Chapter 4

Multiple Regression Analysis: Inference



© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (1 of 37)

- Statistical inference in the regression model
 - Hypothesis tests about population parameters
 - Construction of confidence intervals

Sampling distributions of the OLS estimators

- The OLS estimators are random variables
- We already know their expected values and their variances
- However, for hypothesis tests we need to know their distribution
- In order to derive their distribution we need additional assumptions
- Assumption about distribution of errors: normal distribution

Multiple Regression Analysis: Inference (2 of 37)

Assumption MLR.6 (Normality of error terms)

 $u_i \sim \text{Normal}(0, \sigma^2)$ independently of $x_{i1}, x_{i2}, \ldots, x_{ik}$



It is assumed that the unobserved factors are normally distributed around the population regression function.

The form and the variance of the distribution does not depend on any of the explanatory variables.

It follows that:

 $y|\mathbf{x} \sim \text{Normal}(\beta_0 + \beta_1 x_1 + \ldots + \beta_k x_k, \sigma^2)$

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (3 of 37)

- Discussion of the normality assumption
- The error term is the sum of "many" different unobserved factors.
- Sums of independent factors are normally distributed (CLT).
- Problems:
 - How many different factors? Number large enough?
 - Possibly very heterogenuous distributions of individual factors
 - How independent are the different factors?
- The normality of the error term is an empirical question.
- At least, the error distribution should be "close" to normal.
- In many cases, normality is questionable or impossible by definition.

Multiple Regression Analysis: Inference (4 of 37)

- Discussion of the normality assumption (cont.)
- Examples where normality cannot hold:
 - Wages (nonnegative; also: minimum wage)
 - Number of arrests (takes on a small number of integer values)
 - Unemployment (indicator variable, takes on only 1 or 0)
- In some cases, normality can be achieved through transformations of the dependent variable (e.g. use log(wage) instead of wage).
- Under normality, OLS is the best (even nonlinear) unbiased estimator
 - Important: For the purposes of statistical inference, the assumption of normality can be replaced by a large sample size.

Multiple Regression Analysis: Inference (5 of 37)

Terminology

 $\underline{MLR.1 - MLR.5}$

"Gauss-Markov assumptions"

$\underbrace{MLR.1 - MLR.6}_{MLR.1}$

"Classical linear model (CLM) assumptions"

- Theorem 4.1 (Normal sampling distributions)
- Under assumptions MLR.1 MLR.6:

$$\hat{\beta}_j \sim \operatorname{Normal}(\beta_j, Var(\hat{\beta}_j))$$

The estimators are normally distributed around the true parameters with the variance that was derived earlier

$$\frac{\hat{\beta}_j - \beta_j}{sd(\hat{\beta}_j)} \sim \operatorname{Normal}(0, 1)$$

The standardized estimators follow a standard normal distribution

Multiple Regression Analysis: Inference (6 of 37)

- Testing hypotheses about a single population parameter
- Theorem 4.2 (t-distribution for the standardized estimators)
 - Under assumptions MLR.1 MLR.6



If the standardization is done using the <u>estimated</u> standard deviation (= standard error), the normal distribution is replaced by a t-distribution

Note: The t-distribution is close to the standard normal distribution if n-k-1 is large.

• Null hypothesis (for more general hypotheses, see below)

 $H_0: \ \beta_j = 0$ The population parameter is equal to zero, i.e. after controlling for the other independent variables, there is no effect of x_i on y

Multiple Regression Analysis: Inference (7 of 37)

• t-statistic (or t-ratio)

 $t_{\hat{\beta}_i} \equiv$

The t-statistic will be used to test the above null hypothesis. The farther the estimated coefficient is away from zero, the less likely it is that the null hypothesis holds true. But what does "far" away from zero mean?

This depends on the variability of the estimated coefficient, i.e. its standard deviation. <u>The t-statistic measures how many estimated</u> <u>standard deviations the estimated coefficient is away from zero.</u>

• Distribution of the t-statistic if the null hypothesis is true

$$t_{\widehat{\beta}_j} \equiv \widehat{\beta}_j / se(\widehat{\beta}_j) = (\widehat{\beta}_j - \beta_j) / se(\widehat{\beta}_j) \sim t_{n-k-1}$$

 Goal: Define a rejection rule so that, if it is true, H₀ is rejected only with a small probability (= significance level, e.g. 5%) Multiple Regression Analysis: Inference (8 of 37)

Testing against one-sided alternatives (greater than zero)



- Reject the null hypothesis in favour of the alternative hypothesis if the estimated coefficient is "too large" (i.e. larger than a critical value).
- Construct the critical value so that, if the null hypothesis is true, it is rejected in, for example, 5% of the cases.
- In the given example, this is the point of the tdistribution with 28 degrees of freedom that is exceeded in 5% of the cases.
- Reject if t-statistic is greater than 1.701

Multiple Regression Analysis: Inference (9 of 37)

• Example: Wage equation

• Test whether, after controlling for education and tenure, higher work experience leads to higher hourly wages.

Wage1.dta

$$\widehat{\log}(wage) = .284 + .092 \ educ + 0041 \ exper + .022 \ tenure$$

(.104) (.007) (.0017) (.003)
 $n = 526, \ R^2 = .316$ Standard errors

Test
$$H_0$$
: $\beta_{exper} = 0$ against H_1 : $\beta_{exper} > 0$.

One would either expect a positive effect of experience on hourly wage or no effect at all.

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (10 of 37)

• Example: Wage equation (cont.)

 $t_{exper} = .0041/.0017 \approx 2.41$ \longleftarrow t-statistic

$$df = n - k - 1 = 526 - 3 - 1 = 522 \longleftarrow$$
 begrees of freedom;
here the standard norma approximation applies

 $c_{0.05} = 1.645$ Critical values for the 5% and the 1% significance level (these are conventional significance levels). $c_{0.01} = 2.326$ The null hypothesis is rejected because the t-statistic exceeds the critical value.

• The effect of experience on hourly wage is statistically greater than zero at the 5% (and even at the 1%) significance level.

Multiple Regression Analysis: Inference (11 of 37)

Testing against one-sided alternatives (less than zero)



- Reject the null hypothesis in favour of the alternative hypothesis if the estimated coefficient is <u>"too small</u>" (i.e. smaller than a critical value).
- Construct the critical value so that, if the null hypothesis is true, it is rejected in, for example, 5% of the cases.
- In the given example, this is the point of the tdistribution with 18 degrees of freedom so that 5% of the cases are below the point.
- Reject if t-statistic is less than -1.734

Multiple Regression Analysis: Inference (12 of 37)

• Example: Student performance and school size

• Test whether smaller school size leads to better student performance



Test
$$H_0$$
: $\beta_{enroll} = 0$ against H_1 : $\beta_{enroll} < 0$.

Do larger schools hamper student performance or is there no such effect?

Multiple Regression Analysis: Inference (13 of 37)

• Example: Student performance and school size (cont.)

$$df = n - k - 1 = 408 - 3 - 1 = 404$$

begrees of freedom;
here the standard normal
approximation applies

- $c_{0.05} = -1.65$ Critical values for the 5% and the 15% significance level. The null hypothesis is not rejected because the t-statistic is not smaller than the critical value.
- One cannot reject the hypothesis that there is no effect of school size on student performance (not even for a lax significance level of 15%).

Multiple Regression Analysis: Inference (14 of 37)

- Example: Student performance and school size (cont.)
 - Alternative specification of functional form:

 $\widehat{math10} = -207.66 + 21.16 \log(totcomp) + 3.98 \log(staff) - 1.29 \log(enroll)$ (48.70) $(48.70) + (4.06) \log(totcomp) + 3.98 \log(staff) - 1.29 \log(enroll)$ $n = 408, R^2 = 0.0654 \longleftarrow$ R-squared slightly higher

Test
$$H_0$$
: $\beta_{\log(enroll)} = 0$ against H_1 : $\beta_{\log(enroll)} < 0$.

t-statistic

Multiple Regression Analysis: Inference (15 of 37)

• Example: Student performance and school size (cont.)

 $c_{0.05} = -1.65$ Critical value for the 5% significance level; reject null hypothesis

- The hypothesis that there is no effect of school size on student performance can be rejected in favor of the hypothesis that the effect is negative.
- How large is the effect?

 $t_{\log(enroll)} = -1.29/.69 \approx -1.87$

$$-1.29 = \frac{\Delta \widehat{math10}}{\Delta \log(enroll)} = \frac{\Delta \widehat{math10}}{\frac{\Delta enroll}{enroll}} = \frac{\frac{-1.29}{100}}{\frac{1}{100}} = \frac{-0.0129}{+1\%} \leftarrow +10\% \text{ enrollment }; -0.129 \text{ percentage points students}$$

Multiple Regression Analysis: Inference (16 of 37)

• Testing against two-sided alternatives



- Reject the null hypothesis in favour of the alternative hypothesis if <u>the absolute value</u> of the estimated coefficient is too large.
- Construct the critical value so that, if the null hypothesis is true, it is rejected in, for example, 5% of the cases.
- In the given example, these are the points of the t-distribution so that 5% of the cases lie in the two tails.
- Reject if absolute value of t-statistic is less than
 -2.06 or greater than 2.06

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

GPA1.dta

Multiple Regression Analysis: Inference (17 of 37)

• Example: Determinants of college GPA

$$col\widehat{GPA} = 1.39 + .412 \ hsGPA + .015 \ ACT - .083 \ skipped$$

(.33) (.094) (.011) (.026)

$$n = 141, R^2 = .234$$

For critical values, use standard normal distribution

 $t_{hsGPA} = 4.38 > c_{0.01} = 2.58$ $t_{ACT} = 1.36 < c_{0.10} = 1.645$ $|t_{skipped}| = |-3.19| > c_{0.01} = 2.58$

The effects of hsGPA and skipped are significantly different from zero at the 1% significance level. The effect of ACT is not significantly different from zero, not even at the 10% significance level.

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (18 of 37)

- "Statistically significant" variables in a regression
 - If a regression coefficient is different from zero in a two-sided test, the corresponding variable is said to be "statistically significant".
 - If the number of degrees of freedom is large enough so that the normal approximation applies, the following rules of thumb apply:

$$|t - ratio| > 1.645 \longrightarrow$$
 "statistically significant at 10% level"
 $|t - ratio| > 1.96 \longrightarrow$ "statistically significant at 5% level"
 $|t - ratio| > 2.576 \longrightarrow$ "statistically significant at 1% level"

Introductory Econometrics: A Modern Approach (7e)

郭凯明等:家庭隔代抚养文化、延迟退休年龄与劳动力供给

				表 2	家庭隔代抚养对生育率的影响估计结果			
					(1)	(2)	(3)	(4)
			Grandparenting		0.214 ***	0. 215 ****	0. 228 ***	0. 228 ***
					(0.0292)	(0.0295)	(0.0306)	(0.0308)
				Edu	- 0. 0138 ****	-0.00982***	-0.0135 ***	-0.0109**
					(0. 00324)	(0.00351)	(0.00401)	(0.00430)
				Health	-0.00562	-0.00456	- 0. 00998	-0.00830
				(0.00822)	(0.00832)	(0. 00964)	(0. 0100)	
	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	-		0. 128 ***	0. 0997 ***	0. 136 ***	0. 104 ****	
					(0.0176)	(0.0179)	(0. 0202)	(0.0232)
				Age_sq	-0.00191	-0.00153	- 0. 00204 ****	-0.00160***
			N		(0.000249)	(0.000246)	(0.000284)	(0.000325)
				Edu_f		$-0.00892^{-0.00}$		-0.00715*
2)	(2)			/		0.00465		0.004277
ance	(3)		Health_f			(0,00944)		(0,00987)
15 ***	0. 228 ***					0. 0398 ***		0.0437**
(295)	(0, 0306)			Age_f		(0.00972)		(0.0195)
,2,3,	(0.05007			1 6		- 0. 000522 ****		-0.000582**
0982 ***	-0.0135 ***			Age_f_sq		(0.000121)		(0.000271)
0351)	(0.00401)			Urban	-0.0412	-0.0395	-0.0114	-0.0119
00456	0.00000	<u> </u>			(0.0301)	(0.0307)	(0.0412)	(0.0419)
00456	-0.00998			In finc	0. 0399 ****	0. 0431 ****	0. 0357 **	0. 0381 **
0832)	(0.00964)				(0.0133)	(0.0136)	(0.0159)	(0.0162)
997 ***	0. 136 ***			Job_class				
)179)	(0.0202)			Job_class_f				
0153 ***	-0.00204 ***						1.201.555	
	(0.00204		截距项	-1.143	(0, 325)	-1.394	-1.613	
)0246)	46) (0.000284) ((0.200)	+ 控制	(0.435)	(0.451)	
0892 ***				区 民 因 完 效 応	が知	水江町	拉制	拉制
					4839	4839	3159	3159
				B ²	0 201	0.203	0 184	0.186
				n	0.201	0. 200	0.104	0.100

注:括号内为稳健标准误,*、**和***分别表示10%、5%和1%的显著性水平。下表同

Tables in the paper

				-
	(1)	(2)	(3)	
coefficient	→ 0.214 ***	0. 215 ***	0. 228 ***	
irening	(0.0292)	(0.0295)	(0.0306)	
hu standars	-0.0138***	- 0. 00982 ****	-0.0135 ***	
	(0.00324)	(0.00351)	(0.00401)	
_141	- 0. 00562	- 0. 00456	- 0. 00998	
uin	(0.00822)	(0.00832)	(0.00964)	
	0. 128 ***	0. 0997 ***	0. 136 ***	
ge	(0.0176)	(0.0179)	(0.0202)	
	– 0. 00191 ***	- 0. 00153 ***	- 0. 00204 ***	
$_sq$	(0.000249)	(0.000246)	(0.000284)	
		- 0. 00892 ***		
u_J				I

A REIMINIUN

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

(5)

0. 228 ***

(0.0307)

-0.0108**

(0.00423)

-0.00828

(0.0101)

0. 104 ***

(0.0233)

-0.00160 ***

(0.000327)

-0.00709*

(0.00425)

-0.00694

(0.00976)

0.0438**

(0.0195)

-0.000583**

(0.000271)

-0.0120

(0.0420)

0.0381 **

(0.0165) 0.000453 (0.0761) -0.00565 (0.0340)

-1.613 ****

(0.451)

控制控制

3159 0. 186 Multiple Regression Analysis: Inference (19 of 37)

- Testing more general hypotheses about a regression coefficient
- Null hypothesis

 H_0 : $\beta_j = a_j \leftarrow$ Hypothesized value of the coefficient

• t-statistic

$$t = \frac{(estimate - hypothesized value)}{standard \ error} = \frac{(\hat{\beta}_j - a_j)}{se(\hat{\beta}_j)}$$

• The test works exactly as before, except that the hypothesized value is substracted from the estimate when forming the statistic.

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (20 of 37)

- Example: Campus crime and enrollment CAMPUS.dta
 - An interesting hypothesis is whether crime increases by one percent if enrollment is increased by one percent.

$$\widehat{\log}(crime) = -6.63 + 1.27 \log(enroll)$$

$$(1.03) + (0.11)$$
Estimate is different from one but is this difference statistically significant?

$$H_0: \beta_{\log(enroll)} = 1, \ H_1: \beta_{\log(enroll)} \neq 1$$

 $t = \frac{1.27 - 1}{.11} \approx 2.45 > 1.985 = c_{0.05} \quad \longleftarrow \begin{array}{l} \text{The hypothesis is rejected} \\ \text{at the 5\% level} \end{array}$

Multiple Regression Analysis: Inference (21 of 37)

Computing p-values for t-tests

- If the significance level is made smaller and smaller, there will be a point where the null hypothesis cannot be rejected anymore.
- The reason is that, by lowering the significance level, one wants to avoid more and more to make the error of rejecting a correct H_{0.}
- <u>The smallest significance level at which the null hypothesis is still rejected</u>, is called the p-value of the hypothesis test.
- A small p-value is evidence against the null hypothesis because one would reject the null hypothesis even at small significance levels.
- A large p-value is evidence in favor of the null hypothesis.
- P-values are more informative than tests at fixed significance levels.

Multiple Regression Analysis: Inference (22 of 37)

• How the p-value is computed (here: two-sided test)?



- The p-value is the significance level at which one is indifferent between rejecting and not rejecting the null hypothesis.
- In the two-sided case, the p-value is thus the probability that the t-distributed variable takes on a larger absolute value than the realized value of the test statistic, e.g.

P(|t| > 1.85) = 2*(0.0359) = .0718

- From this, it is clear that a null hypothesis is rejected if and only if the corresponding p-value is smaller than the significance level.
- For example, for a significance level of 5% the t-statistic would not lie in the rejection region.

Multiple Regression Analysis: Inference (23 of 37)

- Guidelines for discussing economic and statistical significance
 - If a variable is statistically significant, discuss the magnitude of the coefficient to get an idea of its economic or practical importance.
 - The fact that a coefficient is statistically significant does not necessarily mean it is economically or practically significant!
 - If a variable is statistically and economically important but has the "wrong" sign, the regression model might be misspecified.
 - If a variable is statistically insignificant at the usual levels (10%, 5%, or 1%), one may think of dropping it from the regression.
 - If the sample size is small, effects might be imprecisely estimated so that the case for dropping insignificant variables is less strong.

Multiple Regression Analysis: Inference (24 of 37)

- Confidence intervals
- Simple manipulation of the result in Theorem 4.2 implies that



- Interpretation of the confidence interval
 - The bounds of the interval are random.
 - In repeated samples, the interval that is constructed in the above way will cover the population regression coefficient in 95% of the cases.

Multiple Regression Analysis: Inference (25 of 37)

• Confidence intervals for typical confidence levels

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.01} \cdot se(\widehat{\beta}_{j})\right) = 0.99$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.05}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.05} \cdot se(\widehat{\beta}_{j})\right) = 0.95$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.10}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.01}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$

• Relationship between confidence intervals and hypotheses tests

$$a_j \notin interval \Rightarrow$$
 reject $H_0 : \beta_j = a_j$ in favor of $H_1 : \beta_j \neq a_j$

Multiple Regression Analysis: Inference (26 of 37)

• Example: Model of firms' R&D expenditures



as the interval is narrow. Moreover, the effect is significantly different from zero because zero is outside the interval.

This effect is imprecisely estimated as the interval is very wide. It is not even statistically significant because zero lies in the interval. Multiple Regression Analysis: Inference (27 of 37)

- Testing hypotheses about a linear combination of the parameters
- Example: Return to education at two-year vs. at four-year colleges



• A possible test statistic would be:



The difference between the estimates is normalized by the estimated standard deviation of the difference. The null hypothesis would have to be rejected if the statistic is "too negative" to believe that the true difference between the parameters is equal to zero. Multiple Regression Analysis: Inference (28 of 37)

• The standard error of the difference in parameters is impossible to with standard regression output

$$se(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{Var(\hat{\beta}_1 - \hat{\beta}_2)} = \sqrt{Var(\hat{\beta}_1) + Var(\hat{\beta}_2) - 2Cov(\hat{\beta}_1, \hat{\beta}_2)}$$

Usually not available in regression output

• An <u>alternative method</u> is to make a substitution in variables. Define $\theta_1 = \beta_1 - \beta_2$ and test $H_0: \theta_1 = 0$ against $H_1: \theta_1 < 0$.

$$log(wage) = \beta_0 + (\theta_1 + \beta_2)jc + \beta_2univ + \beta_3exper + u$$
$$= \beta_0 + \theta_1jc + \beta_2(jc + univ) + \beta_3exper + u$$
Insert into original regression a new regressor (= total years of college)

Multiple Regression Analysis: Inference (29 of 37)

• Estimation results

$$\widehat{\log}(wage) = 1.472 - .0102 \ jc + .0769 \ totcoll + .0049 \ exper (.021) \ (.0069) \ (.0023) \ (.0002)$$

$$n = 6,763, R^2 = .222$$

$$t = -.0102/.0069 = -1.48$$

$$p - value = P(t - ratio < -1.48) = .070$$

 $-.0102 \pm 1.96(.0069) = (-.0237, .0003)$

• This method works always for single linear hypotheses

Multiple Regression Analysis: Inference (30 of 37)

- Testing multiple linear restrictions: The F-test
- Testing exclusion restrictions

Salary of major lea- Years in Average number of gue baseball player the league games per year $\log(salary) = \beta_0 + \beta_1 years + \beta_2 gamesyr$ $+\beta_3 bavg + \beta_4 hrunsyr + \beta_5 rbisyr + u$ MLB1.DTA 1 1 1 Batting average Home runs per year Runs batted in per year $H_0: \beta_3 = 0, \beta_4 = 0, \beta_5 = 0$ against $H_1: H_0$ is not true Test whether performance measures have no effect/can be excluded from regression.

Multiple Regression Analysis: Inference (31 of 37)

• Estimation of the unrestricted model

 $\widehat{\log}(salary) = 11.19 + .0689 \ years + .0126 \ gamesyr$ (0.29) (.0121) (.0026)

None of these variabels is statistically significant when tested individually

$$n = 353, SSR = 183.186, R^2 = .6278$$

Idea: How would the model fit be if these variables were dropped from the regression?

Multiple Regression Analysis: Inference (32 of 37)

Estimation of the restricted model

 $\widehat{\log}(salary) = 11.22 + .0713 \ years + .0202 \ gamesyr \\ (0.11) \ (.0125) \ (.0013)$

$$n = 353, SSR = 198.311, R^2 = .5971$$

The sum of squared residuals necessarily increases, but is the increase statistically significant?



Multiple Regression Analysis: Inference (33 of 37)

• Rejection rule



- A F-distributed variable only takes on positive values. This corresponds to the fact that the sum of squared residuals can only increase if one moves from H₁ to H₀.
- Choose the critical value so that we incorrectly reject the null hypothesis in, for example, only 5% of the cases.

© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or in part, except for use as permitted in a license distributed with a certain product or service or otherwise on a password-protected website or school-approved learning management system for classroom use.

Multiple Regression Analysis: Inference (34 of 37)

• Test decision in example



$$F \sim F_{3,347} \Rightarrow c_{0.01} = 3.78$$

The null hypothesis is overwhelmingly rejected (even at very small significance levels).

- P(F statistic > 9.55) = 0.000
- Discussion
 - The three variables are "jointly significant"
 - They were not significant when tested individually
 - The likely reason is multicollinearity between them

Multiple Regression Analysis: Inference (35 of 37)

• Test of overall significance of a regression

$$y = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \ldots + \beta_k x_{ik} + u$$

 $H_0: \beta_1 = \beta_2 = \ldots = \beta_k = 0$ The null hypothesis states that the explanatory variables are not useful at all in explaining the dependent variable

$$y = \beta_0 + u$$
 \longleftarrow Restricted model (regression on constant)

$$F = \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n-k-1)} = \frac{R^2/k}{(1-R^2)/(n-k-1)} \sim F_{k,n-k-1}$$

- Discussion
 - The three variables are "jointly significant"
 - They were not significant when tested individually
 - The likely reason is multicollinearity between them

Multiple Regression Analysis: Inference (36 of 37)

- Testing general linear restrictions with the F-test
- Example: Test whether house price assessments are rational

The assessed housing value Size of lot Actual house price (before the house was sold) (in square feet) $\log(price) = \beta_0 + \beta_1 \log(assess) + \beta_2 \log(lotsize)$ $+\beta_3 \log(sqrft) + \beta_4 bdrms + u$ Hprice1.dta Square footage Number of bedrooms

 $H_0: \beta_1 = 1, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0$ In addition, other known factors should not influence the price once the assessed value has been controlled for If house price assessments are rational, a 1% change in the assessment should be associated with a 1% change in price.

not influence the price once the assessed value has been controlled for.

Multiple Regression Analysis: Inference (37 of 37)

Unrestricted regression

 $\log(price) = \beta_0 + \beta_1 \log(assess) + \beta_2 \log(lotsize) + \beta_3 \log(sqrft) + \beta_4 bdrms + u$

Restricted regression

 $log(price) = \beta_0 + log(assess) + u$ $log(price) - log(assess) = \beta_0 + u$ The restricted model is actually a regression of log(price) – log(assess) on a constant

• Test statistic

$$F = \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n-k-1)} = \frac{(1.880 - 1.822)/4}{1.822/(88 - 4 - 1)} \approx .661$$

 $F \sim F_{4,83} \Rightarrow c_{0.05} = 2.50 \Rightarrow H_0$ cannot be rejected