

# Applied Regression Analysis Classnotes

Dustin Leatherman

April 17, 2020

## Contents

<b>1</b>	<b>Session 1 - Summary and Review</b>	<b>4</b>
1.1	Relationships . . . . .	4
1.1.1	Functional . . . . .	4
1.1.2	Statistical . . . . .	4
1.2	Basic Concepts . . . . .	4
1.2.1	Construction of Regression models . . . . .	5
1.2.2	Uses of Regression Analysis . . . . .	5
1.3	Simple Linear Regression (SLR) . . . . .	5
1.3.1	Properties of $\epsilon_i$ . . . . .	6
1.3.2	Properties . . . . .	6
1.3.3	Alternative forms of SLR . . . . .	6
1.3.4	Method of Least Squares . . . . .	6
1.3.5	Gauss-Markov Theorem . . . . .	7
1.3.6	Residual . . . . .	8
1.4	Normal Error Regression Model . . . . .	9
<b>2</b>	<b>Session 2 - Inferences in Regression and Correlation Analysis (2019/09/18)</b>	<b>9</b>
2.1	Properties . . . . .	9
2.2	$\beta_1$ . . . . .	9
2.2.1	Inferences . . . . .	9
2.2.2	Sampling Distribution . . . . .	10
2.2.3	PROOF: $b_1$ is a linear combination of Y's . . . . .	10
2.2.4	Properties . . . . .	10
2.3	$\beta_0$ . . . . .	11
2.4	Spacing of X Levels . . . . .	11
2.5	Prediction of new observations . . . . .	11

2.5.1	Interval Estimation of $E(Y_0)$	11
2.5.2	Sampling Distribution	11
2.5.3	Prediction	12
2.6	ANOVA Approach to Regression	12
<b>3</b>	<b>Session 3 - General Linear Testing &amp; Model Selection (2019/09/25)</b>	<b>13</b>
3.1	General Linear Test Approach	13
3.1.1	Reduced Model	13
3.1.2	Coefficients of Determination ( $R^2$ )	14
3.1.3	Coefficient of Correlation: $r = \pm\sqrt{R^2}$	14
3.2	Assessing the Quality of a Model	15
3.2.1	Residuals (observed error)	15
3.2.2	Residual Plots	15
3.2.3	Test of Randomness	17
3.2.4	Constant Variance	19
<b>4</b>	<b>Session 4 - Transformations &amp; Inference (2019/10/02)</b>	<b>19</b>
4.1	Transformations	19
4.1.1	Box-Cox Transformations	20
4.2	Simultaneous Inference	20
4.2.1	Working-Hotelling Procedure	20
4.2.2	Bonferonni Procedure	21
<b>5</b>	<b>Session 5 - Prediction &amp; Linear Algebra in Regression</b>	<b>21</b>
5.1	Simultaneous Intervals	21
5.1.1	Confidence	21
5.1.2	Prediction	21
5.2	Inverse Prediction (“Calibration”)	21
5.3	Linear Algebra in Regression	22
5.3.1	Review	22
5.3.2	Expectations	25
5.3.3	Variance-Covariance Matrix	25
5.3.4	Multivariate Normal Distribution	26
5.3.5	Least Squares Estimation	27
<b>6</b>	<b>Session 6 - Sums of Squares and Multiple Linear Regression</b>	<b>27</b>
6.1	Sum of Squares	27
6.1.1	SSE	28
6.1.2	SSTo	28
6.1.3	SSR	28

6.2	Mean Estimates $\sigma^2$ . . . . .	28
6.2.1	Mean Responses . . . . .	28
6.3	Variance of $\hat{Y}_h$ . . . . .	28
6.4	Multiple Regression Models . . . . .	29
6.4.1	Interpretation . . . . .	29
6.4.2	Aside: Multi-Collinearity . . . . .	29
6.4.3	Matrix Notation . . . . .	29
6.4.4	ANOVA Table . . . . .	30
6.4.5	Omnibus F-Test for Regression Relation . . . . .	30
6.4.6	Coefficient of Multiple Determination . . . . .	30
6.4.7	Coefficient of Multiple Correlation . . . . .	30
6.4.8	Inferences in $\beta_k$ . . . . .	30
<b>7</b>	<b>Session 7 - Multiple Regression &amp; Qualitative\Quantitative Predictors</b>	<b>31</b>
7.1	Multiple Regression . . . . .	31
7.1.1	Extra Sums of Squares . . . . .	31
7.2	Multi-collinearity . . . . .	32
7.2.1	Effects . . . . .	33
7.3	Polynomial Regression Models . . . . .	33
<b>8</b>	<b>Session 8 - Interaction Models &amp; Model Selection</b>	<b>33</b>
8.1	Interaction Regression Models . . . . .	33
8.1.1	Additive Effects . . . . .	34
8.1.2	Qualitative Predictors . . . . .	34
8.2	Model and Variable Selection . . . . .	35
8.2.1	Criterion for Model Selection . . . . .	35
<b>9</b>	<b>Session 9 - Model and Variable Selection &amp; Assessing Diagnostics</b>	<b>38</b>
9.1	Model and Variable Selection . . . . .	38
9.1.1	Automatic Search Procedures . . . . .	38
9.1.2	Model Validation . . . . .	38
9.2	Assessing Diagnostics . . . . .	39
9.2.1	Added-variable Plots . . . . .	39
<b>10</b>	<b>Session 10 - Outliers &amp; Weighted Least Squares</b>	<b>40</b>
10.1	Outliers . . . . .	40
10.1.1	Identifying Outlying Y Observations . . . . .	40
10.1.2	What is an Outlying Y Observation? . . . . .	41

10.1.3	Identifying Outlying X Observations . . . . .	41
10.2	Influential Cases . . . . .	42
10.2.1	Identifying Influential Cases . . . . .	42
10.3	Variance Inflation Factors . . . . .	43
10.4	Weighted Least Squares . . . . .	44
10.4.1	Iteratively Reweighted Least Squares . . . . .	44
<b>11</b>	<b>Extra Curricular - Weighted Least Squares, Ridge, and Robust Regression</b>	<b>44</b>
11.1	Weighted Least Squares . . . . .	44
11.2	OLS with Heteroskedasticity . . . . .	45
11.3	Ridge Regression . . . . .	45
11.4	robust regression . . . . .	46
11.4.1	$u$ . . . . .	47
11.5	Regression Tree (Non-parametric Method) . . . . .	48

## 1 Session 1 - Summary and Review

### 1.1 Relationships

#### 1.1.1 Functional

Mathematical formula

$$Y = f(x)$$

#### 1.1.2 Statistical

Not a perfect relationship

$$Y = f'(x) + \epsilon$$

observations = trials = case

quadratic = curvilinear

Linear models means that the *slope* is not raised to any powers

- $\hat{y} = \beta_0 + \beta_1 X^2$  is linear
- $\hat{y} = \beta_0 + \beta_1^2 X$  is **not** linear

### 1.2 Basic Concepts

- Tendency of Y to vary with X in a *systematic fashion*
- scatter of points around a curve of a statistical relation

- Prob. Distr of Y for each **Level** of X
- means of Y's distr. to vary for each value of X
  - each point on the regression line can be represented as  $\mu_{Y|X_i}$

NOTE: Sir Francis Galton came up with the term “regression”

**Regression function of Y on X:** Means of the prob. distr. have a systematic relation to the Level of X

**Regression curve:** graph of the regression function

### 1.2.1 Construction of Regression models

1. Selection of pred. vars
2. Functional form of the regression relation
  - Summary plots
  - Scatter plots
3. Scope of Model
  - **Scope:** What is the Domain? (range of X's)
  - Making predictions outside the Domain is considered extrapolation and is dangerous

### 1.2.2 Uses of Regression Analysis

1. Description
2. Control
3. Prediction (most abused)

**Association does not imply causation!**

## 1.3 Simple Linear Regression (SLR)

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

- One Predictor
- Linear in the Parameters
- Linear in the Predictor Variables

- SLR = First-Order Model (term from outside of statistics)

$Y_i$ : Value of the response for the  $i$ th trial  $\beta_0$ : Parameters (unknown. estimate these)  $X_i$ : Value of the predictor of the  $i$ th term (known)  $\epsilon_i$ : random error term of the  $i$ th observation

### 1.3.1 Properties of $\epsilon_i$

- $E(\epsilon_i) = 0$
- $\sigma^2(\epsilon_i) = Var(\epsilon_i) = \sigma^2$
- $\epsilon_i$  and  $\epsilon_j$  are uncorrelated

### 1.3.2 Properties

1.  $Y_i$  is the sum of two components. It is **random** because its composed of a random term. constant term:  $\beta_0 + \beta_1 X_i$  random term:  $\epsilon_i$
2.  $E(Y_i) = E(\beta_0 + \beta_1 X_i + \epsilon_i) \rightarrow E(\beta_0 + \beta_1 X_i) + E(\epsilon_i) \rightarrow \beta_0 + \beta_1 X_i$
3.  $Y_i$  falls short of regression function by  $\epsilon_i$
4.  $Var(\epsilon_i) = \sigma^2$  error terms have constant variance
5. Since error terms are uncorrelated, then responses ( $Y_i$  and  $Y_j$  are uncorrelated)

### 1.3.3 Alternative forms of SLR

1.  $Y_i = \beta_0 X_0 + \beta_1 X_i + \epsilon_i$ , where  $X_0 = 1$
2.  $Y_i = \beta_0 + \beta_1(X_i - \bar{x}) + \beta_1 \bar{x} + \epsilon_i \rightarrow (\beta_0 + \beta_1 \bar{x}) + \beta_1(x_i - \bar{x})\epsilon_i \rightarrow \beta_0^* + \beta_1(x_i - \bar{x}) + \epsilon_i$

### 1.3.4 Method of Least Squares

1. Goal Find estimators of  $\beta_0$  and  $\beta_1$

For each  $(X_i, Y_i)$ :  $Y_i - (\beta_0 + \beta_1 X_i)$   $Q = \sum_1^n [Y_i - \beta_0 - \beta_1 X_i]^2$

$\hat{\beta}_0$  and  $\hat{\beta}_1$  are estimators of  $\beta_0$  &  $\beta_1$  that minimize Q for given data  $(X_i, Y_i)$ ,  $i = [1, n]$

### 1.3.5 Gauss-Markov Theorem

1. Proof First, lets find the value of  $b_0$  by taking the partial derivative of  $Q$  with respect to  $\beta_1$

$$\begin{aligned} Q &= \sum_1^n [Y_i - \beta_0 - \beta_1 X_i]^2 \\ \frac{dQ}{d\beta_1} &= -2 \sum_1^n X_i [Y_i - \beta_0 - \beta_1 X_i] \\ &\rightarrow \sum_1^n X_i (Y_i - b_0 - b_1 X_i) = 0 \\ &\rightarrow \sum_1^n X_i Y_i - b_0 \sum_1^n x_i - b_1 \sum_1^n x_i^2 = 0 \\ &\rightarrow \sum_1^n Y_i - nb_0 - b_1 \sum_1^n x_i = 0 \\ &\rightarrow \sum_1^n Y_i - b_1 \sum_1^n x_i = nb_0 \\ &\rightarrow \bar{Y} - b_1 \bar{x} = b_0 \end{aligned} \tag{1}$$

Once  $b_0$  is found, lets use it to find the value of  $b_1$ . Replace values of  $b_0$  with the equation above.

$$\begin{aligned}
& \sum_1^n X_i Y_i - b_0 \sum_1^n x_i - b_1 \sum_1^n x_i^2 = 0 \\
& \rightarrow \sum_1^n X_i Y_i - (\bar{Y} - b_1 \bar{x}) \sum_1^n x_i - b_1 \sum_1^n x_i^2 = 0 \\
& \rightarrow \sum_1^n X_i Y_i - \left( \frac{\sum_1^n Y_i}{n} - b_1 \frac{\sum_1^n x_i}{n} \right) \sum_1^n x_i - b_1 \sum_1^n x_i^2 = 0 \\
& \rightarrow \sum_1^n X_i Y_i - \frac{\sum_1^n x_i \sum_1^n y_i}{n} + b_1 \frac{(\sum_1^n x_i)^2}{n} - b_1 \sum_1^n x_i^2 \\
& \rightarrow \sum_1^n x_i y_i - \frac{\sum_1^n x_i \sum_1^n y_i}{n} = b_1 \left[ \sum_1^n x_i^2 - \frac{(\sum_1^n x_i)^2}{n} \right] \\
& = \dots = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_1^n (x_i - \bar{x})^2}
\end{aligned} \tag{2}$$

## 2. Properties

- (a)  $E(b_0) = \beta_0$  &  $E(b_1) = \beta_1$
- (b)  $b_0$  &  $b_1$  are more precise than any other unbiased estimators of  $\beta_0$  and  $\beta_1$  that are linear functions of  $Y_i$

### 1.3.6 Residual

Difference between the observation and the estimated value  $e_i = Y_i - \hat{Y}_i$ ,  $i = 1, n$

1.  $\sum_i^n e_i = 0$
  2.  $\sum_i^n e_i^2$  is a minimum
  3.  $\sum_i^n Y_i = \sum_i^n \hat{Y}_i$
1. Goal Estimate  $\sigma^2$  know  $E(S^2) = E\left(\frac{\sum(Y_i - \bar{Y})^2}{n-1}\right)$ 
    - numerator == sum of squares
    - $n - 1$  == df
    - $S^2 = \text{Mean Square} = \frac{SS}{df}$



2.  $SSE = \sum (Y_i - \hat{Y}_i)^2 = \sum e_i^2$

- $SSE = \text{Sum of Square Error} = \text{Residual Sums of Squares}$
- $MSE = SSE / n - 2$
- $df \text{ of } SSE = n - 2$
- $E(MSE) = \sigma^2$

## 1.4 Normal Error Regression Model

$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$  where  $\epsilon \approx iidN(0, \sigma^2)$ ,  $i = [1, n]$  so  $Y_i \approx N(\beta_0 + \beta_1 X_i, \sigma^2)$   
 To find MLE's of  $\beta_0$  &  $\beta_1$  i.e.  $\hat{\beta}_0$  &  $\hat{\beta}_1$   $L(\beta_0, \beta_1 | \sigma^2) = \prod pdf$

- MLE of  $\beta_0$ :  $\hat{\beta}_0 = b_0$
- MLE of  $\beta_1$ :  $\hat{\beta}_1 = b_1$

## 2 Session 2 - Inferences in Regression and Correlation Analysis (2019/09/18)

Model =  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

### 2.1 Properties

- $\epsilon_i \approx iidN(0, \sigma^2)$
- $Y_i \approx iidN(\beta_0 + \beta_1 X_i, \sigma^2)$
- $X_i$ : known constant
- $\beta_0$  &  $\beta_1$  are parameters to investigate

### 2.2 $\beta_1$

#### 2.2.1 Inferences

$H_0 : \beta_1 = 0$  (implies no linear association)  $H_1 : \beta_1 \neq 0$

This hypothesis test determines if there is a relationship

### 2.2.2 Sampling Distribution

$$b_1 = \frac{\Sigma((x_i - \bar{x})(y_i - \bar{y}))}{\Sigma(x_i - \bar{x})^2}$$

- $E(b_1) = \beta_1$
- $Var(b_1) = \frac{\sigma^2}{\Sigma(x_i - \bar{x})^2}$

### 2.2.3 PROOF: $b_1$ is a linear combination of Y's

- $b_1 = \frac{\Sigma((x_i - \bar{x})(y_i - \bar{y}))}{\Sigma(x_i - \bar{x})^2}$
- $b_1 = \frac{\Sigma(x_i - \bar{x})y_i - \bar{y}\Sigma(x_i - \bar{x})}{\Sigma(x_i - \bar{x})^2}$
- $b_1 = \frac{\Sigma(x_i - \bar{x})y_i}{\Sigma(x_i - \bar{x})^2}$

Let  $K_i = \frac{x_i - \bar{x}}{\Sigma(x_i - \bar{x})^2}$

**Facts about  $K_i$**

- $\Sigma K_i = \Sigma \frac{x_i - \bar{x}}{\Sigma(x_i - \bar{x})^2} = 0$
- $\Sigma K_i^2 = \Sigma \left( \frac{x_i - \bar{x}}{\Sigma(x_i - \bar{x})^2} \right)^2 = \frac{1}{\Sigma(x_i - \bar{x})^2}$
- $b_1 = \Sigma K_i Y_i$

Therefore  $b_1$  is a linear combination of  $Y_i$

### 2.2.4 Properties

- $E(\hat{\beta}_1) = E(\Sigma K_i Y_i) = \Sigma K_i E(Y_i) = \Sigma K_i (\beta_0 + \beta_1 X_i) = \beta_1 \Sigma K_i X_i = \beta_1$

More detailed proof of  $\Sigma K_i X_i = 1$  exists in notes. It was a sidebar in class.

- $\beta_1 \approx N\left(\beta_1, \frac{\sigma^2}{\Sigma(x_i - \bar{x})^2}\right)$
- $\frac{b_1 - \beta_1}{\sqrt{\frac{\sigma^2}{\Sigma(x_i - \bar{x})^2}}} \approx N(0, 1)$

Recall  $E(MSE) = E\left(\frac{SSE}{n-2}\right) = \sigma^2$

Thus  $\frac{b_1 - \beta_1}{\sqrt{\frac{MSE}{\Sigma(x_i - \bar{x})^2}}} \approx t_{n-2}$  NOTE: a T Distribution is a standard normal distribution divided by a chi-square distribution scaled by its DF

## 1. Solving the Hypothesis Test

Recall

$$H_0 : \beta_1 = 0 \quad H_1 : \beta_1 \neq 0$$

**Test Statistic**

$$t^* = \frac{b_1}{\sqrt{\frac{MSE}{\sum (x_i - \bar{x})^2}}} = \frac{b_1}{SE_{b1}} \approx t_{n-2}$$

Then p-value can be calculated

## 2.3 $\beta_0$

$$b_0 \approx N(\beta_0, \sigma^2[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2}])$$

If  $Y_i$  are not exactly normal,  $b_0$  and  $b_1$  are approx. normal. Thus the t statistic provides some level of confidence.

## 2.4 Spacing of X Levels

- The greater the spread of x, the larger  $\sum (x_i - \bar{x})^2$
- $\text{Var}(b_1)$  and  $\text{Var}(b_0)$  decrease

## 2.5 Prediction of new observations

Let a new observation be defined as  $Y_0$

### 2.5.1 Interval Estimation of $E(Y_0)$

- $X_0$ : level of x we want to estimate the mean response
- $E(Y_0)$ : mean response when  $X = X_0$
- $\hat{Y}_0 = b_0 + b_1 X_0$ : Point estimate of  $E(Y_0)$

### 2.5.2 Sampling Distribution

$$\hat{Y}_0 \approx N(E(Y_0), \sigma^2[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}])$$

$$\hat{Y}_0 \pm t_{\frac{\alpha}{2}, n-2} \sqrt{MSE(\frac{1}{n} + \frac{(X_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2})}$$

NOTE: **confidence interval** == **mean prediction interval** == **single value**

### 2.5.3 Prediction

$\hat{Y}_1$ : predicted individual outcome drawn from the distr. of  $Y$

#### Assumptions

- $E(Y_1)$ : estimated by  $\hat{Y}_1$
- $\text{Var}(Y_1)$ : estimated by MSE

$$\text{Var}(\text{pred}) = \text{Var}(\hat{Y}_1) + \text{Var}(\hat{Y}_0) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

100(1 -  $\alpha$ )% **Prediction Interval**

- $\hat{Y}_1 \pm t_{\frac{\alpha}{2}, n-2} \sqrt{MSE \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right)}$

## 2.6 ANOVA Approach to Regression

Partition the Total Sums of Squares

1. When ignoring the predictor variable, Variation is based on  $Y_i - \bar{Y}$  deviations.

*SSTo*: Total Sums of Squares (or TSS) Therefore,  $SSTo = \sum (Y_i - \bar{y})^2$

2. When using the predictor variable, variation based on  $Y_i - \hat{Y}_i$  deviations. i.e. residuals

*SSE*: Error Sum of Squares Therefore,  $SSE = \sum (Y_i - \hat{Y}_i)^2$

*SSR*: Regression Sum of Squares  $SSR = \sum (Y_i - \bar{y})^2$

**NOTE**:  $SSR = SSTo - SSE$  **OR**  $SSTo = SSR + SSE$ . proof is in notebook.

record here if needed

#### Degrees of Freedom (df)

- *SSTo*:  $n - 1$ .  $Y_i - \bar{y}$
- *SSE*:  $n - 2$ .  $Y_i - \hat{Y}_i$
- *SSR*:  $2 - 1 = 1$ .  $\hat{Y}_i - \bar{y}$

NOTE:

- $E(MSE) = \sigma^2$
- $E(MSR) = \sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2$

Source	SS	df	MS	F Statistic
Regression	SSR	1	$MSR = \frac{SSR}{1}$	$F = \frac{MSR}{MSE}$
Error	SSE	n - 2	$MSE = \frac{SSE}{n-2}$	
<b>Total</b>	SSTo	n - 1		

$F^*$  is the test statistic for

$H_0 : \beta_1 = 0$   $H_1 : \beta_1 \neq 0$

$F^* \approx F_{1,n-2}$  if  $H_0$  is true  $(t^*)^2 = F^*$  if  $F^* \approx F_{1,n-2}$

### 3 Session 3 - General Linear Testing & Model Selection (2019/09/25)

#### 3.1 General Linear Test Approach

Full Model:  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$  where  $\epsilon_i \approx iidN(0, \sigma^2)$

This can be fit by either Least Squares or Maximum Likelihood

**Notes** F = Full Model R = Reduced Model

$$\begin{aligned}
 SSE(F) &= \sum [Y_i - (b_0 + b_1 X_i)]^2 \\
 &= \sum (Y_i - \hat{Y}_i)^2 \\
 &= SSE
 \end{aligned} \tag{3}$$

##### 3.1.1 Reduced Model

$$\begin{aligned}
 H_0 : \beta_1 &= 0 \text{ if } H_0 \text{ then } Y_i = \beta_0 + \epsilon_i \\
 H_A : \beta_1 &\neq 0
 \end{aligned} \tag{4}$$

**Test Statistic:**  $SSE(F) \leq SSE(R)$

The more parameters in the model, the better the fit **thus** smaller deviations around the fitted regression model.

A small diff suggests  $H_0$  holds.  $(SSE(R) - SSE(F))$

$$F^* = \frac{\frac{SSE(R) - SSE(F)}{df_R - df_F}}{\frac{SSE(F)}{df_F}} \tag{5}$$

**Note:** The full model has less variation because the hope is that the predictor (X) helps explain the spread in the response (Y).

p-value =  $P(F_{df_R - df_F, df_F} \geq F^*)$  For SLR and testing the null hypothesis ( $H_0 : \beta_1 = 0$ ),

$$\begin{aligned}
F^* &= \frac{\frac{SST_o - SSE}{(n-1) - (n-2)}}{\frac{SSE}{n-2}} \\
&= \frac{SSR}{MSE} \\
&= \frac{MSR}{MSE}
\end{aligned} \tag{6}$$

This is exactly like the ANOVA table!

### 3.1.2 Coefficients of Determination ( $R^2$ )

**Goal:** Quantify how much variation in the response is explained by the model. **Def:** The proportion of variation in Y explained by regressing Y on X.

$$R^2 = \frac{SSR}{SST_o} = 1 - \frac{SSE}{SST_o}$$

**Properties**

- $0 \leq R^2 \leq 1$
- $R^2 = 1$  indicates a perfect fit
- $R^2 = 0 \rightarrow b_1 = 0$  thus a horizontal line **OR** a non-linear pattern

A high  $R^2$  value does NOT indicate

- useful predictions can be made
- estimated regression line is a good fit
- x and y are related

### 3.1.3 Coefficient of Correlation: $r = \pm\sqrt{R^2}$

A measure of the linear association between Y and X when Y and X are random variables. **Properties**

- $-1 \leq r \leq 1$
- sign of correlation matches sign of slope

### 3.2 Assessing the Quality of a Model

#### Diagnostics for X (predictor variable)

1. Dot Plot
2. Sequence Plot  $X_1, \dots, X_n$ . No pattern is good
3. Stem-and-Leaf plot ( $< 100$  observations)
4. Box Plot
5. Histogram

#### 3.2.1 Residuals (observed error)

$$e_i = Y_i - \hat{Y}_i$$

##### Properties

- $\bar{e} = \frac{\sum e_i}{n} = 0$
- $S_e^2 = \frac{\sum (e_i - \bar{e})^2}{n-2} = \frac{\sum e_i^2}{n-2} = \frac{SSE}{n-2} = MSE$
- $e_i$ 's are **not** independent random variables.
  - If large  $n$ , the dependence of  $e_i$  is relatively unimportant and can be ignored

##### Standardized vs Studentized

- Standardized =  $\frac{Y_i - \bar{y}}{\sigma}$
- Studentized =  $\frac{\frac{Y_i - \mu}{\sigma}}{\frac{1}{n}}$

$$\text{Semi-studentized Residuals } e_i^* = \frac{e_i - \bar{e}}{\sqrt{MSE}} = \frac{e_i}{\sqrt{MSE}}$$

#### 3.2.2 Residual Plots

##### Residual Plot Form

1. Tests
  - (a) Non-linearity of regression function A pattern indicates linear regression not appropriate

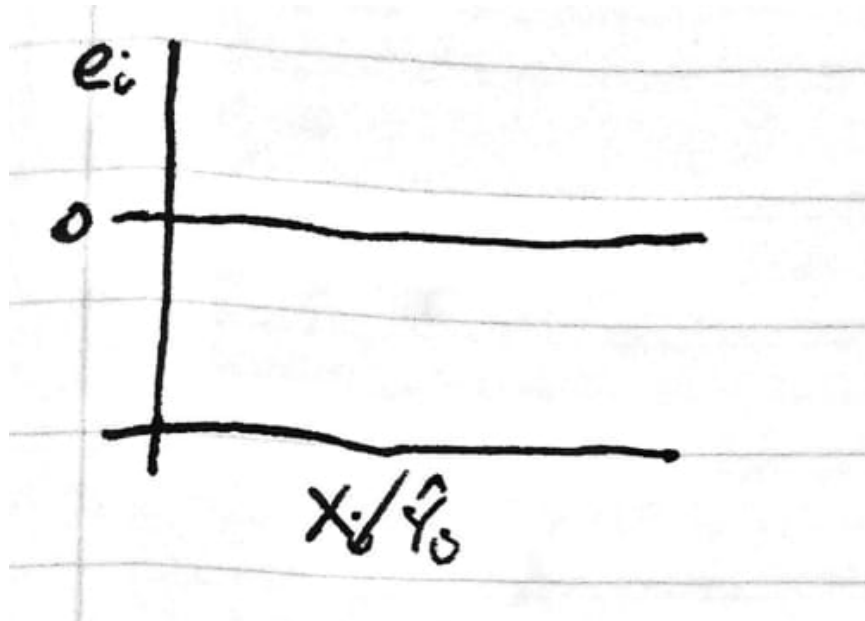


Figure 1: Empty Residual Plot

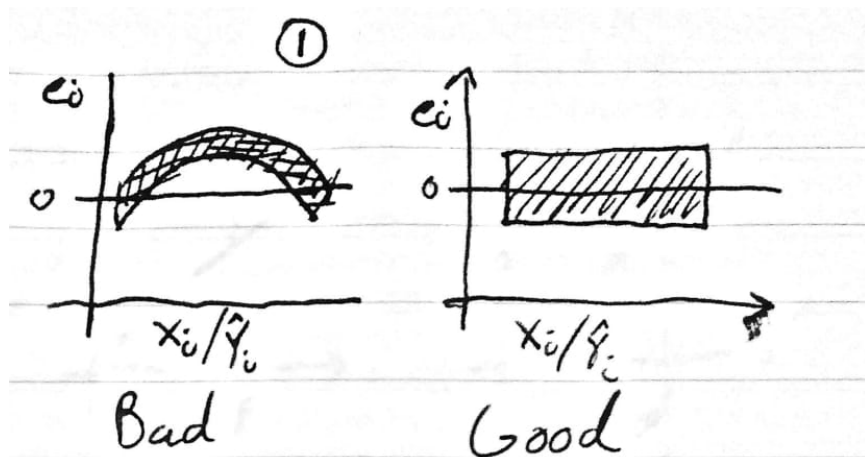


Figure 2: Plots 1



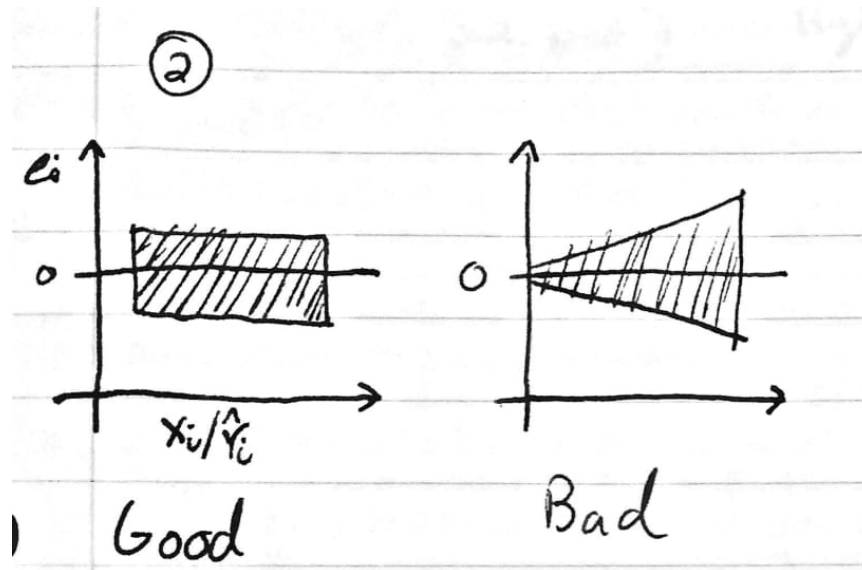


Figure 3: Plots 2

- (a) Non-constancy of error terms Fanning indicates different variances for different values of  $X_i$  or  $\hat{Y}_i$
- (b) Presence of outliers Graph Semi-studentized residuals on a Residual plot **OR** a Box Plot  
if  $|e_i^*| \geq 4$ , outlier
- (c) Non-independence of error terms (more of a concern with time-series) No pattern is good. Error terms safe to assume independent.
- (d) Normality of Error Terms
  - Use a normal probability plot. The closer the points the fall on a straight line, the closer they are to a normal distribution.
- (e) Omission of Important Predictors? A Pattern indicates that there might be a relationship between the residuals and some other predictor. This can be used to determine whether a predictor should be used before modeling it. Probably not as necessary anymore since it is easy to run and compare models.

### 3.2.3 Test of Randomness

1. Durbin-Watson Test

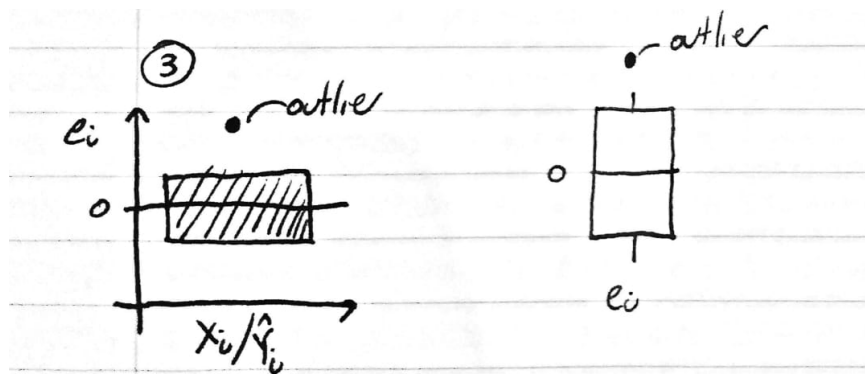


Figure 4: Plots 3

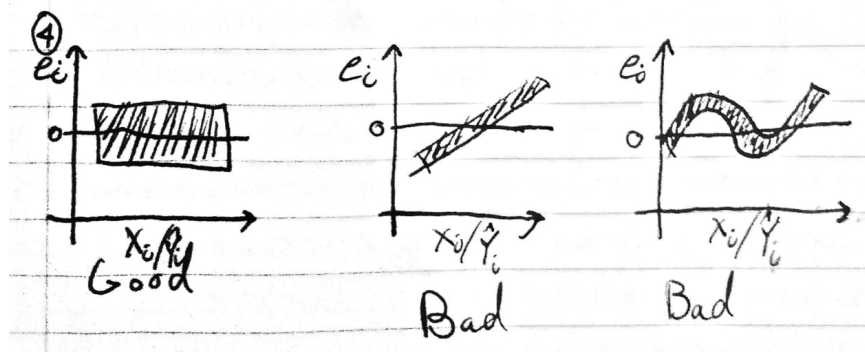


Figure 5: Plots 4

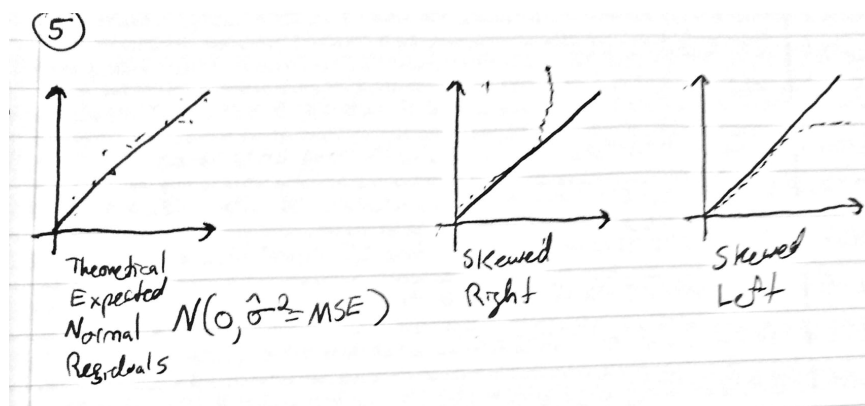


Figure 6: Plots 5

$$\begin{aligned} H_0 : \phi &= 0 \text{ where } \phi \text{ is an autocorrelation coefficient} \\ H_A : \phi &> 0 \text{ most assume positive correlation} \end{aligned} \quad (7)$$

`lmtest::dwtest(modle)`

2. Shapiro-Wilk Test for Normality Not writing much here because I know it already

`shapiro.test()`

### 3.2.4 Constant Variance

1. Brown-Forsyth Test Robust since it uses Median

`lawstat::levene.test()`

2. Breusch-Pagan Test Sensitive to departures from Normality

$$\log(\sigma^2) = \gamma_0 + \gamma_1 x_i$$

$$\begin{aligned} H_0 : \gamma_1 &= 0 \\ H_A : \gamma_1 &\neq 0 \end{aligned} \quad (8)$$

`lmtest::bptest()`

**NOTES:** Heteroscedascity means non-constant variance

## 4 Session 4 - Transformations & Inference (2019/10/02)

### 4.1 Transformations

If non-normality and unequal error variance:

1. Transform Y:  $Y' = f(Y)$
2. Transform X:  $X' = f(X)$

If non-linearity (rarer)

1. Transform X:  $X' = f(X)$

In order to determine which transformation to choose, look at the raw data and make a judgement call.

In Class Example

$$Y'_i = \log(Y_i) = \beta_0 + \beta_1 X_i + \epsilon_i \equiv Y_i = \exp(\beta_0 + \beta_1 X_i + \epsilon_i)$$

A 1 unit increase in X is associated with a  $\exp(\beta_1)$  multiplicative effect on the **geometric** mean. This link explains in detail the impact of log transformed variables.

$$\text{Geometric mean} = (\prod x_i)^{\frac{1}{n}}$$

$$\hat{Y}_i = \log(Y)$$

$$X'_i = \sqrt{x}$$

$$\hat{Y}_i = 4.896 + 4.325X'_i \rightarrow \exp(4.235) = 75.528$$

For each 1 unit increase in  $X'$ , the estimated increase in the geometric mean price is 75.53 times its previous value.

#### 4.1.1 Box-Cox Transformations

There is a value  $\lambda$  that is the optimal transformation to the response for equal variance and normality. It is optimal in the sense that it finds the value of  $\lambda$  which produces the smallest SSE for  $Y_i$ .

$$Y_i^\lambda = \beta_0 + \beta_1 X_i + \epsilon_i \text{ where } \epsilon_i \sim \text{iid } N(0, \sigma^2)$$

$\lambda$	2	0.5	0	-0.5	-1
$Y' = Y^\lambda$	$Y^2$	$\sqrt{Y}$	$\log(Y)$	$\frac{1}{\sqrt{Y}}$	$\frac{1}{Y}$

```
lindia::gg-boxcox(model)
```

## 4.2 Simultaneous Inference

**Goal:** Try to estimate more than one mean response at a time.

$$(0.95)^3 = 0.857375$$

### 4.2.1 Working-Hotelling Procedure

Based on the confidence band for the regression line.

100(1 -  $\alpha$ )% simultaneous confidence limits for g mean responses  $E(Y_h)$

$$Y_h \pm W \sqrt{MSE \left( \frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right)} \text{ where } W^2 = 2F_{1-\alpha, 2, n-2}$$

```
qf(1 - $alpha$, 2, n - 2)
```

#### 4.2.2 Bonferonni Procedure

$$Y_h \pm B \sqrt{MSE(\frac{1}{n} + \frac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2})} \text{ where } B = t_{1-\frac{\alpha}{2g}, n-2}$$

qt(1 - alpha / 2g, n - 2)

## 5 Session 5 - Prediction & Linear Algebra in Regression

### 5.1 Simultaneous Intervals

#### 5.1.1 Confidence

Using the Bonferonni adjustment, The simultaneous confidence interval for mean winning percentage for RunDiff of  $X_h = -100, 0, 100$  has a confidence level  $= 1 - \frac{\alpha}{g}$  where  $\alpha = .05$  and  $g = 3$

This is good for a smaller number of predictors. i.e.  $g < 10$

#### 5.1.2 Prediction

**Bonferroni:**  $\hat{Y}_h \pm t_{1-\frac{\alpha}{2g}, n-2} \sqrt{MSE(1 + \frac{1}{n} + \frac{(x_h - \bar{x})^2}{\sum (x_i - \bar{x})^2})}$  level  $= 1 - \frac{\alpha}{g}$

**Scheffe:**  $\hat{Y}_h \pm S \sqrt{MSE(1 + \frac{1}{n} + \frac{(x_h - \bar{x})^2}{\sum (x_i - \bar{x})^2})}$  where  $S = \sqrt{gF_{1-\alpha, g, n-2}}$

Scheffe is more efficient with a larger  $g$  (i.e.  $g > 10$ ). An in-class example showed that this was not the case so the jury is still out.

### 5.2 Inverse Prediction (“Calibration”)

First, construct a model where  $Y = X$

Goal: Make a prediction of  $X$  that was used to predict a new value of  $Y$ .

$$\begin{aligned} \hat{Y}_i &= \beta_0 + \beta_1 X_i + \epsilon_i \text{ where } \epsilon_i \sim \text{iid } N(0, \sigma^2) \\ \hat{Y} &= b_0 + b_1 x \end{aligned} \tag{9}$$

We are given  $Y_{h(new)}$ , so what is  $X_{h(new)}$ ?

$$X_{h(new)} = \frac{Y_{h(new)} - b_0}{b_1}$$

$$X_{h(new)} \pm t_{1-\frac{\alpha}{2}, n-2} \sqrt{\frac{MSE}{b_1^2} (1 + \frac{1}{n} + \frac{(x_{h(new)} - \bar{x})^2}{\sum (x_i - \bar{x})^2})}$$

investr::calibrate(model, Y, interval = "Wald")

The approximate confidence interval is appropriate if the following quantity is small (i.e.  $< .1$ ):

$$\frac{t_{1-\frac{\alpha}{2}, n-2}^2 MSE}{b_1^2 \Sigma (X_i - \bar{X})^2}$$

## 5.3 Linear Algebra in Regression

### 5.3.1 Review

$$\begin{matrix} \vec{Y} \\ (n \times 1) \end{matrix} = \begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{bmatrix}$$

$$\begin{matrix} \vec{Y}^T \\ (1 \times n) \end{matrix} = [Y_1 \quad \dots \quad Y_n]$$

**Design Matrix**

$$\begin{matrix} X \\ (n \times 2) \end{matrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix}$$

$$\begin{matrix} x^T \\ (2 \times n) \end{matrix} = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}$$

#### 1. Matrix Addition & Subtraction

$$Y_i = E(Y_i) + \epsilon_i$$

$$\vec{Y} = E(\vec{Y}) + \vec{\epsilon}$$

$$E(\vec{Y}) = \begin{bmatrix} E(Y_1) \\ \dots \\ E(Y_n) \end{bmatrix} \tag{10}$$

$$\vec{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \dots \\ \epsilon_n \end{bmatrix}$$

## 2. Matrix Multiplication

$$\begin{aligned}
\vec{Y}^T \vec{Y} &= [Y_1 \quad \dots \quad Y_n] \begin{bmatrix} Y_1 \\ \dots \\ Y_n \end{bmatrix} = \sum_1^n Y_i^2 \\
X^T X &= \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} n & \Sigma X_i \\ \Sigma X_i & \Sigma X_i^2 \end{bmatrix} \\
X^T \vec{Y} &= \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} Y_1 \\ \dots \\ Y_n \end{bmatrix} = \begin{bmatrix} \Sigma Y_i \\ \Sigma X_i Y_i \end{bmatrix}
\end{aligned} \tag{11}$$

3. Special Matrices **Symmetric**:  $A = A^T$  This implies a square matrix.  
i.e.  $n \times n$

**Diagonal:**  $A_{(n \times n)} = \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ 0 & a_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & a_{nn} \end{bmatrix}$

**Identity Matrix:**  $I_{(n \times n)} = \begin{bmatrix} 1 & \dots & 0 \\ \dots & 1 & \dots \\ 0 & \dots & 1 \end{bmatrix}$

**Scalar:**  $gI = \begin{bmatrix} g & \dots & 0 \\ \dots & g & \dots \\ 0 & \dots & g \end{bmatrix}$  where  $g$  is a scalar value

**One vectors**

$$\begin{aligned}
\vec{1}_{(n \times 1)} &= \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix} \\
J_{(n \times n)} &= \begin{bmatrix} 1 & \dots & 1 \\ \dots & 1 & \dots \\ 1 & \dots & 1 \end{bmatrix} \\
\vec{1}^T \vec{1}_{(1 \times n)(n \times 1)} &= n \\
\vec{1} \vec{1}^T_{(n \times 1)(1 \times n)} &= \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} = J
\end{aligned} \tag{12}$$

#### 4. Inverse of a Matrix

$$\begin{aligned} A_{(2 \times 2)} &= \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\ A_{(2 \times 2)}^{-1} &= \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \end{aligned} \quad (13)$$

#### Application to Regression

$$\begin{aligned} (X^T X)_{2 \times 2}^{-1} &= \frac{1}{\det(X^T X)} \begin{bmatrix} \Sigma x_i^2 & -\Sigma x_i \\ -\Sigma x_i & n \end{bmatrix} = \dots = \begin{bmatrix} \frac{\Sigma x_i^2}{n \Sigma (x_i - \bar{x})^2} & -\frac{\Sigma x_i}{n \Sigma (x_i - \bar{x})^2} \\ -\frac{\Sigma x_i}{n \Sigma (x_i - \bar{x})^2} & \frac{n}{n \Sigma (x_i - \bar{x})^2} \end{bmatrix} \\ \det(X^T X) &= n \Sigma x_i^2 - (\Sigma x_i)^2 \\ &= n \Sigma x_i^2 - \frac{n(\Sigma x_i)^2}{n^2} \\ &= n \left[ \Sigma x_i^2 - \frac{(\Sigma x_i)^2}{n} \right] \\ &= n \Sigma (x_i - \bar{x})^2 \end{aligned} \quad (14)$$

#### Side Note

$$\begin{aligned} \Sigma x_i &= n\bar{x} \\ \Sigma (x_i - \bar{x})^2 &= \Sigma x_i^2 - n\bar{x}^2 \\ \Sigma x_i^2 &= \Sigma (x_i - \bar{x})^2 + n\bar{x}^2 \\ (X^T X)^{-1} &= \begin{bmatrix} \frac{1}{n} + \frac{\bar{x}^2}{\Sigma (x_i - \bar{x})^2} & -\frac{\bar{x}}{\Sigma (x_i - \bar{x})^2} \\ -\frac{\bar{x}}{\Sigma (x_i - \bar{x})^2} & \frac{1}{\Sigma (x_i - \bar{x})^2} \end{bmatrix} \end{aligned} \quad (15)$$



## 5. Matrix Rules

$$\begin{aligned}
A + B &= B + A \\
(A + B) + C &= A + (B + C) \\
(AB)C &= A(BC) \\
C(A + B) &= CA + CB \\
(A^T)^T &= A \\
(A + B)^T &= A^T + B^T \\
(AB)^T &= B^T A^T \\
(AB)^{-1} &= B^{-1} A^{-1} \\
(A^{-1})^{-1} &= A \\
(A^T)^{-1} &= (A^{-1})^T
\end{aligned} \tag{16}$$

### 5.3.2 Expectations

$$\begin{aligned}
\vec{Y}_{(n \times 1)} &= \begin{bmatrix} Y_1 \\ \dots \\ Y_n \end{bmatrix} \\
E(\vec{Y})_{(n \times 1)} &= \begin{bmatrix} E(Y_1) \\ \dots \\ E(Y_n) \end{bmatrix} \\
\vec{\epsilon} &= \begin{bmatrix} \epsilon_1 \\ \dots \\ \epsilon_n \end{bmatrix} \\
E(\vec{\epsilon}) &= \vec{0}
\end{aligned} \tag{17}$$

### 5.3.3 Variance-Covariance Matrix

$$\sigma^2(\vec{Y}) = \begin{bmatrix} Var(Y_1) & \dots & Cov(Y_1, Y_n) \\ \dots & \dots & \dots \\ Cov(Y_n, Y_1) & \dots & Var(Y_n) \end{bmatrix} \tag{18}$$

When  $Y_i$  independent, the off diagonals are 0 meaning  $\sigma^2(\vec{Y}) = \sigma^2 I$   
**Aside**

$$\begin{aligned} Var(Y) &= E[(Y - E(Y))^2] \\ \sigma^2(\vec{Y}) &= E[(\vec{Y} - E(\vec{Y}))(\vec{Y} - E(\vec{Y}))^T] \end{aligned} \quad (19)$$

Let  $\vec{W} = \begin{matrix} (p \times 1) \\ (p \times n)(n \times 1) \end{matrix} A\vec{Y}$  where A is a matrix of **constants** and Y is a **random vector**

$$\begin{aligned} E(A) &= A \\ E(\vec{W}) &= AE(\vec{Y}) \\ \sigma^2(\vec{W}) &= E[(\vec{W} - E(\vec{W}))(\vec{W} - E(\vec{W}))^T] \\ &= E[(A\vec{Y} - AE(\vec{Y}))(A\vec{Y} - AE(\vec{Y}))^T] \\ &= E[A(\vec{Y} - E(\vec{Y}))(A(\vec{Y} - E(\vec{Y})))^T] \\ &= E[A(\vec{Y} - E(\vec{Y}))(\vec{Y} - E(\vec{Y}))^T A^T] \\ &= AE[(\vec{Y} - E(\vec{Y}))(\vec{Y} - E(\vec{Y}))^T]A^T \\ &= A\sigma^2(\vec{Y})A^T \end{aligned} \quad (20)$$

#### 5.3.4 Multivariate Normal Distribution

$$\begin{aligned} \vec{Y} &= \begin{bmatrix} Y_1 \\ \dots \\ Y_p \end{bmatrix} \\ \vec{\mu} &= \begin{bmatrix} \mu_1 \\ \dots \\ \mu_p \end{bmatrix} \end{aligned} \quad (21)$$

$\Sigma =$  Variance-Covariance Matrix  
( $p \times p$ )

$$f(\vec{Y}) = \frac{1}{(2\pi)^{\frac{p}{2}} \sqrt{\det(\Sigma)}} \exp(-\frac{1}{2}(\vec{Y} - \vec{\mu})^T \Sigma^{-1}(\vec{Y} - \vec{\mu}))$$

If  $Y_1, \dots, Y_p$  are jointly normally distributed (i.e in the multivariate normal distr.), then  $Y_k \sim N(\mu_k, \sigma_k^2)$  where  $k = [1, p]$

Recall the Linear Regression equation  $Y_i \beta_0 + \beta_1 X_i + \epsilon_i$  where  $\epsilon_i \sim N(0, \sigma^2)$ .

$$\vec{Y} = \begin{matrix} (n \times 1) \\ (n \times 2)(2 \times 1) \end{matrix} X\vec{\beta} + \vec{\epsilon} \text{ where } \begin{matrix} (n \times 1) \\ (n \times 1) \end{matrix} \vec{\epsilon} \sim N_n(\vec{0}, \sigma^2 I)$$

$N_n$  is a dimensions of a multivariate normal

$$\vec{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$E(\vec{Y}) = X\vec{\beta}$$

### 5.3.5 Least Squares Estimation

#### Normal Equations from Week 2

$$\begin{aligned}nb_o + b_1 \Sigma x_i &= \Sigma Y_i \\ b_0 \Sigma X_i + b_1 \Sigma X_i^2 &= \Sigma X_i Y_i\end{aligned}\tag{22}$$

$$X^T X \vec{b} = X^T \vec{Y}$$

So?

**Least Squares Estimator:**  $\vec{b} = (X^T X)^{-1} X^T \vec{Y}$

$$\underset{(n \times 1)}{\vec{Y}} = \underset{(n \times 2)}{X} \underset{(2 \times 1)}{\vec{b}} = (X^T X)^{-1} X^T \vec{Y}$$

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

The Hat Matrix is important for computing diagnostics for the model such as Cook's Distance.

#### Properties

- symmetric ( $H^T = H$ )
- Idempotent ( $HH = H$ )

2. Residuals  $E_i = Y_i - \hat{Y}_i \rightarrow \vec{Y} - \hat{\vec{Y}} = \vec{Y} - X\vec{b} = \vec{Y} - H\vec{Y} = (I - H)\vec{Y}$   
 $\sigma^2(\vec{e}) = \sigma^2(I - H)$

This is estimated by:  $MSE(I - H)$

## 6 Session 6 - Sums of Squares and Multiple Linear Regression

### 6.1 Sum of Squares

$$\underset{(1 \times n)(n \times 1)}{\vec{Y}^T \vec{Y}} = \Sigma Y_i^2\tag{23}$$

**Quadratic Form:** Contains squares of observations **and** their cross products. These are known as second-degree polynomials.

Quadratic forms scaled by  $\sigma^2$  allow us to treat the random variable Y as an observation of  $\chi_{n-1}^2$  distribution.

This is unlike  $\sigma^2(A\vec{Y}) = A\sigma^2(\vec{Y})A^T$  since that is squaring a matrix of **constants** whereas  $\vec{Y}^T \vec{Y}$  squares a matrix of **random variables** i.e. Y

### 6.1.1 SSE

$$\begin{aligned}
 SSE &= \sum e_i^2 \\
 &= \vec{e}^T \vec{e} \\
 &= \vec{Y}^T (I - H) \vec{Y}
 \end{aligned} \tag{24}$$

### 6.1.2 SSTo

$$\begin{aligned}
 SSTo &= \sum (Y_i - \bar{Y})^2 \\
 &= \sum Y_i^2 - \frac{(\sum Y_i)^2}{n} \\
 &= \vec{Y}^T (I - \frac{1}{n} J) \vec{Y}
 \end{aligned} \tag{25}$$

### 6.1.3 SSR

$$\begin{aligned}
 SSR &= \sum (\hat{Y}_i - \bar{Y})^2 \\
 &= \vec{Y}^T (H - \frac{1}{n} J) \vec{Y}
 \end{aligned} \tag{26}$$

## 6.2 Mean Estimates $\sigma^2$

### 6.2.1 Mean Responses

$$\hat{Y}_h = b_0 + b_1 X_h$$

so? we would like

$$\underset{(1 \times 1)}{\hat{Y}_h} = [1 \quad X_h] \underset{(1 \times 1)}{\vec{b}} \tag{27}$$

$$\text{Let } \vec{X}_h = \begin{bmatrix} 1 \\ X_h \end{bmatrix}$$

$$\text{Then, } \hat{Y}_h = \vec{X}_h^T \vec{b}$$

This is an estimate of the mean response!

### 6.3 Variance of $\hat{Y}_h$

$$\begin{aligned}
 \underset{(1 \times 1)}{Var(\hat{Y}_h)} &= Var(\vec{X}_h^T \vec{b}) \\
 &= \vec{X}_h^T Var(\vec{b}) \vec{X}_h \\
 &= \vec{X}_h^T \sigma^2 (X^T X)^{-1} \vec{X}_h \\
 &= \sigma^2 \underset{(1 \times 2)(2 \times 2)(2 \times 1)}{X_h^T (X^T X)^{-1} X_h}
 \end{aligned} \tag{28}$$

## 6.4 Multiple Regression Models

$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_{p-1} X_{i,p-1} + \epsilon_i$  where  $\epsilon_i \sim iidN(0, \sigma^2)$

$$E(Y_i) = \beta_0 + \beta_1 X_{i1} + \dots + \beta_{p-1} X_{i,p-1}$$

$$Y_i \sim indepN(E(Y_i), \sigma^2).$$

The parameters of this model are  $\{\beta_0, \dots, \beta_p\}$ . Thus there are **p** regression coefficients.

### 6.4.1 Interpretation

Using the model,  $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2}$

let's interpret the coefficients.

$\beta_0$ : The mean response of Y when  $X_1 = 0, X_2 = 0$

$\beta_1$ : For a fixed value of  $X_2$ , the associated increase in mean response in Y is  $\beta_1$  for every 1 unit increase in  $X_1$ . **This is known as a partial effect**

$\beta_2$ : For a fixed value of  $X_1$ , the associated increase in mean response in Y is  $\beta_2$  for every 1 unit increase in  $X_2$ .

$\beta_k$ : Associated change in mean response of Y for every 1 unit increase in  $X_k$ , given all other predictors are held constant.

### 6.4.2 Aside: Multi-Collinearity

**Multicollinearity** occurs when two or more predictors are highly correlated.

- Standard Errors blow up which makes test statistic small, which makes p-values high. This affects the ability for us to make **inferences**
- Multicollinearity is acceptable when using models for **prediction** but not when using them for **inference**.

### 6.4.3 Matrix Notation

$$\begin{aligned} \vec{Y} &= X\vec{\beta} + \vec{\epsilon} \\ \begin{matrix} (n \times 1) & (n \times p)(p \times 1) & (n \times 1) \end{matrix} & \end{aligned} \tag{29}$$
$$\begin{matrix} Var(\vec{\epsilon}) = \sigma^2 I \\ (n \times n) \end{matrix}$$

1. Fitted Values  $\hat{Y}_i = b_0 + b_1 X_{i1} + \dots + b_{p-1} X_{i,p-1}$

**Residuals:**  $e_i = Y_i - \hat{Y}_i$

2. Least Squares Estimators

$$\begin{matrix} \vec{b} \\ (p \times 1) \end{matrix} = \begin{matrix} (X^T X)^{-1} & X^T \vec{Y} \\ (p \times n)(n \times p) & (p \times n)(n \times 1) \end{matrix}$$

#### 6.4.4 ANOVA Table

Source	SS	DF	MS	F	p-value
Regression	$SSR = \Sigma(\hat{Y}_i - \bar{Y})^2$	p - 1	$MSR = \frac{SSR}{p-1}$	$F^* = \frac{MSR}{MSE}$	$P(F_{p-1, n-p} \geq F^*)$
Error	$SSE = \Sigma(Y_i - \hat{Y}_i)^2$	n - p	$MSE = \frac{SSE}{n-p}$		
Total	$SSto = \Sigma(Y_i - \bar{Y})^2$	n - 1			

#### 6.4.5 Omnibus F-Test for Regression Relation

$$\begin{aligned} H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0 \\ H_A : \text{at least one } \beta_k \neq 0 \end{aligned} \quad (30)$$

Test statistic:  $F^* = \frac{MSR}{MSE}$ . If  $H_0$  is true,  $F^* \sim F_{p-1, n-p}$

#### 6.4.6 Coefficient of Multiple Determination

$$R^2 = 1 - \frac{SSE}{SSto}$$

The issue with  $R^2$  is that it increases with the number of predictors **irrespective** of the predictor improving the model.

$$R^2_{adj} = 1 - \frac{\frac{SSE}{n-p}}{\frac{SSto}{n-1}}$$

#### 6.4.7 Coefficient of Multiple Correlation

$$R = \sqrt{R^2}$$

#### 6.4.8 Inferences in $\beta_k$

$$\begin{aligned} H_0 : \beta_k = 0 \\ H_A : \beta_k \neq 0 \end{aligned} \quad (31)$$

**Test Statistic:**  $t^* = \frac{b_k}{SE_{b_k}}$   
If  $H_0$  is true, then  $t^* \sim t_{n-p}$   
p-value =  $2P(t_{n-p} \geq |t|)$

2 \* (1 - pt(abs(t.star), n - p))

$$100(1 - \alpha)SE_{b_k}$$

## 7 Session 7 - Multiple Regression & Qualitative\Quantitative Predictors

### 7.1 Multiple Regression

#### 7.1.1 Extra Sums of Squares

**Def:** The marginal reduction in SSE when one or several predictors are added to the regression model, **given** other predictors are already in the model.

$$SSR(X_2|X_1) = SSE(X_1) - SSE(X_1, X_2)$$

$$SSR(X_2|X_1) = SSR(X_1, X_2) - SSR(X_1)$$

These are equivalent because any reduction in SSE implies an increase in SSR per the ANOVA definition:  $SSTo = SSR + SSE$

#### 1. Multiple Predictors

$$SSR(X_3|X_1, X_2) = SSE(X_1, X_2) - SSE(X_1, X_2, X_3)$$

$$SSR(X_3|X_1, X_2) = SSR(X_1, X_2, X_3) - SSR(X_1, X_2)$$

$$SSR(X_1, X_2) = SSR(X_2) + SSR(X_1|X_2)$$

Source	SS	df	MSE
Regression	$SSR(X_1, X_2, X_3)$	3	$MSR(X_1, X_2, X_3)$
$X_1$	$SSR(X_1)$	1	$MSR(X_1)$
$X_2 X_1$	$SSR(X_2 X_1)$	1	$MSR(X_2 X_3)$
$X_3 X_1, X_2$	$SSR(X_3 X_1, X_2)$	1	$MSR(X_3 X_1, X_2)$
Error	$SSE(X_1, X_2, X_3)$	n - 4	$MSE(X_1, X_2, X_3)$
Total	SSTo	n - 1	

#### 2. Hypothesis Test - $\beta_k = 0$ $H_0 : \mu_k = 0$

$$H_A : \mu_k \neq 0$$

This is the  $\mu_k X_k$  dropped from the model.

**Test Statistic:**  $t^* = \frac{b_k}{SE_{b_k}} \quad df = n - p$

(a) Full model  $Y_i = \mu_0 + \mu_1 X_1 + \dots + \mu_{p-1} X_{i,p-1} + \epsilon_i$

“p - 1” predictor variables

$$SSE(F) = SSE(X, \dots, X_{p-1})$$

(b) Reduced Model  $Y_i = \mu_0 + \mu_1 X_1 + \dots + \mu_{p-2} X_{i,p-2} + \epsilon_i$

“p - 2” predictor variables

$$SSE(R) = SSR(X, \dots, X_{k-1}, X_{p-1})$$

$$F^* = \frac{SSE(R) - SSE(F)}{\frac{df_R - df_F}{SSE(F)}} = \frac{\frac{SSE(X_1, \dots, X_{k-1}, X_k, \dots, X_{p-1}) - SSE(X, \dots, X_{p-1})}{n - (p-1) - (n-p)}}{\frac{SSE(X, \dots, X_{p-1})}{n-p}}$$

3. Hypothesis Test -  $\beta_0 = \dots = \beta_k = 0$

(a) Reduced Model  $Y_i = \mu_0 + \mu_1 X_{i1} + \dots + \mu_{k-1} X_{i,k-1} + \mu_k X_{ik} + \dots + \mu_{p-1} X_{i,p-1} + \epsilon_i$

“p - g - 1” predictors **OR** “p - g” regression coefficients

$$F^* = \frac{\frac{SSE(X_1, \dots, X_{k-1}, X_k, \dots, X_{p-1}) - SSE(X, \dots, X_{p-1})}{n - (p-g) - (n-p)}}{\frac{SSE(X, \dots, X_{p-1})}{n-p}} = \frac{SSR(X_k, \dots, X_{k+(g-1)} | X_k, \dots, X_{k+g}, \dots, X_{p-1})}{MSE(X_1, \dots, X_{p-1})}$$

If  $H_0$  is true,  $F^* \sim F_{g, n-p}$

4.  $R^2$ : Coefficient of multiple determination

- proportion of variation in Y explained by the regression of Y on  $X_1, \dots, X_{p-1}$

Ex

$$Y_i = \mu_0 + \mu_1 X_{i,1} + \mu_2 X_{i,2}$$

$SSE(X_2)$ : variation when only  $X_2$  is in the model.

$SSE(X_1, X_2)$ : variation when both  $X_1, X_2$  are in the model.

Marginal reduction in variation when  $X_1$  is added to the model?

$$\frac{SSE(X_2) - SSE(X_1, X_2)}{SSE(X_2)}$$

$$R^2_{Y_1/Y_2} = \frac{SSR(X_1|X_2)}{SSE(X_2)}$$

$$R^2_{Y_2/Y_1} = \frac{SSR(X_2|X_1)}{SSE(X_1)}$$

3 predictors

$$R^2_{Y_{3|2,1}} = \frac{SSR(X_3|X_1, X_2)}{SSE(X_1, X_2)}$$

Recipe for correlation coefficient:

- Take sqrt of partial  $R^2$
- Sign of partial correlation = sign of correlation corresponding coefficient

## 7.2 Multi-collinearity

Predictors that are highly correlated with each other. **10 N values per predictor**



### 7.2.1 Effects

1. There is no unique sum of squares that can be assigned to the predictor variable
2. May inflate standard error of  $b_k$  least square error.

It does not greatly impact the value of predictions.  
 $ETA^2$  tells  $R^2$  given the previously given variable  $R^2$

### 7.3 Polynomial Regression Models

- true curvilinear response
- true curvilinear response is unknown but a polynomial function provides a good approximation to the true function.

One prediction variable **and** second order:

$$Y_i = \mu_0 + \mu_1 X + \mu_2 X^2 + \epsilon_i \text{ where } X_i = x_i - \bar{x}$$

$$E(Y) = \mu_0 + \mu_1 X_1 + \mu_2 X^2$$

Two parameters **and** second order:

$$x_{i,1} = x_{i,1} - \bar{x}_1$$

$$X_{i,2} = x_{i,2} - \bar{x}_2$$

$$Y_i = \beta_0 \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,1}^2 + \beta_4 X_{i,2}^2 + \beta_5 X_{i,1} X_{i,2} + \epsilon_i$$

Strategy? Fit higher order models and compare to reduced models.

`summary(model)`

$$\hat{Y} = b_0 + b_1 x + b_2 x^2 \quad \hat{Y} = b'_0 + b'_1 x + b'_2 x^2$$

$$b'_0 = b_0 - b_1 \bar{x} - b_2 \bar{x}^2$$

$$b'_1 = b_1 - 2b_2 \bar{x}$$

$$b'_2 = b_2$$

Why do this? Solving a regression model with a non-linear  $E(Y_i)$

## 8 Session 8 - Interaction Models & Model Selection

### 8.1 Interaction Regression Models

- p: # of regression coefficients. i.e. parameters
- p - 1: predictor variables

### 8.1.1 Additive Effects

$$E(Y) = \sum_{i=1}^{p-1} f_K(x_k) \quad (32)$$

but  $E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2$  is **not** additive since  $X_1 X_2$  is an interaction term

Consider the following:

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon_i \text{ where } \epsilon_i \sim iidN(0, \sigma^2)$$

A one-unit **increase** in  $X_2$  for a fixed value of  $X_1$ , results in an associated change of  $\beta_1 + \beta_3 X_1$  units in mean response  $Y$ .

- $\beta_3 = 0 \Rightarrow$  additive model
- $\beta_3 > 0 \Rightarrow$  reinforcement or synergistic interaction\*
- $\beta_3 < 0 \Rightarrow$  interference or antagonistic interaction\*

\*if  $\beta_1$  and  $\beta_2$  are negative, these terms flip

parallel lines indicate **additive** terms, otherwise interactive

**Aside**

To avoid multicollinearity between predictors, center variables!

$$X_{ik} = X_{ik} - \bar{X}_k$$

Does Standardizing also help reduce multicollinearity? Yes, but makes interpretation more difficult. This is done in PCA and as I've seen, interpreting PCA can be hairy or a best guess.

- Try to identify possible interactions ahead of time prior to fitting the model.
- When looking at removing **one** term, the  $t$  statistic is sufficient to rule out a parameter.

### 8.1.2 Qualitative Predictors

Qualitative Predictor with two classes. i.e. two values This is sometimes called: Indicators, Binary, dummy variables

For representing  $C$  classes, use  $C - 1$  indicator variables.

Example

$$\text{Let } C = 4 \quad C = \begin{bmatrix} A \\ B \\ C \\ D \end{bmatrix}$$

$$\begin{aligned}
Y_i &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \epsilon_i \text{ where} \\
X_2 &= \begin{cases} 1, & A \\ 0, & \text{else} \end{cases} \\
X_3 &= \begin{cases} 1, & B \\ 0, & \text{else} \end{cases} \\
X_4 &= \begin{cases} 1, & C \\ 0, & \text{else} \end{cases}
\end{aligned} \tag{33}$$

if  $X_2 = X_3 = X_4 = 0$ , indicates effect of  $C = D$  on mean  $Y_i$

$$\begin{aligned}
A : E(Y) &= (\beta_0 + \beta_2) + \beta_1 X_i \\
B : E(Y) &= (\beta_0 + \beta_3) + \beta_1 X_i \\
C : E(Y) &= (\beta_0 + \beta_4) + \beta_1 X_i \\
D : E(Y) &= (\beta_0) + \beta_1 X_i
\end{aligned} \tag{34}$$

D is considered the **baseline category**

1 Qualitative variable and 1 Quantitative variable in the same model is known as **ancova**: Analysis of Covariance. ANCOVA assumes that each group has the same slope.

Interpret  $\beta_0$ : The diff in mean response of  $Y$  between  $A$  and  $D$  group for a given value of  $X_1$

**Estimate**  $\beta_3 - \beta_4$

1.  $b_3 - b_4$
2.  $Var(b_3 - b_4) = var(b_3) + var(b_4) - 2cov(b_3, b_4)$

If doing time series, one can use indicator variables to model time periods

## 8.2 Model and Variable Selection

### 8.2.1 Criterion for Model Selection

p: # of parameters (regression coefficients)

1.  $R_p^2$  or  $SSE_p$  criterion. Both indicate the same thing.

$$R_p^2 = 1 - \frac{SSE_p}{SST_o}$$

**Look for**

- *High*  $R_p^2$
- *Small*  $SSE_p$

2.  $R_{a,p}^2$  or  $MSE_p$  criterion

$$R_p^2 = 1 - \frac{\frac{SSE_p}{n-p}}{\frac{SST_o}{n-1}} = 1 - \frac{MSE_p}{\frac{SST_o}{n-1}}$$

**Look For:**

- *High*  $R_{a,p}^2$
- *Small*  $MSE_p$

3. Mallows'  $C_p$  Criterion

$$C_p = \frac{SSE_p}{MSE(X_1, \dots, X_{p-1})} - (n - 2p)$$

$MSE(X_1, \dots, X_{p-1})$ : MSE for the model with **all** potential predictors of interest.

For largest possible value of  $P$ ,  $C_p = p$

proof

$$\begin{aligned} MSE &= \frac{SSE}{n-p} \\ \frac{SSE_p}{\frac{SSE_p}{n-p}} &= n - p - (n - 2p) = p \end{aligned} \tag{35}$$

**Look for**

- *Small*  $C_p$  or  $C_P \leq p$ . This means the model has a small amount of bias.

Recall  $MSE(Y) = Bias^2(Y) + Var(Y)$

1.  $AIC_p$  or  $SBC_p$  Criterion

- $AIC_p$ : Akaike's Information Criterion:  $n \ln(SSE_p) - n \ln(n) + 2p - 2 \loglik + 2p$
- $SBC_p$ : Schwartz' Bayesian Information Criterion -  $n \ln(SSE_p) - n \ln(n) + p \ln(n) - 2 \loglik + 2 \ln(p)$

**Look for**

- *Small  $SSE_p$*
- *Small  $AIC_p$  and/or  $SBC_p$*

## 2. $PRESS_p$ Criterion

Prediction Sum of Squares

$$PRESS_P = \sum_1^n (Y_i - \hat{Y}_{i(i)})^2$$

$\hat{Y}_{i(i)}$

- Ignore the  $i$ th case
- Fit model on remaining  $n - 1$  cases
- Find Fitted value based on deleted  $i$ th case

This is **not** the same as bootstrapping, mostly because there is no resampling going on.

```
leaps::regsubsets(formula, data, method="exhaustive", nbest=30)
```

```
#+NAME: fortify_leaps
```

```
fortify.regsubsets <- function(model, data, ...){
```

```
  require(plyr)
```

```
  stopifnot(model$intercept)
```

```
  models <- summary(model)$which
```

```
  rownames(models) <- NULL
```

```
  model_stats <- as.data.frame(summary(model)[c("bic","cp","rss","rsq","adjr2")])
```

```
  dfs <- lapply(coef(model, 1:nrow(models)), function(x) as.data.frame(t(x)))
```

```
  model_coefs <- plyr::rbind.fill(dfs)
```

```
  model_coefs[is.na(model_coefs)] <- 0
```

```
  model_stats <- cbind(model_stats, model_coefs)
```

```
  # terms_short <- abbreviate(colnames(models))
```

```
  terms_short <- colnames(models)
```

```
  model_stats$model_words <- aapply(models, 1, function(row) paste(terms_short[1:
```

```
  model_stats$size <- rowSums(summary(model)$which)
```

```
  model_stats
```

```
}
```

```
get_model_coefs <- function(model){
```

```
  models <- summary(model)$which
```

```
  dfs <- lapply(coef(model, 1:nrow(models)), function(x) as.data.frame(t(x)))
```

```
  model_coefs <- plyr::rbind.fill(dfs)
```

```

    model_coefs[is.na(model_coefs)] <- 0
    model_coefs
  }

```

## 9 Session 9 - Model and Variable Selection & Assessing Diagnostics

### 9.1 Model and Variable Selection

#### 9.1.1 Automatic Search Procedures

1. Backward Selection

Full Model -> reduce parameters to “smallest” AIC

```
step(model.full, direction = "backward")
```

2. Forward Selection

Intercept-only model -> add parameters to “smallest” AIC

```
step(lm.null, scope = list(lower, upper), direction = "forward")
```

3. Step-wise

Intercept-only model -> add one -> subtract\add one for “smallest” AIC

```
step(lm.null, scope = list(upper), direction = "both")
```

#### 9.1.2 Model Validation

1. Collect new data to check the model and its predictive validity

$$MSPR = \frac{\sum_1^{n^*} (Y_i - \hat{Y}_i)^2}{n^*}$$

If MSPR approximately your Model’s MSE, then your model is not necessarily biased. If the difference is large, MSPR is a good indicator on how well it predicts.

#### Definitions

- MSPR: Mean Square Prediction Error

- $Y_i$ : Value of the response variable in the  $i$ th validation case
  - $\hat{Y}_i$ : Predicted value of the  $i$ th validation case using the model you previously built.
  - $n^*$ : number of cases in the validation dataset.
2. Compare results with theoretical expectations empirical results, and simulation results.
  3. Use a holdout sample to check the model and its predictive ability. This is standard practice for predictive models

## 9.2 Assessing Diagnostics

### 9.2.1 Added-variable Plots

Also known as:

- Partial Regression Plots
- Adjusted Variable Plots

These plots show:

- Marginal Importance of this variable in reducing residual variability
- May provide info about the nature of the marginal regression relation for predictor variable  $X_k$  under consideration

1. Example  $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$

**Goal:** What is  $X_i$ 's effect given that  $X_2$  is in the model?

$$\begin{aligned}\hat{Y}_i(X_2) &= b_0 + b_2 X_{i2} \\ e_i(Y|X_2) &= Y_i - \hat{Y}_i(X_2)\end{aligned}\tag{36}$$

*fitted values + residuals from the model with only  $X_2$*

$$\begin{aligned}\hat{X}_{i1}(X_2) &= b_0^* + b_2^* X_{i2} \\ e_i(X_1|X_2) &= X_{i1} - \hat{X}_{i1}(X_2)\end{aligned}\tag{37}$$

fitted value + residuals from the model with  $X_1$  as the response and  $X_2$  as the predictor.

## 2. Reading Plots

```
car::avPlots(model)
```

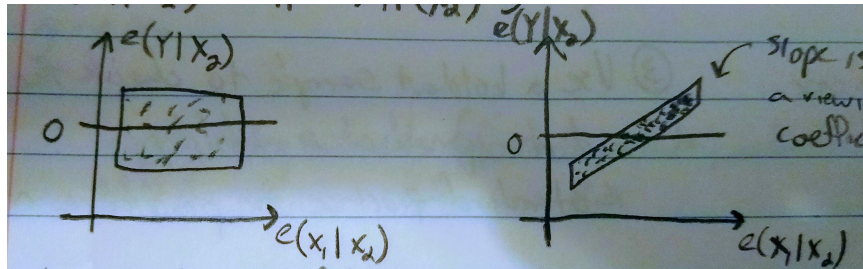


Figure 7: Partial Residuals vs Fitted Values

### (a) Partial Residuals $X_1$ vs Fitted Values

Notice the even distribution of residuals around  $y = 0$ .  $X_1$  provides no useful information given  $X_2$  is in the model.

### (b) Partial Residuals $X_2$ vs Fitted Values

Notice the pattern.  $X_1$  may be a good addition to the model given  $X_2$  is already in the model.

**Goal:** Identify outlying Y observations. i.e. which Y observations are influential on our own regression model?

- Residuals:  $e_i = Y_i - \hat{Y}_i$
- Semi-studentized Residuals:  $e^* = \frac{e_i}{\sqrt{MSE}}$
- Studentized Residuals:  $R_i = \frac{e_i}{\sqrt{MSE(1-h_{ii})}}$   
 $h_{ii}$ : the  $i$ th diagonal value from the hat matrix H

```
rstandard(model)
```

## 10 Session 10 - Outliers & Weighted Least Squares

### 10.1 Outliers

#### 10.1.1 Identifying Outlying Y Observations

- Use Studentized Deleted Residuals to identify *outlying Y Observations*



**Residuals:**  $e_i = Y_i - \hat{Y}_i$

**Semi-studentized Residuals:**  $e_i^* = \frac{e_i}{\sqrt{MSE}}$

**Studentized Residuals:**  $r_i = \frac{e_i}{\sqrt{MSE(1-h_{ii})}}$

$h_{ii}$ : Standard Error of  $e_i$ . aka Standard Error of the  $i$ th residual

**Deleted Residuals:**  $d_i = Y_i - \hat{Y}_{i(i)} = \frac{e_i}{1-h_{ii}}$

**Studentized Deleted Residuals** (rstudent):

$$\begin{aligned} t_i &= \frac{d_i}{SE_{d_i}} \\ &= \frac{e_i}{\sqrt{MSE_{(i)}(1-h_{ii})}} \\ &= e_i \sqrt{\frac{n-p-1}{SSE(1-h_{ii}) - e_i^2}} \end{aligned} \quad (38)$$

### 10.1.2 What is an Outlying Y Observation?

$$|t_i| > t_{1-\frac{\alpha}{2n}, n-p-1}$$

qt(1 - alpha / 2n, n - p - 1)

- The “- 1” is the residual that is being deleted

### 10.1.3 Identifying Outlying X Observations

- Use leverage values. i.e. “hat matrix leverage values”

$h_{ii}$ : leverage (in terms of X values)

1.  $0 \leq h_{ii} \leq 1, i = [1, n]$
2.  $\sum_1^n h_{ii} = p$  (number of parameters in the model)

Recall:  $Var(e_i) = MSE(1 - h_{ii})$

- The larger  $h_{ii}$ ,  $Var(e_i)$  **decreases**, thus making  $\hat{Y}_i$  close to  $Y_i$

How large is a large  $h_{ii}$ ?

- if  $h_{ii} > s\bar{h} = \frac{2p}{n}$ , the cases are outlying cases in terms of X.

## 10.2 Influential Cases

How influential are “new” cases?

$$h_{new,new} = X_{new,new}^T (X^T X)^{-1} X_{new,new}$$

If  $h_{new,new}$  is much larger than  $h_{ii}$ , there may be some extrapolation. There are no set guidelines for this.

### 10.2.1 Identifying Influential Cases

1. Influence of the  $i$ th case on a single fitted value, .
  - Use DFFITS (Difference of Fits)

$$\begin{aligned} DFFITS_i &= \frac{\hat{Y}_i - Y_{i(i)}}{\sqrt{MSE_{(i)} h_{ii}}} \\ &= e_i \sqrt{\frac{n-p-1}{SSE(1-h_{ii}) - e_i^2}} \sqrt{\frac{h_{ii}}{1-h_{ii}}} \\ &= t_i \sqrt{\frac{h_{ii}}{1-h_{ii}}} \end{aligned} \quad (39)$$

#### Notes

- $MSE_{(i)}$ : calculated with the  $i$ th case removed
- $t_i$ : Studentized Deleted Residuals

What is influential?

- Small - Med Dataset:  $|DFFITS_i| > 1$
- Large Dataset:  $|DFFITS_i| > 2\sqrt{\frac{p}{n}}$

2. Influence of the  $i$ th case on all fitted values

- Cooks Distance

$$\begin{aligned} D_i &= \frac{\sum_{j=1}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{pMSE} \\ &= \frac{e_i^2}{pMSE} \left[ \frac{h_{ii}}{(1-h_{ii})^2} \right] \end{aligned} \quad (40)$$

#### Notes

- $Y_{j(i)}$ : fitted value when the  $i$ th case is left out

What is an influential case? Compare  $D_i$  to  $F_{p,n-p}$

- If  $P(F_{p,n-p} \leq D_i) < 0.1, 0.2$ , the  $i$ th case has very little influence.
- If  $P(F_{p,n-p} \leq D_i) > 0.5$ , the  $i$ th case has major influence.

### 3. Influence of the $i$ th case on the regression coefficients

- DFBETAS

$$(DFBETAS)_{k(i)} = \frac{b_k - b_{k(i)}}{\sqrt{MSE_{(i)} C_{kk}}} \quad \text{Notes:}$$

- $C_{kk}$ : Diagonal term of  $(X^T X)^{-1}$
- $Var(\vec{b}) = \sigma^2 (X^T X)^{-1} = \sigma^2 C_{kk}$

## 10.3 Variance Inflation Factors

- used to assess Multicollinearity

$$VIF = \frac{1}{1 - R_k^2}$$

**Notes**

- $R_k^2$  is  $R^2$  from 'lm( $X_k \sim X_1 + \dots + X_{(k-1)} + X_{(k+1)} + \dots + X_{(p-1)}$ )'

\* This is a mishmash of math and R

$$\min VIF_k = 1 \quad \max VIF = \infty$$

- Sometimes (rarely) signs flip
- multicollinearity causes increase variance

### Interpretation

- $VIF > 4$ , mild/moderate multicollinearity
- $VIF > 10$ , severe multicollinearity
- Ideal?  $VIF$  close to 1

If experiencing high multicollinearity, check for correlation between response and each predictor.

## 10.4 Weighted Least Squares

- Good use if Variance is Unequal

Possible Weight:  $W_i = \frac{1}{\sigma_i^2}$

### 10.4.1 Iteratively Reweighted Least Squares

1. Fit regular least squares model and analyze results
2. Estimate the variance function or the standard deviation function by regressing  $e_i^2$  or  $|e_i|$  on the predictors.
3. Use the fitted values from the estimated  $Var(\hat{V}_i)$  or estimate std. dev ( $\hat{S}_i$ ) function to obtain weights  $w_i$ .
4. Estimate regression coefficients use the weights. So?
  - $e_i^2$  estimates  $\sigma_i^2$
  - $|e_i|$  estimates  $\sigma_i$

$W_i = \frac{1}{(\hat{S}_i)^2}$  using  $|e_i|$

OR

$W_i = \frac{1}{\hat{V}_i}$  using  $e_i^2$

## 11 Extra Curricular - Weighted Least Squares, Ridge, and Robust Regression

### 11.1 Weighted Least Squares

- Useful for models with heteroskedasticity (non-constant variance)

$$\begin{aligned}\vec{b} &= (X^T X)^{-1} X^T Y \\ \vec{b}_w &= (X^T W X)^{-1} X^T W Y \\ &_{(n \times n)}\end{aligned}\quad (41)$$
$$W = \begin{bmatrix} w_1 & 0 & \dots & 0 \\ 0 & w_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & w_n \end{bmatrix}$$

- OLS is a special case of WLS where  $W = J = 1$ .

- $w_i = k(\frac{1}{\sigma_i^2})$ . if error variances known (rare)
- $w_i = \frac{1}{(\hat{s}_i)^2}$ . if using fitted standard error
- $w_i = \frac{1}{(\hat{v}_i)}$ . if using fitted variance

Using the weights to estimate regression coefficient is called *Iteratively Reweighted Least Squares*. Typically done until coefficients have stabilized.

#### Notes

- $R^2$  does not have a clearcut meaning for WLS.

### 11.2 OLS with Heteroskedasticity

OLS can still be used with unequal error variances via White's Estimator. This leverages something called the *Robust Covariance Matrix*.

$$\begin{aligned}\sigma^2(b) &= (X^T X)^{-1} (X^T \sigma^2(e) X) (X^T X)^{-1} \\ S^2(b) &= (X^T X)^{-1} (X^T S_0 X) (X^T X)^{-1} \\ S_0 &= \begin{bmatrix} e_1^2 & 0 & \dots & 0 \\ 0 & e_2^2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & \dots & e_n^2 \end{bmatrix} \end{aligned} \quad (42)$$

$e_i$ : OLS estimator of the residuals squared.

### 11.3 Ridge Regression

- Useful for cases with severe Multicollinearity.

What to do when you have multicollinearity?

1. If only estimating and no conf intervals, nothing
2. Center predictor variables
3. Drop Predictors
  - Downside: some predictors not accounted for and there is a relationship affecting the response that is not being represented in the model.
4. Add cases that break multicollinearity.
5. PCA

Definition: Modifies OLS to allow biased estimators to lower variance.

Recall  $MSE = Var(Y) + (Bias(Y))^2$

$E(b^R - \beta)^2 = \sigma^2(b^R) + (E(b^R) - \beta)^2$  where  $b^R$ : biased estimator

Least Squares Normal Equations given by:  $r_{XX}b = r_{XY}$

$r_{XX}$ : correlation matrix of X variables

$r_{XY}$ : Vector of coefficients of simple correlation variables between Y and each X Variable

Ridge Standardized Regression:  $(r_{XX} + cI)b^R = r_{XY}$

$$b^R_{(p-1) \times 1} = \begin{bmatrix} b_1^R \\ \dots \\ b_{p-1}^R \end{bmatrix}$$

$$b^R = (r_{XX} + cI)^{-1}r_{XY}$$

A **biasing** constant  $c \geq 0$  can be chosen.

- bias increases as c increases. likewise variance decreases
- There is always some value  $c$  where  $b^R$  has a smaller MSE than OLS  $b$ .

– Optimal Values  $c$  varies by application and is unknown

**Ridge Trace**: Method often used to determine  $c$ . This is combined with

VIF

**look for**

- spots where the line smooths out
- where least change in  $b_k^r$  happens

finding  $c$  is a bit of an art.

this formula can be used to convert standardized coefficients to unstandardized coefficients.

$$b_k = \left(\frac{s_y}{s_k}\right)b_k^r$$

$s_y$ : standard dev of y  $s_k$ : standard error of  $b_k$

## 11.4 robust regression

- reduce influential cases

uses **iteratively reweighted least squares** where  $w_i$  dampens influential cases instead of heteroskedasticity.

$u$ : Scaled residual 0.345|4.685: tuning constants that are robust for 95% of normal data. Huber:

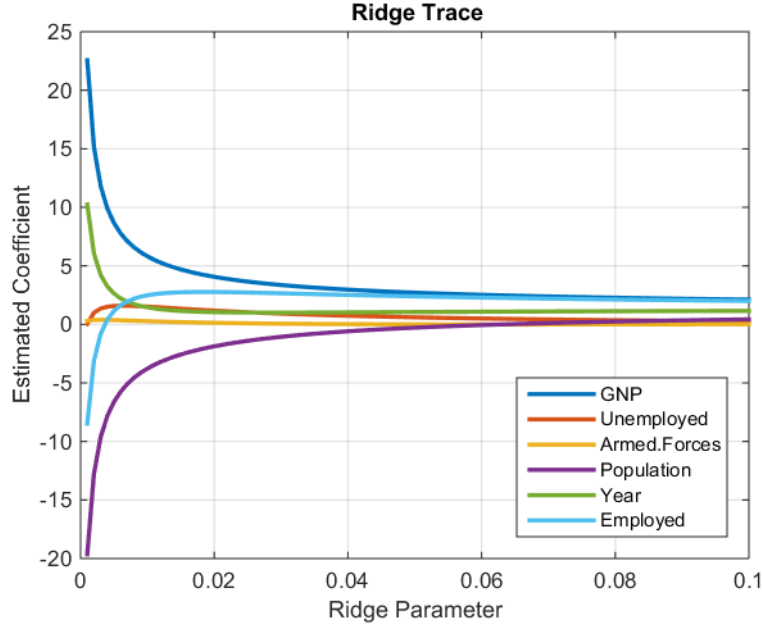


Figure 8: Ridge Trace Example

$$w = \begin{cases} 1 & |u| \leq 1.345 \\ \frac{1.345}{|u|} & |u| > 1.345 \end{cases} \quad (43)$$

Bisquare:

$$w = \begin{cases} [1 - (\frac{u}{4.685})^2]^2 & |u| \leq 4.685 \\ 0 & |u| > 4.685 \end{cases} \quad (44)$$

Huber is often used to obtain starting weights for Bisquare.

#### 11.4.1 $u$

- Semi-studentized residuals could be used but they are not resistant to outliers
- Mean Absolute Deviation (MAD) often used.

$$MAD = \frac{1}{0.6745} med(|e_i - med(e_i)|)$$

$$u_i = \frac{e_i}{MAD}$$

0.6745 is used to make this an unbiased estimate for  $\sigma$  from a normal distribution.

### 11.5 Regression Tree (Non-parametric Method)

- Split  $X$ 's into distinct regions  $r$  and run a regression on each region.
- “Growing a tree” is finding the number of regions  $r$  and the boundaries/split points between them.
- If the variance of the residuals in each region seem constant, splitting may not be necessary.
- The best split point minimizes  $SSE = \sum_{k=1}^r SSE(R_{rk})$
- Once the optimal  $r$  is chosen, each  $r$  is subdivided to find the most optimal SSE.
- The chosen number of regions is done through validation studies, such as choosing the tree that minimizes MSPR