

Multicollinearity

A crucial condition for the application of least squares is that the explanatory variables are not perfectly linearly co-related

$$i.e. \quad r_{x_i x_j} \neq 1$$

If there is no linear relationship between the regression of EV then they are said to be orthogonal. When the regressors are orthogonal inferences are like -

- A) Identifying relative effects of the regressor variables.
- B) Prediction and/or estimation, and
- C) Selection of an appropriate set of variables for the model.

Can be made relatively very easily. Unfortunately in most applications of regression regressors are not orthogonal. However in some situation the regressors are ~~not~~ nearly perfectly linearly correlated and in such cases, the inferences based on the regression model can be misleading or erroneous. When there are near interrelationship or dependence

among the EV, then there exist the problem of multicollinearity. Hence, the term multicollinearity is used to denote the presence of linear relationship or linear relationship among the EV.

If the EV are perfectly linearly correlated i.e. if correlation coefficient of these variables is equal to unity then the parameters become ~~interdet~~ indeterminate and it becomes impossible to obtain numerical values of each parameter separately and the method of least square breaks down.

When $|r_{xixj}| = 1$ then there is complete violation of the assumption, while $|r_{xixj}|$ lies between 0 and 1, then there is near correlation.

- * $0 < |r_{xixj}| < 1 \rightarrow$ near correlation
- $|r_{xixj}| = 1 \rightarrow$ Perfectly linearly correlated
- $|r_{xixj}| = 0 \rightarrow$ no correlation
- $|r_{xixj}| \neq 1 \rightarrow$ (assumption) not linearly perfectly correlated.

In practice neither of the extreme cases, there is some degree of i.e. of orthogonal X 's is often met. In most cases, there is some degree of inter correlation among the explanatory variables due to the interdependence of many economic magnitudes overtime. In this cases the simple correlation coefficient for each pair of explanatory variable will have a value between 0 and 1.

Multi co-linearity is not a condition that either exist or does not exist in economic functions but it is rather a phenomena inherent in most relationships due to the nature of economic magnitudes. There is no conclusive evidence concerning the degree of co-linearity, which if present will effect seriously the parameter estimates. Intuitively, when two explanatory variables are changing in nearly the same way it becomes extremely difficult to establish the influence of each one of the regressors on Y separately eg, assume that consumption expenditure of an individual depends on his income and liquid asset change by the same proportion then, the influence of one of these variables on consumption may be erroneously attributed to the other and the effects of these variables are

consumption can not be sensitively investigated due to high. Inter correlated multicollinearity is thus a situation where because of strong interrelationship among the explanatory variables, it becomes very difficult to disentangle their separate effects on the dependent variable, i.e. it becomes difficult to obtain precise estimates of the regression coefficients i.e.

$$\hat{\beta} = (x'x)^{-1} x'y$$

does not exist or we do not have an exact solution. Thus the term multicollinearity which was due to Ragnar Frisch originally meant "existence of a perfect or exact or linear relationship among some or all the explanatory variables of a regression model".

Today however multicollinearity is used in a broader sense to include the case of perfect multicollinearity as well as the case where the explanatory variables are inter correlated but not perfectly inter correlated.

Causes of multicollinearity / sources :

① The data collection method can lead to multicollinearity problems when an analyst samples only a sub-space of the region of the regressions defined.

② Constraints on the model or in the population being sampled can cause multicollinearity.
eg - Suppose that an electric utility is investigating the effect of family income (x_1) and house size (x_2) on residential electricity consumption, the levels of the two regressor obtain in the sample data show that if the data lie approximately along the straight line thus indicating potential multicollinearity problem.

In this example of physical constraint in the population has caused this phenomena i.e families with higher income generally have larger homes than families with lower. when physical constraint such regardless of the sampling method employed, constraints after occur in problems involving production or chemical process whether regressor are the component of a product and these components add to a constraint.

③ Multicollinearity may be induced by the choice of model. Moreover, adding polynomial terms to the regression model can result in significant multicollinearity.

④ An over (adding unnecessary explanatory variables) defined model with more regressor variables than observations may be a source of multicollinearity. These models are sometimes encountered in medical and behavioural research where there may be small number of sample units, eg - the small number of patients, and information is collected for a large number of regressors on each subject. The usual approach is dealing with multicollinearity in this context is to eliminate some of the regressor variables from consideration.

⑤ There is a tendency of economic variables to move together overtime which is a cause of multicollinearity. Economic magnitudes are influenced by the same factors and in consequences once these determining factors become operative, the economic variables show the same broad pattern of

behaviour over time eg - In periods of boom or rapid economic growth, the basic economic magnitudes grow although some tend to lag behind others. Thus income, consumption, saving, investment, prices, employment tend to rise in the periods of economic expansion and decrease in periods of recession. Growth and trend factors in time series are the most serious cause of multicollinearity.

⑥ The use of lagged values of some explanatory variables as separate independent factors in the model is another cause of multicollinearity. eg - In consumption function, it is customary to include, among the explanatory variables, past as well as present values of income. Similarly, in investment functions distributed lags concerning past levels of economic activities are introduced as separate explanatory variables. Naturally, the successive values of these variables are inter-correlated and hence multicollinearity is almost certain to exist in distributed lag models.

Consequences of multicollinearity

(A) Perfect multicollinearity : If the interco-relation between the explanatory variables is perfect, i.e. $k_{x_1 x_2} = 1$, then consequences are as follows -

(a) The estimates of the coefficients are indeterminate.

(b) The standard errors of these estimates becomes infinitely large.

Proof of (a)

Suppose that the relation be estimated is $y = b_0 + b_1 x_1 + b_2 x_2 + u$ and that x_1 and x_2 are related with the exact relation, $x_2 = k x_1$, where k is any arbitrary constant.

The formula for the estimation of the coefficients \hat{b}_1 and \hat{b}_2 are -

$$\hat{b}_1 = \frac{(\sum x_1 y) (\sum x_2^2) - (\sum x_2 y) (\sum x_1 x_2)}{(\sum x_1^2) (\sum x_2^2) - (\sum x_1 x_2)^2} \quad \text{--- (1)}$$

$$\hat{b}_2 = \frac{(\sum x_2 y) (\sum x_1^2) - (\sum x_1 y) (\sum x_1 x_2)}{(\sum x_1^2) (\sum x_2^2) - (\sum x_1 x_2)^2} \quad \text{--- (2)}$$

Now putting, $x_2 = kx_1$ in eqⁿ ①

$$\begin{aligned} \hat{b}_1 &= \frac{(\sum x_1 y) \sum (kx_1)^2 - (\sum kx_1 y) (\sum x_1 kx_1)}{(\sum x_1^2) \sum (kx_1)^2 - (\sum x_1 kx_1)^2} \\ &= \frac{k^2 (\sum x_1 y) (\sum x_1^2) - k^2 (\sum x_1 y) (\sum x_1^2)}{k^2 (\sum x_1^2) (\sum x_1^2) - k^2 (\sum x_1^2)^2} \\ &= \frac{k^2 (\sum x_1 y) (\sum x_1^2) - k^2 (\sum x_1 y) (\sum x_1^2)}{k^2 (\sum x_1^2)^2 - k^2 (\sum x_1^2)^2} \\ &= \frac{0}{0} = 0 \end{aligned}$$

Now putting $x_2 = kx_1$ in eqⁿ ②

$$\begin{aligned} \hat{b}_2 &= \frac{(\sum x_2 y) (\sum x_1^2) - (\sum x_1 y) (\sum x_1 x_2)}{(\sum x_1^2) (\sum x_2^2) - (\sum x_1 x_2)^2} \\ &= \frac{(\sum kx_1 y) (\sum x_1^2) - (\sum x_1 y) (\sum x_1 kx_1)}{(\sum x_1^2) \sum (kx_1)^2 - (\sum kx_1 x_1)^2} \\ &= \frac{k (\sum x_1 y) (\sum x_1^2) - k (\sum x_1 y) (\sum x_1^2)}{k (\sum x_1^2)^2 - k (\sum x_1^2)^2} \\ &= \frac{0}{0} = 0 \end{aligned}$$

Therefore the parameters are indeterminate and there is no way of finding separate values of these coefficients.

Proof of (b)

If $kx_1x_2 = 1$, then the standard errors of the estimates become infinitely large.

In a two variable model, $Y = b_0 + b_1x_1 + b_2x_2 + u$

$$\text{var}(\hat{b}_1) = \sigma_u^2 \frac{\sum x_2^2}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2} \quad (1)$$

$$\text{var}(\hat{b}_2) = \sigma_u^2 \frac{\sum x_1^2}{\sum x_1^2 \sum x_2^2 - (\sum x_1 x_2)^2} \quad (2)$$

Now substituting, $x_2 = kx_1$, in eqⁿ (1) & (2)

$$\text{var}(\hat{b}_1) = \sigma_u^2 \frac{\sum (kx_1)^2}{\sum x_1^2 \sum (kx_1)^2 - (\sum x_1 kx_1)^2}$$

$$= \sigma_u^2 \frac{k^2 (\sum x_1)^2}{\sum x_1^2 k^2 x_1^2 - k^2 (\sum x_1^2)^2}$$

$$= \sigma_u^2 \frac{k^2 \sum x_1^2}{k^2 (\sum x_1^2)^2 - k^2 (\sum x_1^2)^2}$$

$$= \sigma_u^2 \frac{k^2 \sum x_1^2}{k^2 (\sum x_1^2)^2 - k^2 (\sum x_1^2)^2}$$

$$= \sigma_u^2 \frac{k^2 \sum x_1^2}{0} = \infty$$

Therefore the parameter are indeterminate and there is no way of finding separate values of these coefficients.

$$\begin{aligned} \text{var}(\hat{\beta}_2) &= \sigma_u^2 \frac{\sum x_1^2}{\sum (kx_1)^2 \cdot \sum x_1^2 - (\sum x_1 kx_1)^2} \\ &= \sigma_u^2 \frac{\sum x_1^2}{k^2 (\sum x_1^2)^2 - k^2 (\sum x_1^2)^2} \\ &= \sigma_u^2 \frac{\sum x_1^2}{0} = \infty \end{aligned}$$

thus the variances of the estimates became infinite unless $\sigma_u^2 = 0$.

(B) Near multicollinearity or (High multicollinearity)

(i) Large variance and co-variance of OLS estimators

Although BLUE, the OLS estimators have large variances and making estimation difficult. Say, in three variable linear regression model—

$$Y_i = \beta_2 x_{2i} + \beta_3 x_{3i} + u_i \quad (1)$$

$$\text{var}(\hat{\beta}_2) = \frac{\sigma^2}{\sum x_{2i}^2 (1 - r_{23}^2)} \quad (2)$$

$$\& \text{var}(\hat{\beta}_3) = \frac{\sigma^2}{\sum x_{3i}^2 (1 - r_{23}^2)} \quad (3)$$

$$\text{and. cov}(\hat{\beta}_2, \hat{\beta}_3) = \frac{-r_{23} \sigma^2}{(1 - r_{23}^2) \sqrt{\sum x_{2i}^2 \sum x_{3i}^2}} \quad (4)$$

Where, r_{23} is correlation coefficient between x_2 and x_3 .

It is now apparent that from (2) and (3) as $r_{23} \xrightarrow{\text{tends to}} 1$, i.e. If collinearity increases, the variance of the two estimators increases and ~~for~~ the co-variance of the two estimators also increases in absolute value.

(2) Wider confidence interval : Because of large standard errors the confidence interval for the relevant population parameter also tends to be very large. Therefore in cases of high multicollinearity the sample data may be compatible with ~~divers~~ diverse set of hypothesis and hence the problem of accepting of all hypothesis increases.

(3) Insignificant t-ratios : To test a null hypothesis we use the t-ratio, i.e. $\frac{\hat{\beta}}{SE(\hat{\beta})}$ and compare the estimated t-value with the critical t value from t table. But in case of high multicollinearity the ~~estima-~~ standard error become dramatically large, thereby making the t-value smaller.

therefore in such cases there is the probability of accepting of false hypothesis.

④ A high R^2 but a few significant t-ratio:

If we consider a K variable linear regression model — $Y_i = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_K x_K + u$ — ① and if there is high and near multicollinearity defined as one and more partial slope coefficient individually statistically insignificant on the basis of 't' test. However R^2 means such situations may be so high say in excess of 0.9 that on the basis of F-test we may reject the hypothesis to zero i.e. $\beta_2 = \beta_3 = \dots = \beta_K = 0$. In fact insignificant t-ratios and a very high R^2 is one of the signals of multicollinearity.

⑤ As long as multicollinearity is not perfect, estimation of the regression coefficient is possible but their estimates and their standard errors become very sensitive to even a small change

in data.

Test / methods of multicollinearity

(1) Test based on Frischo's congruence analysis:

The seriousness of the effects of multicollinearity since the seems to depend on the effects of degree of inter-correlation i.e. r_{xixj} as well the overall co-relation co-efficient R^2 . The standard error, the partial correlation co-efficient and total correlation coefficient may be used for testing multi-collinearity. However, none of these, by itself is a satisfactory indicator of multicollinearity, because —

(i) Large standard error do not always appear with multicollinearity. Large standard error may arise for various other reasons.

(ii) A high r_{xixj} is sufficient condition but not a necessary condition of multicollinearity.

(iii) R^2 may be high and yet the estimates may not be significant and imprecise.

However, the all these criteria should help the detection of multicollinearity and has been suggested by Frischo's test.

Frisch's test

In this test it is required to estimate all the possible regressions between the variables which are present in relationship taking each variable successively as the dependent variable and considering all possible regression of each variables and all others which are introduced into the analysis.

(2) The Farrow - Glauber test of multicollinearity:

A statistical test of multicollinearity has been developed by Farrow - Glauber and it is a set of 3 tests.

(i) Firstly, chi-square test is used to detect the presence and severity of multicollinearity in a function of several explanatory variables. The basic hypothesis here are -

H_0 : The X's are orthogonal

H_a : The X's are not orthogonal.

It is observed that ^{if the} chi-square is greater than the theoretical chi-square value, ^{with $\frac{1}{2}$} $k(k-1)$ ^{degrees of freedom}, then we reject the assumption of orthogonality and we say ^{that} there is multicollinearity in the

function. If the observed chi-square is less than theoretical chi-square ~~is less~~ than theoretical chi-square we accept the assumption of orthogonality and say there is no multicollinearity.

(ii) Secondly, an F test is also used for the location of multicollinearity. To locate the factors which are multicollinear, Farrer and Gleuser computed the multiple correlation coefficients amongst the explanatory variables. ($R^2_{x_1, x_2, x_3, \dots, x_k}$, $R^2_{x_2, x_1, x_3, \dots, x_k}$) and in general $R^2_{x_i, x_1, x_2, \dots, x_k}$ and then test the statistical significance of these multiple correlation coefficient with an F test. For each multiple correlation coefficient, we compute the observed F^* .

$$i.e. F^* = \frac{(R^2_{x_i, x_1, x_2, \dots, x_k}) / (k-1)}{(1 - R^2_{x_i, x_1, x_2, \dots, x_k}) / (n-k)}$$

where, $n =$ size of the sample

$k =$ Number of explanatory variables.

The hypotheses are $H_0: R^2_{x_i, x_1, x_2, \dots, x_k} = 0$

$H_a: R^2_{x_i, x_1, x_2, \dots, x_k} \neq 0$

The observed F^* is compared with theoretical F value with degrees of freedom $U_1 = (K-1)$
 $U_2 = (n-1)$

If $F^* > F$, we reject H_0 and say that there is multicollinearity.

And if $F^* < F$, we accept H_0 and say that there is no multicollinearity.

(iii) Thirdly, a t -test is used to detect the variables which cause multicollinearity. To find which variables are responsible for the multicollinearity, we compute the partial correlation coefficient among the explanatory variables and test their statistical significance with the t -statistics.

Remedial measures

1. A priori information:

Suppose we consider the model —

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + U_i \quad \text{--- (1)}$$

where, Y = Consumption

X_2 = Income

X_3 = Wealth

As noted income and wealth variables are highly collinear but suppose a priori we believe that $\beta_3 = 10\beta_2$, i.e. the rate of change of consumption with respect to wealth is $\frac{1}{10}$ the corresponding rate with respect to income. We can then run the following regression; —

$$Y_i = \beta_1 + \beta_2 X_{2i} + 10\beta_2 X_{3i} + U_i$$

$$= \beta_1 + \beta_2 X_i + U_i$$

where, $X_i = X_{2i} + 10X_{3i}$

Once we obtain $\hat{\beta}_2$ we can estimate $\hat{\beta}_3$ from the postulated relationship between β_2 and β_3 } However a priori information comes from previous empirical work in which the collinearity problem happens to be less serious or from the relevant field of study.

2. Combining cross-section data and time-series data :

Pooling the data i.e. combining the cross-section data and time series data is another remedial measure for the problem of multicollinearity. Suppose we want to study demand for the auto-mobiles in country

and assumed we have time series data on the numbers of cars sold, average prices of the cars and consumer income. Suppose we have here, —

$$\log Y_t = \beta_1 + \beta_2 \log P_t + \beta_3 \log I_t + U_t$$

where, Y = Number of cars sold

P = Average price

I = Income

T = Time

Now our aim is to estimate β_2 the price elasticity and β_3 the income elasticity.

Now in time series data the price and the income variables tend to be collinear. Therefore if we run the preceding regression we will be faced with the problem of multi-collinearity. Way out of this is that if we have a cross section data, we can obtain a fairly reliable estimate of the income elasticity of β_3 because in such data we are at one point of time and the prices do not change much.

3. Dropping a variable :

When faced with severe multicollinearity one of the simplest way to drop one of the multicollinear variables, But however in dropping from

a model, we may be committing specific
-ation error which arises from the
incorrect specification of the model used
in the analysis.

4. Transformation of variables :

Suppose we have time-series data on
consumption (Y), expenditure, income and wealth.
One reason for high multicollinearity
between Y and w in such data is that
over time both the variables tend to
move in the same direction. If the
relation,

$$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + U_t \quad \text{--- (1)}$$

holds at time t , it must also hold at
time, $(t-1)$ because the origin of t is
arbitrary any way. Therefore, we have

$$Y_{t-1} = \beta_1 + \beta_2 X_{2,t-1} + \beta_3 X_{3,t-1} + U_{t-1} \quad \text{--- (2)}$$

If we subtract (2) from (1), we have

$$Y_t - Y_{t-1} = \beta_2 (X_{2t} - X_{2,t-1}) + \beta_3 (X_{3t} - X_{3,t-1}) + V_t \quad \text{--- (3)}$$

where, $V_t = U_t - U_{t-1}$

Eqⁿ (3) is known as the first differential form because we run the regression not only on original variables, but on the differential of successive values of the variables.

The 1st differential regression model reduces the severity of multicollinearity, because although the levels of X_2 and X_3 are highly co-related, there is no a priori reason to believe their differential will also be highly collinear.

5. Additional / new data :

Since multicollinearity is a sample feature, it is possible that in another sample involving the same variables, collinearity may not be as serious as in the 1st sample. Sometimes simply increase size of the sample, we may get rid of collinearity problem.