STA121: Applied Regression Analysis Linear Regression Analysis - Chapters 3 and 4 in Dielman

Artin Armagan

Department of Statistical Science

September 15, 2009

イロト イポト イヨト イヨト

ъ

Armagan

Outline

Simple Linear Regression Analysis

- Using simple regression to describe a linear relationship
- Inferences From a Simple Regression Analysis
- Assessing the Fit of the Regression Line
- Prediction with a Sample Linear Regression Equation

2 Multiple Linear Regression

- Using Multiple Linear Regression to Explain a Relationship
- Inferences From a Multiple Regression Analysis
- Assessing the Fit of the Regression Line
- Comparing Two Regression Models
- Multicollinearity

ヘロト ヘアト ヘビト ヘビト

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

・ロト ・回ト ・ヨト ・ヨト

Purpose and Formulation

- Regression analysis is a statistical technique used to describe relationships among variables.
- In the simplest case where bivariate data are observed, the simple linear regression is used.
- The variable that we are trying to model is referred to as the *dependent* variable and often denoted by *y*.
- The variable that we are trying to explain *y* with is referred to as the *independent* or *explanatory* variable and often denoted by *x*.
- If a linear relationship between *y* and *x* is believed to exist, this relationship is expressed through an equation for a line:

$$y=b_0+b_1x$$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロト ヘアト ヘヨト ヘ

Purpose and Formulation

- Above equation gives an *exact* or a *deterministic* relationship meaning there exists no randomness.
- In this case recall that having only two pairs of observations (x, y) would suffice to construct a line.
- However many things we observe have a random component to it which we try to understand through various probability distributions.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

Example



y = 1 + 2x

 $\hat{y} = -0.2 + 2.2x$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イロト イヨト イヨト

Least Squares Criterion to Fit a Line

- We need to specify a method to find the "best" fitting line to the observed data.
- When we pass a line through the the observations, there will be differences between the actual observed values and the values predicted by the fitted line. This difference at each x value is called a *residual* and represents the "error".
- It is only sensible to try to minimize the total error we make while fitting the line.
- The least squares criterion minimizes the sum of squared errors to fit a line, i.e. min $\sum_{i=1}^{n} (y_i \hat{y}_i)^2$.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

◆□ > ◆□ > ◆豆 > ◆豆 > -

Least Squares Criterion to Fit a Line

This is a simple minimization problem and results in the following expressions for b_0 and b_1 :

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

These are simply obtained by differentiating $\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ $(\hat{y}_i = b_0 + b_1 x_i)$ with respect to b_0 and b_1 and setting them equal to zero at the solution which leaves us with two equations and two unknowns.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ъ

Pricing Communication Nodes



Armagan

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

A B A B A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

э

Estimating Residential Real Estate Values



SIZE

Armagan

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

・ロト ・ 『 ト ・ ヨ ト

Assumptions

 It is assumed that there exists a linear deterministic relationship between x and the mean of y, μ_{y|x}:

$$\mu_{\mathbf{y}|\mathbf{x}} = \beta_0 + \beta_1 \mathbf{x}$$

Since the actual observations deviate from this line, we need to add a noise term giving

$$y_i = \beta_0 + \beta_1 x_i + \boldsymbol{e}_i.$$

- The expected value of this error term is zero: $E(e_i) = 0$.
- The variance of each e_i is equal to σ²_θ. This assumptions suggests a constant variance along the regression line.
- The *e_i* are normally distributed.
- The *e_i* are independent.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イロト イヨト イヨト

Inferences about β_0 and β_1

- The point estimates of β₀ and β₁ are justified by the least squares criterion such that b₀ and b₁ minimize the sum of squared errors for the observed sample.
- It should be also noted that, under the assumptions made earlier, the *maximum likelihood estimator* for β_0 and β_1 is identical to the least squares estimator.
- Recall that a statistic is a function of a sample (which is a realization of a random variable), thus is a random variable itself. b₀ and b₁ have sampling distributions.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Sampling Distribution of b_0

•
$$E(b_0) = \beta_0$$

•
$$Var(b_0) = \sigma_e^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x}^2)^2} \right)$$

• The sampling distribution of b_0 is normal.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

Sampling Distribution of b₁

•
$$E(b_1) = \beta_1$$

•
$$Var(b_0) = \frac{\sigma_e^2}{\sum_{i=1}^n (x_i - \bar{x}^2)^2}$$

• The sampling distribution of *b*₁ is normal.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロト ヘアト ヘヨト ヘ

Properties of b_0 and b_1

- b_0 and b_1 are unbiased estimators for β_0 and β_1
- b_0 and b_1 are consistent estimators for β_0 and β_1
- b₀ and b₁ are minimum variance unbiased estimators for β₀ and β₁. That said, they have smaller sampling errors than any other unbiased estimator for β₀ and β₁.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Estimating σ_e^2

- The sampling distributions of b_0 and b_1 are normal when σ_e^2 is known.
- In realistic cases we won't know σ_2^2 .
- An unbiased estimate of σ_e^2 is given by

$$s_{e}^{2} = rac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{n-2} = rac{SSE}{n-2} = MSE$$

where $\hat{y}_i = b_0 + b_1 x_i$.

• Substituting s_e for σ_e earlier, s_{b_0} and s_{b_1} can be obtained.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Constructing Confidence Intervals for β_0 and β_1

- Now that σ_e is not known, the sampling distributions of b_0 and b_1 are t, i.e. $\frac{b_0 \beta_0}{s_{b_0}} \sim t_{n-2}$ and $\frac{b_1 \beta_1}{s_{b_1}} \sim t_{n-2}$.
- (1 α)100% confidence intervals then can be constructed as

$$\begin{array}{rcl} (b_0 - t_{\alpha/2,n-2} s_{b_0} &, & b_0 + t_{\alpha/2,n-2} s_{b_0}) \\ (b_1 - t_{\alpha/2,n-2} s_{b_1} &, & b_1 + t_{\alpha/2,n-2} s_{b_1}). \end{array}$$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

Hypothesis tests about β_0 and β_1

- Conducting a hypothesis test is no more involved than constructing a confidence interval. We make use of the same pivotal quantity, $\frac{b_1 \beta_1}{s_{b_1}}$ which is *t* distributed.
- Since we often include the intercept in our model anyway, a hypothesis test on β_0 may be redundant. Our main goal is to see whether there exists a linear relationship between the two variables which is implied by the slope, β_1 .
- We first state the null and alternative hypotheses:

$$H_0 : \beta_1 = (\geq, \leq)\beta_1^*$$

$$H_a : \beta_1 \neq (<, >)\beta_1^*$$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロン 人間 とくほど くほどう

Hypothesis tests about β_0 and β_1

- To test this hypothesis, a *t* statistic is used, $t = \frac{b_1 \beta_1^*}{s_{b_1}}$.
- A significance level, α, is specified to decide whether or not reject the null hypothesis.
- Possible alternative hypotheses and corresponding decision rules are Alternative Decision Rule
 - $\begin{array}{l} \overline{H_a: \beta_1 \neq \beta_1^*} & \text{Reject } H_0 \text{ if } |t| > t_{\alpha/2, n-2} \\ H_a: \beta_1 < \beta_1^* & \text{Reject } H_0 \text{ if } t < -t_{\alpha, n-2} \\ H_a: \beta_1 > \beta_1^* & \text{Reject } H_0 \text{ if } t > t_{\alpha, n-2} \end{array}$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

😐 💀 🚵 🗙 🗅 🚍 🌒 🖹 🗎 🖨	R Console		
~/Desktop/STA121/Chap3&4	\supset	Q	
<pre>dat = read.csv("comnode3.txt")</pre>			6
dat			
COST NUMPORTS			
52388 68			
51/61 52			
2600F 22			
27500 16			
57988 56			
54475 56			
33969 28			
31309 24			
0 23444 24			
1 24269 12			
2 53479 52			
3 33543 20			
4 33056 24			
<pre>summary(im(CUSI~NUMPORIS,data=dat))</pre>			
111			
m(formula = COST ~ NUMPORTS, data = dat)			
actor autor = cost = non onts, ducu = ducy			
tesiduals:			
Min 10 Median 30 Max			
8753.7 -873.9 681.0 2675.4 5019.9			
oefficients:			
Estimate Std. Error t value Pr(>ItI)			
Intercept) 16593.65 2687.05 6.175 4.76e-05			
UMPORTS 650.17 66.91 9.717 4.88e-07			
light+, codes: 0 000 0.001 0.01 0.05 . 0.1	1		
lesidual standard error: 4307 on 12 dearees of freedom			
Aultiple R-squared: 0.8872, Adjusted R-squared: 0.8778			
-statistic: 94.41 on 1 and 12 DF, p-value: 4.882e-07			
confint(obj)			n -
2.5 % 97.5 %			
Tekenseek) 18720 8692 22440 226			

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロト ヘアト ヘビト ヘビト

3



Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Lab Assignment #1

- Work on Exercises 6 and 7 in Chapter 3. You are encouraged to use R but may use JMP if you feel more comfortable.
- Due date is September 9.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

ъ



Analysis of Variance Sum of Source DF Squares Mean Square F Ratio Model 1 1751268376 1.7513e+994.4105 Error 12 222594146 18549512 Prob > FC. Total 13 1973862522 <.0001*

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロト ヘアト ヘヨト ヘ

The Coeficient of Determination

- In an exact or deterministic relationship, SSR=SST and SSE=0. This would imply that a straight line could be drawn through each observed value.
- Since this is not the case in real life, we need a a measure of how well the regression line fits the data.
- The coefficient of determination gives the proportion of total variation explained in the response by the regression line and is denoted by R².

$$R^2 = \frac{SSR}{SST}$$

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

The Correlation Coefficient

- For simple linear regression the correlation coefficient is $r = \pm \sqrt{R^2}$.
- This does not apply to multiple linear regression.
- If the sign of *r* is positive, then the relationship between the variables is direct, otherwise is inverse.
- r ranges between -1 and 1.
- A correlation of 0 merely implies no linear relationship.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロン イロン イヨン イヨン

The F Statistic

• An additional measure of how well the regression line fits the data is provided by the *F* statistic, which tests whether the equation $\hat{y} = b_0 + b_1 x$ provides a better fit to the data than the equation $\hat{y} = \bar{y}$.

$$F = \frac{MSR}{MSE}$$

where MSR = SSR/1 and MSE = SSE/(n-2).

 The degrees of freedom corresponding to SSR and SSE add up to the total degrees of freedom, n − 1.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロン イロン イヨン イヨン

The F Statistic

- To formalize the use of *F* statistic, consider the hypotheses $H_0: \beta_1 = 0$ vs. $H_a: \beta_1 \neq 0$.
- We reject H_0 if $F > F_{\alpha,1,n-2}$.
- For simple linear regression, $F = \frac{MSR}{MSE} = t^2$.
- Since both a *t*-test and an *F*-test will yield the same conclusions, it doesn't matter which one we use.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

ヘロト ヘアト ヘビト ヘビト

3



Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

э



Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation



Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

What Makes a Prediction Interval Wider?

• The difference arises from the difference between the variation in the mean of *y* and the variation in one individual *y* value.

•
$$\mathbb{V}(\bar{y}_f) = \sigma_e^2\left(\frac{1}{n} + \frac{(x_f - \bar{x})^2}{(n-1)s_x^2}\right)$$

•
$$\mathbb{V}(\mathbf{y}_f) = \sigma_e^2 \left(1 + \frac{1}{n} + \frac{(x_f - \bar{x})^2}{(n-1)s_x^2}\right)$$

• Replace σ_e^2 by s_e^2 when the error variance is not known and is to be estimated.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Assessing the Quality of Fit

• The mean square deviation is used commonly.

$$MSD = \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n_h}$$

where n_h is the size of the hold-out sample.

Using simple regression to describe a linear relationship Inferences From a Simple Regression Analysis Assessing the Fit of the Regression Line Prediction with a Sample Linear Regression Equation

イロト イポト イヨト イヨト

Lab Assignment #2

- Work on Exercises 8, 9, 10, 11, 12 and 13 in Chapter 3.
 You are encouraged to use R but may use JMP if you feel more comfortable.
- Due date is September 18.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト



 If a linear relationship between y and a set of xs is believed to exist, this relationship is expressed through an equation for a plane:

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_p x_p$$

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

Meddicorp Sales



Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト



 Assumptions are the same with simple linear regression model. Thus the population regression equation is written as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_p x_{ip} + e_i$$

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト

> lm.med=lm(SALES~ADV+BONUS+MKTSHR+COMPET.data=dat)

Inferences About the Population Regression Coefficients

							Call:							
Summar	y of Fit						lm(formula	= SALES	$\sim ADV$	+ BONUS	+ MKTSHR	<pre>k + COMPET,</pre>	, data	= dat)
RSquare RSquare Ad Root Mean Mean of Re Observation	lj Square Error sponse ns (or Sum Wg	0.85 0.831 93.769 1269. gts)	92 04 72 02 25				Residuals: Min -186.98 -7	1Q M 73.97	ledian 16.95	3Q 55.62	Max 125.52			
Analysis	of Varianc	e					Coefficient	ts:						
Source Model Error C. Total	DF 4 10 20 1 24 124	Sum of Squares M 73118.5 75855.2 48973.7	ean Squar 26828 879	e F Rat 0 30.51 3 Prob > <.000	io 14 F 1*		(Intercept) ADV BONUS MKTSHR	Estin) -593.5 2.5 1.9 2.6	0059	d. Error 259.1959 0.3143 0.7424 4.6357	-2.290 7.997 2.567 0.572	0.0330 1.17e-07 0.0184 0.5738	•	
Paramet	er Estimate	s					COMPET	-0.1	1207	0.3718	-0.325	0.7488		
Term Intercept ADV BONUS MKTSHR COMPET	Estimate -593.5375 2.513138 1.905948 2.651007 -0.120731	Std Error 259.1959 0.314275 0.742386 4.635655 0.371815	t Ratio -2.29 8.00 2.57 0.57 -0.32	Prob> t 0.0330* <.0001* 0.0184* 0.5738 0.7488	Lower 95% -1134.211 1.8575708 0.3573588 -7.018801 -0.896324	Upper 95% -52.86438 3.1687052 3.4545372 12.320815 0.654861	Signif. coo Residual st Multiple R- F-statistic	des: 0 tandard -squared :: 30.51	****' error: d: 0.85 L on 4	0.001 '* 93.77 o 92, Adju and 20 D	*'0.01' n 20 degr sted R-sq F, p-val	*' 0.05'. rees of fre uared: 0.8 ue: 2.937e	' 0.1 redom 331 e-08	''1

> summary(lm.med)

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

ヘロト ヘワト ヘビト ヘビト



- An R^2 value is computed as in the case of simple linear regression.
- Although it has a nice interpretation, it also has a drawback in the case of multiple linear regression.
- Intuitively, R² will never decrease as we add more independent variable into the model disregarding the fact that the variables being thrown in to the model may be explaining an insignificant portion of the variation in y. And as far as we can tell, the closer R² is to 1, the better.
- This means, we have to somehow account for how many variables we include in our model. In other words, we need to somehow "penalize" for the number of variables included in the model.
- Always remember that, the simpler the model we come up with, the better.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

(日)



• The adjusted *R*² does not suffer from this limitation and it accounts for the number of variables included in the model

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

 Note that, in this case, if a variable is causing an insignificant amount of decrease in SSE, the denominator in the above equation may actually be increasing, leading to a smaller R²_{adj} value.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イ理ト イヨト イヨト

F Statistic

- And *F* statistic is computed in a similar fashion to that of simple linear regression case.
- A different use of the *F* statistic will come into play with the comparison of nested models in the multiple linear regression case.
- This helps us compare a larger model to a reduced model which comprises a subset of variables included in the full model.
- This statistic is computed for each variable in the model if one uses anova() in R or looks at "Effect tests" in JMP.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

F Statistic

Effect 1	Test	s				
				Sum of		
Source	Np	arm	DF	Squares	F Ratio	Prob > F
ADV		1	1	562259.56	63.9457	<.0001*
BONUS		1	1	57954.64	6.5912	0.0184*
MKTSHR		1	1	2875.57	0.3270	0.5738
COMPET		1	1	927.07	0.1054	0.7488
Response:	SAL	ES				
	Df	Sum Sa	Mean So	F value	Pr(>E)	
ADV	1	1012408	1012408	8 115.1411	9.546e-10 **	*
BONUS	1	55389	55389	6.2994	0.02079 *	
MKTSHR	1	4394	4394	4 0.4997	0.48777	
COMPET	1	927	927	7 0.1054	0.74877	
Residuals	20	175855	8793	3		
Signif. co	odes	: 0 '*	**' 0.00	01 '**' 0.0	0.05 (**)	.' 0.1 ' ' 1

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト 不得 とくほと くほとう

F Statistic

- $H_0: \beta_{l+1} = ... = \beta_p = 0$
 - H_a : At least one of $\beta_{l+1}, ..., \beta_p$ is not equal to zero.
- This implies, under the reduced model we have

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_p x_{il} + e_i.$$

• If we want to compare this reduced model to the full model and find out if the reduction was reasonable, we have to compute an *F* statistic:

$$F = \frac{(SSE_R - SSE_F)/(p-I)}{SSE_F/(n-p-1)}$$

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト

Meddicorp

*H*₀: β₃ = ... = β₄ = 0
 H_a: At least one of β₃, β₄ is not equal to zero.

			Analysis of variance						
Source Model Error C. Total	DF 4 20 24	Sum of Squares 1073118.5 175855.2 1248973.7	Mean Square 268280 8793	F Ratio 30.5114 Prob > F <.0001*	Source Model Error C. Total	DF 2 22 24	Sum of Squares 1067797.3 181176.4 1248973.7	Mean Square 533899 8235	F Ratio 64.8306 Prob > F <.0001*

- $F = \frac{(181176 175853)/2}{175855/20} = 0.303$
- 3.49 is the 5% F critical value with 2 numerator and 20 denominator degrees of freedom. Thus we accept H₀.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト

Using Conditional Sums of Squares

```
> anova(lm.med)
Analysis of Variance Table
Response: SALES
              Sum Sq Mean Sa
          Df
                              F value
                                          Pr(>F)
             1012408 1012408 115.1411 9.546e-10 ***
ADV
           1
BONUS
               55389
                       55389
                               6.2994
                                         0.02079 *
           1
MKTSHR
                4394
                        4394
                               0.4997
                                         0.48777
           1
COMPET
           1
                 927
                         927
                               0.1054
                                         0.74877
Residuals 20
                        8793
              175855
_ _ _
Signif. codes:
                        0.001
                               (**)
                                   0.01 (*' 0.05 (.' 0.1 (' 1
```

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト

Lab Assignment #3

- Work on Exercises 1, 2, 3 and 4 in Chapter 4. You are encouraged to use R but may use JMP if you feel more comfortable.
- Due date is September 25.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

・ コ ト ・ 四 ト ・ 回 ト ・

Multicollinearity

- When explanatory variables are correlated with one another, the problem of multicollinearity is said to exist.
- The presence of a high degree multicollinearity among the explanatory variables result in the following problem:
 - The standard deviations of the regression coefficients are disproportionately large leading to small *t*-score although the corresponding variable may be an important one.
 - The regression coefficient estimates are highly unstable. Due to high standard errors, reliable estimation is not possible.

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト イポト イヨト イヨト

Detecting multicollinearity

- Pairwise correlations.
- Large F, small t.
- Variance inflation factor (VIF).

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

・ロト ・回ト ・ヨト ・ヨト

Example - Hald's data

Montgomery and Peck (1982) illustrated variable selection techniques on the Hald cement data and gave several references to other analysis. The response variable v is the heat evolved in a cement mix. The four explanatory variables are ingredients of the mix, i.e., x1: tricalcium aluminate, x2: tricalcium silicate, x3: tetracalcium alumino ferrite, x4: dicalcium silicate. An important feature of these data is that the variables x1 and x3 are highly correlated (corr(x1,x3)=-0.824), as well as the variables x2 and x4 (with corr(x2,x4)=-0.975). Thus we should expect any subset of (x1,x2,x3,x4) that includes one variable from highly correlated pair to do as any subset that also includes the other member.

Multicollinearity

イロン 不同 とくほ とくほ とう

ъ

Example - Hald's data

[1,]

[2,]



Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

Example - Hald's data

Call:				
lm(formula	-	y.hald	~	x.hald)

Residuals:

Min	10	Median	3Q	Max	
-3.1750	-1.6709	0.2508	1.3783	3,9254	

Coefficients:

	Estimate	Std. Error	t value	Pr(>ltl)			
(Intercept) 62.4054	70.0710	0.891	0.3991			
x.hald1	1.5511	0.7448	2.083	0.0708			
x.hald2	0.5102	0.7238	0.705	0.5009			
x.hald3	0.1019	0.7547	0.135	0.8959			
x.hald4	-0.1441	0.7091	-0.203	0.8441			
Signif. co	des: 0 **	**' 0.001 '	**' 0.01	'*' 0.05	·.' 0.1	· * *	1

Residual standard error: 2.446 on 8 degrees of freedom Multiple R-squared: 0.9824, Adjusted R-squared: 0.9736 F-statistic: 111.5 on 4 and 8 DF, p-value: 4.756e-07

1	Summa	ry of Fit							
	RSquare			0.9	82376				
	RSquare A	٨dj		0.9	73563				
	Root Mea	n Square E	ror	2.4	46008				
	Mean of R	lesponse		95.	42308				
	Observati	ons (or Sur	n Wg	ts)	13				
1	Analys	is of Vari	ianc	e					
				Sum of					
	Source	DF		Squares	Mear	Squa	re	F Rat	io
	Model	4	266	7.8994		666.97	5 11	1.47	92
	Error	8	- 4	7.8636		5.98	33 P	rob >	F
	C. Total	12	271	5.7631				<.000	1*
٦	Parame	eter Estin	nate	s					
	Term	Estim	ate	Std Err	or t	Ratio	Prob	> t	
	Intercept	62.405	369	70.070	96	0.89	0.39	91	
	X1	1.5511	026	0.744	77	2.08	0.07	'08	
	X2	0.5101	676	0.7237	88	0.70	0.50	09	
	X3	0.1019	094	0.7547	09	0.14	0.89	59	
	X4	-0.144	061	0.7090	52	-0.20	0.84	41	
1	Effect 1	Tests							
				S	um of				
	Source	Nparm	DF	Sq	uares	E F	Ratio	Pro	b > F
	X1	1	1	25.95	0911	4.	3375	0.0	708
	X2	1	1	2.97	2478	0.	4968	0.5	009
	X3	1	1	0.10	9090	0.	0182	0.8	959
	X4	1	1	0.24	6975	0.	0413	0.8	441

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

イロト 不得 とくほ とくほう

3

Example - Hald's data

Paramet	ter Estimate	S			
Term	Estimate	Std Error	t Ratio	Prob> t	VIF
Intercept	62.405369	70.07096	0.89	0.3991	
X1	1.5511026	0.74477	2.08	0.0708	38.496211
X2	0.5101676	0.723788	0.70	0.5009	254.42317
X3	0.1019094	0.754709	0.14	0.8959	46.868386
X4	-0.144061	0.709052	-0.20	0.8441	282.51286

Using Multiple Linear Regression to Explain a Relationship Inferences From a Multiple Regression Analysis Assessing the Fit of the Regression Line Comparing Two Regression Models Multicollinearity

3

Example - Hald's data

Summa	ary of Fi	t						
RSquare RSquare Root Mea Mean of F Observati	Adj In Square Response ions (or Si	Error um Wg	0.9786 0.9744 2.4063 95.423 ts)	578 14 335 308 13				
Analys	is of Va	rianco	e					
Source Model Error C. Total	DF 2 10 12	265 5 271	Sum of Squares M 7.8586 7.9045 5.7631	ean Sq 132	uare 8.93 5.79	F Rat 229.503 Prob > <.000	io 37 F 1*	
Parame	eter Esti	imate	S					
Term Intercept X1 X2	Esti 52.57 1.468 0.662	mate 7349 3057 2505	Std Error 2.286174 0.121301 0.045855	t Rat 23.0 12.1 14.4	io P 00 < 10 < 14 <	rob> t <.0001* <.0001* <.0001*	1.055 1.055	VIF 129 129
Effect	Tests							
Source X1 X2	Nparm 1 1	DF 1 1	Sum Squai 848.43 1207.78	of res 19 14 23 20	F Rat 46.522	io Prol 27 <.0 18 <.0	b > F 1001* 1001*	