

Lecture 4

Multiple Linear Regression II

Reading: Forecasting, Time Series, and Regression (4th edition) by Bowerman, O'Connell, and Koehler: Chapter 4

MATH 4070: Regression and Time-Series Analysis

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Whitney Huang
Clemson University

- 1 General Linear F -Test
- 2 Prediction
- 3 Multicollinearity
- 4 Model Selection
- 5 Model Diagnostics
- 6 Non-Constant Variance & Transformation

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- t -test: Testing one predictor

- 1 **Null/Alternative Hypotheses:** $H_0 : \beta_j = 0$ vs. $H_a : \beta_j \neq 0$

- 2 **Test Statistic:** $t^* = \frac{\hat{\beta}_j - 0}{\text{se}(\hat{\beta}_j)}$

- 3 Reject H_0 if $|t^*| > t_{1-\alpha/2, n-p}$

- Overall F -test: Test of all the predictors

- 1 $H_0 : \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$

- 2 $H_a : \text{at least one } \beta_j \neq 0, 1 \leq j \leq p-1$

- 3 **Test Statistic:** $F^* = \frac{\text{MSR}}{\text{MSE}}$

- 4 Reject H_0 if $F^* > F_{1-\alpha, p-1, n-p}$

Both tests are special cases of **General Linear F -test**

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- Comparison of a “full model” and “reduced model” that involves a **subset of full model predictors**
- Consider a full model with k predictors and reduced model with ℓ predictors ($\ell < k$)
- Test statistic: $F^* = \frac{(SSE_{\text{reduce}} - SSE_{\text{full}})/(k - \ell)}{SSE_{\text{full}}/(n - k - 1)} \Rightarrow$ Testing H_0 that the regression coefficients for the extra variables are all zero
 - Example 1: x_1, x_2, \dots, x_{p-1} vs. intercept only \Rightarrow Overall F -test

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- Comparison of a “full model” and “reduced model” that involves **a subset of full model predictors**
- Consider a full model with k predictors and reduced model with ℓ predictors ($\ell < k$)
- Test statistic: $F^* = \frac{(SSE_{\text{reduce}} - SSE_{\text{full}})/(k - \ell)}{SSE_{\text{full}}/(n - k - 1)} \Rightarrow$ Testing H_0 that the regression coefficients for the extra variables are all zero
 - Example 1: x_1, x_2, \dots, x_{p-1} vs. intercept only \Rightarrow Overall F -test
 - Example 2: $x_j, 1 \leq j \leq p - 1$ vs. intercept only $\Rightarrow t$ -test for β_j

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- Comparison of a “full model” and “reduced model” that involves a **subset of full model predictors**
- Consider a full model with k predictors and reduced model with ℓ predictors ($\ell < k$)
- Test statistic: $F^* = \frac{(SSE_{\text{reduce}} - SSE_{\text{full}})/(k - \ell)}{SSE_{\text{full}}/(n - k - 1)} \Rightarrow$ Testing H_0 that the regression coefficients for the extra variables are all zero
 - Example 1: x_1, x_2, \dots, x_{p-1} vs. intercept only \Rightarrow Overall F -test
 - Example 2: $x_j, 1 \leq j \leq p - 1$ vs. intercept only $\Rightarrow t$ -test for β_j
 - Example 3: x_1, x_2, x_3, x_4 vs. $x_1, x_3 \Rightarrow H_0 : \beta_2 = \beta_4 = 0$

General Linear F -Test

Prediction

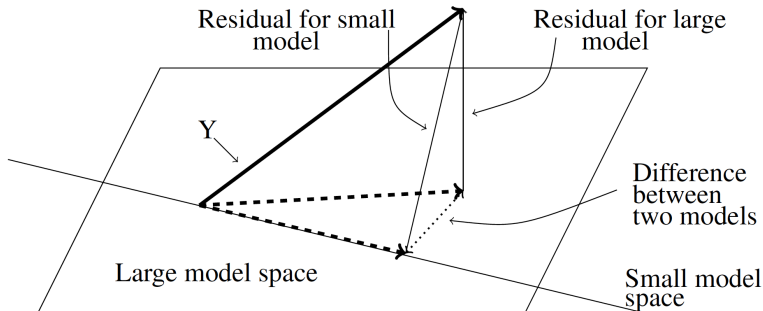
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Geometric Illustration of General Linear F -Test



Source: Faraway, *Linear Models with R*, 2014, p.34

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Species Diversity on the Galapagos Islands: Full Model

```
> summary(gala_fit2)
```

Call:

```
lm(formula = Species ~ Elevation + Area)
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|---------|---------|-------|---------|
| -192.619 | -33.534 | -19.199 | 7.541 | 261.514 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|------------|
| (Intercept) | 17.10519 | 20.94211 | 0.817 | 0.42120 |
| Elevation | 0.17174 | 0.05317 | 3.230 | 0.00325 ** |
| Area | 0.01880 | 0.02594 | 0.725 | 0.47478 |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 79.34 on 27 degrees of freedom

Multiple R-squared: 0.554, Adjusted R-squared: 0.521

F-statistic: 16.77 on 2 and 27 DF, p-value: 1.843e-05

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Species Diversity on the Galapagos Islands: Reduce Model

```
> summary(gala_fit1)
```

Call:

```
lm(formula = Species ~ Elevation)
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|---------|---------|-------|---------|
| -218.319 | -30.721 | -14.690 | 4.634 | 259.180 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 11.33511 | 19.20529 | 0.590 | 0.56 |
| Elevation | 0.20079 | 0.03465 | 5.795 | 3.18e-06 *** |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 78.66 on 28 degrees of freedom

Multiple R-squared: 0.5454, Adjusted R-squared: 0.5291

F-statistic: 33.59 on 1 and 28 DF, p-value: 3.177e-06

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Performing a General Linear F -Test

● $H_0 : \beta_{\text{Area}} = 0$ vs. $H_a : \beta_{\text{Area}} \neq 0$

● $F^* = \frac{(173254 - 169947)/(2-1)}{169947/(30-2-1)} = 0.5254$

● P-value: $P[F > 0.5254] = 0.4748$, where $F \sim F_{\underbrace{1}_{k-\ell}, \underbrace{27}_{n-k-1}}$

```
> anova(gala_fit1, gala_fit2)
```

Analysis of Variance Table

Model 1: Species ~ Elevation

Model 2: Species ~ Elevation + Area

| | Res.Df | RSS | Df | Sum of Sq | F | Pr(>F) |
|---|--------|--------|----|-----------|--------|--------|
| 1 | 28 | 173254 | | | | |
| 2 | 27 | 169947 | 1 | 3307 | 0.5254 | 0.4748 |

General Linear F -Test

Prediction

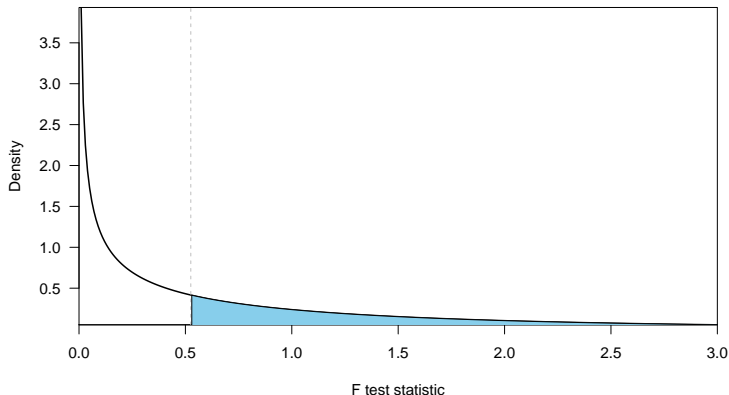
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Visualizing p -value



p -value is the shaded area under the density curve of the null distribution

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Another Example of General Linear F -Test

```
> full <- lm(Species ~ Area + Elevation + Nearest + Scrutz + Adjacent,  
  data = gala)  
> anova(full)
```

Analysis of Variance Table

Response: Species

| | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|-----------|----|--------|---------|---------|---------------|
| Area | 1 | 145470 | 145470 | 39.1262 | 1.826e-06 *** |
| Elevation | 1 | 65664 | 65664 | 17.6613 | 0.0003155 *** |
| Nearest | 1 | 29 | 29 | 0.0079 | 0.9300674 |
| Scrutz | 1 | 14280 | 14280 | 3.8408 | 0.0617324 . |
| Adjacent | 1 | 66406 | 66406 | 17.8609 | 0.0002971 *** |
| Residuals | 24 | 89231 | 3718 | | |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```
> reduced <- lm(Species ~ Elevation + Adjacent)  
> anova(reduced)
```

Analysis of Variance Table

Response: Species

| | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|-----------|----|--------|---------|---------|---------------|
| Elevation | 1 | 207828 | 207828 | 56.112 | 4.662e-08 *** |
| Adjacent | 1 | 73251 | 73251 | 19.777 | 0.0001344 *** |
| Residuals | 27 | 100003 | 3704 | | |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- Null and alternative hypotheses:

$$H_0 : \beta_{\text{Area}} = \beta_{\text{Nearest}} = \beta_{\text{Scruz}} = 0$$

$$H_a : \text{at least one of the three coefficients} \neq 0$$

- $F^* = \frac{(100003 - 89231)/(5-2)}{89231/(30-5-1)} = 0.9657$

- $p\text{-value: } P[F > 0.9657] = 0.425, \text{ where } F \sim F_{3,24}$

```
> anova(reduced, full)
```

Analysis of Variance Table

Model 1: Species ~ Elevation + Adjacent

Model 2: Species ~ Area + Elevation + Nearest + Scruz + Adjacent

| | Res.Df | RSS | Df | Sum of Sq | F | Pr(>F) |
|---|--------|--------|----|-----------|--------|--------|
| 1 | 27 | 100003 | | | | |
| 2 | 24 | 89231 | 3 | 10772 | 0.9657 | 0.425 |

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Given a new set of predictors, $\mathbf{x}_0 = (1, x_{0,1}, x_{0,2}, \dots, x_{0,p-1})^T$, the predicted response is

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{0,1} + \hat{\beta}_2 x_{0,2} + \dots + \hat{\beta}_{p-1} x_{0,p-1}.$$

Again, we can use matrix representation to simplify the notation

$$\hat{y}_0 = \mathbf{x}_0^T \hat{\boldsymbol{\beta}},$$

where $\mathbf{x}_0^T = (1, x_{0,1}, x_{0,2}, \dots, x_{0,p-1})$

We will use this formula to carry out two different kinds of predictions

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

There are two kinds of predictions can be made for a given x_0 :

- **Predicting a future response:**

Based on MLR, we have $y_0 = x_0^T \beta + \varepsilon$. Since $E(\varepsilon) = 0$, therefore the predicted value is

$$\hat{y}_0 = x_0^T \hat{\beta}$$

- **Predicting the mean response:**

Since $E(y_0) = x_0^T \beta$, there we have the predicted mean response

$$\widehat{E(y_0)} = x_0^T \hat{\beta},$$

the same predicted value as predicting a future response

Next, we need to assess their **prediction uncertainties**, and then we will identify the differences in terms of these uncertainties

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

From page 30 of slides 3, we have $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$.
Therefore we have

$$\text{Var}(\hat{y}_0) = \text{Var}(\mathbf{x}_0^T \hat{\beta}) = \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0$$

We can now construct $100(1 - \alpha)\%$ CI for the two kinds of predictions:

- **Predicting a future response y_0 :**

$$\mathbf{x}_0^T \hat{\beta} \pm t_{1-\alpha/2, n-p} \times \hat{\sigma} \sqrt{\underbrace{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}_{\text{accounting for } \varepsilon}}$$

- **Predicting the mean response $E(y_0)$:**

$$\mathbf{x}_0^T \hat{\beta} \pm t_{1-\alpha/2, n-p} \times \hat{\sigma} \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}$$

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Example: Predicting Body Fat (Faraway 2014 Chapter 4.2)

```
lm(formula = brozek ~ age + weight + height + neck + chest +  
    abdom + hip + thigh + knee + ankle + biceps + forearm + wrist,  
    data = fat)
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|---------|--------|--------|-------|-------|
| -10.264 | -2.572 | -0.097 | 2.898 | 9.327 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|-----------|------------|---------|-------------|
| (Intercept) | -15.29255 | 16.06992 | -0.952 | 0.34225 |
| age | 0.05679 | 0.02996 | 1.895 | 0.05929 . |
| weight | -0.08031 | 0.04958 | -1.620 | 0.10660 |
| height | -0.06460 | 0.08893 | -0.726 | 0.46830 |
| neck | -0.43754 | 0.21533 | -2.032 | 0.04327 * |
| chest | -0.02360 | 0.09184 | -0.257 | 0.79740 |
| abdom | 0.88543 | 0.08008 | 11.057 | < 2e-16 *** |
| hip | -0.19842 | 0.13516 | -1.468 | 0.14341 |
| thigh | 0.23190 | 0.13372 | 1.734 | 0.08418 . |
| knee | -0.01168 | 0.22414 | -0.052 | 0.95850 |
| ankle | 0.16354 | 0.20514 | 0.797 | 0.42614 |
| biceps | 0.15280 | 0.15851 | 0.964 | 0.33605 |
| forearm | 0.43049 | 0.18445 | 2.334 | 0.02044 * |
| wrist | -1.47654 | 0.49552 | -2.980 | 0.00318 ** |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.988 on 238 degrees of freedom

Multiple R-squared: 0.749, Adjusted R-squared: 0.7353

F-statistic: 54.63 on 13 and 238 DF, p-value: < 2.2e-16

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

What is our prediction for the future response of a “typical” (e.g., each predictor takes its median value) man?

Example: Predicting Body Fat Cont'd

- 1 Calculate the median for each predictor to get x_0
- 2 Compute the predicted value $\hat{y}_0 = x_0^T \hat{\beta}$
- 3 Quantify the prediction uncertainty

```
> X <- model.matrix(lmod)
> (x0 <- apply(x, 2, median))
(Intercept)      age      weight      height      neck      chest      abdom
      1.00      43.00      176.50      70.00      38.00      99.65      90.95
      hip      thigh      knee      ankle      biceps      forearm      wrist
      99.30      59.00      38.50      22.80      32.05      28.70      18.30
> (y0 <- sum(x0 * coef(lmod)))
[1] 17.49322
> predict(lmod, new = data.frame(t(x0)))
      1
17.49322
> predict(lmod, new = data.frame(t(x0)), interval = "prediction")
      fit      lwr      upr
1 17.49322  9.61783 25.36861
> predict(lmod, new = data.frame(t(x0)), interval = "confidence")
      fit      lwr      upr
1 17.49322 16.94426 18.04219
```

General Linear F -Test

Prediction

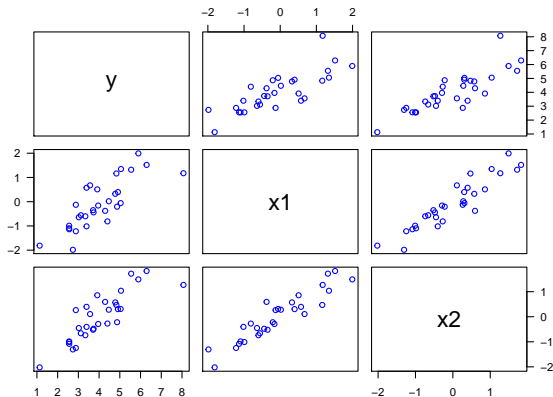
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Multicollinearity



```
> cor(sim1)
```

| | y | x1 | x2 |
|----|-----------|-----------|-----------|
| y | 1.0000000 | 0.7987777 | 0.8481084 |
| x1 | 0.7987777 | 1.0000000 | 0.9281514 |
| x2 | 0.8481084 | 0.9281514 | 1.0000000 |

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Multicollinearity is a phenomenon of high inter-correlations among the predictor variables

- Numerical issue \Rightarrow the matrix $\mathbf{X}^T \mathbf{X}$ is nearly singular
- Statistical issues/consequences
 - β 's are not well estimated \Rightarrow spurious regression coefficient estimates
 - R^2 and predicted values are usually okay even with multicollinearity

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Suppose the true relationship between response y and predictors (x_1, x_2) is

$$Y = 4 + 0.8x_1 + 0.6x_2 + \varepsilon,$$

where $\varepsilon \sim N(0, 1)$ and x_1 and x_2 are positively correlated with $\rho = 0.9$. Let's fit the following models:

- Model 1: $Y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \varepsilon_1$
This is the true model with parameters unknown
- Model 2: $Y = \beta_0 + \beta_1x_1 + \varepsilon_2$
This is the wrong model because x_2 is omitted

General Linear F -Test

Prediction

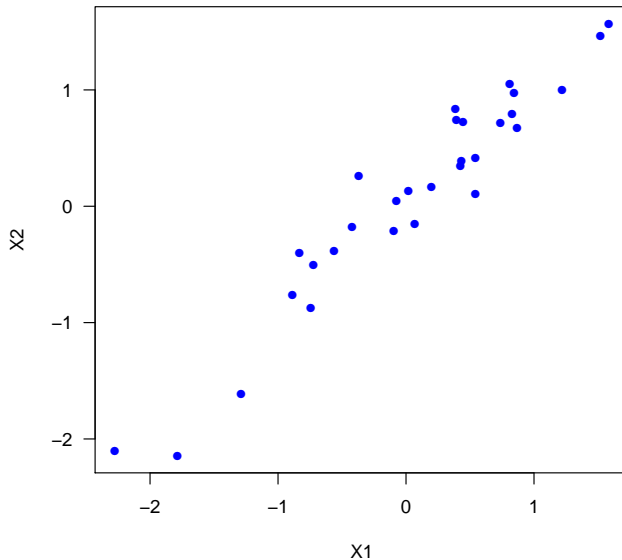
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Scatter Plot: x_1 vs. x_2



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Call:

```
lm(formula = Y ~ X1 + X2)
```

Residuals:

| | Min | 1Q | Median | 3Q | Max |
|--|----------|----------|---------|---------|---------|
| | -1.91369 | -0.73658 | 0.05475 | 0.87080 | 1.55150 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 4.0710 | 0.1778 | 22.898 | < 2e-16 *** |
| X1 | 2.2429 | 0.7187 | 3.121 | 0.00426 ** |
| X2 | -0.8339 | 0.7093 | -1.176 | 0.24997 |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9569 on 27 degrees of freedom

Multiple R-squared: 0.673, Adjusted R-squared: 0.6488

F-statistic: 27.78 on 2 and 27 DF, p-value: 2.798e-07

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Call:

```
lm(formula = Y ~ X1)
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|----------|----------|----------|---------|---------|
| -2.09663 | -0.67031 | -0.07229 | 0.87881 | 1.49739 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 4.0347 | 0.1763 | 22.888 | < 2e-16 *** |
| X1 | 1.4293 | 0.1955 | 7.311 | 5.84e-08 *** |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9634 on 28 degrees of freedom

Multiple R-squared: 0.6562, Adjusted R-squared: 0.644

F-statistic: 53.45 on 1 and 28 DF, p-value: 5.839e-08

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Takeaways

Model 1 fit:

```
Call:
lm(formula = Y ~ X1 + X2)
```

Residuals:

| | Min | 1Q | Median | 3Q | Max |
|--|----------|----------|---------|---------|---------|
| | -1.91369 | -0.73658 | 0.05475 | 0.87080 | 1.55150 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 4.0710 | 0.1778 | 22.898 | < 2e-16 *** |
| X1 | 2.2429 | 0.7187 | 3.121 | 0.00426 ** |
| X2 | -0.8339 | 0.7093 | -1.176 | 0.24997 |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9569 on 27 degrees of freedom
Multiple R-squared: 0.673, Adjusted R-squared: 0.6488
F-statistic: 27.78 on 2 and 27 DF, p-value: 2.798e-07

Recall the true model:

$$Y = 4 + 0.8x_1 + 0.6x_2 + \varepsilon,$$

where $\varepsilon \sim N(0, 1)$, x_1 and x_2 are positively correlated with $\rho = 0.9$

Summary:

- β 's are not well estimated in model 1
- Spurious regression coefficient estimates
- In model 2, R^2 and predicted values are OK compared to model 1

Model 2 fit:

```
Call:
lm(formula = Y ~ X1)
```

Residuals:

| | Min | 1Q | Median | 3Q | Max |
|--|----------|----------|----------|---------|---------|
| | -2.09663 | -0.67031 | -0.07229 | 0.87881 | 1.49739 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|--------------|
| (Intercept) | 4.0347 | 0.1763 | 22.888 | < 2e-16 *** |
| X1 | 1.4293 | 0.1955 | 7.311 | 5.84e-08 *** |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9634 on 28 degrees of freedom
Multiple R-squared: 0.6562, Adjusted R-squared: 0.644
F-statistic: 53.45 on 1 and 28 DF, p-value: 5.839e-08

Variance Inflation Factor (VIF)

We can use the **variance inflation factor (VIF)**

$$\text{VIF}_i = \frac{1}{1 - R_i^2}$$

to quantifies the severity of multicollinearity in MLR, where R_i^2 is the **coefficient of determination** when X_i is regressed on the remaining predictors

R example code

```
> library(faraway)
> vif(sim1[, 2:3])
      x1      x2
7.218394 7.218394
```

$\sqrt{\text{VIF}}$ indicates how much larger the standard error increases compared to if that variable had 0 correlation to other predictor variables in the model.

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Multiple Linear Regression Model:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \varepsilon, \quad \varepsilon \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Basic Problem: how to choose between competing linear regression models?

- **Model too “small”:** underfit the data; poor predictions; high **bias**; low **variance**
- **Model too big:** “overfit” the data; poor predictions; low **bias**; high **variance**

In the next few slides we will discuss some commonly used model selection criteria to choose the “right” model to balance bias and variance

General Linear F -Test

Prediction

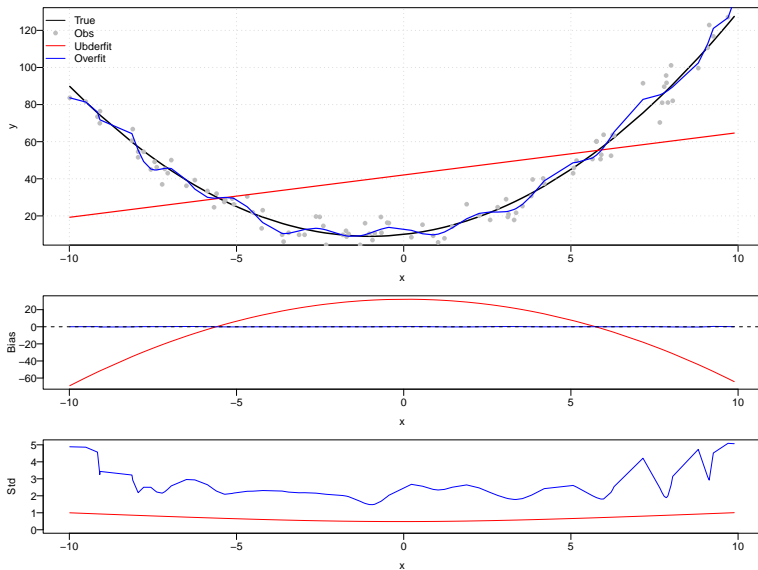
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

An Example of Bias and Variance Tradeoff



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Balancing Bias And Variance: Mallows' C_p Criterion

A good model should balance **bias** and **variance** to get good predictions

$$\begin{aligned}
 (\hat{Y}_i - \mu_i)^2 &= (\hat{Y}_i - \mathbb{E}(\hat{Y}_i) + \mathbb{E}(\hat{Y}_i) - \mu_i)^2 \\
 &= \underbrace{(\hat{Y}_i - \mathbb{E}(\hat{Y}_i))^2}_{\sigma_{\hat{Y}_i}^2 \text{ Variance}} + \underbrace{(\mathbb{E}(\hat{Y}_i) - \mu_i)^2}_{\text{Bias}^2},
 \end{aligned}$$

where $\mu_i = \mathbb{E}(Y_i | X_i = x_i)$

- Mean squared prediction error (MSPE):

$$\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (\mathbb{E}(\hat{Y}_i) - \mu_i)^2$$

- C_p criterion measure:

$$\begin{aligned}
 \Gamma_p &= \frac{\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (\mathbb{E}(\hat{Y}_i) - \mu_i)^2}{\sigma^2} \\
 &= \frac{\sum \text{Var}_{\text{pred}} + \sum \text{Bias}^2}{\text{Var}_{\text{error}}}
 \end{aligned}$$

General Linear F^2 -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

C_p statistic:

$$C_p = \frac{\text{SSE}}{\text{MSE}_F} + 2p - n$$

- When model is correct $E(C_p) \approx p$
- When plotting models against p
 - Biased models will fall above $C_p = p$
 - Unbiased models will fall around line $C_p = p$
 - By definition: C_p for full model equals p

We desire models with small p and C_p around or less than p . See R session for an example

Adjusted R^2 , denoted by R_{adj}^2 , attempts to take account of the phenomenon of the R^2 automatically and spuriously increasing when extra explanatory variables are added to the model.

$$R_{\text{adj}}^2 = 1 - \frac{\text{SSE}/(n - p - 1)}{\text{SST}/(n - 1)}$$

- Choose model which maximizes R_{adj}^2
- Same approach as choosing model with smallest MSE

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Information criteria are statistical measures used for model selection. Commonly used information criteria include:

- Akaike's information criterion (AIC)

$$n \log\left(\frac{\text{SSE}_k}{n}\right) + 2k$$

- Bayesian information criterion (BIC)

$$n \log\left(\frac{\text{SSE}_k}{n}\right) + k \log(n)$$

Here k is the number of the parameters in the model.

These criteria balance the goodness of fit of a model with its complexity

General Linear F^2 -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- **Forward Selection:** begins with no predictors and then adds in predictors one by one using some criterion (e.g., p -value or AIC)
- **Backward Elimination:** starts with all the predictors and then removes predictors one by one using some criterion
- **Stepwise Search:** a combination of backward elimination and forward selection. Can add or delete predictor at each stage
- **All Subset Selection:** Comparing all possible models using a selected criterion. Impractical for “large” number of predictors

General Linear F^2 -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Model:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \varepsilon, \quad \varepsilon \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

We make the following **assumptions**:

- Linearity:

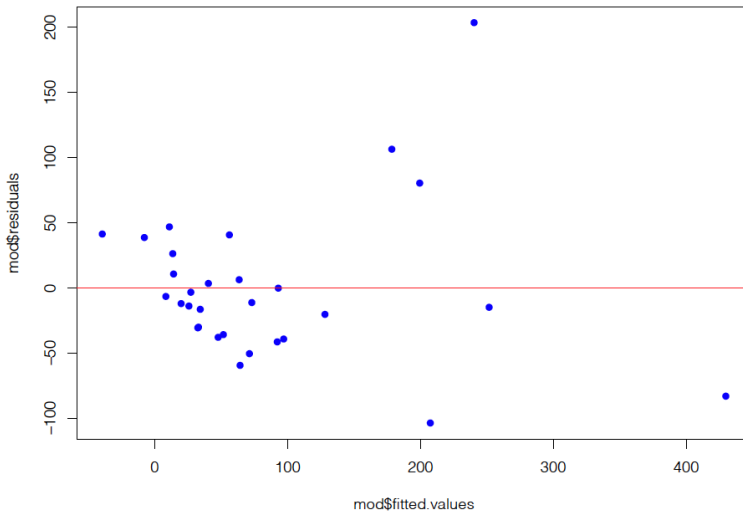
$$E(Y|x_1, x_2, \dots, x_{p-1}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1}$$

- Errors have constant variance, are independent, and normally distributed

$$\varepsilon \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Residuals versus Fits Plot

```
plot(mod$fitted.values, mod$residuals, pch = 16, col = "blue")  
abline(h = 0, col = "red")
```



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

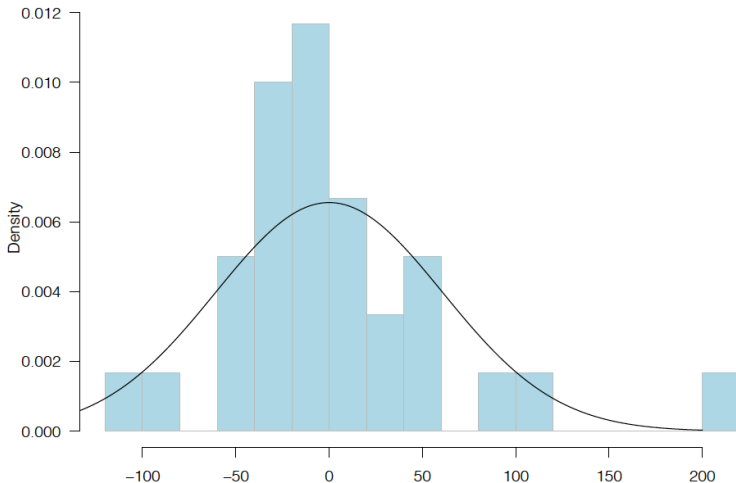
Non-Constant
Variance &
Transformation

We will revisit this in the end of the lecture

Assessing Normality of Residuals: Histogram

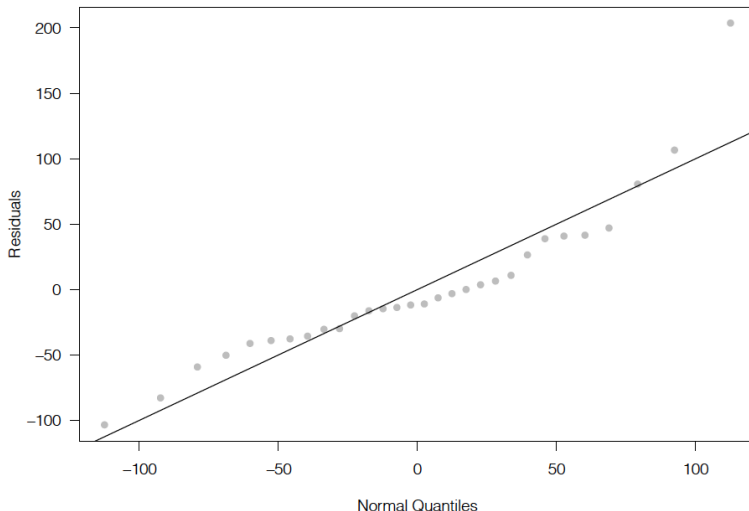
```
par(las = 1)
hist(mod$residuals, 12, prob = T,
     col = "lightblue", border = "gray")
xg <- seq(-200, 200, 1)
yg <- dnorm(xg, 0, 60.86)
lines(xg, yg)
```

Histogram of mod\$residuals



Assessing Normality of Residuals: QQ Plot

```
plot(qnorm(1:30 / 31, 0, 60.86), sort(mod$residuals), pch = 16,  
     col = "gray", xlab = "Normal Quantiles", ylab = "Residuals")  
abline(0, 1)
```



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

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

- The leverage value 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 $\bar{h} = \sum_{i=1}^n \frac{h_i}{n} = \frac{p}{n} \Rightarrow$ a “rule of thumb” is that leverages greater than $\frac{2p}{n}$ should be examined more closely

General Linear F -Test

Prediction

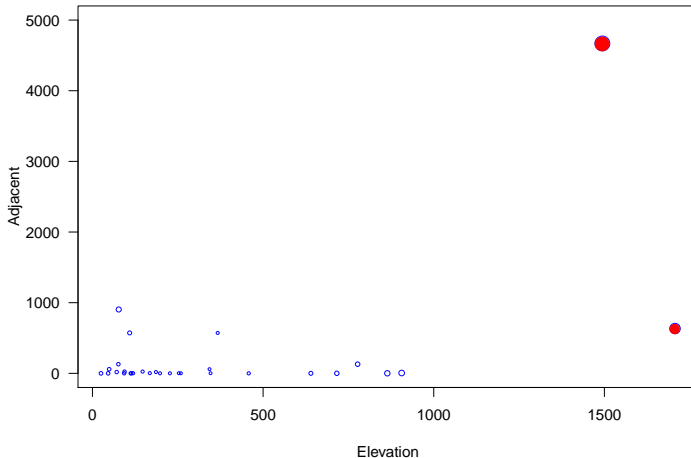
Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Leverage Values of Species ~ Elev + Adj



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

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 **standardized 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 $\text{Corr}(r_i, r_j)$ tends to be small

General Linear F -Test

Prediction

Multicollinearity

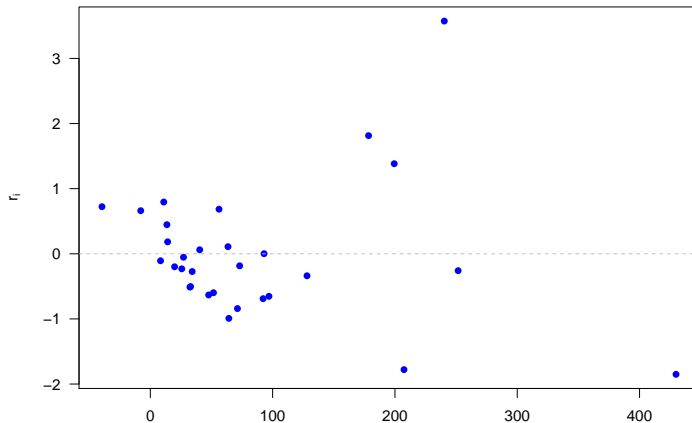
Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Standardized Residuals of Species ~ Elev + Adj

Studentized Residuals



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

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”
- Can show
$$\text{Var}(\hat{y}_{i(i)} - y_i) = \sigma_{(i)}^2 \left(1 + \mathbf{x}_i^T (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right) = \sigma_{(i)}^2 (1 - h_i)$$
- Define the **Studentized (Jackknife) Residuals** as

$$t_i = \frac{\hat{y}_{i(i)} - y_i}{\sqrt{\hat{\sigma}_{(i)}^2 (1 - h_i)}} = \frac{\hat{y}_{i(i)} - y_i}{\sqrt{\text{MSE}_{(i)} (1 - h_i)}}$$

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

General Linear F^* -Test

Prediction

Multicollinearity

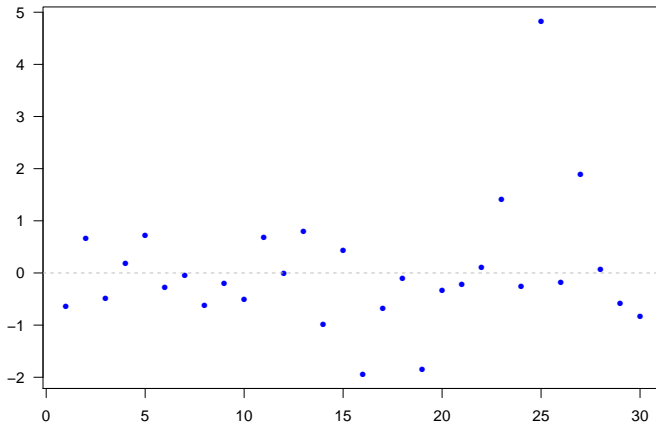
Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Studentized (Jackknife) Residuals of Species ~ Elev + Adj

Jackknife Residuals



General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Identifying Influential Observations: Cook's Distance

Cook's Distance quantifies how much the predicted values change when a particular observation is excluded from the analysis.

- Cook's distance measure (D_i) is defined as:

$$D_i = \frac{(y_i - \hat{y}_i)^2}{p \times \text{MSE}} \left(\frac{h_i}{(1 - h_i)^2} \right)$$

- Cook's Distance considers both leverage and residual, providing a broader measure of influence
- Here are the guidelines commonly used:
 - ① If $D_i > 0.5$, then the i^{th} data point is worthy of further investigation as it may be influential
 - ② If $D_i > 1$, then the i^{th} data point is quite likely to be influential

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

Cook's Distance of Species ~ Elev + Adj

General Linear F -Test

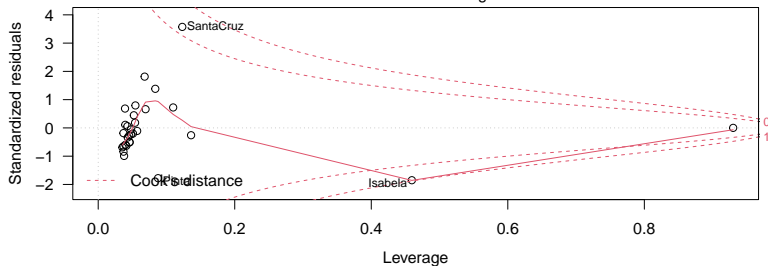
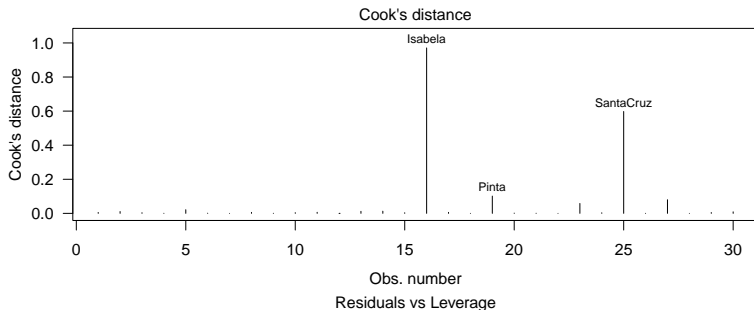
Prediction

Multicollinearity

Model Selection

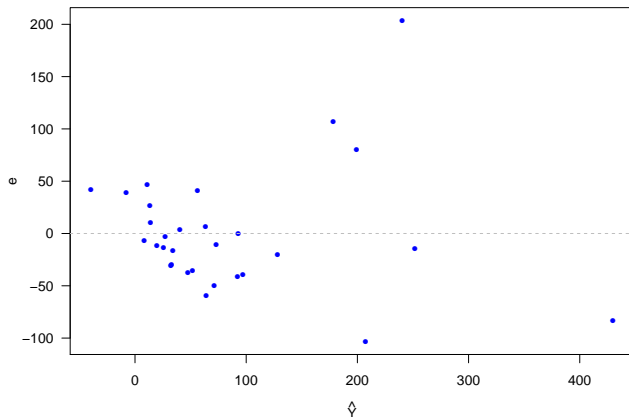
Model Diagnostics

Non-Constant
Variance &
Transformation



Residual Plot of Species ~ Elev + Adj

Residuals



Such a residual plot suggests a violation of constant variance

General Linear F -Test

Prediction

Multicollinearity

Model Selection

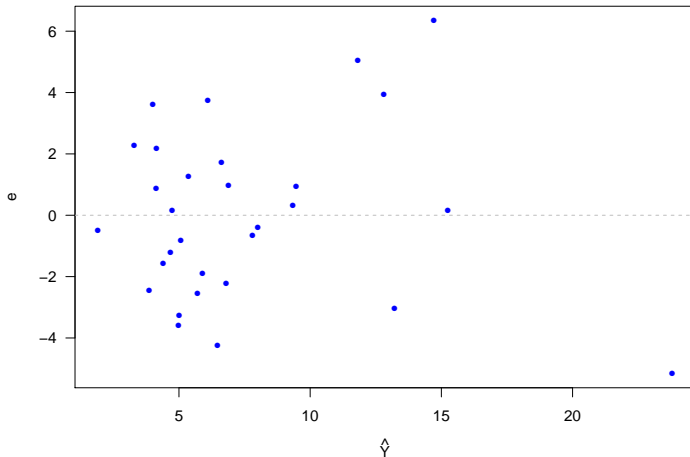
Model Diagnostics

Non-Constant
Variance &
Transformation

Residual Plot After Square Root Transformation

$$\sqrt{\text{Species}} \sim \text{Elev} + \text{Adj}$$

Residuals



General Linear F -Test

Prediction

Multicollinearity

Model Selection

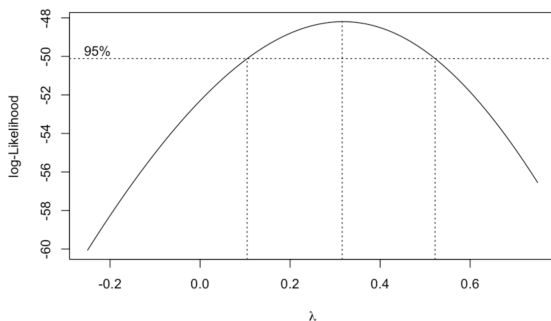
Model Diagnostics

Non-Constant
Variance &
Transformation

Box-Cox Transformation

The Box-Cox method [Box and Cox, 1964] is a powerful way to determine if a transformation on the response is needed

$$g_{\lambda}(y) = \begin{cases} \frac{y^{\lambda}-1}{\lambda} & \text{if } \lambda \neq 0; \\ \log(y) & \text{if } \lambda = 0. \end{cases}$$



In R, we can use the `boxcox` function from the `MASS` package to perform a Box-Cox transformation. The plot suggests a cube root may be needed

These slides cover:

- General Linear F -Test provides a unifying framework for hypothesis tests
- Making predictions and quantifying prediction uncertainty
- Multicollinearity and its implications for MLR
- Model/variable selection can be done via some criterion-based methods to balance bias and variance
- Model diagnostics is crucial to ensure valid statistical inference
- Box-Cox Transformation can be used to transform the response in order to correct model violations

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation

- `anova` for model comparison based on F -test
- `predict`: obtain predicted values from a fitted model
- `vif` under the `faraway` library: computes the variance inflation factors
- `regsubsets` in the `leaps` library and `step` for model selection
- `influence.measures` includes a suite of functions (`hatvalues`, `rstandard`, `rstudent`, `cooks.distance`) for computing regression diagnostics
- `boxcox` in the `MASS` library for performing a **Box-Cox transformation**

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant
Variance &
Transformation