

## Chapter 3: Simple Linear Regression

Regression Model:  $Y = B_0 + B_1X + e$

Regression Equation:  $y = b_0 + b_1X$

The regression equation is used to estimate the regression model. The "e" represents the random error term that always occurs when building a model.

### Hypothesis Testing

#### Steps for conducting a Hypothesis Test:

##### 1) State Hypothesis

Ho: This is what we assume to be true. This statement contains the = sign.

Ha: This is what we are trying to prove. This statement never contains the = sign.

##### 2) State Critical Value/Decision Rule

Use a) the alpha level and b) the degrees of freedom to look up the value used as a comparison point in the appropriate table. This value is the Critical Value.

##### 3) Calculate Test Statistic

Use the information in the problem to calculate the Test Statistic value for the test by using the appropriate formula.

##### 4) Calculate the p-value

Find the probability region beyond the calculated test value.

OR Find the Test Statistic in the table while looking only at the corresponding df row.

Identify the two values within which the Test Statistic falls in between and look up at the top of the table to identify the two corresponding alpha level. Express the p-value as a range.

##### 5) Conclude

Show how the Test Statistic compares to the Critical Value and draw conclusions.

#### There are three possible hypothesis testing questions for simple linear regression:

##### 1) Is the model significant/valid?

Use the T-test to test the slope of the regression equation.

Ho:  $B_1 = 0$

Ha:  $B_1 \neq 0$

CV =  $T^*(\alpha/2, df) = T^*(\alpha/2, n - 2)$  where "n" is the number of observations

TS =  $T = (b_1 - 0) / sb_1$

Test Stat:  $T_{stat} = \frac{b_1 - 0}{sb_1}$

$T_{critical} = T^*(\frac{\alpha}{2}, n-2)$

## 2) Is the relationship positive/negative?

Use the T-test to test if the slope of the regression equation is bigger/smaller than zero.

$$H_0: B_1 \leq 0 \quad \text{OR} \quad B_1 \geq 0$$

$$H_a: B_1 > 0 \quad \text{OR} \quad B_1 < 0$$

$CV = T^*(\alpha, df) = T^*(\alpha, n - 2)$  where "n" is the number of observations

$$TS = T = (b_1 - 0) / sb_1$$

$$T_{\text{test stat}} = \frac{b_1 - 0}{sb_1} \quad T^*(\alpha, n-2)$$

## 3) Can we claim that each X is more/less than A?

Use the T-test to test if the slope of the regression equation is bigger/smaller than A.

$$H_0: B_1 \leq A \quad \text{OR} \quad B_1 \geq A$$

$$H_a: B_1 > A \quad \text{OR} \quad B_1 < A$$

$CV = T^*(\alpha, df) = T^*(\alpha, n - 2)$  where "n" is the number of observations

$$TS = T = (b_1 - A) / sb_1$$

$$T_{\text{test stat}} = \frac{b_1 - A}{sb_1} \quad T^*(\alpha, n-2)$$

## Chapter 4: Multiple Linear Regression

Regression Model:  $Y = B_0 + B_1X_1 + B_2X_2 + (\dots) + e$

Regression Equation:  $y = b_0 + b_1X_1 + b_2X_2$

*Note: The regression equation is also called the least squares prediction equation*

### Hypothesis Testing

#### 1) Is the the model significant? (Test the overall significance/validity of the model)

Use the F-test and build the ANOVA table to calculate the F Statistic.

$$H_0: B_1 = B_2 = B_3 = \dots = B_k = 0$$

$$H_a: \text{At least one } B_i \neq 0$$

$CV = F^*(\alpha, d1, d2) = F^*(\alpha, k, n - k - 1)$  where k is the number of X variables

$$TS = F = (SSR / k) / (SSE / n - k - 1) \quad T_{\text{test stat}} = F = \frac{SSR/k}{SSE/(n-k-1)} \quad F^*(\alpha, k, n-k-1)$$

#### 2) Is Xi a significant variable in the model? (Test the partial slope of Xi)

Use the T-test for the slope of the individual variable and use the associated p-value to determine its significance based on the alpha level.

$$H_0: B_i = 0$$

$$H_a: B_i \neq 0$$

$CV = T^*(\alpha/2, df) = T^*(\alpha/2, n - 2)$  where "n" is the number of observations

$$TS = T = (b_i - 0) / sb_i$$

If p-value < alpha, the variable is significant and we reject  $H_0$ .

$$T_{\text{stat}} = \frac{b_i - 0}{sb_i} \quad T^*: (\alpha/2, n-2)$$

### 3) Should we keep $X_i$ and $X_j$ in the model?

Use the Partial F-Test to determine if we should keep the full model or accept the reduced model.

$$H_0: B_i = B_j = 0$$

$$H_a: \text{At least one } B \neq 0$$

$CV = F^*(\alpha, d_1, d_2) = F^*(\alpha, f - r, n - k - 1)$  where  $f$  is the number of  $X$  variables in the full model and  $r$  is the number of variables in the reduced model.

$$TS = F = (SSE_r - SSE_f / f - r) / (SSE_f / n - k - 1)$$

$$T_{stat} = F = \frac{\frac{SSE_{red} - SSE_{full}}{f - r}}{\frac{SSE_{full}}{n - k - 1}} \quad F^*(\alpha, f - r, n - k - 1)$$

#### Multicollinearity

This phenomenon occurs when the  $X$  variables are correlated with each other. The prediction ability of the model becomes redundant.

Consequences:

- The standard errors of the slopes ( $s_{b_i}$ ) are larger than they should be
  - ✓ This is apparent when all the individual  $T$ -values are near 0 and  $F$  is quite large.
- It is hard to get a good estimate for the model's slope coefficients ( $B_i$ )
  - ✓ This is apparent when the slope ( $b_i$ ) is the wrong sign or when one of the slopes ( $b_i$ ) changes significantly when another variable is dropped.

How to measure:

$$VIF = 1 / (1 - R^2_j)$$

- If individual  $VIF > 10$ , this variable is a problem
- If average  $VIF > 1$ , SSE is inflated
- If the  $VIF$  for the entire model is larger than an individual  $VIF$ , multicollinearity is not a problem

#### Table of Correlations

Identify the variable pairs with a correlation larger than 0.50, these variables are related in some way and should be investigated.

## Chapter 6: Assessing Assumptions

Note: Error = Residual = Disturbance

Assumptions:

- 1) The errors are independent
- 2) The errors have constant variance
- 3) The errors are normally distributed
- 4) The average of the errors "e" = 0 and the regression line passes through y-bar

\* If one assumption is not respected, we cannot proceed with the regression model as it is will be biased. The error terms must respect all these conditions.

How to perform a Residual Analysis:  $\neq$  ANOVA ASSUMPTIONS

### 1) Check for Independence

Ho: Residuals are independent

Ha: Residuals are not independent

Comment on the distribution of the error terms by looking at the errors plot. If there is no distinct pattern in the plot, we do not reject Ho and assume the errors to be independent.

### 2) Check for Constant Variance

Ho: Variance is constant

Ha: Variance is not constant

Run a Heteroschedasticity test (First & Second Moment Specification).

If the p-value is smaller than 0.05, do not reject Ho and assume constant variance.

### 3) Check for Normality

Ho: Studentized Residual is normal

Ha: Studentized Residual is not normal

If the majority of the p-values are smaller than 0.05, do not reject Ho and assume normality of residuals.

### 4) Check for Linearity

Use the *Pure-Error Lack-of-Fit Test* OR the *Data Subsetting Test*. Unfortunately, only SAS coding yields the first test and only Minitab can provide the p-value table for the second test. Do not test this assumption unless the appropriate table is provided to you.

## How to identify outliers in Residual Analysis:

### 1) Leverage

Flag value if it is larger than  $2 \cdot (k + 1) / n$

### 2) DFIT

Flag value if it is larger than  $2 \cdot \text{SQRT}((k + 1) / n)$

### 3) Cook's Distance

Flag value if it is larger than  $F(k + 1, n - k - 1, \alpha)$

### 4) Studentized Residual

Flag value if it is larger than 2 (absolute)

\*An outlier fails at least 2 tests and must be investigated further.

## Chapter 7: Indicators and Interactions

### Indicators variables

These are used to introduce qualitative data into the model.

### Interaction variables

These are used to try to build a better model by multiplying two already existing variables together into one variable.

$$Y = B_0 + B_1X_1 + B_2X_2 + B_3X_1X_2$$

$$Y = B_0 + (B_1 + B_3X_2) X_1 + B_2X_2$$

### Find the BEST model

The BEST model is found by eliminating the least significant variables and running the regression again. We can conclude whether the new model is better than the old model by conducting a partial F-test (Chapter 4).

How to identify the least significant variables:

- Variables that have p-values that are greater than alpha
- Variables that have a p-value relatively far from zero

### Sum of Squares Type 1

By adding variable X to the model, SSR will increase by the value of SS1 associated with variable X.

- This is a weak assessment because the sequence is set by the order in which the variables were entered into the model.

### Sum of Squares Type 2

By adding variable X to the model, in addition to all other variables, SSR will increase by the value of SS2 associated with variable X.

- This is a better assessment because it is what you have gained from adding a variable to the model.