value

- P-value for test that observation with largest studentized residual is an outlier

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(max(abs.ti), bodyfat.lm$df - 1))
```

P-Value

- ▶ P-value for test that observation with largest studentized residual is an outlier

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(max(abs.ti), bodyfat.lm$df - 1))
```

- ▶ Issues with multiple comparisons if we compare each p-value to $\alpha = 0.05$

P-Value

how many
tests?

hypothesis
n-tests

- P-value for test that observation with largest studentized residual is an outlier

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(max(abs.ti), bodyfat.lm$df - 1))
```

- Issues with multiple comparisons if we compare each p-value to $\alpha = 0.05$
- Bonferroni compares p-values to α/n

Bonferonni Correction & Multiple Testing

- ▶ H_1, \dots, H_n are a family of hypotheses and p_1, \dots, p_n their corresponding p-values

Bonferonni Correction & Multiple Testing

- ▶ H_1, \dots, H_n are a family of hypotheses and p_1, \dots, p_n their corresponding p-values
- ▶ n_0 of the n are true

Bonferroni Correction & Multiple Testing

- ▶ H_1, \dots, H_n are a family of hypotheses and p_1, \dots, p_n their corresponding p-values
- ▶ n_0 of the n are true
- ▶ The **familywise error rate** (FWER) is the probability of rejecting at least one true H_i (making at least one type I error).

$$t(i) = 0.$$

how indep. is $t_{(1)}$ $t_{(2)}$ h_{11}

$$\text{FWER} = P \left\{ \bigcup_{i=1}^{n_0} \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq \sum_{i=1}^{n_0} \left\{ P \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq n_0 \frac{\alpha}{n} \leq n \frac{\alpha}{n}$$

$= \alpha$

$P \left(\bigcup_{i=1}^{n_0} \left(p_i \leq \frac{\alpha}{n} \right) \right) \leq \sum_{i=1}^{n_0} P \left(p_i \leq \frac{\alpha}{n} \right)$

equality if all n_{i1} are big $\leq n_0$ $\frac{\alpha}{n} \leq \frac{\alpha}{n}$

Bonferonni Correction & Multiple Testing

- ▶ H_1, \dots, H_n are a family of hypotheses and p_1, \dots, p_n their corresponding p-values
- ▶ n_0 of the n are true
- ▶ The **familywise error rate** (FWER) is the probability of rejecting at least one true H_i (making at least one type I error).

$$\begin{aligned} \text{FWER} &= P \left\{ \bigcup_{i=1}^{n_0} \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq \sum_{i=1}^{n_0} \left\{ P \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq n_0 \frac{\alpha}{n} \leq n \frac{\alpha}{n} \\ &= \alpha \end{aligned}$$

- ▶ This does not require any assumptions about dependence among the p-values or about how many of the null hypotheses are true.

Bonferonni Correction & Multiple Testing

- ▶ H_1, \dots, H_n are a family of hypotheses and p_1, \dots, p_n their corresponding p-values
- ▶ n_0 of the n are true
- ▶ The **familywise error rate** (FWER) is the probability of rejecting at least one true H_i (making at least one type I error).

$$\begin{aligned} \text{FWER} &= P \left\{ \bigcup_{i=1}^{n_0} \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq \sum_{i=1}^{n_0} \left\{ P \left(p_i \leq \frac{\alpha}{n} \right) \right\} \leq n_0 \frac{\alpha}{n} \leq n \frac{\alpha}{n} \\ &= \alpha \end{aligned}$$

- ▶ This does not require any assumptions about dependence among the p-values or about how many of the null hypotheses are true.
- ▶ Link https://en.wikipedia.org/wiki/Bonferroni_correction

Bonferroni Correction

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(abs.ti, bodyfat.lm$df - 1))  
min(pval) < .05/nrow(bodyfat)
```

```
## [1] TRUE
```

```
sum(pval < .05/nrow(bodyfat))
```

```
## [1] 1
```

- ▶ Bonferroni multiplicity adjustment compare each p-value to α/n and reject null (point is not an outlier) if the p-value is less than α/n

Bonferroni Correction

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(abs.ti, bodyfat.lm$df - 1))  
min(pval) < .05/nrow(bodyfat)
```

```
## [1] TRUE
```

```
sum(pval < .05/nrow(bodyfat))
```

```
## [1] 1
```

- ▶ Bonferroni multiplicity adjustment compare each p-value to α/n and reject null (point is not an outlier) if the p-value is less than α/n
- ▶ Start with max absolute value of t_i (or min p-value)

Bonferroni Correction

```
abs.ti = abs(rstudent(bodyfat.lm))  
pval= 2*(1- pt(abs.ti, bodyfat.lm$df - 1))  
min(pval) < .05/nrow(bodyfat)
```

```
## [1] TRUE
```

```
sum(pval < .05/nrow(bodyfat))
```

```
## [1] 1
```

- ▶ Bonferroni multiplicity adjustment compare each p-value to α/n and reject null (point is not an outlier) if the p-value is less than α/n
- ▶ Start with max absolute value of t_i (or min p-value)
- ▶ Case 39 would be considered an outlier based on Bonferroni or other multiplicity adjustments. no other outliers

Cook's Distance

$$\hat{\mathbf{y}}_{(i)} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_n \end{pmatrix} \quad \mathbf{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

- Measure of influence of case i on predictions

$$D_i = \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}_{(i)}\|^2}{\hat{\sigma}^2 p}$$

after removing the i th case

$$\|\mathbf{Y} - \hat{\mathbf{Y}}_{(i)}\|^2$$

Cook's Distance

- ▶ Measure of influence of case i on predictions

$\hat{\sigma}_{(i)}^2$ \gg $\hat{\sigma}^2$

after removing the i th case

$$D_i = \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}_{(i)}\|^2}{\hat{\sigma}^2 p}$$

why
scale
by

$e_i^2, \hat{\sigma}_{(i)}^2, \hat{\sigma}^2, p$

- ▶ Easier way to calculate

$$D_i = \frac{e_i^2}{\hat{\sigma}^2 p} \left[\frac{h_{ii}}{(1 - h_{ii})^2} \right], \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}_{(i)}\|^2}{\sum_{j=1}^n (y_j - \hat{y}_j^{(i)})^2}$$

$$\hat{\sigma}^2 = \text{Var}(e_i)$$

$$= \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \hat{\sigma}_{(i)}^2$$

Model Assessment

- ▶ Always look at residual plots!

Model Assessment

- ▶ Always look at residual plots!
- ▶ Check constant variance, outliers, influence, normality assumption

Model Assessment

- ▶ Always look at residual plots!
- ▶ Check constant variance, outliers, influence, normality assumption
- ▶ Treat e_i as “new data” - look at structure, other predictors
avplots

Model Assessment

- ▶ Always look at residual plots!
- ▶ Check constant variance, outliers, influence, normality assumption
- ▶ Treat e_i as “new data” - look at structure, other predictors
avplots
- ▶ Case 39 looks an influential outlier!

Model Assessment

- ▶ Always look at residual plots!
- ▶ Check constant variance, outliers, influence, normality assumption
- ▶ Treat e_i as “new data” - look at structure, other predictors
avplots
- ▶ Case 39 looks an influential outlier!
- ▶ Impact on predictions?

Predictions with Case 39

```
predict(bodyfat.lm, newdata=bodyfat[39,],  
        se=T, interval="prediction")
```

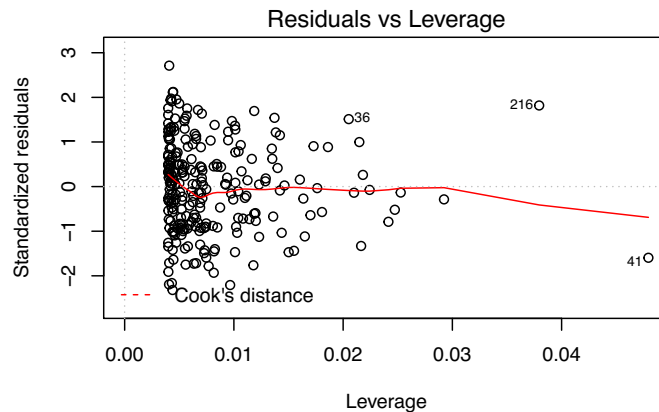
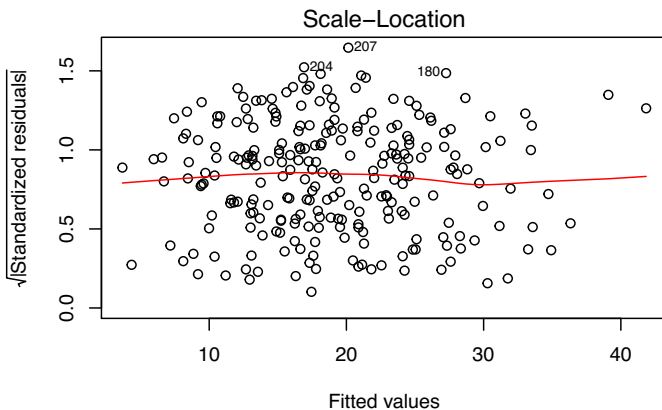
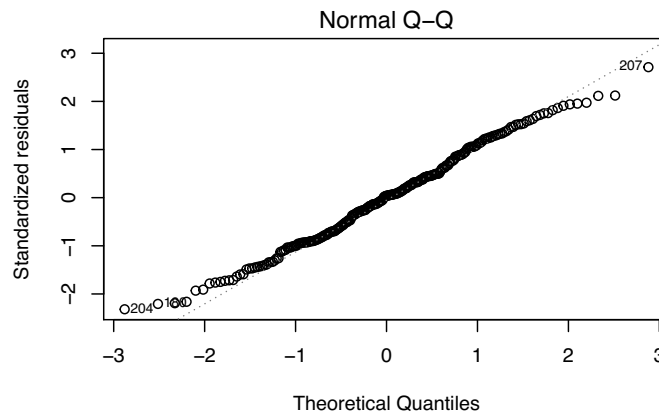
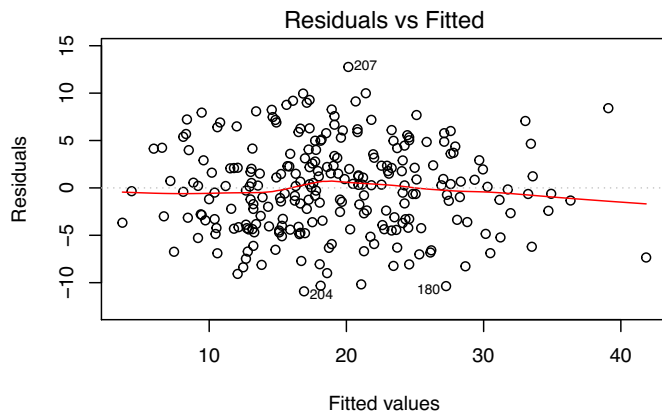
```
## $fit  
##           fit      lwr      upr  
## 39 54.21599 44.0967 64.33528  
##  
## $se.fit  
## [1] 1.615311  
##  
## $df  
## [1] 250  
##  
## $residual.scale  
## [1] 4.877484
```

Predictions without Case 39

```
bodyfatsub.lm = lm(Bodyfat ~ Abdomen, data=bodyfat,  
                   subset=c(-39))  
predict(bodyfatsub.lm, newdata=bodyfat[39,],  
        se=T, interval="prediction")
```

```
## $fit  
##           fit           lwr           upr  
## 39 56.55856 46.71172 66.40541  
##  
## $se.fit  
## [1] 1.655744  
##  
## $df  
## [1] 249  
##  
## $residual.scale  
## [1] 4.717441
```

Residual Checks without Case 39



How should we proceed?

- ▶ Reproducible Research - Document removing a case

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc
 - ▶ Influential points (not outliers): report analysis with & without

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc
 - ▶ Influential points (not outliers): report analysis with & without
- ▶ If we remove Case 39, are there other outliers or influential points?

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc
 - ▶ Influential points (not outliers): report analysis with & without
- ▶ If we remove Case 39, are there other outliers or influential points?
- ▶ Model Uncertainty (more later)

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc
 - ▶ Influential points (not outliers): report analysis with & without
- ▶ If we remove Case 39, are there other outliers or influential points?
- ▶ Model Uncertainty (more later)
- ▶ Robust Models (more later)

How should we proceed?

- ▶ Reproducible Research - Document removing a case
 - ▶ Adjust for multiple testing!
 - ▶ Remove statistically significant outliers if you cannot confirm other data entry errors, etc
 - ▶ Influential points (not outliers): report analysis with & without
- ▶ If we remove Case 39, are there other outliers or influential points?
- ▶ Model Uncertainty (more later)
- ▶ Robust Models (more later)
- ▶ Next: Transformations