

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

Whitney Huang
Clemson University

Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
CLEMSON UNIVERSITY

General Linear F -Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant Variance & Transformation

4.1

Notes

Notes

Notes

Agenda

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

Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
CLEMSON UNIVERSITY

General Linear F -Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant Variance & Transformation

4.2

Review: t -Test and F -Test in Linear Regression

- 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

Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
CLEMSON UNIVERSITY

General Linear F -Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant Variance & Transformation

4.3

General Linear F-Test

- 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: $F^* = \frac{(SSE_{\text{reduced}} - SSE_{\text{full}})/(k - \ell)}{SSE_{\text{full}}/(n - k - 1)} \Rightarrow$ 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. $x_1, x_3 \Rightarrow H_0 : \beta_2 = \beta_4 = 0$

Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT

General Linear F-Test

Prediction

Multicollinearity

Model Selection

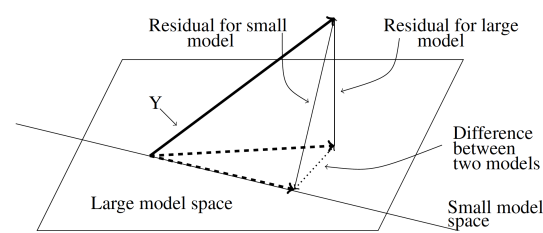
Model Diagnostics

Non-Constant Variance & Transformation

4.4

Notes

Geometric Illustration of General Linear F-Test



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

Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.5

Notes

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
```

Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.6

Notes

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
```

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.7

Notes


Performing a General Linear F-Test

- $H_0 : \beta_{Area} = 0$ vs. $H_a : \beta_{Area} \neq 0$
- $F^* = \frac{(173254 - 169947)/(2-1)}{169947/(30-2-1)} = 0.5254$
- 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
```

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

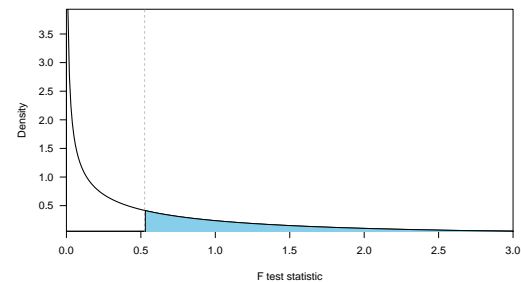
Model Diagnostics

Non-Constant Variance & Transformation

4.8


Notes

Visualizing p-value



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

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.9

Notes

Another Example of General Linear F-Test

```
> full <- lm(Species ~ Area + Elevation + Nearest + Scrub + Adjacent,
data = gale)
> anova(full)
Analysis of Variance Table

Response: Species
      Df Sum Sq Mean Sq F value    Pr(>F)    
Area    1 145478   145478 39.1262 1.826e-06 ***
Elevation 1  65664    65664 17.6613 0.0003155 ***
Nearest   1    29         29  0.0079 0.9300674    
Scrub     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
```

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.10

Notes

Performing a General Linear F-Test

- Null and alternative hypotheses:
 $H_0 : \beta_{Area} = \beta_{Nearest} = \beta_{Scrub} = 0$
 $H_a : \text{at least one of the three coefficients} \neq 0$
- $F^* = \frac{(100003-89231)/(5-2)}{89231/(30-5-1)} = 0.9657$
- p-value: $P[F > 0.9657] = 0.425$, where $F \sim F_{3,24}$

```
> anova(reduced, full)
Analysis of Variance Table

Model 1: Species ~ Elevation + Adjacent
Model 2: Species ~ Area + Elevation + Nearest + Scrub + Adjacent
      Res.Df  RSS Df Sum of Sq    F Pr(>F)
1         27 100003
2         24  89231   3    10772 0.9657  0.425
```

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.11

Notes

Multiple Linear Regression Prediction

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

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.12

Notes

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      abdomen
      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
```

Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

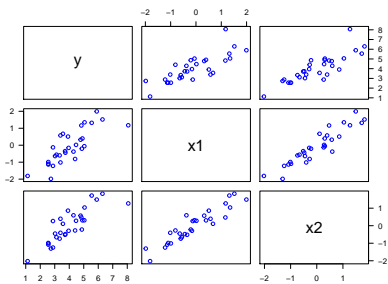
Model Diagnostics

Non-Constant Variance & Transformation

4.16


Notes

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
```

Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.17


Notes

Multicollinearity Cont'd

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

- Numerical issue \Rightarrow the matrix $X^T X$ 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

Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.18

Notes

An Simulated Example

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

Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
FACULTY OF MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT OF APPLIED AND COMPUTATIONAL MATHEMATICS

General Linear F-Test

Prediction

Multicollinearity

Model Selection

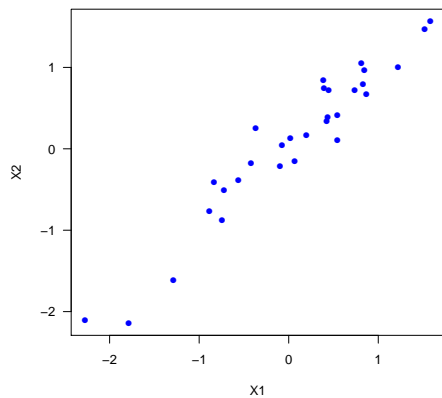
Model Diagnostics

Non-Constant Variance & Transformation

4.19

Notes

Scatter Plot: x_1 vs. x_2



Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
FACULTY OF MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT OF APPLIED AND COMPUTATIONAL MATHEMATICS

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.20

Notes

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
```

Multiple Linear Regression II

UNIVERSITY OF CAMBRIDGE
FACULTY OF MATHEMATICAL AND STATISTICAL SCIENCES
DEPARTMENT OF APPLIED AND COMPUTATIONAL MATHEMATICS

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.21

Notes

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
```

Multiple Linear Regression II

UNIVERSITY OF OXFORD
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.22

Notes

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  0.00456 **
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

Multiple Linear Regression II

UNIVERSITY OF OXFORD
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.23

Notes

Variance Inflation Factor (VIF)

We can use the [variance inflation factor \(VIF\)](#)

$$VIF_i = \frac{1}{1 - R_i^2}$$

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

Multiple Linear Regression II

UNIVERSITY OF OXFORD
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.24

Notes

Model Selection in Multiple Linear Regression

Multiple Linear Regression Model:

Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{p-1} x_{p-1} + \epsilon, \epsilon \sim^{i.i.d.} 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

Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

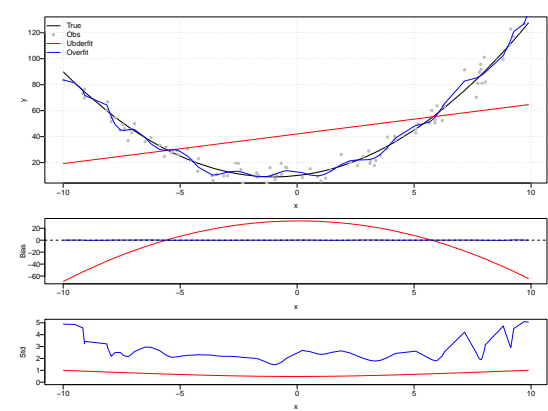
Model Diagnostics

Non-Constant Variance & Transformation

4.25

Notes

An Example of Bias and Variance Tradeoff



Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.26

Notes

Balancing Bias And Variance: Mallows' Cp Criterion

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


(\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}

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
- Cp criterion measure:

\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}}}

Multiple Linear Regression II

UNIVERSITY OF MISSOURI
MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.27

Notes

C_p Criterion

C_p statistic:

$$C_p = \frac{SSE}{MSE_F} + 2p - n$$

- 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

Multiple Linear Regression II

Journal of MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.28

Notes

Adjusted R^2 Criterion

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

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

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

Multiple Linear Regression II

Journal of MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.29

Notes

Information criteria

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

- Akaike's information criterion (AIC)

$$n \log\left(\frac{SSE_k}{n}\right) + 2k$$

- Bayesian information criterion (BIC)

$$n \log\left(\frac{SSE_k}{n}\right) + k \log(n)$$

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

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

Multiple Linear Regression II

Journal of MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F -Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation


4.30

Notes

Automatic Search Procedures

- **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

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.31

Notes

Model Assumptions

Model:


$$Y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \cdots + \beta_{p-1}x_{p-1} + \varepsilon, \quad \varepsilon \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

We make the following assumptions:

- Linearity:
$$E(Y|x_1, x_2, \cdots, x_{p-1}) = \beta_0 + \beta_1x_1 + \beta_2x_2 + \cdots + \beta_{p-1}x_{p-1}$$
- Errors have constant variance, are independent, and normally distributed

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

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

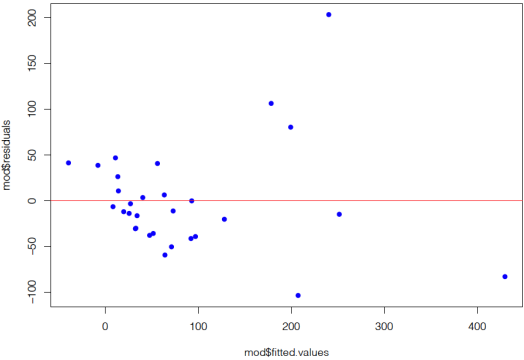
Non-Constant Variance & Transformation

4.32

Notes


Residuals versus Fits Plot

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



We will revisit this in the end of the lecture

Multiple Linear Regression II

UNIVERSITY OF MELBOURNE
SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear F-Test

Prediction

Multicollinearity

Model Selection

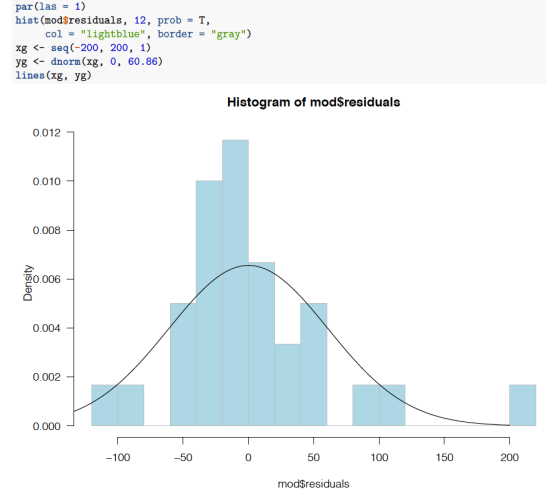
Model Diagnostics

Non-Constant Variance & Transformation


4.33

Notes

Assessing Normality of Residuals: Histogram



Multiple Linear Regression II

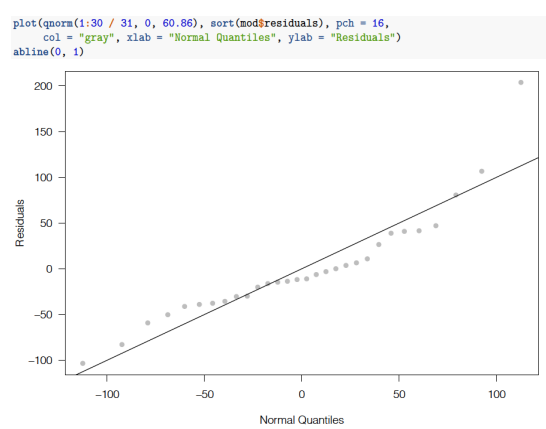
 SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.34

Notes

Assessing Normality of Residuals: QQ Plot



Multiple Linear Regression II

 SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.35

Notes

Leverage: Detecting “Extreme” Predictor Values

Recall in MLR that $\hat{y} = (X^T X)^{-1} X^T y = H y$ 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
- $\frac{1}{n} \leq h_i \leq 1, \quad 1 \leq i \leq n$ and $\bar{h} = \sum_{i=1}^n \frac{h_i}{n} = \frac{p}{n} \Rightarrow$ a “rule of thumb” is that leverages greater than $\frac{2p}{n}$ should be examined more closely

Multiple Linear Regression II

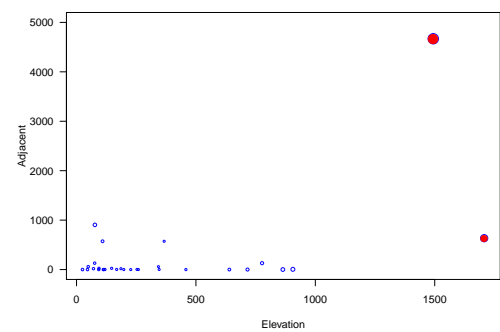
 SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES

General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.36

Notes

Leverage Values of Species ~ Elev + Adj



Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
UNIVERSITY OF OKLAHOMA

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.37


Notes

Standardized 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 **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

Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
UNIVERSITY OF OKLAHOMA

General Linear F-Test

Prediction

Multicollinearity

Model Selection

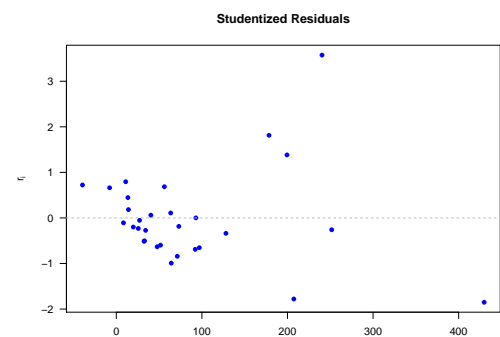
Model Diagnostics

Non-Constant Variance & Transformation


4.38

Notes

Standardized Residuals of Species ~ Elev + Adj



Multiple Linear Regression II

SCHOOL OF MATHEMATICAL AND STATISTICAL SCIENCES
UNIVERSITY OF OKLAHOMA

General Linear F-Test

Prediction

Multicollinearity

Model Selection

Model Diagnostics

Non-Constant Variance & Transformation

4.39

Notes

Studentized (Jackknife) Residuals

- 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 $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

Multiple Linear Regression II

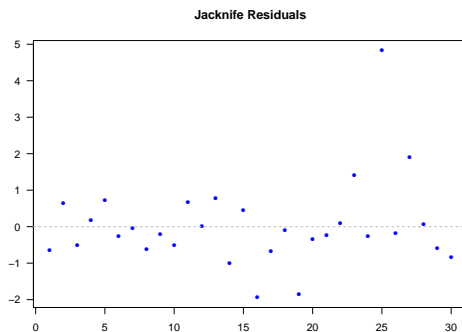


General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.40

Notes

Studentized (Jackknife) Residuals of Species ~ Elev + Adj



Multiple Linear Regression II



General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.41

Notes

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

Multiple Linear Regression II

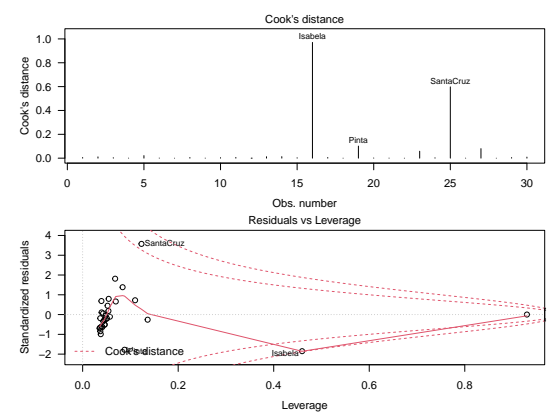


General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.42

Notes

Cook's Distance of Species ~ Elev + Adj



Multiple Linear Regression II



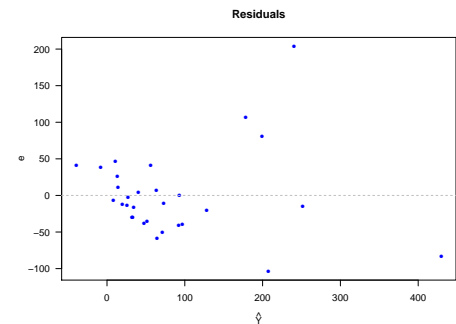
UNIVERSITY OF
MATHEMATICAL AND
STATISTICAL SCIENCES

General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.43


Notes

Residual Plot of Species ~ Elev + Adj



Such a residual plot suggests a violation of constant variance

Multiple Linear Regression II



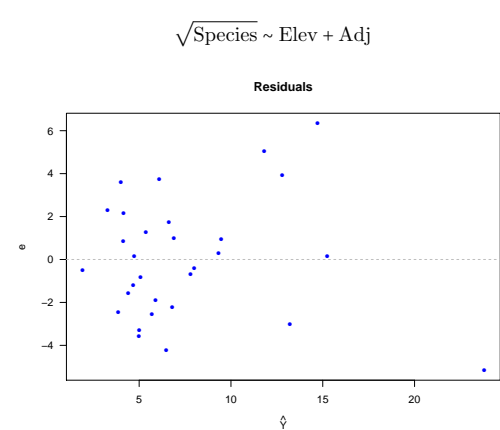
UNIVERSITY OF
MATHEMATICAL AND
STATISTICAL SCIENCES

General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation


4.44

Notes

Residual Plot After Square Root Transformation



Multiple Linear Regression II



UNIVERSITY OF
MATHEMATICAL AND
STATISTICAL SCIENCES

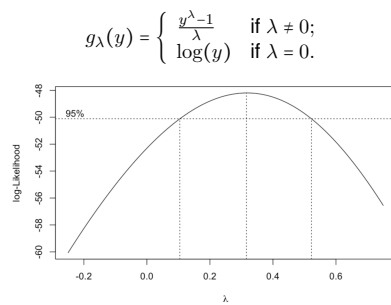
General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.45

Notes

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



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

Multiple Linear Regression II



General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.46

Notes

Summary

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

Multiple Linear Regression II



General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.47

Notes

R Functions to Know

- `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 `stepAIC` 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

Multiple Linear Regression II



General Linear
F-Test
Prediction
Multicollinearity
Model Selection
Model Diagnostics
Non-Constant
Variance &
Transformation

4.48

Notes
