

## Chapter 4

# Multiple Regression Analysis: Inference



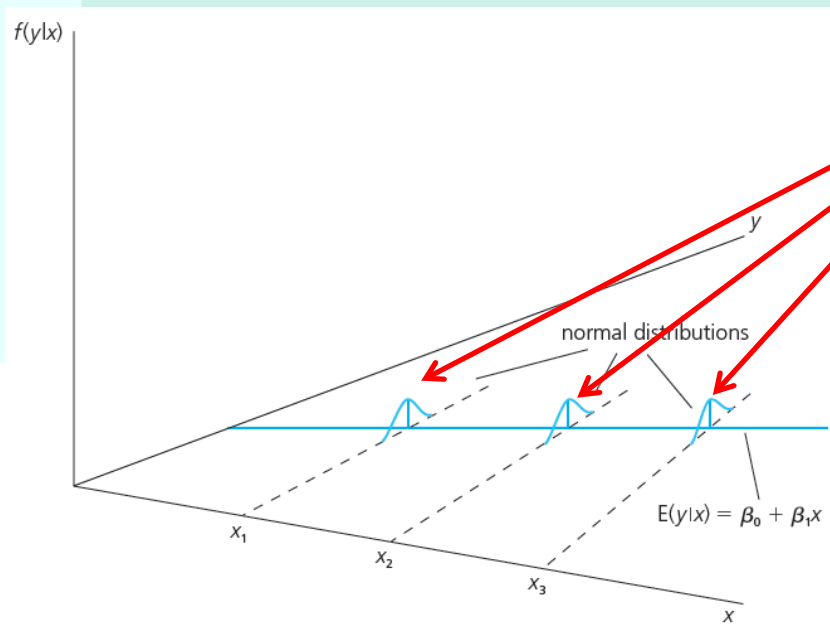
# Multiple Regression Analysis: Inference

- **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
  - 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

- Assumption MLR.6 (Normality of error terms)

$u_i \sim \text{Normal}(0, \sigma^2)$  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.

It follows that:

$$y|x \sim \text{Normal}(\beta_0 + \beta_1x_1 + \dots + \beta_kx_k, \sigma^2)$$



# Multiple Regression Analysis: Inference

- **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

- **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

- Terminology

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

"Gauss-Markov assumptions"

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

"Classical linear model (CLM) assumptions"

- Theorem 4.1 (Normal sampling distributions)

Under assumptions MLR.1 – MLR.6:

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

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

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

↑  
The standardized estimators follow a standard normal distribution

# Multiple Regression Analysis: Inference

- Testing hypotheses about a single population parameter
- Theorem 4.2 (t-distribution for the 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

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

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

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**

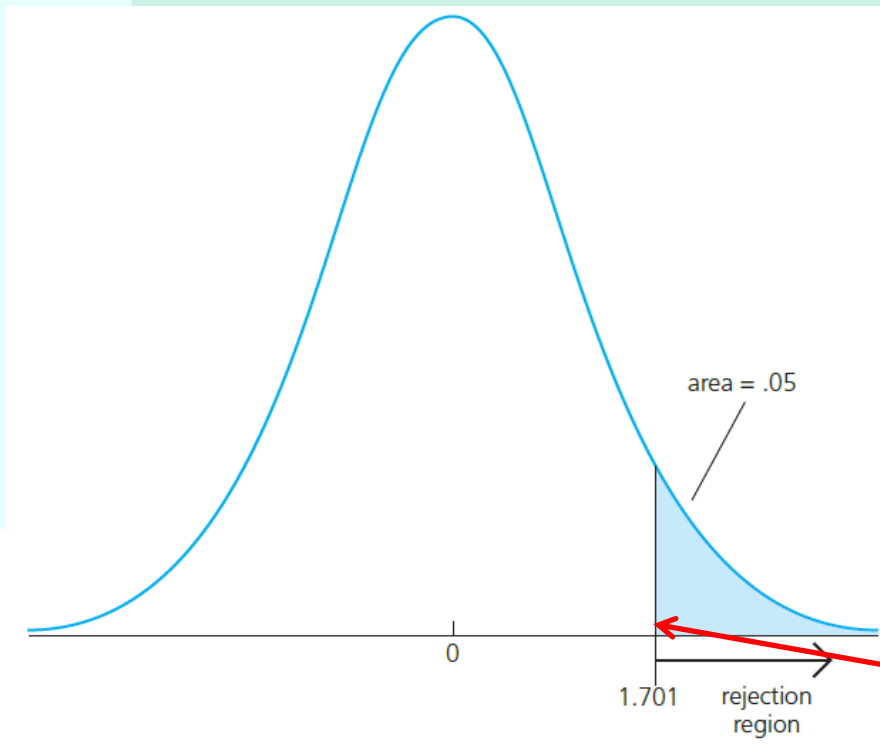
$$t_{\hat{\beta}_j} \equiv \hat{\beta}_j / se(\hat{\beta}_j) = (\hat{\beta}_j - \beta_j) / se(\hat{\beta}_j) \sim t_{n-k-1}$$

- **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 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 is greater than 1.701

# Multiple Regression Analysis: Inference

- **Example: Wage equation**

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

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

(.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

- **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).

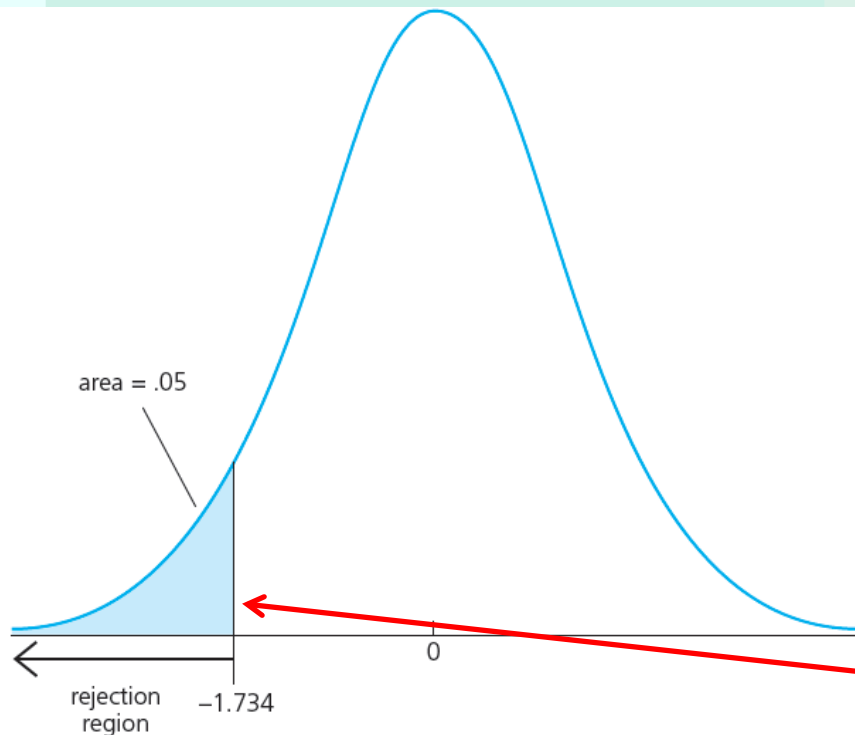
$$c_{0.01} = 2.326$$

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 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 is less than -1.734

# Multiple Regression Analysis: Inference

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

Percentage of students passing maths test

Average annual teacher compensation

Staff per one thousand students

Student enrollment (= school size)

$$\widehat{math10} = + 2.274 + .00046 \text{ totcomp} + .048 \text{ staff} - .00020 \text{ enroll}$$

(6.113)    (.00010)    (.040)    (.00022)

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

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

- 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

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

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

$n = 408, R^2 = .0654$  ← R-squared slightly higher

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

# Multiple Regression Analysis: Inference

- 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

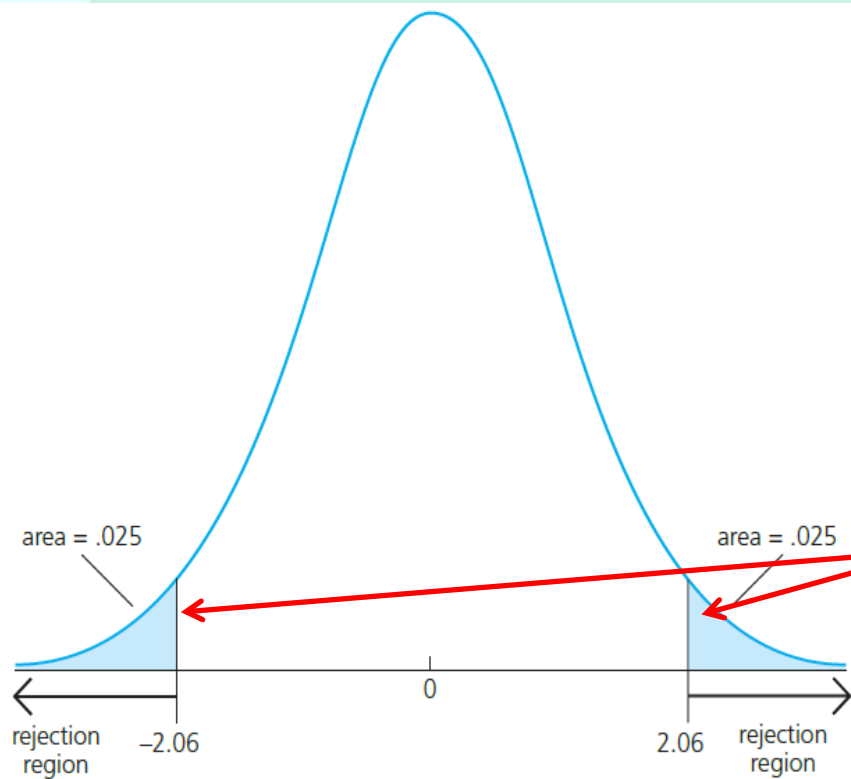
$$-1.29 = \frac{\Delta \widehat{math10}}{\Delta \log(enroll)} = \frac{\Delta \widehat{math10}}{\frac{\Delta enroll}{enroll}} = \frac{\frac{-1.29}{100}}{\frac{1}{100}} = \frac{-0.0129}{+1\%}$$

(small effect)



# Multiple Regression Analysis: Inference

- 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 absolute value of t-statistic is less than -2.06 or greater than 2.06

# Multiple Regression Analysis: Inference

- Example: Determinants of college GPA

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



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

Lectures missed per week  



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

For critical values, use standard normal distribution

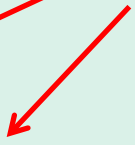
$$t_{hsGPA} = 4.38 > c_{0.01} = 2.58$$

$$t_{ACT} = 1.36 < c_{0.10} = 1.645$$



$$|t_{skipped}| = |-3.19| > c_{0.01} = 2.58$$



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

- **“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

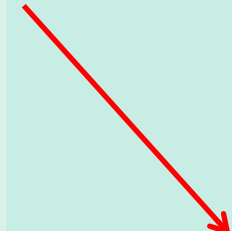
- Testing more general hypotheses about a regression coefficient
- Null hypothesis

$$H_0 : \beta_j = a_j$$

Hypothesized value of the coefficient



- t-statistic

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


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

# Multiple Regression Analysis: Inference

- **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

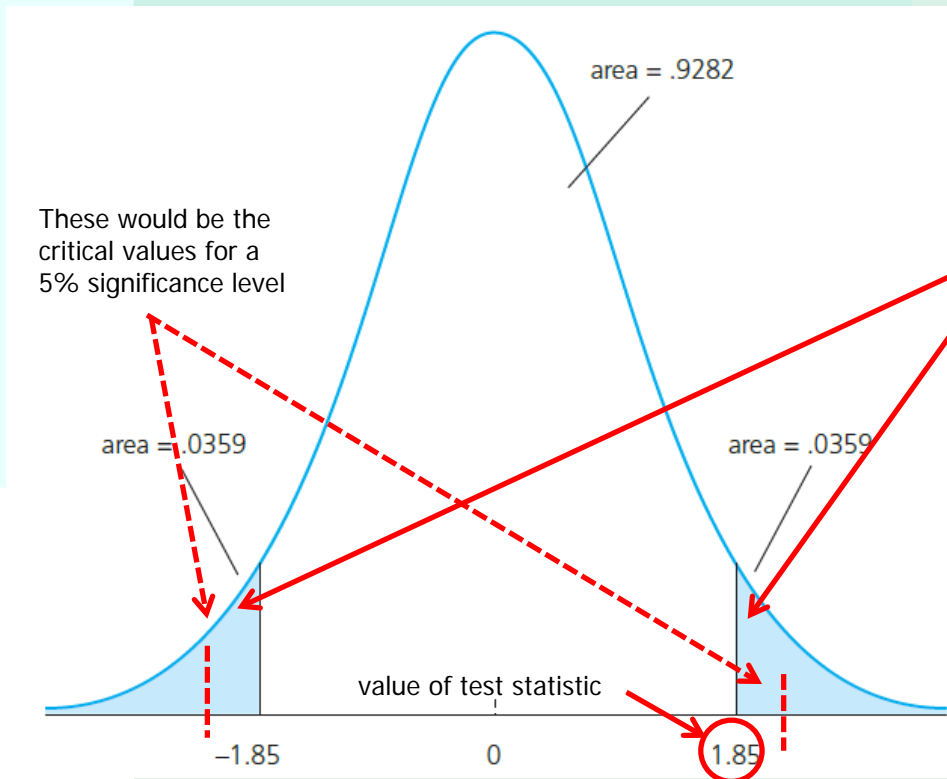


# Multiple Regression Analysis: Inference

- **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  $H_0$
  - 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

- 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.:

$$P(|t - ratio| > 1.85) = 2(.0359) = .0718$$

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

- **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

- **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 of two-sided test

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 for typical confidence levels

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

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

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

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

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

# Multiple Regression Analysis: Inference

- Example: Model of firms' R&D expenditures

Spending on R&D

Annual sales

Profits as percentage of sales

$$\widehat{\log(rd)} = -4.38 + 1.084 \log(sales) + .0217 \text{ profmarg}$$

(.47)            (.060)            (.0128)

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

$$1.084 \pm 2.045(.060) \quad .0217 \pm 2.045(.0218)$$

$$= (.961, 1.21)$$

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

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

Years of education  
at two year  
colleges

Years of  
education at four  
year colleges

Months in the  
workforce

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

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

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

- 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

$$\widehat{\log(wage)} = 1.472 + \underbrace{-.0102}_{(.0069)}jc + \underbrace{.0769}_{(.0023)}\overset{\text{Total years of college}}{totcoll} + \underbrace{.0049}_{(.0002)}exper$$

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

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

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

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

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 baseball 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

$$\begin{aligned}\widehat{\log(salary)} = & 11.19 + .0689 \text{ years} + .0126 \text{ gamesyr} \\ & (0.29) \quad (.0121) \quad (.0026) \\ & + .00098 \text{ bavg} + .0144 \text{ hrunsyr} + .0108 \text{ rbisyr} \\ & (.00110) \quad (.0161) \quad (.0072)\end{aligned}$$

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

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

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

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

- Test statistic

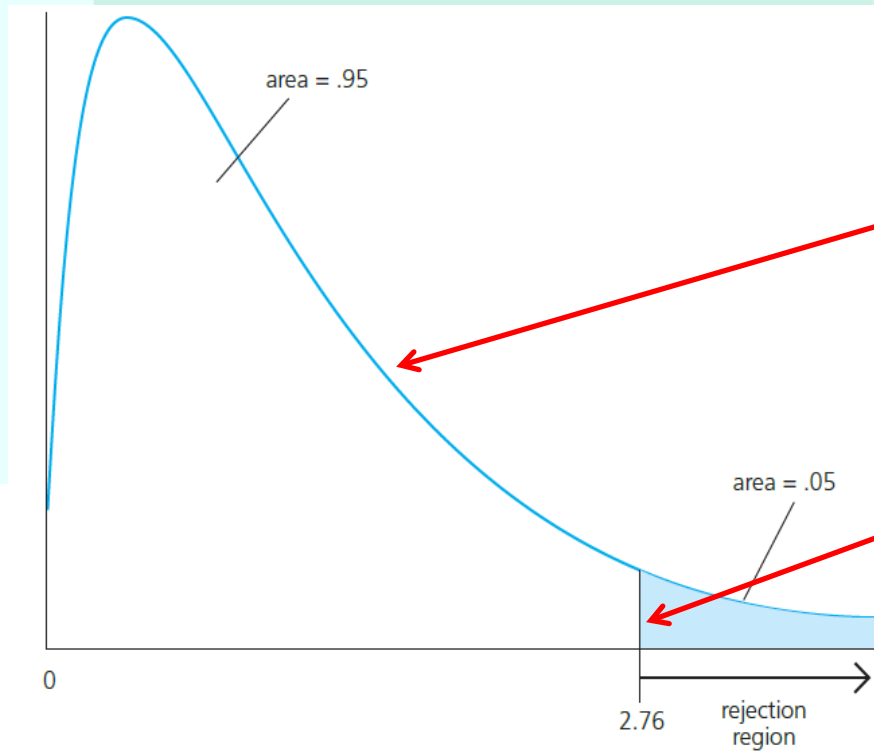
Number of restrictions

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

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

# Multiple Regression Analysis: Inference

- **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.

# Multiple Regression Analysis: Inference

- **Test decision in example**

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

Number of restrictions to be tested

Degrees of freedom in the unrestricted model

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

$$P(F - statistic > 9.55) = 0.000$$

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 square feet)

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

   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**

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

- **Restricted regression**

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

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

- **Test statistic**

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

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



# Model specification and Diagnostic testing





# Agenda

- **Model selection criteria**
- **Types of specification errors & their consequences**
- **Tests of specification errors**
- **Errors of measurement**
- **Nested and non nested hypothesis**
- **Model selection criteria**



## What are the criteria in choosing a model for empirical analysis ?

**Model should satisfy following criteria:**

- **Be data admissible; that is.....**  
**Be consistent with the theory; that is**
- **Have weakly exogenous regressors; that is .....**
- **Exhibit parameter constancy; that is .....**
- **Exhibit data coherency; that is .....**
- **Be encompassing.....**



# Types of specification errors

- **Omission of a relevant variable (s)**
- **Inclusion of unnecessary models**
- **Adopting the wrong functional form**
- **Errors of measurement**



# Types of the specification errors

- Model specification errors. We have in mind the “true” model but somehow we do not estimate the correct model
- Model misspecification errors. We do not know what the true model is to begin with.



# Types of specification errors

## **Underfitting a model**

- **Omission of a relevant variable (s)**
- **Inclusion of unnecessary models**

## **Overfitting the model**

- **Adopting the wrong functional form**
- **Errors of measurement**



# Consequences of Model specification errors

# Omitting a relevant variable

Suppose the true model has the following form:

$$Y_i = \beta_1 + \beta_2 X_i + \beta_3 X_i^2 + \beta_4 X_i^3 + e_{1i}$$

Suppose researcher decides to use

$$Y_i = \alpha_1 + \alpha_2 X_{2i} + \beta_3 X_{3i}^2 + e_{2i}$$

Since the first model is true, adopting second model will lead to specification error, error constituting omitting a relevant variable  $X_i^3$ . Therefore, the error term in second model is in fact:

$$e_{2i} = \beta_1 X_i^3 + e_{1i}$$

# Consequences of omitting a relevant Variable

- If  $r_{23} \neq 0$  then  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  are both biased and inconsistent estimates.  $E(\hat{\alpha}_1) \neq \beta_1$   $E(\hat{\alpha}_2) \neq \beta_2$
- If  $r_{23} = 0$  then  $\hat{\alpha}_1$  biased  $\hat{\alpha}_2$  unbiased
- The disturbance variance  $\sigma^2$  incorrectly estimated
- Conventionally measured variance of  $\hat{\alpha}_2 (= \sigma^2 / \sum x_{2i}^2)$  is a biased estimator of the variance of true estimator  $\beta_2^{\text{est}}$
- Misleading conclusions about the statistical significance of estimated parameters
- Unreliable forecasts

Once a model is formulated on the basis of relevant theory one is ill advised to drop a variable from such a model.

## ➤ Inclusion of an irrelevant variable

**Suppose the true model has the following form:**

$$Y_i = \alpha_1 + \alpha_2 X_{2i} + e_{2i}$$

**Suppose researcher decides to use**

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + e_{1i}$$


**Since the first model is true, adopting second model will lead to specification error, error constituting inclusion of an irrelevant variable  $X_{3i}$ .**



# Consequences of inclusion of an irrelevant Variable

- OLS estimates are unbiased and consistent
- The error variance  $\sigma^2$  correctly estimated
- Usual confidence interval and test of hypothesis remain valid
- Estimated  $\alpha$ 's are generally inefficient, that is their variances will be larger than those of  $\beta$

(the implication of this finding is that inclusion of unnecessary variable  $X_3$  makes the variance of  $\hat{\alpha}_2$  larger than necessary, thereby making  $\hat{\alpha}_2$  less precise)



## Is it better to include irrelevant variables than omit relevant ones?

- **Addition of variables will lead to loss in efficiency of estimators and may also lead to problem of multicollinearity, not mentioning the loss in the degrees of freedom**
- **The best approach is to include only explanatory variables, that on the theoretical grounds, directly influence the dependent variable and that are not accounted for by other included variables.**

# Tests of specification errors

- Data mining approach
- Durbin Watson D statistic
- Ramsey`s RESET test
- Lagrangean Multiplier (LM) test
- Tests of incorrect functional form and omitted variable

# Data mining approach

## Detecting the presence of unnecessary variables

Suppose we develop a model where we are not sure about  $X_{ki}$

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + e_{1i}$$

To find that we look at:

1. Test the significance of  $\beta_k$ ,
2. Or we if we are not sure about two random variables then we look at F test

However t test and F test should not be used for model building. In the above approach we assume that we know the model.

The approach of starting with a smaller model and step by step expansion of that model is called data mining





# Tests of incorrect functional form and omitted variable

- In practice we never sure that the model adapted is the true model. We look at
- Signs of estimated coefficients
- Statistical significance
- R squared and adjusted R squared
- F test
- DW statistics

# Examination of the residuals

- Examination of the residuals is a good visual diagnostic to detect autocorrelation or heteroscedasticity and also to detect omitted variable and incorrect functional form

For instance consider cubic cost function

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i}^2 + \beta_4 X_{4i}^3 + e_{1i}$$

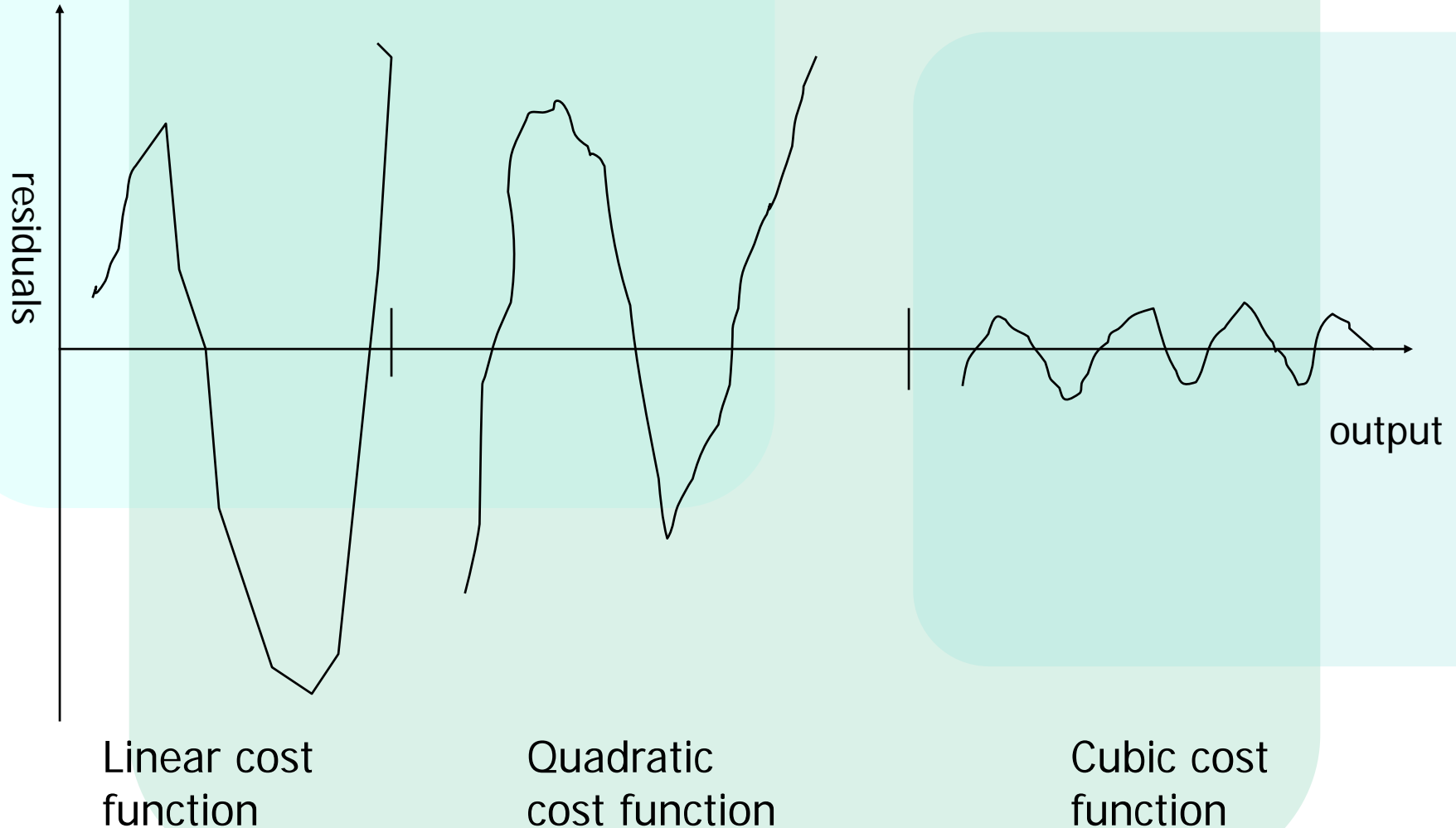
Suppose researcher estimates

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i}^2 + e_{1i}$$

An other researcher estimates

$$Y_i = \beta_1 + \beta_2 X_{2i} + e_{1i}$$

# Examination of the residuals



# Durbin Watson D statistic

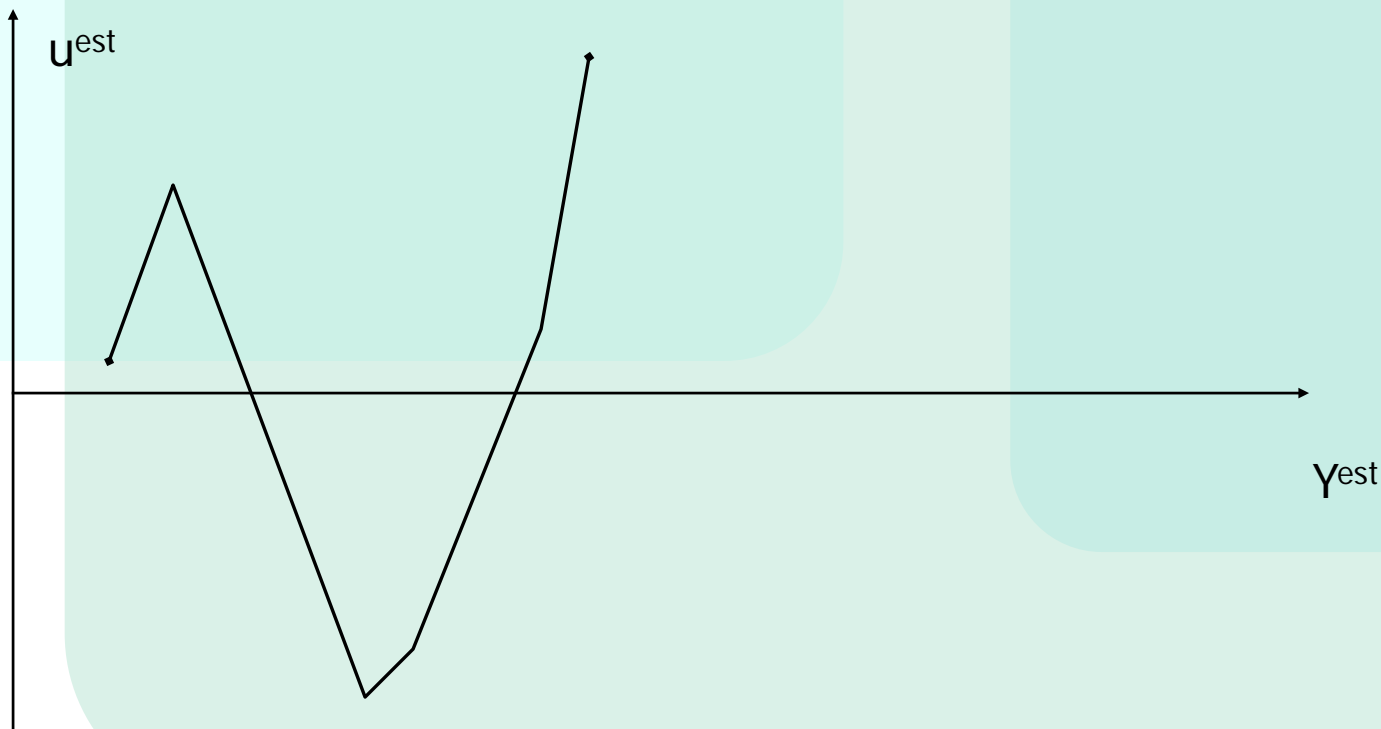
1. From the observed data obtain OLS estimates
2. If it is believed that the model is mis-specified because it excludes a relevant explanatory variable, say Z, order the residuals according increasing values of Z
3. Compute d statistic
$$d = \frac{\sum_{t=2}^n (u_t^{\text{est}} - u_{t-1}^{\text{est}})^2}{\sum_{t=1}^n (u_t^{\text{est}})^2}$$
4. From the Durbin Watson tables find out whether estimated d values is statistically significant, then one can accept the hypothesis of model mis-specification



# Ramsey's RESET (regression specification error test)

Assume cost function is linear in output where Y total cost and X is the output

$$Y_i = \lambda_1 + \lambda_2 X_{2i} + e_{1i}$$





# The idea behind RESET test

**From the graph one can see that the mean of error term systematically changes with  $Y$ . this suggests that if we introduce estimated  $Y$  in the form of the regressor then it should increase  $R$  squared. And if on the basis of  $F$  test increase in  $R$  squared is statistically significant then it would suggest that the model is misperceived.**

# Steps of Reset test

- From the chosen model obtain estimated  $Y_i$
- Rerun the regression introducing the estimated  $Y_i$  thus we run

$$Y_i = \lambda_1 + \lambda_2 X_{2i} + \lambda_3 (Y^{\text{est}})^2 + \lambda_4 (Y^{\text{est}})^3 + e_{1i}$$

- let  $R^2$  obtained from the above model be  $R^2_{\text{new}}$  and from previous model  $R^2_{\text{old}}$ . Then we use F test

$$F = \frac{(R^2_{\text{new}} - R^2_{\text{old}})}{\text{number of new regressors}} : \frac{(1 - R^2_{\text{new}})}{(n - \text{number of parameters in new model})}$$

- If the computed F value is significant, say, at 5 percent level, one can accept the hypothesis that the model is mis-specified

# LM test for adding the variables

- Estimate the restricted regression by OLS and obtain the residuals,  $u^{est}$
- If in fact the unrestricted regression is the true regression, the residuals obtained from the restricted regression should be related to the squared and cubed output terms, that is  $X^2_i$  and  $X^3_i$
- Regress the  $e^{est}$  on all regressors

$$e^{est}_{1i} = \lambda_1 + \lambda_2 X_{2i} + \lambda_3 X^2 + \lambda_4 X^3 + v$$

- For large sample
- $nR^2 \sim$  chi squared distribution with (number of restrictions)
- If the Chi squared value obtained exceeds critical chi squared value at the chosen level of significance we reject the restricted regression



# Errors in measurement

# Errors in measurement of dependent variable

$$Y^* = \alpha + \beta X_i + e_i$$

$Y^*$  - permanent income

$X$ - current income

$$Y^* = Y_i + u_i$$

$$Y^* = (\alpha + \beta X_i + e_i) + u_i$$

$$v = e_i + u_i$$

Errors in measurement of dependent variable

give unbiased estimates however with larger

# Errors in independent variable

$$Y_i = \alpha + \beta X_i^* + e_i$$

$Y_i$  – current consumption expenditure

$X_i^*$  – permanent income

$$X_i = X_i^* + w_i$$

$$Y^* = (\alpha + \beta X_i + e_i) + w_i$$

$$v = e_i + \beta w_i$$

**Errors in measurement of independent variable give both biased estimates and inefficient estimates**

# Model selection criteria

- R Square:**  $R^2 = \frac{ESS}{TSS} = \frac{\sum (\hat{Y} - \bar{Y})^2}{\sum (Y - \bar{Y})^2}$   $TSS = ESS + RSS$

- Adjusted R square**  $\bar{R}^2 = 1 - \frac{RSS / (n - k)}{TSS / (n - 1)} = 1 - (1 - R^2) \frac{n - 1}{n - k}$

$$AIC = e^{2k/n} \frac{\sum u_i^2}{n} = e^{2k/n} \frac{RSS}{n}$$

- AIC**

$$SIC = n^{k/n} \frac{\sum u_i^2}{n} = n^{k/n} \frac{RSS}{n}$$





# Reading

- **Gujarati chapter 13**