

Chapter # 04

→ Linear Regression with one regressor

→ Regression (Simple)

Simple linear regression is a regression model that estimates the relationship between one independent variable and one dependent variable.

Simple linear regression gets its adjective "Simple", because it concerns the study of only one predictor variable. In contrast multiple linear regression, which we will study later in this course, gets its adjective "multiple" because it concerns the study of two or more predictor variables.

In simple regression model we talked about dependent (y) and independent variables (x) for which we have some examples such as:

X (Independent)

- i) Income \rightarrow expenditure
- ii) Study hours \rightarrow Marks
- iii) Height of parents \rightarrow Height of children

Y (dependent)

to follow and to follow

Income

50K

55K

75K

82K

Here

we can be questioned
that if a professional
is 81K so what will be its
expenditure which is surely dependent
on income.

So we can estimate or predict
our expenditure from our income.

So from the above given examples
we can define simple regression
such as in which we can

check the dependence of a
dependent variable on another
independent variable.

⇒ The Linear regression model or equations

$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

i) Y_i = Dependent variable
ii) β_0 = Intercept (is the value of
a dependent variable when independent)

Expenditure

40K

42K

50K

?

Variable is zero) \Rightarrow present only theoretically.
But, note that "When there is no predictor don't predict."

$$Y_i = \beta_0 + \beta_1 X_i$$

↪ independent variable, $X_i = 0$ so.

$$Y_i = \beta_0 + \beta_1(0)$$

↪ intercept.

iii) β_1 = Regression coefficient or slope. (. tells the rate of change in dependent variable with unit change in independent variable.)

$$\beta_{\text{class size}} = \frac{\text{Change in test score } \Delta Y}{\text{Change in class size } \Delta X}$$

$$\beta_1 = \frac{\Delta Y}{\Delta X} \rightarrow \begin{matrix} \text{dependent.} \\ \text{independent.} \end{matrix}$$

iv) X_i = Independent variable

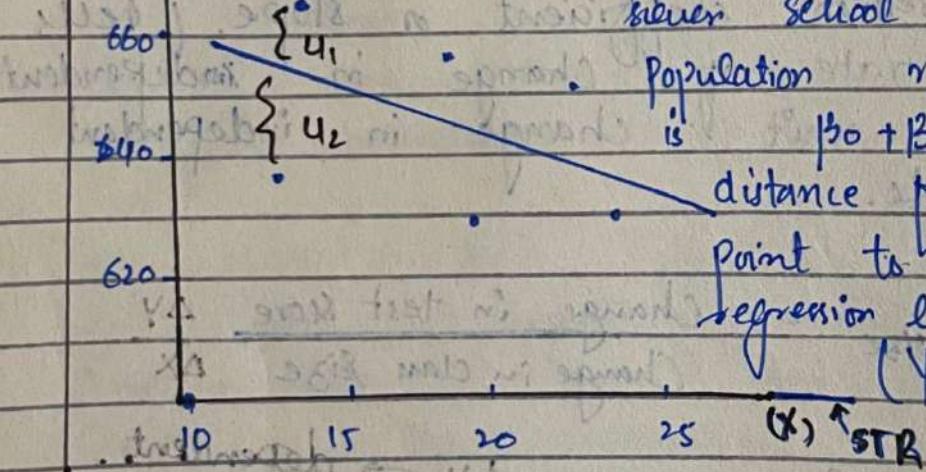
v) $U_{i,n}$ = Error term.

According to the equation, if you knew β_0 and β_1 , not only would you be able to determine the change in dependent variable (test score) at a district associated with a change in independent

variable (class size), but you would also be able to predict the average test score itself for a given class size.

02 Scatterplot of Test Scores vs. Student - Teacher Ratio

Test score (y_i)



The scatterplot shows hypothetical observations for seven school districts. The population regression line is $\beta_0 + \beta_1 x$. The vertical distance from the i th point to the population regression line is $(y_i - \beta_0 - \beta_1 x)$,

which is the population error term u_i for the i th observation.

In the about scatterplot the error term u_1 is positive because the test score in district # 1 is better than the predicted by the population regression line. And y_2 is below the population regression line.

so the test score for that district were worse than predicted, and U_{210} means it is negative.

→ The linear regression model general equation is

$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

It is a general equation model based on population data and when we apply this equation on sample data then we will have fitted model.

Before that we should know the difference b/w parameter and estimator.

When we calculate any value from population it is called parameter.

Whereas when we calculate any value from sample then it is called estimator.

Fitted models

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{U}_i$$

→ \hat{Y} on X equation.

$$\hat{Y}_i = \text{Estimated value of } Y_i$$

→ observed value.

\hat{B}_0 = Estimator of B_0
 \hat{B}_1 = Estimator of B_1 ,
 X = Independent variable.

In some cases we have dependent and independent variables in the equation $y = \beta_0 + \beta_1 x + \epsilon$. If we know the values of β_0 and β_1 , then we can easily calculate the value of the dependent variable for that. We should have the values of β_0 and β_1 .

→ Estimating the coefficients of the linear regression model.

In a practical situation such as the application to class size and test scores, the intercept is β_0 and slope β_1 of the population regression line are unknown. Therefore, we must use data to estimate the unknown slope and intercept of the population regression line.

This estimation problem is similar to others you have faced in statistics, the same extends to the linear regression model. We do not know the population value of β_1 , the

slope of regression line relating X (class size) and Y (test scores). But just as it was possible to learn about the population mean using a simple sample of data drawn from that population, so is it possible to learn about the population slope (β_1) using a sample of data.

As we know that we have the values of (β_0) intercept and slope (β_1) then we can easily calculate the value of dependent variable using the simple linear regression model.

So, for the β_0 and β_1 we can have the following formulas or we can use the OLS estimator to estimate the values of coefficient (β_0 and β_1)

$$\hat{\beta}_0 = \frac{n \sum y - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\hat{\beta}_1 = \frac{\sum xy - n \bar{x} \bar{y}}{\sum x^2 - n \bar{x}^2}$$

$$\text{iii) } \beta_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

and for β_0 we have formula.

$$\text{ii) } \beta_0 = \bar{y} - b \bar{x}$$

$$\text{i) } \beta_0 = \bar{y} - b \bar{x}$$

To understand these formulas we have solve a question and how will come to know that how they are calculated.

	X	y	XY	X^2	$\hat{y} = 2.001 + 0.387(x)$	$e = y - \hat{y}$
1	2.1	2.1	2.1	1	2.001 + 0.387(1) = 2.388	2.001 - 2.388 = -0.388
1	2.5	2.5	2.5	1	2.775	2.388 0.112
2	3.1	6.2	6.2	4	3.162	2.775 0.325
3	3.3	9	9	9	3.549	3.162 -0.162
4	3.8	15.2	15.2	16	3.549	3.549 0.251
4	3.2	12.8	12.8	16	3.549	-0.349
5	4.3	21.5	21.5	25	3.936	0.364
6	3.9	23.4	23.4	36	4.323	-0.423
6	4.4	26.4	26.4	36	4.323	0.077
7	4.8	33.6	33.6	49	4.71	0.09
					$\sum e =$	
					$\sum (y - \hat{y}) =$	
	$\sum x =$	$\sum y =$	$\sum xy =$	$\sum x^2 =$	0.003	
	39	35.1	152.7	193	≈ 0	

As we have the simple regression model such as.

$$\hat{Y} = \beta_0 + \beta_1 X_i$$

so know calculating β_1 such as using the formula

$$\hat{\beta}_1 = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$\hat{\beta}_1 = \frac{10(152.7) - (39)(35.1)}{10(193) - (39)^2}$$

$$\hat{\beta}_1 = \frac{10(152.7) - (39)(35.1)}{1930 - (39)^2}$$

$$\hat{\beta}_1 = \frac{158.1}{409}$$

$$\hat{\beta}_1 = 0.387$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_0 = 35.1 - 0.387(39)$$

$$\hat{\beta}_0 = 20.001$$

so now we calculate the value of dependent variables as in column 25.

of the table using different values of x .

→ If $x = 10$ then

$$\hat{y} = 2.001 + 0.387(10)$$
$$\hat{y} = 5.871$$

→ if $x = 0$ then

$$\hat{y} = 2.001$$

⇒ Ordinary Least Squares Estimator (OLS Estimator)

In data analysis, we use OLS estimator for estimating the unknown parameters in linear regression model.

The goal is to minimize the differences between the calculated observations in some arbitrary dataset and the responses predicted by the linear approximation of the data.

The OLS estimator chooses the regression coefficients so that the estimated regression line is as close as possible to the

observed measured mistakes data, where closeness is made in predicting Y given X .

The sum of these squared prediction mistakes over all observations is.

$$Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)$$

$$Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i$$

so it becomes as

$$\sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \rightarrow \text{min}$$

The estimators of the intercept and slope that minimize the sum of squared mistakes in equation (i) are called the ordinary least squares (OLS) estimators of β_0 and β_1 .

You could compute the OLS estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ by trying different values of β_0 and β_1 repeatedly.

Find those that minimize the total squared mistakes in equation or expression (i). They are the least squares estimators.

Because it is BLUE and consistent

⇒ Why Use the OLS estimators?

There are both practical and theoretical reasons to use the OLS estimators $\hat{\beta}_0$ and $\hat{\beta}_1$. Because OLS is the dominant method used in practice, it has become the common language for regression analysis throughout economics, finance, and the social sciences more generally.

The OLS formulas are built into virtually all spreadsheet and statistical software packages, making OLS easy to use.

→ Under the assumptions of the simple linear regression model, the OLS estimator is unbiased and consistent. The OLS estimator is also efficient among a certain class of unbiased estimators.

→ Measures of Fit & Having estimated a linear regression, you might wonder how well that regression line describes the data. Does the account for much

or for little of the variation in the dependent variable? Are the observations tightly clustered around the regression line, or are they spread out.

The R^2 and the standard error of regression measures how well the OLS regression line fits the data. The R^2 ranges from 0 and 1 and measures the fraction of the variance of y_i that is explained by x_i . The standard error of regression measures how far is the y_i from its predicted value.

⇒ The R^2 is a statistical measure that represents the proportion of a variance for a dependent variable that is explained by an independent variable or variables in a regression model.

Mathematically, the R^2 can be written as the ratio of the explained sum of squares (ESS) to the total sum of squares.

of squares (TSS). The explained sum of squares (ESS) is the sum of squared deviations of the predicted value \hat{Y}_i , from its average, and the total sum of squares (TSS) is the sum of squared deviations of Y_i from its average.

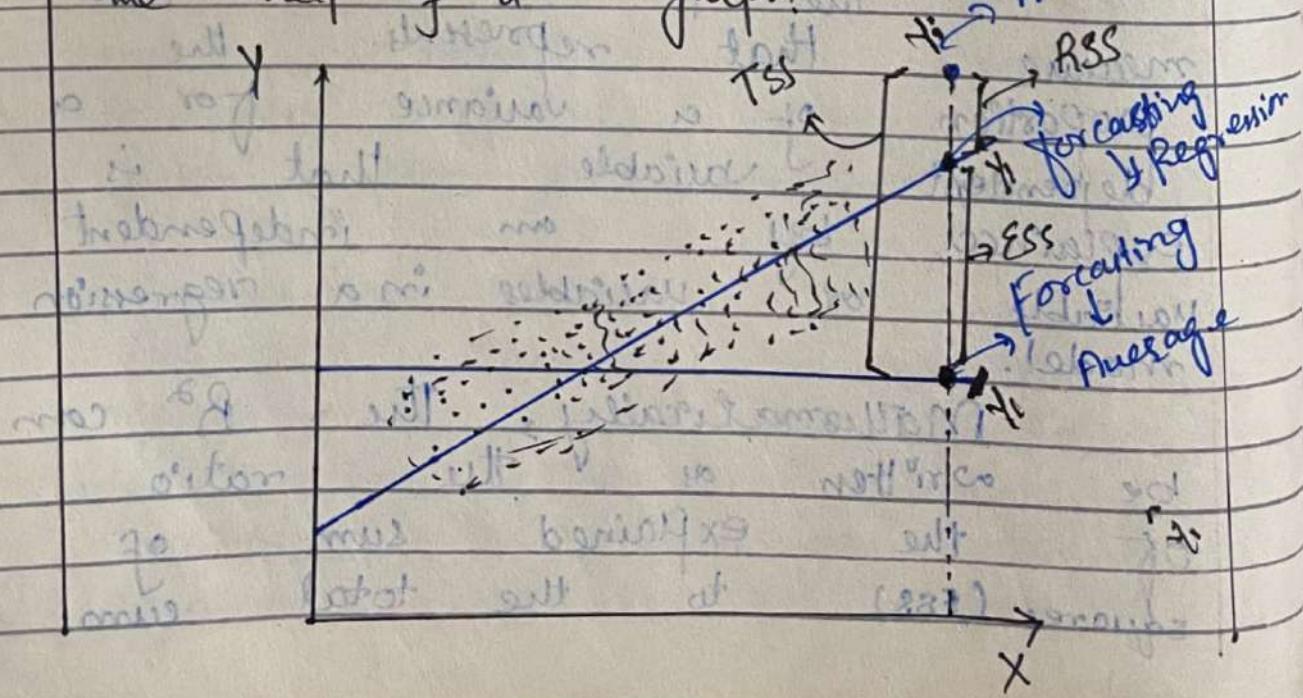
$$ESS = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Actual value
from regression

Forecasting from Average

Now we will explain with the help of a graph.



Predicting

In the graph we are finding the value of y given x .
And here we are predicting the value of y through forecasting, (i) Forecasting by average.
(ii) and then forecasting through a regression line.

When we forecasted through average then the difference b/w the predicted value and the actual value is called Total sum of square.

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

And after that we forecasted through the regression line and we found that there is a little variation from the y_i that in forecasting through average. so some of the portion of the TSS is explained by the regression line so we said it explained sum of square.

$$ESS = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

And The difference between regression line and the observed value is called Residual sum of squares. (RSS)

So,

$$TSS = ESS + RSS$$

So, when we call about the ratio of ESS in TSS, it is basically called R^2 .

$$R^2 = \frac{ESS}{TSS}$$

So as R^2 is greater the model will be considered fit.

R^2 ranges b/w 0 and 1

$$0 \leq R^2 \leq 1$$

R^2 is also called coefficient of determination.

$$R^2 = \frac{ESS}{TSS}$$

$$\text{As } TSS = ESS + RSS$$

$$\text{So } ESS = TSS - RSS$$

equation becomes

$$R^2 = \frac{TSS - RSS}{TSS}$$

$$R^2 = \frac{TSY - RSS}{TSS}$$

$$R^2 = 1 - \frac{RSS}{TSS}$$

$$\therefore R^2 = \frac{ESS}{TSS} \Rightarrow R^2 = 1 - \frac{RSS}{TSS}$$

\Rightarrow The standard error of the regression
The standard error of the regression will basically measures the magnitude of a typical regression residual in the units of y .

The SER measures the spread of the distribution of u . The standard error of regression is almost the sample standard deviation of the OLS residuals.

$$SER = \sqrt{\frac{SSR}{n-2}}$$

OR it can also be defined as

The Standard error of regression also known as the standard error of the estimate, represents the average distance that the

Observed regression values fall from the line. Conveniently, it tells you how wrong the regression model is on average using the units of the response variable.

⇒ What should we make of low R^2 and large SER

The fact that if the R^2 of any regression is low and the SER is large does not, by itself, imply that this regression is either "good" or "bad".

What the low R^2 does tell us is that other important factors influence test score. These factors could include differences in the student body across school districts, differences in quality unrelated to the student teacher ratio, or luck on the test.

The low R^2 and high SER do not tell us what these factors are, but they indicate that the student teacher ratio alone explains only a small part.

of the variation in test scores
in these data.

⇒ The Least Squares Assumptions

This section presents a set of three assumptions on the linear regression model and the sampling scheme under which OLS provides an appropriate estimator of the unknown regression coefficients β_0 and β_1 .

(1) ⇒ Assumption # 01 & The independent variable and the error term should not be correlated.

(*) ⇒ The conditional distribution of u_i given x_i has a mean of zero.

This assumption is a formal mathematical statement about the "other factors" contained in u_i and asserts that these other factors are unrelated to x_i in the sense that, given a value of x_i , the mean of the distribution of these other factors is zero.

So, the error term u_i has a conditional mean zero given x_i : $E(u_i | x_i) = 0$.

So, if the conditional mean of one random variable given another is zero, then the two random variables have zero covariance and thus are uncorrelated. Thus the conditional mean assumption $E(u_i | X_i) = 0$ implies that X_i and u_i are uncorrelated or $\text{corr}(X_i | u_i) = 0$. However, if X_i and u_i are correlated, then it must be the case that $E(u_i | X_i) \neq 0$.

So, if X_i and u_i are correlated, then the conditional mean assumption is violated.

(2) Assumption # 02:

(X_i, Y_i) , $i = 1, \dots, n$, are independently and identically distributed.

This assumption is a statement about how the sample is drawn.

If the observations are drawn by simple random sampling from a single large population, then

(X_i, Y_i) , $i = 1, \dots, n$, are i.i.d.

For example, let X be the age of the worker and Y be his or

her earnings, and imagine drawing a person at random from the population of workers. That randomly drawn person will have a certain age and earnings. If I draw a sample of n workers is drawn from this population, then (X_i, Y_i) , $i = 1, \dots, n$, necessarily have the same distribution. If they are drawn at random they are also distributed independently from one observation to the next; that is, they are i.i.d.

The i.i.d. assumption is a reasonable one for many data collection schemes. For example, survey data from a randomly chosen population typically can be treated as i.i.d.

(3) → Large Outliers are Unlikely & Squares
The third least squares assumption is that large outliers — that is, observations with values of X_i, Y_i , or that both are far outside the usual range of data — are unlikely. Large outliers can make OLS regression

results misleading.

Another way to state this assumption is that X and Y have finite kurtosis.

⇒ Sampling distributions of the OLS Estimators

Because the OLS estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ are computed from a randomly drawn sample, the estimators themselves are random variables with a probability distribution — the sampling distribution that describes the values they could take over different possible random samples.

This section presents these sampling distributions.

In small samples, these distributions are complicated, but in large samples, they are approximately normal.

⇒ Completed

Exercises

Chapter No 53

⇒ Regression with a single regressor;
Hypothesis Tests and Confidence
Intervals

This Chapter continues the treatment of a linear regression with a single regressor.

⇒ Testing hypotheses about one of the regression Coefficients

Your client, the superintendent, calls you with a problem. She has an angry taxpayer who claims that cutting class size will not help boost test scores, so reducing them is a waste of money. Class size, the taxpayer claims, has no effect on test scores.

The taxpayer is asserting that the population regression line is flat — that is, the slope ' B_1 ' of the population regression line is zero, so the superintendent asks whether there is any evidence that you can reject the hypothesis that $B_1 = 0$ or should accept it, at least tentatively.

pending further new evidence?

So, this section discusses tests of hypotheses about the slope β_1 and intercept β_0 of the population regression line. We start by discussing two-sided tests of the slope β_1 in detail, then turn to one-sided tests and to tests of hypotheses regarding the intercept β_0 .

→ Two-sided hypothesis test concerning β_1 &

The general approach to test hypotheses about the coefficient β_1 is the same as to test hypotheses about the population mean.

Testing hypotheses about the population mean.

First of all we should know about the null and alternative hypotheses.

A null hypothesis is a statement in which there is no relationship between two variables.

An alternative hypothesis is a statement in which there is some statistical significance.

between two measured phenomena.

So,

Null :- The effect equals zero

Alternative :- The effect doesn't equal zero

Recall that the null hypothesis is that the mean value of Y is a specific value $\mu_{Y,0}$. Can be written as $H_0 : E(Y) = \mu_{Y,0}$, and the two-sided alternative is $H_1 : E(Y) \neq \mu_{Y,0}$.

To test the null hypothesis H_0 against the two-sided alternative proceeds as follows in three steps.

(i) The first is to compute the standard error ($SE(\bar{Y})$)

(ii) The second step is to compute the t-statistic.

$$\text{The t-Statistic} = \frac{(\bar{Y} - \mu_{Y,0})}{SE(\bar{Y})}$$

(iii) The third step is to compute the p-value, which is the smallest significance level at which the null hypothesis could be rejected, based on

test statistic has actually been observed.

Alternatively, the third step can be replaced by simply comparing the t-statistic to the critical value, appropriate for the test with the desired significance level.

For example, a two-sided test with a 5% significance level would reject the null hypotheses if $t > 1.96$. In this case, the population mean is said to be statistically significant different from the hypothesized value at the 5% significance level.

→ Testing hypotheses about the slope $\beta_{1,0}$

The null and alternative hypotheses need to be tested precisely before they can be tested.

The "angry taxpayer" hypothesis is that under the true population some specific value, $\beta_{1,0}$. Under the two-sided alternative, β_1 does not equal $\beta_{1,0}$.

That is, the null hypothesis and the two-sided alternative hypothesis are:

$$H_0: \beta_1 = \beta_{1,0} \text{ vs } H_1: \beta_1 \neq \beta_{1,0}$$

To test the null hypothesis H_0 , we follow the same three steps as (i) for the population mean.

(i) Compute the standard error of $\hat{\beta}_1$.

In applications the $SE(\hat{\beta}_1)$ is computed by regression software so that it is easy to use in practice.

(ii) The 2nd step is to compute the t-statistic

$$\text{t-statistic} = \frac{\text{Estimator} - \text{hypothesized value}}{SE(\text{estimator})}$$

$$\text{t-statistic} = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)}$$

(iii) The 3rd step is to compute the p-value.

Reject the hypothesis at 5% significance level if the p-value is less than 5%, or

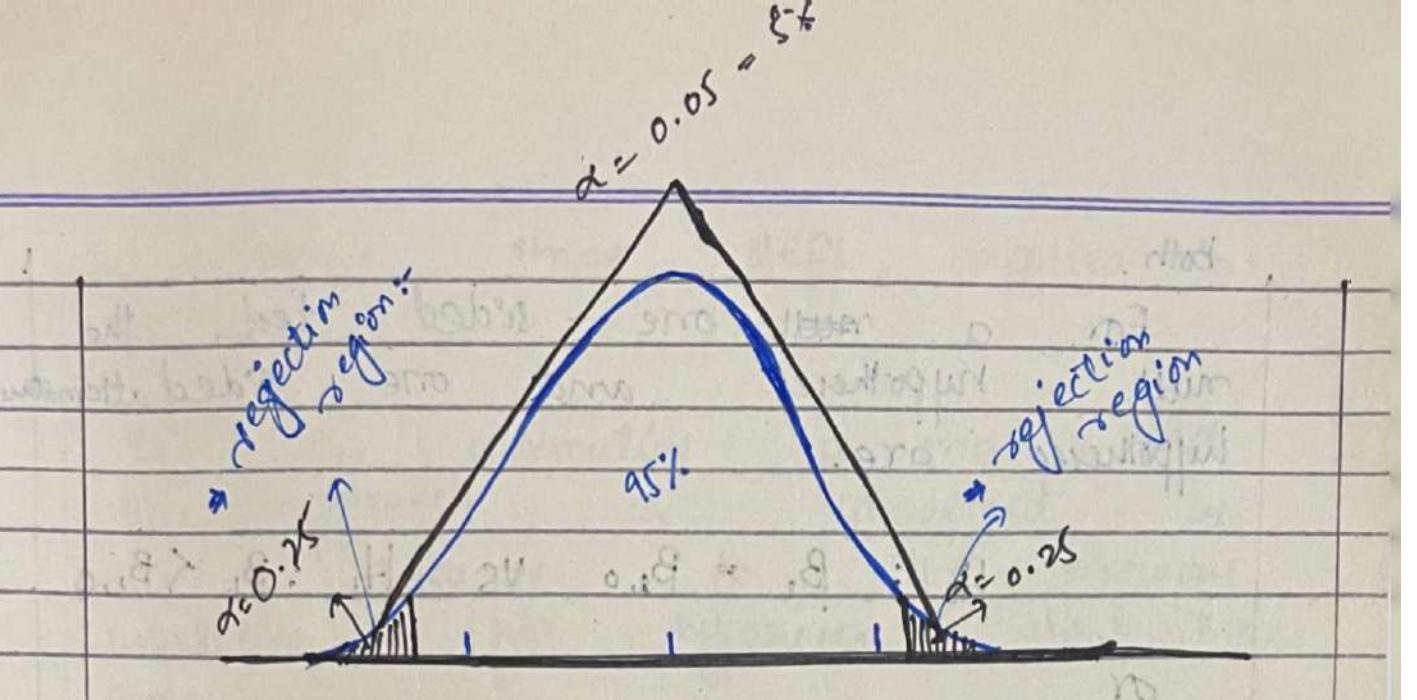
equivalently, if $|t^{act}| > 1.96$.

The standard error and t-statistic and p-value for testing $H_0: \beta_1 = 0$ are computed automatically by regression software.

A p-value of less than 5% provides evidence against the null hypothesis in the sense that, under the null hypotheses, the probability of obtaining a value of $\hat{\beta}_1$ at least as far from the null as that actually observed is less than 5%. So the null hypothesis is rejected at a 5% significance level.

Alternatively, the hypothesis can be rejected tested at the 5% significance level simply by comparing the absolute value of the t-statistic to 1.96, the critical value of for a two-sided test, and rejecting the null hypothesis at a 5% level if $|t^{act}| > 1.96$.

Graphically,



If that sample t-statistic value falls in any of the rejection regions that means we can reject our null hypotheses.

→ One-sided hypotheses
In one-sided hypotheses the null and alternative hypotheses are.

Null & The effect is greater than or equal to zero

Alternative: The effect is less than zero

A one sided hypothesis test is a statistical test in which the critical area of a distribution is one-sided so that is either greater than or less than a certain value, but not

both.

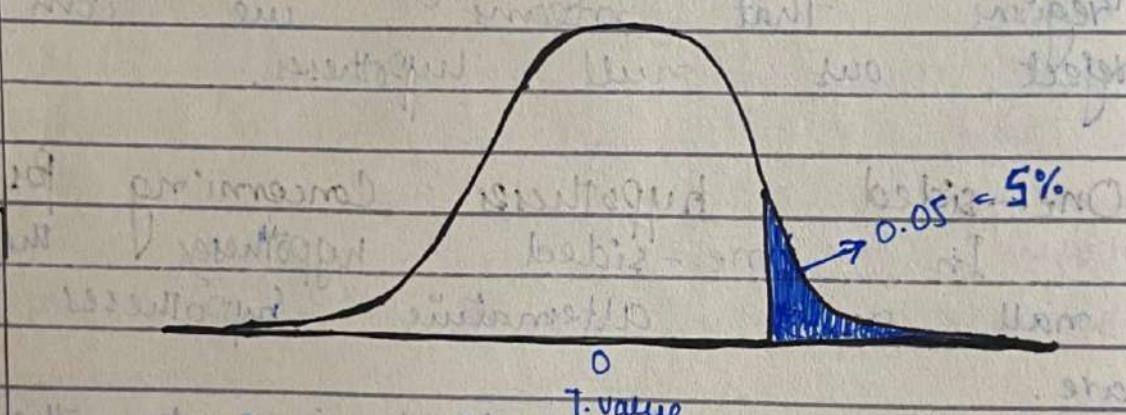
For a ~~null~~ one-sided test, the null hypothesis and one sided alternative hypotheses are.

$$H_0 : \beta_1 = \beta_{1,0} \text{ vs } H_1 : \beta_1 < \beta_{1,0}$$

or

$$H_0 : \beta_1 = \beta_{1,0} \text{ vs } H_1 : \beta_1 > \beta_{1,0}$$

(One sided alternative).



→ Testing hypotheses about the intercept β_0 .

The null hypotheses concerning the intercept and the two-sided alternative are.

$$H_0 : \beta_0 = \beta_{0,0}, H_1 : \beta_0 \neq \beta_{0,0}$$

The general approach to testing this null hypothesis consists of

the same three steps applied to

β_1

If the alternative is one-sided, this approach is modified as was discussed in the previous subsection hypotheses about the slope.

→ Hypotheses Test are useful if you have a specific null hypothesis in mind. Being able to accept or reject this null hypothesis based on the statistical evidence.

Yet there are many times that no single hypotheses about a regression coefficient is dominant, and instead one would like to know a range of values of coefficients that are consistent with the data. This is called for construction a confidence interval.

→ Confidence Intervals for a regression

Coefficients

Because any statistical estimate of the slope β_1 necessarily has sampling uncertainty, we cannot determine the true value of

B_1 exactly from the sample of data. It is possible, however, to use the OLS estimator and its standard error to construct a confidence interval for a slope B_1 or for the intercept B_0 .

→ Confidence interval for B_1 &

Confidence interval means that

for example as we take 100 samples of size 420, then 95 samples will contain true population parametric value and 5 samples will not contain (Note:- if B_1 is normally distributed).

Recall that a 95% confidence is, equivalently:

↳ A set of points that cannot be rejected at a 5% significance level.

↳ A set-valued function of a data that contains the true parameters value 95% of the time in repeated samples.

Because the t-statistic for B_1 is $N(0,1)$ in large samples,

construction of a 95% confidence β_1 is just like the case in the sample mean.

95% confidence interval for β_1 =

$$= \{ \hat{\beta}_1 \pm 1.96 \times SE(\hat{\beta}_1) \}$$

Confidence interval example :-

Test score and STR

Estimated regression line :-

$$TS = 698.9 - 2.28 \times STR$$

$$SE(\hat{\beta}_0) = 10.4 \quad SE(\hat{\beta}_1) = 0.52$$

95% confidence interval for $\hat{\beta}_1$:-

$$\{ \hat{\beta}_1 \pm 1.96 \times SE(\hat{\beta}_1) \} = \{ -2.28 \pm 1.96 \times 0.52 \} \\ = \{ -3.30, -1.26 \}$$

So, in the above example the $\beta_1 \neq 0$ because it does cover zero so we reject the null hypothesis.

The following two statements are equivalent.

i) The 95% confidence interval

does not include zero;
(iii) Null hypothesis $B_1 = 0$ is rejected
at the 5% level.

Now:- implication of this example of our hypothetical superintendent is contemplating reducing the STR by 2. Because the 95% confidence interval for B_1 is $[-3.30, -1.26]$, the effect of reducing the STR by 2 could be as great as $-3.30 \times -2 = 6.60$ or as little as $-1.26 \times -2 = 2.52$.

Thus decreasing the student-teacher ratio by 2 is predicted to increase test scores by between 2.52 and 6.60 points, with a 95% confidence level.

⇒ Regression When X is a Binary Variable:
Regression analysis can also be used when the regressor is binary - that is, when it takes on only two values, 0 or 1.

For example:-

- 4. $X = 1$ if small class size, = 0 if not.
- 4. $X = 1$ if female, = 0 if male.

• $x = 1$ if Rural, $= 0$ if Urban

Binary regressors are sometimes called "Dummy" variables.

→ Interpretation of the regression Coefficients

The mechanics of regression with a binary regressor are the same as if it was/is continuous. The interpretation of β_1 , however, is different, and it turns out that regression with a binary variable is equivalent to performing a difference of means analysis.

Simple regression model $= Y_i = \beta_0 + \beta_1 X_i + u_i$, where X is binary ($X_i = 0$ or 1).

(i) When $X_i = 0$

$$Y_i = \beta_0 + u_i$$

that is, $E(Y_i | X_i = 0) = \beta_0$

(ii) When $X_i = 1 \rightarrow Y_i = \beta_0 + \beta_1 + u_i$

↳ The mean of Y_i is $\beta_0 + \beta_1$

↳ That is, $E(Y_i | X_i = 1) = \beta_0 + \beta_1$.

Here we know that X_i is not continuous; so, it is not useful to.

think of β_1 as a slope; indeed, because x_i can take on only two values, there is no "line", so it makes no sense to talk about the slope. Thus we will not refer to β_1 as a slope in equation $y_i = \beta_0 + \beta_1 x_i + u_i$ (when x_i is binary); instead we will simply refer to β_1 as the coefficient multiplying x_i .

Q: β_1 is not a slope so what is it? The best way to interpret β_0 and β_1 in a regression with a binary regressor is to consider, one at a time, the two possible cases, $x_i = 0$ and $x_i = 1$.
So we know that

(P) when

$$x_i = 0.$$

$$\text{then } y_i = \beta_0$$

(P) And when $x_i = 1$

$$\text{then } y_i = \beta_0 + \beta_1$$

So, the difference between them conditional expectation of y_i when $x_i = 1$ and when $x_i = 0$ or

$$\beta_1 = (E(Y_i | X_i \neq 0)) - (E(Y_i | X_i = 0)).$$

So $(\beta_0 + \beta_1) - \beta_0 = \beta_1$ is the difference b/w these two means.

→ Heteroskedasticity and Homoskedasticity
 Our only assumption about the distribution of U_i conditional on X_i is that it has mean of zero. Furthermore, the variance of this conditional distribution does not depend on X_i , then the error term is said to be homoskedastic, otherwise the error term is heteroskedastic.

→ If $\text{Var}(U|X = x)$ is constant — that is, if the variance of the conditional distribution of U given X does not depend on X — then U is said to be homoskedastic, otherwise, U is heteroskedastic.

Equal group Variances = Homoskedasticity.
 Unequal group Variances = Heteroskedasticity.

When &

$$E(U|X_i) = 0$$

↳ The variance of \hat{U} does not depend on X .

⇒ We now have two formulas for Standard error for $\hat{\beta}_1$.

↳ Homoskedasticity - only standard errors.

↳ These are valid only for if errors are homoskedastic.

↳ The usual standard errors - to differentiate the two, it is conventional to call these heteroskedasticity-robust Standard error, because they are valid whether or not the errors are heteroskedastic.

↳ The main advantage of the homoskedasticity-only standard error is that the formula is simpler. But the disadvantage is that the formula is only correct if the errors are homoskedastic.

↳ If the errors are either homoskedastic or heteroskedastic and you use heteroskedastic-robust standard

errors, you are OK.

- ↳ If the error is heteroskedastic and you use the homoskedasticity-only formula for standard error, your standard error will be wrong.
- ↳ The two formulas coincide (when n is large) in the special case of homoskedasticity.
- ↳ So, you should always account for heteroskedasticity errors.

⇒ The extended least squares Assumptions.

- (i) $E(U|X_i) = 0$
- (ii) $(X_i, Y_i), i = 1, \dots, n$, are i.i.d.
- (iii) Large outliers are unlikely.
- (iv) U is homoskedastic
- (v) U is distributed $N(0, \sigma^2)$.

⇒ The Gauss - Markov Theorem says
If the three OLS assumptions hold
and if errors are homoskedastic,
then the OLS estimator is
Best (most efficient) Linear conditionally
Unbiased Estimator (BLUE)



Chapter # 06 &

Linear regression with multiple Regressors &

This Chapter explains how to estimate the coefficients of the multiple linear regression model.

(As we know that school districts with low STR tends to have higher test scores in California data set, perhaps students from districts with small classes have other advantages that help them perform well on standardized tests. Could this have produced misleading results, and, if so what can be done?)

Omitted factor, such as student characteristics, can, in fact, make the OLS estimator of the effect of class size or more precisely biased.

This chapter explains the "omitted variable bias", ~~here~~ and introduces multiple regression, a method that can eliminate omitted variable bias.

→ Definition of Omitted Variable Bias.

→ Omitted variable bias &

If the regressor (the STR) is correlated with a variable that has been omitted from the analysis (The el-pct) and that determines in part, the dependent variable (Test score), then the OLS estimator will have omitted variable bias.

Omitted variable bias occurs when two conditions are true:

- (i) When the omitted variable is correlated with the included regressor, and
- (ii) When the omitted variable is a determinant of the dependent variable.

Both conditions must hold for the omission of Z to result in omitted variable bias.

In the test score example,

- 1) English language ability (whether the student has English as a 2nd language) plausibly affects the standardized test scores, so

(el-pct) or Z

U is determinant of Y

- ② Immigrant communities tend to be less affluent and thus have smaller school budgets - and higher STR, so Z is correlated with X .

→ so, this results into omitted variable bias.

→ Omitted variable bias and the first least squares assumption. 8

Omitted variable bias means that the 1st least squares assumption - that $E(u_i | X_i) = 0$ is incorrect. To see why, recall that the error term u_i in the linear regression model, with a single regressor represents all factors other than X_i , that are determinants of Y_i . If one of these other factors is correlated with X_i this means that the error term (which contains this factor) is correlated with X_i . Because u_i and X_i are correlated, the conditional mean of u_i given X_i is nonzero. This correlation therefore violates the first least

squares assumption, and the consequence is serious. The OLS estimator is biased. This bias does not vanish even in very large samples, and the OLS estimator is inconsistent.

⇒ The Formula for Omitted Variable Bias
 The discussion about the omitted variable bias can be summarized mathematically by a formula for this bias.
 Let the correlation b/w X and u_i be $\text{corr}(X_i, u_i) = \rho_{Xu}$.

Suppose that the 2nd and 3rd OLS assumption hold, but the 1st does not because ρ_{Xu} is nonzero. Then the OLS estimator has the limit.

$$\hat{\beta}_1 \leftarrow \beta_1 + \frac{\rho_{Xu} \sigma_u}{\sigma_X}$$

This formula summarizes several of the ideas discussed about omitted variable bias.

④ Omitted variable bias is a problem whether the sample

size of the sample large or small, because $\hat{\beta}_1$ does not converge in probability to the true value β_1 , $\hat{\beta}_1$ is biased and inconsistent. So, the term $Pxu \sigma_u^2$ in formula is the bias in $\hat{\beta}_1$ that persists even in large sample.

2) Whether this bias is large or small in practice depends on the correlation between the regressor and the error term. The larger $|Pxu|$ is, the larger the bias.

3) The direction of bias in $\hat{\beta}_1$ depends on whether x and u are positively or negatively correlated. For example the percentage of students learning English has a negative effect on district test score, so the percentage of english learners enter the error term with a negative sign, thus the student teacher ratio (x) would be negatively correlated with the error term (u) so $Pxu < 0$ and the coefficients β_1 on the student-teacher ratio would be biased towards a negative number.

⇒ The multiple regression model
Consider the case of two regressors,

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, \dots, n$$

Y = Dependent variable.

X_1, X_2 = two independent variables.

β_0 = population intercept.

β_1 = $\frac{\Delta Y}{\Delta X_1}$, holding X_2 constant.

β_2 = $\frac{\Delta Y}{\Delta X_2}$, holding X_1 constant.

u_i = the regression error term.

⇒ The OLS Estimator in multiple Regression

This section describes how the coefficients of the multiple regression model can be estimated using OLS.

The OLS estimators $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ are the values of b_0, b_1, \dots, b_k that minimizes the sum of squared mistakes $\sum_{i=1}^n (Y_i - b_0 - b_1 X_{1i})^2$.

The OLS predicted values and residuals are:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \dots + \hat{\beta}_k X_{ki}, \quad i = 1, \dots, n,$$

and

$$\hat{u}_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n.$$

⇒ Measures of Fit in multiple Regression

Three commonly used summary statistics in multiple regression are the standard error of the regression, the regression R^2 , and the adjusted R^2 (also known as \bar{R}^2). All three statistics measure how well the OLS estimate of the multiple regression line describes, or "fits" the data.

→ The standard error of regression (SER)

The SER estimates the standard deviation of the error term U_i . Thus the SER is a measure of spread of the distribution of Y around the regression line. In multiple regression, the SER is

$$SER = \sqrt{\frac{SSR}{n-k-1}}$$

The only difference b/w SER in single and multiple regression is that here the divisor is $n-k-1$ rather than $n-2$.

If there is a single regressor, $k=1$ then both the formulas becomes the same. When n is large, the effect in both formulas become negligible.

$$SER = \sqrt{\frac{SSR}{n-2}}$$

$$SER = \sqrt{\frac{SSR}{n-k-1}}$$

when n is large it becomes negligible.

→ \hat{R}^2 is the regression R^2 is the fraction of the sample variance of y_i explained by the regressors. Equivalently R^2 is

$$R^2 = \frac{ESS}{TSS}, \text{ or } R^2 = 1 - \frac{RSS}{TSS}$$

In multiple regression when R^2 increases, whenever a regressor is added unless the estimated coefficient on the added regressor is exactly zero, so in general the RSS will decrease when a new variable is added. But this means that R^2 generally increases when a new regressor is added.

→ "The Adjusted R^2 " (\bar{R}^2) is

Because the R^2 increases when a new variable is added, an increase in R^2 does not mean that adding a variable actually improves the fit of the model. In this sense, the R^2 gives an inflated estimate of how well the regression fits the data. One way to correct for this is to deflate or reduce the R^2 by some factor, and this is what the adjusted R^2 , or \bar{R}^2 does.

The adjusted R^2 is the modified version of the R^2 that does not necessarily increase when a regressor is added. The \bar{R}^2 is

$$\bar{R}^2 = 1 - \frac{n-1}{n-k-1} \left(\frac{SSR}{TSS} \right)$$

- (i) $(m-1)/m-k-1$ is always greater than 1, so R^2 is always less than \bar{R}^2 .
- (ii) Adding new regressors will decrease SSR, which increase \bar{R}^2 . Whether \bar{R}^2 ↑ or ↓ depends on which of their effect is strong.
- (iii) The \bar{R}^2 can be negative, when SSR↑ in small amount and $\frac{n-1}{m-k-1}$ increases in large amount.

→ The Least Squares Assumptions in Multiple Regression

The multiple regression model is written as.

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

- (1) The conditional distribution of u given the X 's has mean zero, that is, $E(u|X_1 = x_1, \dots, X_k = x_k) = 0$.
- (2) $(X_{1i}, \dots, X_{ki}, Y_i)$, $i = 1, \dots, n$, are i.i.d.
- (3) Large outliers are unlikely.
- (4) There is no perfect multicollinearity.

⇒ Perfect multicollinearity is perfect multicollinearity when one of the regressors is an exact linear function of the other regressor.

For example

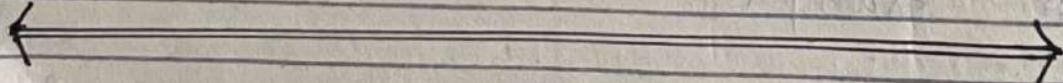
- In the previous regression, β_1 is the effect on Testscore of a unit change holding STR constant, so

it makes no sense, so if you have perfect multicollinearity your statistical software will let you know - either by crashing or giving an error message or by "dropping" one of the variables arbitrarily.

→ The solution to perfect multicollinearity is to modify your list of regressors so that you no longer have perfect multicollinearity.

→ Imperfect multicollinearity

Imperfect multicollinearity occurs when the two or more regressors are very highly correlated but not perfectly correlated. Imperfect multicollinearity does not pose any problem for the theory of OLS estimators. So, the imperfect multicollinearity is not an error.



Chapter # 07 "Hypothesis test and confidence intervals in Multiple Regression."

As discussed in Chapter 6, multiple regression analysis provides a way to mitigate the problem of omitted variable bias by including the additional regressors, thereby controlling the effects of those additional regressors. The coefficients in the multiple regression model can be estimated by OLS. Like all estimators, the OLS estimator has some sampling uncertainty because its value differs from one sample to the next.

This chapter presents methods for quantifying the sampling uncertainty of the OLS estimator through the use of standard errors, statistical hypothesis tests, and confidence intervals. One new possibility that arises in multiple regression is a hypothesis that simultaneously involves two or more coefficients of a regression. The general approach to testing such "joint" hypotheses involves a new test statistic, the F-statistic.

⇒ Hypothesis Tests and Confidence intervals for a single Coefficient β

This section describes how to test hypotheses, and how to construct confidence intervals for a single coefficient in a multiple regression model.

4) → Hypothesis tests for a single Coefficient β

Suppose that you want to test the hypothesis that a change in the student teacher ratio has no effect on test scores, holding the PCTEL in the district.

This corresponds to hypothesizing that the true coefficient β_2 on the STR in population regression of test scores on STR and PctEL is zero ($\beta_2 = 0$).

So here the null hypothesis is $H_0 = \beta_2 = 0$, and the alternative hypothesis is $H_1 = \beta_2 \neq 0$ (two sided alternative).

The procedure for testing this

null hypothesis when there is a single regression is done in three steps

(i) Compute the standard error of $\hat{\beta}_1$, $SE(\hat{\beta}_1)$

(ii) Compute the t-statistic,

$$t = \frac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)}$$

(iii) Compute the p-value,

$$p\text{-value} = 2\Phi(-|t^{\text{act}}|),$$

Where t^{act} is the value of the t-statistics actually computed. Reject the hypothesis at the 5% significance level if the p-value is less than 0.05 or, equivalently,

(iv) Confidence intervals for a single coefficient

The method for constructing a confidence interval in the multiple regression model is also the same as in the single-regressor model.

The 95% confidence interval is

(i) 95% confidence interval for β_1 =

$$\hat{\beta}_1 \pm 1.96 \times SE(\hat{\beta}_1) \quad (\text{A}) 38$$

(ii) A 90% confidence interval is obtained by replacing 1.96 in above equation with 1.64
(A) 38

The method for conducting a hypothesis test and the method for conducting a confidence interval rely on large samples.

→ An angry taxpayer asserts that the population values of both the coefficient on STR (β_1) and coefficient on spending per pupil (β_2) are zero; that is, he hypothesizes that both $\beta_1 = 0$ and $\beta_2 = 0$. So,

→ the tax payer's hypothesis is a joint hypothesis, and we need a new F. Statistic.

→ Tests of joint hypotheses &

This section describes how to formulate multiple joint regression hypothesis on how to test coefficient and them using an F-statistic.

Testing hypothesis on two or more Coefficients

Joint Null Hypothesis

Consider the regression

$$\widehat{TS} = 649.6 - 0.29 \text{ STR} + 3.87 \text{ Expn} - 0.656 \text{ PctEL}$$

(15.5) (0.48) (1.59) (0.032)

Our angry taxpayer hypothesizes that neither test nor expenditure on English have an effect on test score, once we control for the percentage of learners.

So, we can write this null hypothesis mathematically as

$$H_0: \beta_1 = 0 \text{ and } \beta_2 = 0 \text{ vs } H_1: \beta_1 \neq 0 \text{ and/or } \beta_2 \neq 0$$

Here the null hypothesis imposes two restrictions on multiple regression model : $\beta_1 = 0$ and $\beta_2 = 0$

In General the joint hypothesis is a hypothesis that imposes two or more restrictions on the regression coefficients.

→ If any one (or more than one) of the equalities under the null hypothesis H_0 is false, then the joint null hypothesis itself is false. Thus, the alternative hypothesis is that at least one of the equalities in the null hypothesis H_0 does not hold.

→ Why can't I just test the individual coefficients one at a time?

Although it seems it is possible to test a joint hypothesis by using the usual t-statistics to test the restriction one at a time, but the calculations show that it is unreliable. When the regressors are correlated,

the situations is even more complicated, so it becomes difficult to test the null hypothesis on "one at a time" using the t-statistics.

→ We can test the joint hypothesis that is more powerful, especially when the regressors are highly correlated. That approach is based on the F-statistics.

⇒ The F-statistic

The F-statistic is used to test joint hypothesis about regression coefficients.

A simple F-statistic formula that is easy to understand. (It is only valid if the errors are homoskedastic, but it might help intuition).

→ The homoskedasticity-only F-statistic
When the errors are homoskedastic, there is a simple formula for computing the "homoskedasticity-only" F-statistic:

→ Run the two regressions, one under the null hypothesis (the

"restricted" regression) and one under the alternative hypothesis (the "unrestricted" regression)

The Homoskedasticity-only F-statistic is given by the formula.

$$F = \frac{(\text{SSR}_{\text{restricted}} - \text{SSR}_{\text{unrestricted}}) / q}{\text{SSR}_{\text{unrestricted}} / (n - k_{\text{unrestricted}} - 1)}$$

q = number of restrictions under the null hypothesis and

k = number of regressors in the unrestricted regression.

An alternative equivalent formula for the homoskedasticity-only F-statistic is based on the R^2 of the two regressions.

$$F = \frac{(R^2_{\text{unrestricted}} - R^2_{\text{restricted}}) / q}{(1 - R^2_{\text{unrestricted}})(n - k_{\text{unrestricted}} - 1)}$$

⇒ Example

(i) Restricted regression:

$$\hat{T}_S = 644.7 - 0.671 \text{ PctEL}, \quad R^2_{\text{restricted}} = 0.4149$$

→ (ii) Unrestricted regression

$$\hat{TS} = 649.6 - 0.29 \text{STR} + 3.87 \text{Expn} - 0.656 \text{PctEL}$$

$$R^2_{\text{un}} = 0.4366, k_m = 3, q_k = 2 \\ n = 420.$$

so

$$F = \frac{(0.4366 - 0.4149)/2}{(1 - 0.4366)(420 - 3 - 1)}$$

$$F = 8.01$$

so it is higher than the critical value so that null hypothesis is rejected that $\beta_1 = 0$ and $\beta_2 = 0$. because 8.01 exceeds the 1% critical value of 4.61, the hypothesis is rejected at 1% level using the homoskedasticity-only test.

→ Testing Single Restrictions involving

Multiple

coefficients

Sometimes economic theory suggests

a single restriction that involves two or more regression coefficients.

for example, theory might suggest a null hypothesis of the form

$\beta_1 = \beta_2$; that is the effect of

the first and the ^{2nd} regressor are the same. In this case the task is to test this null hypothesis against the alternative that the two coefficients differ.

$$H_0 : \beta_1 = \beta_2 \quad \text{vs} \quad H_1 : \beta_1 \neq \beta_2$$

This null hypothesis has a single restriction, so $g = 1$, but the restriction involves multiple coefficients (β_1 and β_2)

→ Confidence sets for multiple coefficients

The method is conceptually similar to the method using the t-statistic, except that the confidence set for multiple coefficients is based on F-statistics.

→ Model specification for multiple regression

The job of determining which variables to include in the regression - that is the problem

of choosing a regression specification can be quite challenging, and no single rule applies in all situations.

It is important to rely on your expertise knowledge of the empirical problem and to focus on obtaining an unbiased estimate of the causal effect of interest; do not rely solely on purely statistical measures of fit as R^2 or \bar{R}^2 .

→ Omitted Variable Bias in Multiple Regressions

Omitted variable bias is the bias in the OLS estimator that arises when one or more included regressors are correlated with an omitted variable. For omitted variable bias to arise, two things must be true:

- (i) At least one of the included regressors must be correlated with the omitted variable.
- (ii) The omitted variable bias must be a determinant of the dependent variable Y .

→ Interpreting the R^2 and the Adjusted R^2 in Practice

The R^2 and \bar{R}^2 tell you whether the regressors are good at predicting, or "explaining", the values of the dependent variable in the sample of data on hand.

If the R^2 (or \bar{R}^2) is nearly 1, then the regressors produce good predictions of the dependent variable in that sample, in the sense that the variance of the OLS residual is small compared to the variance of the dependent variable. If the R^2 or (\bar{R}^2) is nearly 0, the opposite is true.

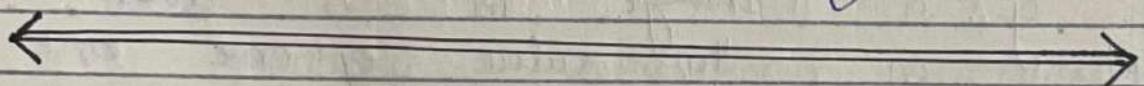
The R^2 and \bar{R}^2 do not tell you whether:

(1) An increase in the R^2 or \bar{R}^2 does not necessarily mean that an added variable is statistically significant.

(2) A high R^2 or \bar{R}^2 does not mean that the regressors are true cause of the dependent variable.

(3) A high R^2 or \bar{R}^2 does not mean that there is no omitted variable bias.

(4) A high \bar{R}^2 or R^2 does not necessarily mean that you have the most appropriate set of regressors, nor does a low R^2 or \bar{R}^2 necessarily mean that you have an inappropriate set of regressors.



chapter # 08

"Nonlinear Regression Functions"

In chapter 4 through 7, the population regression function was assumed to be linear. In other words, the slope of the population regression function was constant, so the effect of Y of a unit change in X does not itself depend on the value of X . But what if the effect on Y of a change in X does depend on the value of one or more of the independent variables? If so, the population regression function is nonlinear.

A non linear function is a function whose line is not constant. A non linear regression function with a slope that is not constant.

→ A General strategy for modeling non linear regression functions

→ Test scores and District incomes

The relationship between income and test scores is not a straight line. Rather it is a nonlinear.

A non linear function is a function with a slope that is not constant. One way to approximate such a curve mathematically is to model the relationship as a quadratic function. That is, we could model test scores as a function of income and the square of income.

A quadratic population regression model relating test scores and income is written mathematically as

$$\text{Test Score}_i = \beta_0 + \beta_1 \text{income}_i + \beta_2 \text{income}_i^2 + u_i$$

It is a quadratic function of the independent variable

The quadratic function captures in the scatterplot: It is steep for low values of district income but flattens out when district income is high. In short, the quadratic regression function seems to fit the data better than the linear one.

This quadratic regression model becomes linear when $\beta_0 = 0$. So it becomes linear

When (null hypothesis $\beta_2 = 0$)

→ Non Linear Regression Population Regression
functions. General ideas &

ID a relation between Y and X is nonlinear.

→ The effect on Y of a change in X depends on the value of X - that is, the marginal effect of X is not constant.

→ The effect on Y of a change in X in nonlinear specifications

The expected effect on Y of a change in X_1 in the non linear regression model:

The expected change in Y, ΔY , associated with the change in X_1 , ΔX_1 , holding X_2, \dots, X_k constant, is the difference b/w the value of the population regression function before changing X_1 , holding X_2, \dots, X_k constant that is the expected change in Y after changing X_1 .

$$\Delta Y = f(X_1 + \Delta X_1, X_2, \dots, X_k) - f(X_1, X_2, \dots, X_k)$$

The estimator of this unknown population difference is the difference b/w the predicted values for these two cases. Let $\hat{f}(X_1, X_2, X_3, \dots, X_k)$ be the predicted values of Y based on the estimator \hat{f} of the population regression function the predicted change in Y is.

$$\Delta \hat{Y} = \hat{f}(X_1 + \Delta X_1, X_2, \dots, X_k) - \hat{f}(X_1, X_2, \dots, X_k)$$

→ Nonlinear Independent Function of a single variable

This section provides two methods for modeling a non linear regression function that involves only one independent variable

We have two ways to capture nonlinearity in variables
Both are complementary approaches

1) Polynomials

Polynomials is an expression of more than two algebraic terms especially the same that contain different

powers of the same variables

A population regression function is approximated by a quadratic, cubic, or higher-degree polynomial

Approximate the population regression function by a polynomial

$$Y_i = \beta_0 + \beta_1 X_i^r + \beta_2 X_i^{2r} + \dots + \beta_r X_i^{(r)} + u_i$$

This is just the linear multiple regression model - except that the regressors are powers of X_i .

Estimation, hypothesis testing, etc. proceed as in the multiple regression model using OLS.

When $r=2$ \Rightarrow quadratic regression model

When $r=3$ \Rightarrow cubic regression model

→ Testing null hypothesis that the population regression function is linear:

According to the null hypothesis H_0 that the regression is linear and the alternative (H_1) that it is polynomial of

degree r corresponds to

$H_0 : \beta_2 = 0, \beta_3 = 0, \dots, \beta_r = 0$ vs $H_1:$ at least one $\beta_j \neq 0, j = 2, \dots, r.$

→ Which degree polynomial should I use?
This is summarized in following steps

1) Pick a maximum value of x and estimate the polynomial regression for that $x.$

2) Use the t-statistic to test the hypothesis that the coefficient on x^r [$\beta_r = 0$] is zero. If you reject this hypothesis, then x^r belongs in the regression, so use all the polynomial of degree $r.$

3) If you do not reject $\beta_r = 0$ in step 2, eliminate x^r from the regression and estimate a polynomial of degree $r-1.$ Test whether the coefficient on x^{r-1} is zero - If you reject, use the polynomial of degree $r-1$.

4) If you do not reject β_{r-1}

$= 0$ in step 3, continue this procedure until the coefficient on the highest power is statistically significant.

⇒ Interpretation of coefficients in polynomial regression modeling

The coefficients in Polynomial regressions do not have a simple regression interpretation. The best way to interpret Polynomial regression is to plot the estimated regression function and calculate the estimated effect on Y associated with a change in X for one or more values of X .

⇒ Logarithms

Usually log of a variable may make non normal data normal, non linear data linear, heteroskedastic data homoskedastic.

Logarithmic variables convert changes in into percentage changes

The logarithmic function has the following useful properties

$$\begin{aligned}\ln(1/x) &= -\ln(x) \\ \ln(ax) &= \ln a + \ln x \\ \ln(x/a) &= \ln x - \ln a \\ \ln(x^a) &= a \ln x\end{aligned}$$

→ Logarithms and percentages
The link between logarithms
and percentages relies on a
key fact:

When Δx is small, the difference
between the logarithms of
 $x + \Delta x$ and the logarithms
of x is approximately
divided by $\frac{\Delta x}{x}$. That is,

$$\ln(x + \Delta x) - \ln(x) \approx \frac{\Delta x}{x} \quad (\text{when } \frac{\Delta x}{x} \text{ is small})$$

The three logarithmic regression models
There are three different cases
in which logarithms might be used.

Logarithms can be used to transform
the dependent variable y , an
independent variable x or both
(but the variable being transformed
must be positive). The following

Summarizes these three cases and the interpretation of the regression coefficient β_1 . In each case OLS can be estimated using the logarithm of the dependent and independent variable.

Regression Specification

Interpretation of β_1

$$1) Y_i = \beta_0 + \beta_1 \ln(x) + \epsilon_i \quad (\text{linear log model})$$

A 1% change in x is associated with a change of β_1 in Y .

$$2) \ln(Y_i) = \beta_0 + \beta_1 x_i + \epsilon_i \quad (\text{log linear model})$$

A change in x by one unit ($\Delta x = 1$) is associated with $100\beta_1\%$ change in Y .

$$3) \ln(Y_i) = \beta_0 + \beta_1 \ln(x) + \epsilon_i \quad (\text{log-log model})$$

A 1% change in x is associated with a $\beta_1\%$ change in Y , so β_1 is the elasticity of Y with respect to x .

\Rightarrow Interactions between An interaction

Independent Variables & occurs when an

independent variable has a different effect on the dependent variable (outcome), depending on the values of another independent variable.

We consider three cases

(i) When both independent variables are binary.

(ii) When one is binary and the other is continuous.

(iii) And when both are continuous.

→ Interaction between two binary variables

The Population Regression of Y_i on there binary variables is

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + u_i$$

→ There could be an interaction between the two binary variables.

If it is easy to modify the specification so that it does by introducing another regressor, the product of the two binary variables, $D_{1i} \times D_{2i}$. The resulting regression is

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + \beta_3 (D_{1i} \times D_{2i}) + u_i$$

→ This model interaction is called population regression binary variable regression model.

A method for interpreting coefficients in regressions with binary variables first compute the expected values of Y for each possible case described by the set of binary variables. Next compare these expected values. Each coefficient can be expressed either as an expected value or as the difference b/w two or more expected values.

→ Interactions between a continuous and a binary variables through the use of interaction term $X_i \times D_i$, the population regression line relating Y_i and the continuous variable X_i can have a slope that depends on the binary variables D_i . There are three possibilities

(i) different intercept, same slope

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + u_i$$

ii) Different intercept and slope.

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + \beta_3 (X_i \times D_i) + u_i$$

iii) Same intercept and different slope.

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 (X_i \times D_i) + u_i$$

(The coefficients of all three can be estimated by OLS)

3) Interaction between Two continuous Variables &

This interaction is modeled by augmenting the linear regression model with an interaction term that is product of X_i and X_{2i} .

$$Y_i = \beta_0 + \beta_1 \underline{X_{1i}} + \beta_2 \underline{X_{2i}} + \beta_3 (X_{1i} \times X_{2i}) + u_i$$

⇒ Interactions in Multiple regressions
The interaction term between two independent variables X_1 and X_2 is their product $X_{1i} \times X_{2i}$. Including this interaction terms allows the effect on Y ~~why~~ of a change

in X_1 to depend on the values of X_2 and, conversely, allows the effect of a change in X_2 to depend on the values of X_1 .

The coefficient on X_1 and X_2 is the effect of a one-unit increase in X_1 and X_2 , above and beyond the sum of the individual effects of a unit increase in X_1 alone and a unit increase in X_2 alone. This is true whether X_1 and/or X_2 are continuous or binary.

⇒ Non linear Effects on Test scores of a student - Teacher ratio &

Briefly explained in Book with the help of a table

