

Multiple Regression Analysis: Inference



Chapter 3 (Chapter 4 from the textbook)

Wooldridge: Introductory Econometrics:
A Modern Approach, 5e

Multiple Regression Analysis: Inference - Introduction

- **Statistical inference in the regression model**
 - Construction of confidence intervals for a coefficient
 - Hypothesis tests about population parameters
- **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 the distribution of errors: normal distribution

Multiple Regression Analysis: Inference - Introduction



■ Examples of test of hypothesis

birth weight \rightarrow $bwght = \beta_0 + \beta_1 cigs + \beta_2 educ + \beta_3 npvis + \beta_4 age + u$

cigarettes smoked per day while pregnant

total number of prenatal visits

Age of the mother

- Is the partial effect of age relevant after controlling for cigs, education and npvis?

$$H_0 : \beta_4 = 0 \quad H_1 : \beta_4 \neq 0$$

Individual statistical significance

- Is the effect of smoking 10 cigarettes canceled by the effect of one more prenatal visit?

$$H_0 : 10\beta_1 + \beta_3 = 0 \quad H_1 : 10\beta_1 + \beta_3 \neq 0$$

Single linear combination of the parameters

Multiple Regression Analysis: Inference - Introduction



■ Examples of test of hypothesis:

birth weight \rightarrow $bwght = \beta_0 + \beta_1 cigs + \beta_2 educ + \beta_3 npvis + \beta_4 age + u$

cigarettes smoked per day while pregnant

total number of prenatal visits

Age of the mother

- Are the partial effect of age, education and npvis jointly irrelevant after controlling for the number of cigarettes smoked?

$$H_0 : \beta_2 = 0, \beta_3 = 0, \beta_4 = 0 \quad H_1 : \beta_2 \neq 0 \vee \beta_3 \neq 0 \vee \beta_4 \neq 0$$

jointly statistical significance
Exclusion restrictions

- Is there any variable in the equation relevant to explain the birth weight?

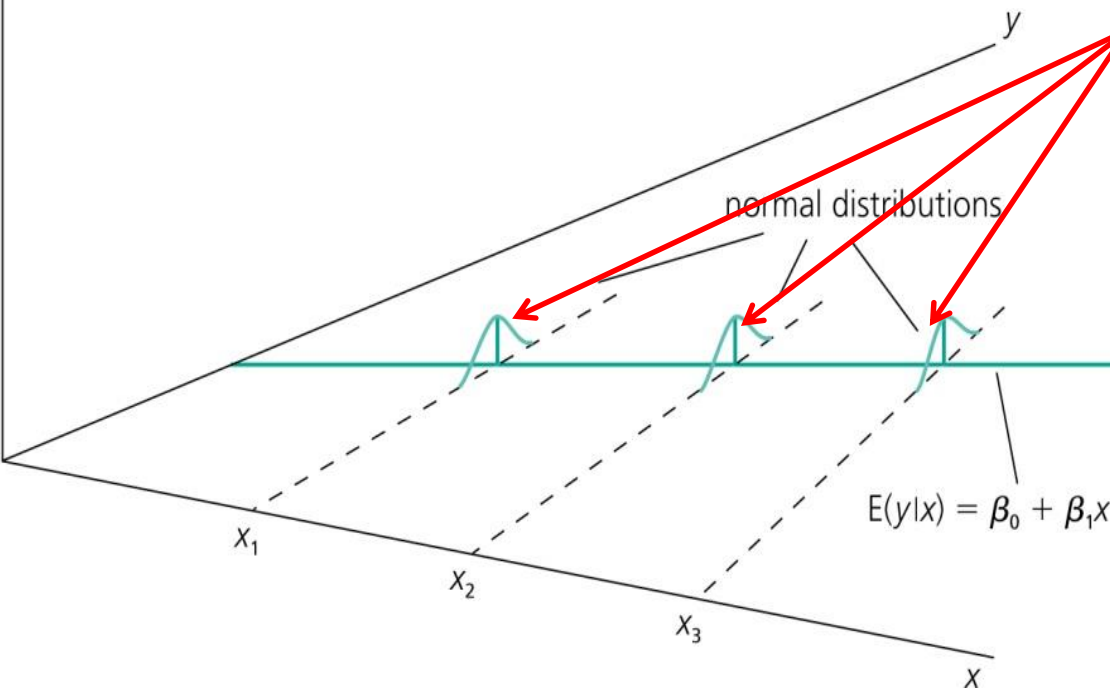
$$H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0 \quad H_1 : \beta_1 \neq 0 \vee \beta_2 \neq 0 \vee \beta_3 \neq 0 \vee \beta_4 \neq 0$$

Overall significance of
the regression

Multiple Regression Analysis: Inference – Sampling Distribution of OLS

- **Assumption MLR.6 (Normality of error terms)**

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



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.

It follows that:

$$y|x \sim N(\beta_0 + \beta_1x_1 + \dots + \beta_kx_k, \sigma^2)$$

Multiple Regression Analysis: Inference – Sampling Distribution of OLS

- **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 heterogeneous 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 – Sampling Distribution of OLS

- **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 $\log(\text{wage})$** 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 – Sampling Distribution of OLS

- **Terminology**

MLR.1 – MLR.5

„Gauss-Markov assumptions“

MLR.1 – MLR.6

„Classical linear model (CLM) assumptions“

- **Theorem 4.1 (Normal sampling distributions)**

Under assumptions MLR.1 – MLR.6:

$$\hat{\beta}_j \sim N(\beta_j, \text{Var}(\hat{\beta}_j))$$

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

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

The standardized estimators follow a standard normal distribution

Multiple Regression Analysis: Inference – Sampling Distribution of OLS

- **Theorem 4.2 (t-distribution for standardized estimators)**

Under assumptions MLR.1 – MLR.6:

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

← If the standardization is done using the estimated 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.

Multiple Regression Analysis: Inference – Confidence Intervals

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

$$P \left(\underbrace{\hat{\beta}_j - c_{0.05} \cdot se(\hat{\beta}_j)}_{\text{Lower bound of the Confidence interval}} \leq \beta_j \leq \underbrace{\hat{\beta}_j + c_{0.05} \cdot se(\hat{\beta}_j)}_{\text{Upper bound of the Confidence interval}} \right) = 0.95$$

Critical value

Confidence level

- **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 – Confidence Intervals

- **Confidence intervals for typical confidence levels**

$$P\left(\hat{\beta}_j - c_{0.01} \cdot se(\hat{\beta}_j) \leq \beta_j \leq \hat{\beta}_j + c_{0.01} \cdot se(\hat{\beta}_j)\right) = 0.99$$

$$P\left(\hat{\beta}_j - c_{0.05} \cdot se(\hat{\beta}_j) \leq \beta_j \leq \hat{\beta}_j + c_{0.05} \cdot se(\hat{\beta}_j)\right) = 0.95$$

$$P\left(\hat{\beta}_j - c_{0.10} \cdot se(\hat{\beta}_j) \leq \beta_j \leq \hat{\beta}_j + c_{0.10} \cdot se(\hat{\beta}_j)\right) = 0.90$$

If n large 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 $H_1 : \beta_j \neq a_j$

Multiple Regression Analysis: Inference



■ Example: Wage equation

- Make a 95% confidence interval for the return of education, β_1

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

(.104) (.007) (.0017) (.003)

$$n = 526, R^2 = .316$$

Standard errors

95% confidence $\rightarrow C_{0.05} = 1.96$

$\beta_1 \in (0.092 - 0.007 \times 1.96 ; 0.092 + 0.007 \times 1.96)$ with 95% of confidence

$\beta_1 \in (0.078 ; 0.106)$ with 95% of confidence

Multiple Regression Analysis: Inference – Confidence Intervals

■ Example: Model of firms' R&D expenditures

Spending on R&D Annual sales Profits as percentage of sales

$$\widehat{\log(rd)} = -4.38 + 1.084 \log(sales) + .0218 \text{ profmarg}$$

(.47) (.060) (.0217)

$$n = 32, R^2 = .918, df = 32 - 2 - 1 = 29 \Rightarrow c_{0.05} = 2.045$$

Confidence intervals at 95%:

$$1.084 \pm 2.045(.060)$$

$$= (.961, 1.21)$$

$$.0217 \pm 2.045(.0218)$$

$$= (-.0045, .0479)$$

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.

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **Testing hypotheses about a single population parameter**
- **Theorem 4.2 (t-distribution for standardized estimators)**

Under assumptions MLR.1 – MLR.6:

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

If the standardization is done using the estimated 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

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **t-statistic (or t-ratio) under H_0**

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

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. The t-statistic measures how many estimated standard deviations the estimated coefficient is away from zero.

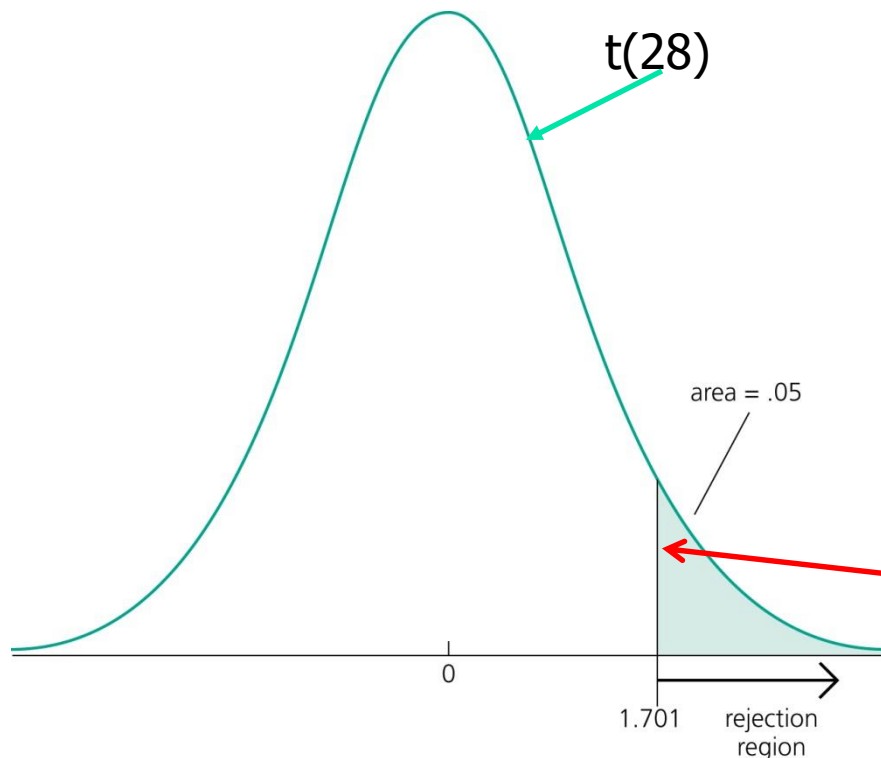
- **Distribution of the t-statistic if the null hypothesis is true**

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

- **Goal: Define a rejection rule so that, if it is true, H_0 is rejected only with a small probability (= significance level, e.g. 5%)**

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **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 „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 t-distribution with 28 degrees of freedom that is exceeded in 5% of the cases.

! Reject if t-statistic greater than 1.701

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

■ Example: Wage equation

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

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

(.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.

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

■ Example: Wage equation (cont.)

$$t_{exper} = .0041 / .0017 \approx 2.41$$

t-statistic

$$df = n - k - 1 = 526 - 3 - 1 = 522$$

Degrees of freedom;
here the standard normal
approximation applies

$$c_{0.05} = 1.645$$

Critical values for the 5% and the 1% significance level (these are conventional significance levels).

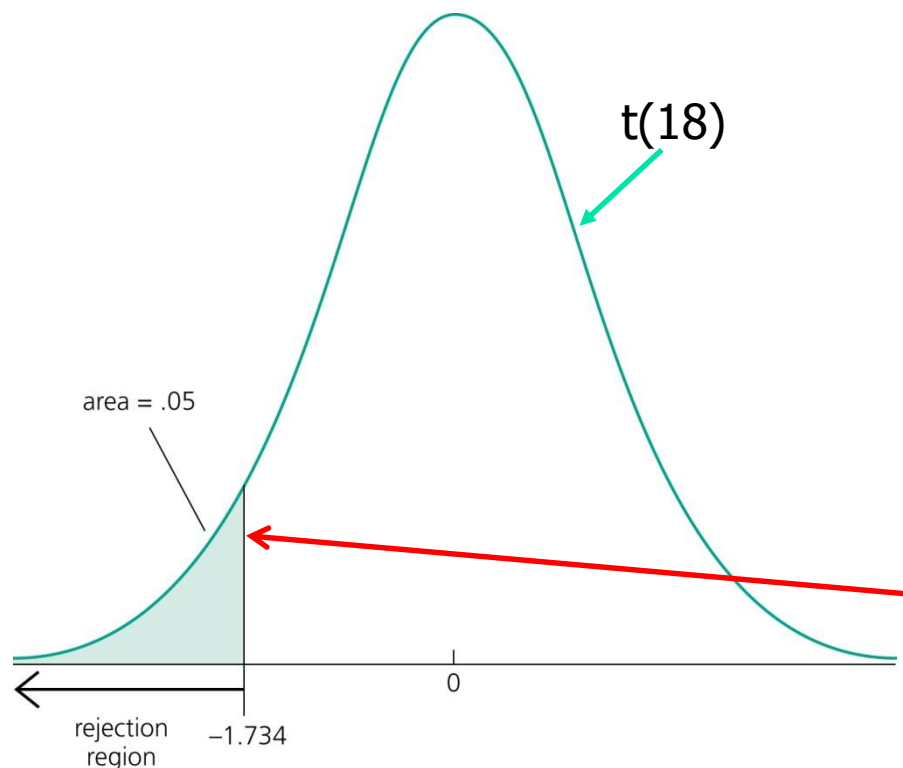
$$c_{0.01} = 2.326$$

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 – Testing hyp. about a single coeff.

■ 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 „too small“ (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 t-distribution with 18 degrees of freedom so that 5% of the cases are below the point.

! Reject if t-statistic less than -1.734

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **Example: Student performance and school size**
 - Test whether smaller school size leads to better student performance

Percentage of students passing maths test	Average annual tea- cher compensation	Staff per one thou- sand students	School enrollment (= school size)
↓	↓	↓	↓
$\widehat{math10} = + 2.274 + .00046 \text{ totcomp} + .048 \text{ staff} - .00020 \text{ enroll}$ <p style="text-align: center;">$(6.113) \quad (.00010) \quad (.040) \quad (.00022)$</p>			

$$n = 408, R^2 = .0541$$

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 – Testing hyp. about a single coeff.

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

$$t_{enroll} = -.00020 / .00022 \approx -.91$$

t-statistic

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

Degrees of freedom; here the standard normal approximation applies

$$c_{0.05} = -1.65$$

Critical values for the 5% and the 15% significance level.

$$c_{0.15} = -1.04$$

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 – Testing hyp. about a single coeff.

- **Example: Student performance and school size (cont.)**

- Alternative specification of functional form:

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

$$n = 408, R^2 = .0654 \leftarrow \text{R-squared slightly higher}$$

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

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **Example: Student performance and school size (cont.)**

$$t_{\log(enroll)} = -1.29 / .69 \approx -1.87 \quad \leftarrow \text{t-statistic}$$

$$c_{0.05} = -1.65 \quad \leftarrow \text{Critical value for the 5\% significance level ! reject null hypothesis}$$

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?

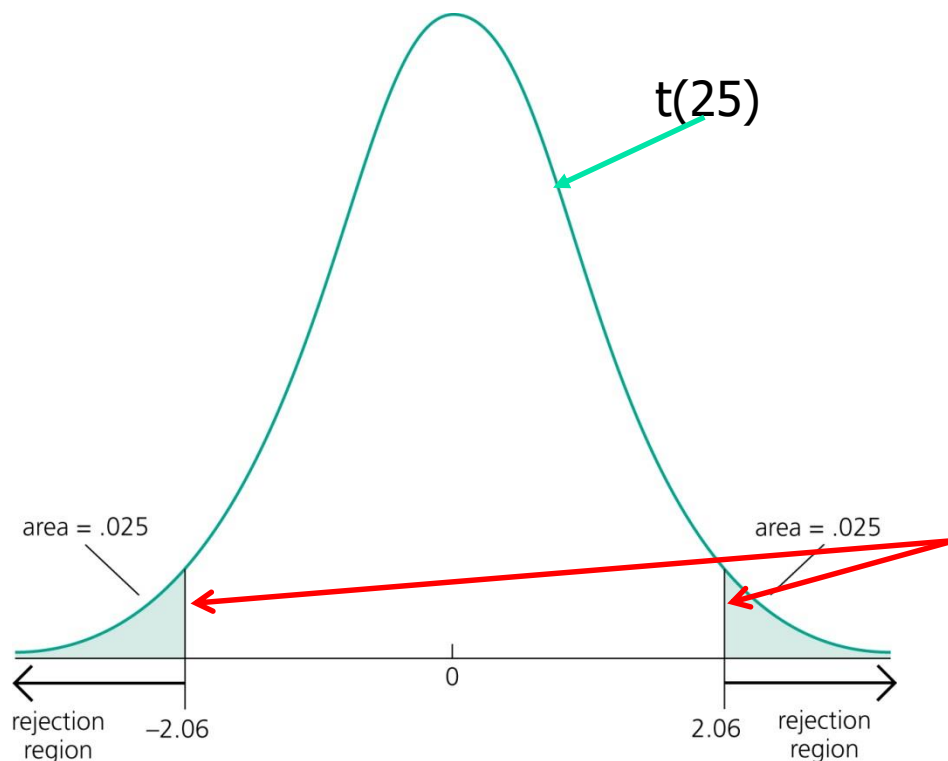
+ 10% enrollment ! -0.129 percentage points
students pass test

$$-1.29 = \frac{\partial math10}{\partial \log(enroll)} = \frac{math10}{\frac{\partial enroll}{enroll}} = \frac{\frac{-1.29}{100}}{\frac{1}{100}} = \frac{-0.0129}{+1\%}$$

(small effect)

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

■ 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 the absolute value 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 t-statistic is less than -2.06 or greater than 2.06, that is, if absolute value of t-statistic is greater than 2.06

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

■ Example: Determinants of college GPA

Lectures missed per week

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

(.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.

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **„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“

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **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%, 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 – Testing hyp. about a single coeff.

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

$$H_0 : \beta_j = a_j$$

Hypothesized value of the coefficient

- **t-statistic**

$$t = \frac{(\text{estimate} - \text{hypothesized value})}{\text{standard error}} = \frac{(\hat{\beta}_j - a_j)}{se(\hat{\beta}_j)} \sim t(n - k - 1)$$

- **The test works exactly as before, except that the hypothesized value is subtracted from the estimate when forming the statistic**

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

■ Example: Campus crime and enrollment

- An interesting hypothesis is whether crime increases by one percent if enrollment is increased by one percent

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

$(1.03) \quad (0.11)$

$$n = 97, R^2 = .585$$

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

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

Estimate is different from one but is this difference statistically significant?

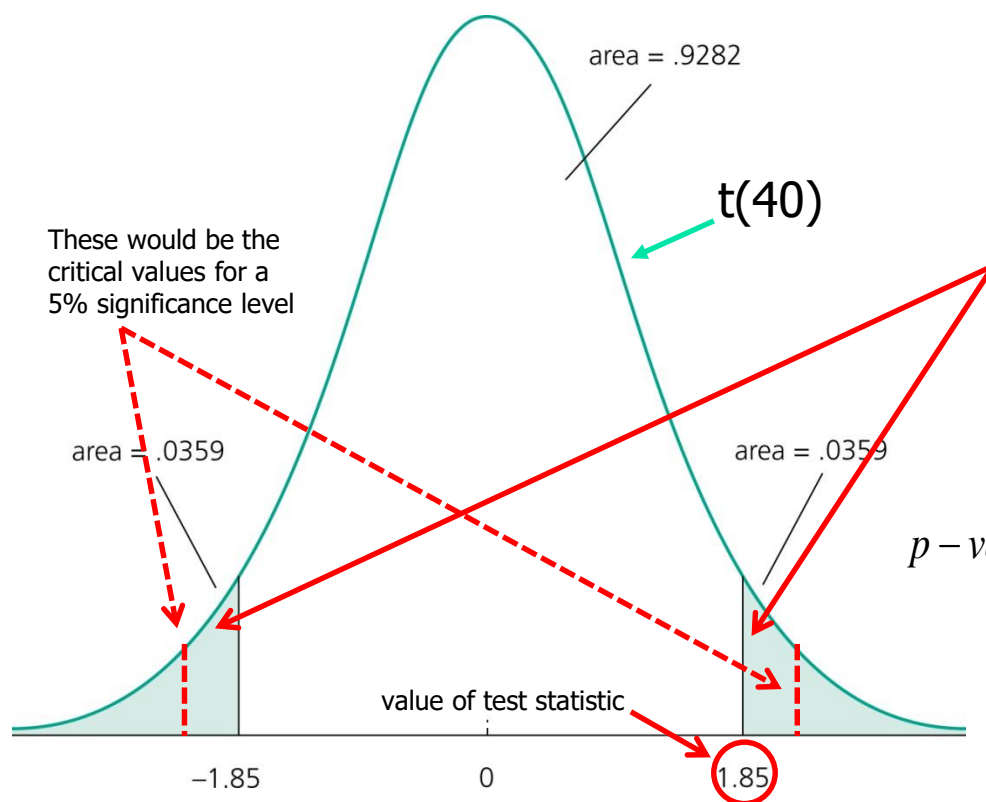
The hypothesis is rejected at the 5% level

Multiple Regression Analysis: Inference – Testing hyp. about a single coeff.

- **Computing p-values for t-tests**
 - **The smallest significance level at which the null hypothesis is still rejected, 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 – Testing hyp. about a single coeff.

■ 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., suppose that the realized t-ratio is $t = -1.85$, then,

$$p\text{-value} = 2P(t(n-k-1) > |-1.85|)$$

$$p\text{-value} = 2P(t(40) > |-1.85|) = 2 \times 0.0359 = 0.0718$$

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

- **Testing hypotheses about a linear combination of parameters**
- **Example: Return to education at 2 year vs. at 4 year colleges**

Years of education
at 2 year colleges

Years of education
at 4 year colleges

$$\log(\text{wage}) = \beta_0 + \beta_1 jc + \beta_2 univ + \beta_3 exper + u$$

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

A possible test statistic would be:

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_2}{se(\hat{\beta}_1 - \hat{\beta}_2)}$$

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

- Impossible to compute with standard regression output because

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

- **Alternative method**

Usually not available in regression output

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

$$\log(wage) = \beta_0 + (\theta_1 + \beta_2)jc + \beta_2univ + \beta_3exper + u$$

$$= \beta_0 + \theta_1jc + \beta_2(jc + univ) + \beta_3exper + u$$

Insert into original regression

a new regressor (= total years of college)

Multiple Regression Analysis: Inference

- **Estimation results**

$$\widehat{\log}(wage) = 1.472 + .0102 jc + .0769 totcoll + .0049 exper$$

jc + univ ↗

(0.021) (0.0069) (0.0023) (0.0002)

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

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

$$p\text{-value} = P(t\text{-ratio} < -1.48) = .070$$

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

Hypothesis is rejected at 10% level but not at 5% level

- **This method works always for single linear hypotheses**

Multiple Regression Analysis: Inference

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

Salary of major league base ball player

Years in the league

Average number of games per year

$$\log(\text{salary}) = \beta_0 + \beta_1 \text{years} + \beta_2 \text{gamesyr}$$

$$+ \beta_3 \text{bavg} + \beta_4 \text{hrunsyr} + \beta_5 \text{rbisyr} + u$$

Batting average

Home runs per year

Runs batted in per year

$$H_0 : \beta_3 = 0, \beta_4 = 0, \beta_5 = 0 \quad \text{against} \quad H_1 : H_0 \text{ is not true}$$

Test whether performance measures have no effect/can be excluded from regression.

Multiple Regression Analysis: Inference

- Estimation of the unrestricted model

$$\widehat{\log}(\text{salary}) = 11.19 + .0689 \text{ years} + .0126 \text{ gamesyr} \\ + .00098 \text{ bavg} + .0144 \text{ hrunsyr} + .0108 \text{ rbisyr}$$

(0.29) (.0121) (.0026)
(.00110) (.0161) (.0072)

None of these variables 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

- **Estimation of the restricted model**

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

(0.11) (.0125) (.0013)

$$n = 353, \text{ SSR} = 198.311, R^2 = .5971$$

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

- **Test statistic**

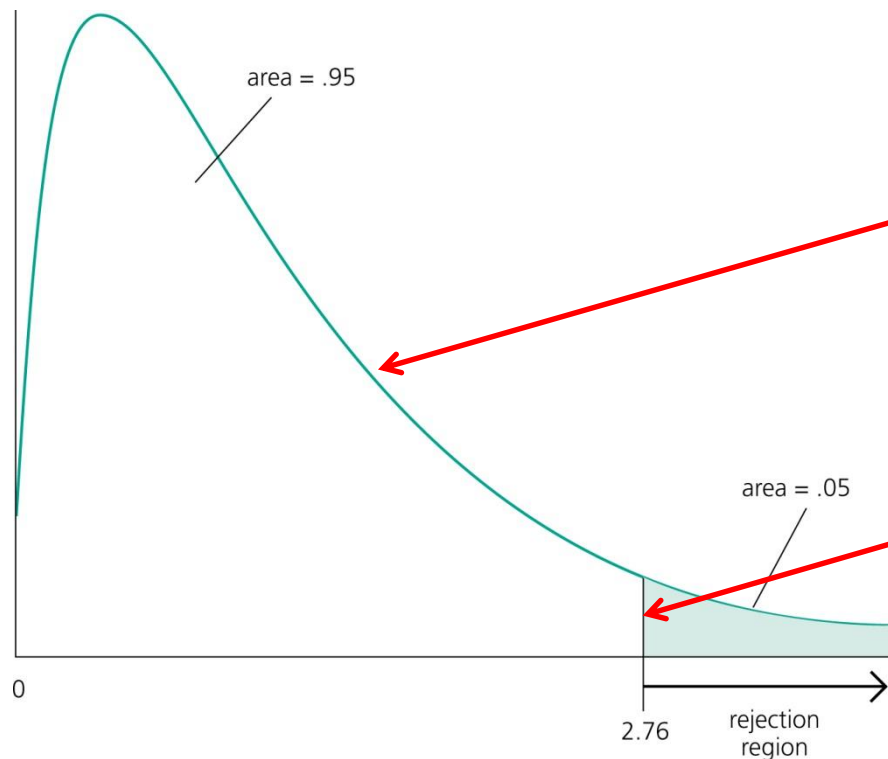
Number of restrictions

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

The relative increase of the sum of squared residuals when going from H_1 to H_0 follows a F-distribution (if the null hypothesis H_0 is correct)

Multiple Regression Analysis: Inference

■ Rejection rule (Figure 4.7)



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.

Multiple Regression Analysis: Inference

■ Test decision in example

$$F = \frac{(198.311 - 183.186) / 3}{183.186 / (353 - 5 - 1)} \approx 9.55$$

Number of restrictions to be tested

Degrees of freedom in the unrestricted model

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

$$p\text{-value} = P(F(3,347) > 9.55) = 0.000$$

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

■ 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

- **Test of overall significance of a regression**

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

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

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

$$y = \beta_0 + u$$

← Restricted model
(regression on constant)

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

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

Multiple Regression Analysis: Inference

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

Actual house price The assessed housing value (before the house was sold) Size of lot (in feet)

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{assess}) + \beta_2 \log(\text{lotsize})$$
$$+ \beta_3 \log(\text{sqr ft}) + \beta_4 \text{bdrms} + u$$

 Square footage Number of bedrooms

$$H_0 : \beta_1 = 1, \beta_2 = 0, \beta_3 = 0, \beta_4 = 0$$

If house price assessments are rational, a 1% change in the assessment should be associated with a 1% change in price.

In addition, other known factors should not influence the price once the assessed value has been controlled for.

Multiple Regression Analysis: Inference

- **Unrestricted regression**

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_3 x_3 + \beta_4 x_4 + u$$

- **Restricted regression**

The restricted model is actually a regression of $[y - x_1]$ on a constant

$$y = \beta_0 + x_1 + u \Rightarrow [y - x_1] = \beta_0 + u$$

- **Test statistic**

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

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

Multiple Regression Analysis: Inference

- **Regression output for the unrestricted regression**

$$\widehat{\log(price)} = .264 + 1.043 \log(assess) + .0074 \log(lotsize) \\ - .1384 \log(sqrft) + .0338 bdrms$$

(.570) (.151) (.0386) (.1032) (.0221)

When tested individually, there is also no evidence against the rationality of house price assessments

$$n = 88, SSR = 1.822, R^2 = .773$$

- **The F-test works for general multiple linear hypotheses**
- **For all tests and confidence intervals, validity of assumptions MLR.1 – MLR.6 has been assumed. Tests may be invalid otherwise.**