# Lecture 9
## Multiple Linear Regression V
Reading: Chapter 13

*STAT 8020 Statistical Methods II*
September 17, 2020

Whitney Huang
Clemson University

# Agenda

**1** **Model Diagnostics: Influential Points**

**2** **Non-Constant Variance & Transformation**

**3** **Regression with Both Quantitative and Qualitative Predictors**

**4** **Polynomial Regression**

# Leverage

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors
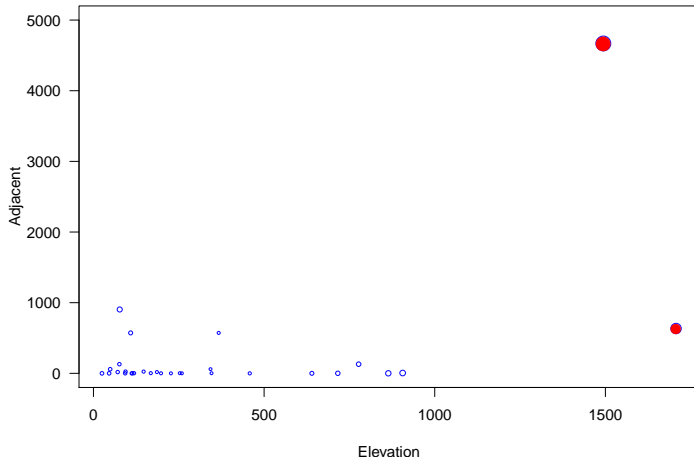
Polynomial Regression

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

- The leverage value for the $i_{\text{th}}$ observation is defined as:

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

- 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 $i_{\text{th}}$ observation

- $\frac{1}{n} \leq h_i \leq 1, \quad 1 \leq i \leq n$ and $\bar{h} = \sum_{i=1}^{n} \frac{h_i}{n} = \frac{p}{n} \Rightarrow$ a "rule of thumb" is that leverages of more than $\frac{2p}{n}$ should be looked at more closely

# Leverage Values of `Species ~ Elev + Adj`

# Studentized Residuals

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 `Species` $\sim$ `Elev` + `Adj`

**Studentized Residuals**

# Studentized Deleted Residuals

CLEMSON
U N I V E R S I T Y

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
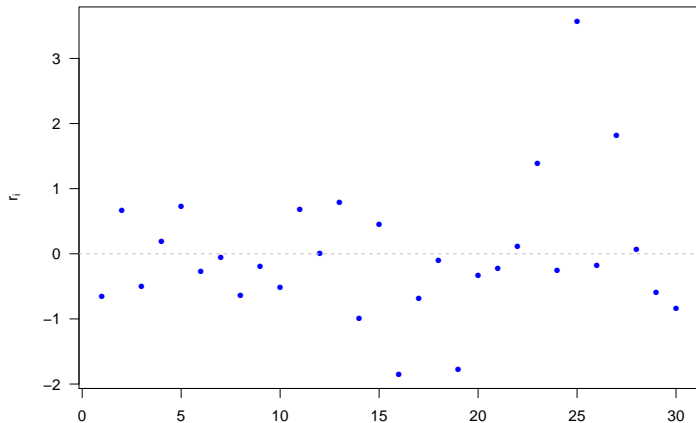Qualitative Predictors

Polynomial Regression

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

- The observation $i$ is an outlier if $\hat{Y}_{i(i)} - Y_i$ is "large"

- Can show
$$\mathsf{Var}(\hat{Y}_{i(i)} - Y_i) = \sigma_{(i)}^2 \left( 1 + \boldsymbol{x}_i^T (\boldsymbol{X}_{(i)}^T \boldsymbol{X}_{(i)})^{-1} \boldsymbol{x}_i \right) = \frac{\sigma_{(i)}^2}{1 - h_i}$$
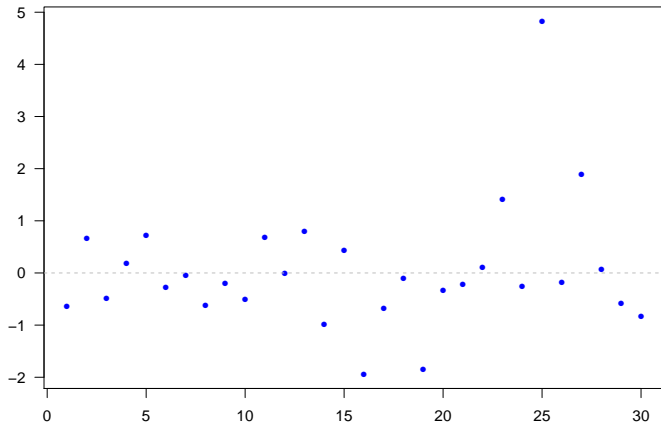
- Define the **Studentized Deleted Residuals** as

$$t_i = \frac{\hat{Y}_{i(i)} - Y_i}{\hat{\sigma}_{(i)}^2 / 1 - h_i} = \frac{\hat{Y}_{i(i)} - Y_i}{\mathsf{MSE}_{(i)}(1 - h_i)^{-1}}$$

which are distributed as a $t_{n-p-1}$ if the model is correct and $\varepsilon \sim \mathrm{N}(\boldsymbol{0}, \sigma^2 \boldsymbol{I})$

# Jackknife Residuals of `Species ~ Elev + Adj`

Multiple Linear Regression V

CLEMS🐾N
U N I V E R S I T Y

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

Jacknife Residuals

# Influential Observations

**Multiple Linear Regression V**

**CLEMSON**
UNIVERSITY

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

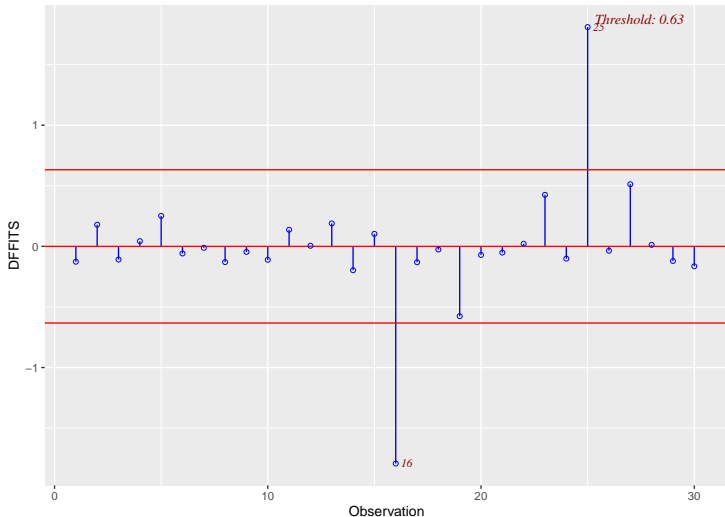Polynomial Regression

## DFFITS

- Difference between the fitted values $\hat{Y}_i$ and the predicted values $\hat{Y}_{i(i)}$

- $\text{DFFITS}_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{\text{MSE}_{(i)} h_i}}$

- 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

Model Diagnostics:
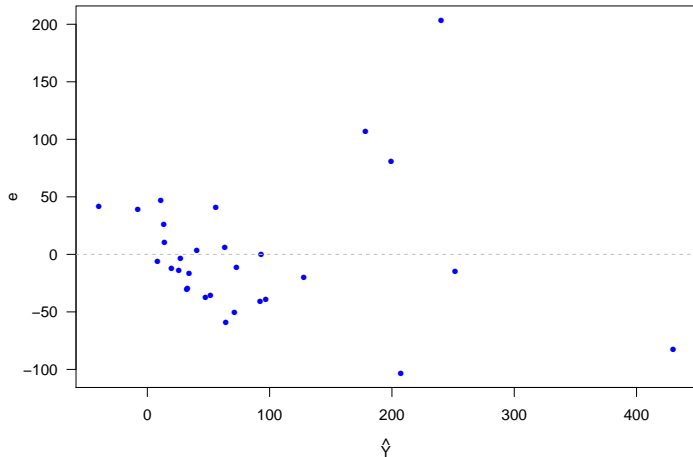Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

# Residual Plot of `Species ~ Elev + Adj`

Model Diagnostics:
Influential Points

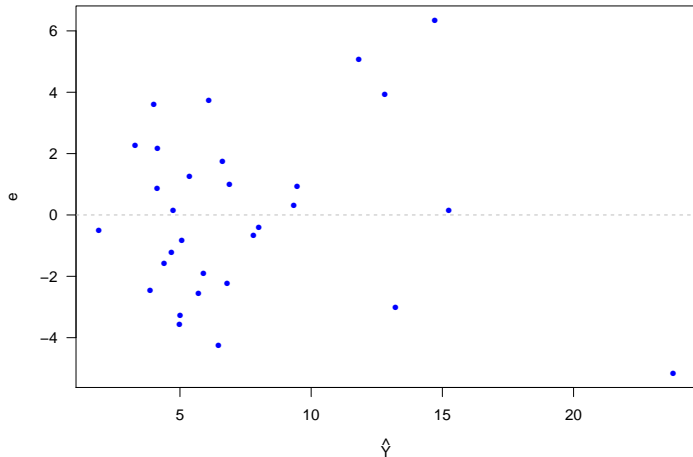Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

**Residuals**

# Residual Plot After Square Root Transformation

Multiple Linear
Regression V

CLEMS😺N
U N I V E R S I T Y

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

Residuals

**Regression with Both Quantitative and Qualitative Predictors**

Multiple Linear Regression V

**CLEMSON**
U N I V E R S I T Y

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

**Multiple Linear Regression**

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_{p-1} X_{p-1} + \varepsilon, \quad \varepsilon \sim \mathrm{N}(0, \sigma^2)$$

$X_1, X_2, \cdots, X_{p-1}$ are the predictors.

**Question**: What if some of the predictors are qualitative (categorical) variables?

$\Rightarrow$ We will need to create **dummy (indicator) variables** for those categorical variables

**Example**: We can encode `Gender` into $1$ (Female) and $0$ (Male)

# Salaries for Professors Data Set

The 2008-09 nine-month academic salary for Assistant Professors, Associate Professors and Professors in a college in the U.S. The data were collected as part of the on-going effort of the college's administration to monitor salary differences between male and female faculty members.

```
> head(Salaries)
      rank discipline yrs.since.phd yrs.service  sex salary
1     Prof          B            19          18 Male 139750
2     Prof          B            20          16 Male 173200
3  AsstProf          B             4           3 Male  79750
4     Prof          B            45          39 Male 115000
5     Prof          B            40          41 Male 141500
6 AssocProf          B             6           6 Male  97000
```

# Predictors

```
> summary(Salaries)
      rank        discipline yrs.since.phd   yrs.service
 AsstProf : 67   A:181      Min.   : 1.00   Min.   : 0.00
 AssocProf: 64   B:216      1st Qu.:12.00   1st Qu.: 7.00
 Prof     :266              Median :21.00   Median :16.00
                            Mean   :22.31   Mean   :17.61
                            3rd Qu.:32.00   3rd Qu.:27.00
                            Max.   :56.00   Max.   :60.00

     sex           salary
 Female: 39   Min.   : 57800
 Male  :358   1st Qu.: 91000
              Median :107300
              Mean   :113706
              3rd Qu.:134185
              Max.   :231545
```

We have three categorical variables, namely, `rank`, `discipline`, and `sex`.

## Dummy Variable

For binary categorical variables:

$$X_{\text{sex}} = \begin{cases} 0 & \text{if } \texttt{sex} = \text{male}, \\ 1 & \text{if } \texttt{sex} = \text{female}. \end{cases}$$

$$X_{\text{discip}} = \begin{cases} 0 & \text{if } \texttt{discip} = \text{A}, \\ 1 & \text{if } \texttt{discip} = \text{B}. \end{cases}$$

For categorical variable with more than two categories:

$$X_{\text{rank1}} = \begin{cases} 0 & \text{if } \texttt{rank} = \text{Assistant Prof}, \\ 1 & \text{if } \texttt{rank} = \text{Associated Prof}. \end{cases}$$

$$X_{\text{rank2}} = \begin{cases} 0 & \text{if } \texttt{rank} = \text{Associated Prof}, \\ 1 & \text{if } \texttt{rank} = \text{Full Prof}. \end{cases}$$

# Design Matrix

```
> head(X)
  (Intercept) rankAssocProf rankProf disciplineB yrs.since.phd
1           1             0        1           1            19
2           1             0        1           1            20
3           1             0        0           1             4
4           1             0        1           1            45
5           1             0        1           1            40
6           1             1        0           1             6
  yrs.service sexMale
1          18       1
2          16       1
3           3       1
4          39       1
5          41       1
6           6       1
```

With the design matrix $X$, we can now use method of least squares to fit the model $Y = X\beta + \varepsilon$
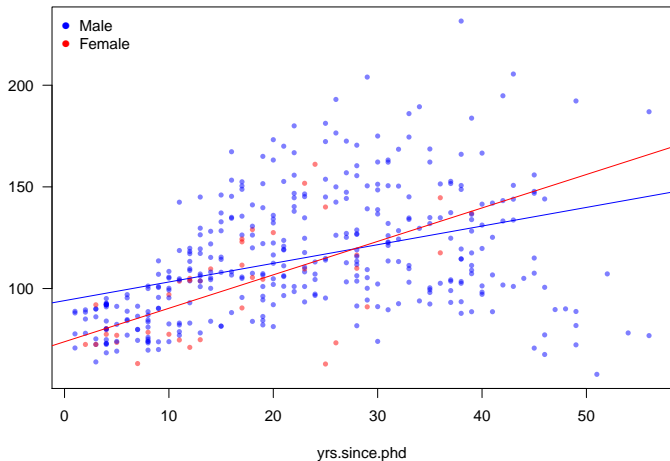
# Model Fit

Model Diagnostics: Influential Points

Non-Constant Variance & Transformation

**Regression with Both Quantitative and Qualitative Predictors**

Polynomial Regression

```
Coefficients:
               Estimate Std. Error t value Pr(>|t|)
(Intercept)     70738.7     3403.0  20.787  < 2e-16 ***
rankAssocProf   12907.6     4145.3   3.114  0.00198 **
rankProf        45066.0     4237.5  10.635  < 2e-16 ***
disciplineB     14417.6     2342.9   6.154 1.88e-09 ***
yrs.since.phd     535.1      241.0   2.220  0.02698 *
yrs.service      -489.5      211.9  -2.310  0.02143 *
sexFemale       -4783.5     3858.7  -1.240  0.21584
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 22540 on 390 degrees of freedom
Multiple R-squared:  0.4547,    Adjusted R-squared:  0.4463
F-statistic:  54.2 on 6 and 390 DF,  p-value: < 2.2e-16
```

**Question**: Interpretation of these dummy variables (e.g. $\hat{\beta}_{\texttt{rankAssocProf}}$)?

# $\texttt{lm(salary} \sim \texttt{sex} * \texttt{yrs.since.phd)}$

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

**Regression with Both
Quantitative and
Qualitative Predictors**
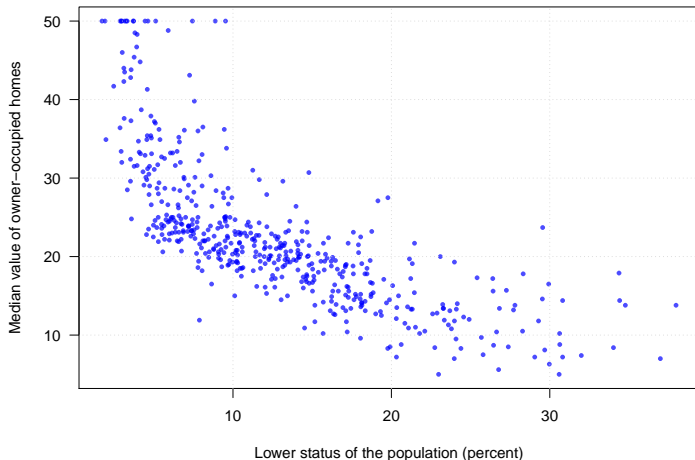
Polynomial Regression

# Polynomial Regression

Suppose we would like to model the relationship between response $Y$ and a predictor $X$ as a $p_{\text{th}}$ degree polynomial in $X$:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \cdots + \beta_p X_i^p + \varepsilon$$

We can treat polynomial regression as a special case of multiple linear regression. In specific, the design matrix takes the following form:

$$\boldsymbol{X} = \begin{pmatrix} 1 & X_1 & X_1^2 & \cdots & X_1^p \\ 1 & X_2 & X_2^2 & \cdots & X_2^p \\ \vdots & \cdots & \ddots & \vdots & \vdots \\ 1 & X_n & X_n^2 & \cdots & X_n^p \end{pmatrix}$$

# Housing Values in Suburbs of Boston Data Set
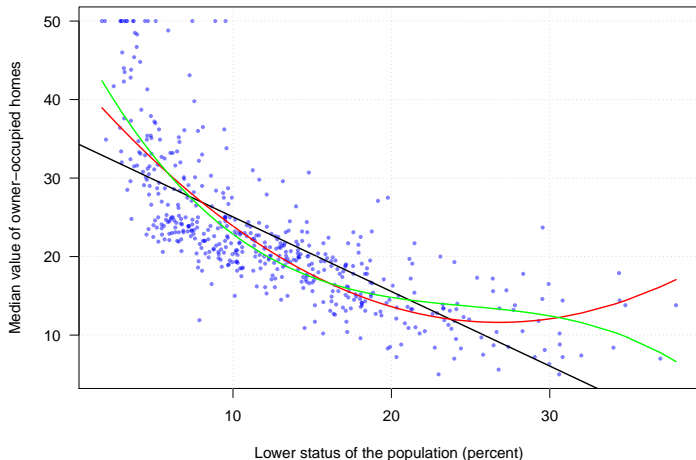
CLEMS❖N
U N I V E R S I T Y

Model Diagnostics:
Influential Points

Non-Constant
Variance &
Transformation

Regression with Both
Quantitative and
Qualitative Predictors

Polynomial Regression

# Polynomial Regression Fits