Multiple Regression Analyses: <u>Statistical</u> <u>Inference</u>

March 10, 2023

Dali Laxton

Today's Lecture

- We are going to discuss how hypotheses about coefficients can be tested in regression models
- We will explain what significance of coefficients mean
- We will learn how to read regression output

- Wooldridge Chapter 4;
- Studenmund Chapter 5.1-5.4

Multiple Regression Analyses: Inference

- 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
 - For hypothesis testing we need to know their <u>distribution</u>

<u>Inference</u>: Sampling distributions of the OLS Estimators

Assumption 6 (Normality of error terms)

 $u_i \sim N(0, \sigma^2)$ independently of $x_{i1}, x_{i2}, \dots, x_{ik}$



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.

Show normality of the error terms in GRETL

- Open GRETL load sample data "Engel"
- Run regression ols foodexp income const
- Generate residuals:

series exphat=\$yhat
genr resid=foodexp-exphat

or

genr resid =foodexp-(\$coeff(const) + \$coeff(income)*income)

• Display distribution of residuals:

freq resid --plot=display

<u>Inference</u>: Sampling distributions of the OLS Estimators

• 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? Observations 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

<u>Inference</u>: Sampling distributions of the OLS Estimators

- Discussion of the normality assumption (cont.)
 - Examples where normality cannot hold:
 - Wages (nonnegative; also: minimum wage)
 - Unemployment (indicator variable, takes on only 1 or 0)
 - In some cases, normality can be achieved through transformations of the dependent variable
 - Under normality, OLS is the best (even nonlinear) unbiased estimator
 - <u>Important</u>: For the purposes of statistical inference, the assumption of normality can be replaced by a large sample size (CLT)

Inference: Sampling distributions of the OLS Estimators CLT



Source: https://statisticsbyjim.com/basics/central-limit-theorem/

Multiple Regression Analyses: *Hypothesis Testing*

- We cannot prove that a given hypothesis is "correct" using hypothesis testing
- All we can do is to state that a particular sample conforms to a particular hypothesis
- We can often reject a given hypothesis with a certain degree of confidence
- In such a case, we conclude that it is very unlikely the sample result would have been observed if the hypothesized theory were correct

Multiple Regression Analyses: *Hypothesis Testing*

- Step 1: state explicitly the hypothesis to be tested
- Null hypothesis: statement of the range of values of the regression coefficient that would be expected to occur if the researcher's theory were not correct
- Alternative hypothesis: specification of the range of values of the coefficient that would be expected to occur if the researcher's theory were correct
- In other words, we define the null hypothesis as the result we do not expect

Type I and Type II Errors

- It would be unrealistic to think that conclusions drawn from regression analysis will always be right
- There are two types of errors we can make:
 - Type I: we reject a true null hypothesis
 - Type II: We fail to reject a false null hypothesis

Type I and Type II Errors

Example:

- H₀: The defendant is innocent
- H_A: The defendant is guilty
 - Type I error: sending an innocent person to jail
 - Type II error: freeing a guilty person
- Lowering the probability of Type I error means increasing the probability of Type II error;
- In hypothesis testing, we focus on Type I error and we ensure that its probability is not unreasonably large

Type I error (false positive)



Type II error (false negative)



- <u>Testing hypotheses about a single population parameter</u>
- Theorem (t-distribution for standardized estimators)

Under assumptions 1 - 6:

$$\frac{\hat{\beta}_j - \beta_j}{se(\hat{\beta}_j)} \sim t_{n-k-1}$$

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_j on y

• t-statistic (or t-ratio)

 $t_{\hat{\beta}_j} = \frac{1}{se(\hat{\beta}_j)}$

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

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

<u>Goal</u>: Define a rejection rule so that, if it is true, H₀ is rejected only with a small probability (= significance level, e.g. 5%)

Testing against one-sided alternatives (greater than zero)



Test
$$H_0: \beta_j = 0$$
 against $H_1: \beta_j > 0$

Reject the null hypothesis in favour of the alternative hypothesis if the estimated coefficient is <u>"too large"</u> (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 greater than 1.701

• Example: Wage equation

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

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

$$(.104) \quad (.007) \qquad (.0017) \qquad (.003)$$

$$n = 526, \ R^2 = .316 \qquad \text{Standard errors}$$
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.



 $c_{0.05} = 1.645$ Critical values for the 5% and the 1% significance level (these are conventional significance levels).

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."

 $c_{0.01} = 2.326$

Testing against one-sided alternatives (less than zero)



Test H_0 : $\beta_j = 0$ against H_1 : $\beta_j < 0$

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 less than -1.734

• Example: Student performance and school size

• Test whether smaller school size leads to better student performance



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



One cannot reject the hypothesis that there is no effect of school size on student performance (not even for a larger significance level of 15%).

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

• Alternative specification of functional form:

 $\widehat{math10} = -207.66 + 21.16 \log(totcomp) \\ (48.70) + (4.06) + 3.98 \log(staff) - 1.29 \log(enroll) \\ (4.19) + (0.69)$

 $n = 408, R^2 = .0654 \longleftarrow$ R-squared slightly higher Test $H_0: \beta_{\log(enroll)} = 0$ against $H_1: \beta_{\log(enroll)} < 0$.

Testing against two-sided alternatives



Test H_0 : $\beta_j = 0$ against H_1 : $\beta_j \neq 0$

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.

<u>Project if absolute value of t-statistic is less than -</u> 2.06 or greater than 2.06



• "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 <u>"statistically significant"</u>
- 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"

• 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
- <u>The fact that a coefficient is statistically significant does not necessarily mean</u> <u>it is economically or practically significant!</u>
- 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%, 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

- <u>Testing more general hypotheses about a regression coefficient</u>
- Null hypothesis



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

• Example: Campus crime and enrollment

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

$$t = (1.27 - 1)/.11 \approx 2.45 > 1.96 = c_{0.05}$$

<u>Computing p-values for t-tests</u>

- 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 H₀
- The smallest significance level at which the null hypothesis is still rejected, is called the <u>p-value</u> 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

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 - ratio| > 1.85) = 2(.0.359) = .0718

From this, it is clear that <u>a null hypothesis is</u> rejected if and only if the corresponding p-value is smaller than the significance level.

For example, for a significance level of 5% the tstatistic would not lie in the rejection region.

Inference: Confidence Intervals



- 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

Inference: Confidence Intervals

• Confidence intervals for typical confidence levels

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

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

$$P\left(\widehat{\beta}_{j} - \overbrace{c_{0.10}}^{\bullet} \cdot se(\widehat{\beta}_{j}) \leq \beta_{j} \leq \widehat{\beta}_{j} + c_{0.10} \cdot se(\widehat{\beta}_{j})\right) = 0.90$$
Use rules of thumb $c_{0.01} = 2.576, c_{0.05} = 1.96, c_{0.10} = 1.645$

• Relationship between confidence intervals and hypotheses tests

$$a_j \notin interval \Rightarrow$$
 reject $H_0 : \beta_j = a_j$ in favor of $H_1 : \beta_j \neq 0$

Inference: Confidence Intervals

Example: Model of firms' R&D expenditures



The effect of sales on R&D is relatively precisely estimated 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.

<u>Inference</u>: Testing hypotheses about a linear combination of parameters

• Example: Return to education at 2 year vs. at 4 year colleges



Test $H_0: \beta_1 - \beta_2 = 0$ against $H_1: \beta_1 - \beta_2 < 0$.

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.

Inference: Testing hypotheses about a linear combination of parameters

Impossible to compute with standard regression output because

$$se(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{\widehat{Var}(\hat{\beta}_1 - \hat{\beta}_2)} = \sqrt{\widehat{Var}(\hat{\beta}_1) + \widehat{Var}(\hat{\beta}_2) - 2\widehat{Cov}(\hat{\beta}_1, \hat{\beta}_2)}$$
Usually not available in regression output

Define $\theta_1 = \beta_1 - \beta_2$ and test H_0 : $\theta_1 = 0$ against H_1 : $\theta_1 < 0$.

$$og(wage) = \beta_0 + (\theta_1 + \beta_2)jc + \beta_2univ + \beta_3exper + u$$
$$= \beta_0 + \theta_1jc + \beta_2(jc + univ) + \beta_3exper + u$$
Insert into original regression

<u>Inference</u>: Testing hypotheses about a linear combination of parameters



• This method works <u>always</u> for single linear hypotheses

<u>Testing multiple linear restrictions: The F-test</u>

• Testing exclusion restrictions



• Estimation of the unrestricted model

$$\widehat{\log}(salary) = 11.19 + .0689 \ years + .0126 \ gamesyr \\ (0.29) \ (.0121) \ (.0026) \\ + .00098 \ bavg + .0144 \ hrunsyr + .0108 \ rbisyr \\ (.00110) \ (.0161) \ (.0072) \\ \text{None of these variabels are statistically significant when tested individually}$$

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

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

Estimation of the restricted model

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

1

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

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

• Test statistic

Number of restrictions

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

The relative increase of the sum of squared residuals when going from H₁ to H₀ follows a F-distribution (if the null hypothesis H₀ is correct)

• 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 H_1 to H_0 .

Choose the critical value so that the null hypothesis is rejected in, for example, 5% of the cases, although it is true.



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

<u>Test of overall significance of a regression</u>

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

(regression on constant)

$$H_0: \beta_1 = \beta_2 = \ldots = \beta_k = 0 \blacktriangleleft$$

 $y = \beta_0 + u$ \leftarrow Restricted model

• The null hypothesis states that the explanatory variables are not useful at all in explaining the dependent variable

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

• The test of overall significance is reported in most regression packages; the null hypothesis is usually overwhelmingly rejected