



7. Multiple Regression Analysis: The Problem of Analysis

Three-Variable Model: Notation and Assumptions

Let us consider the following three-variable PRF as:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i \quad (k=3)$$
Regressand
Regressors

where

- ✓ Y_i is the dependent variable (regressand)
- ✓ X_2 and X_3 are the regressors or the explanatory variables
- ✓ u_i is the stochastic disturbance term

Remark: the subscript i is denoted the observation i from our sample data.

In case our data are time series, the subscript t will denote the t observation.

β_1 means the average value of Y when X_2 and X_3 are set equal to zero

β_2 and β_3 are called the partial regression coefficients.

We will talk about the meaning of β_1 and β_2 shortly after knowing the assumptions of the classical linear regression model (CLRM)

$Y_i \quad X_{2i} \quad X_{3i}$

	1		
	2		
	3		
	4		
	5		
	6		
	7		
	8		
	9		
	10		
	11		
	12		
	13		
	14		
	15		
	16		
	17		
	18		
	19		
	20		
	21		
	22		
	23		
	24		
	25		
	26		
	27		
	28		
	29		
	30		
	31		
	32		
	33		
	34		
	35		
	36		
	37		
	38		
	39		
	40		

$Y_t = \beta_1 + \beta_2 X_{2t} + \beta_3 X_{3t} + u_t$

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + u_i$$
(k parameters)

of X_i = $k-1$ terms

Chapter 7. Multiple Regression Analysis: The Problem of Analysis

Under the CLRM, we assume:

1. Zero mean value of u_i

$$E(u_i | X_{2i}, X_{3i}) = 0$$
Mean value of Residuals = 0

2. No serial correlation

$$\text{Cov}(u_i, u_j) = 0 \quad \text{for all } i \neq j$$

3. Homoscedasticity

$$\text{Var}(u_i) = \sigma^2 \rightarrow \text{constant variance of } u_i$$

4. Zero covariance between u_i and each X variable, or

$$\text{Cov}(u_i, X_{2i}) = 0, \quad \text{Cov}(u_i, X_{3i}) = 0$$

5. No specification bias or

The model is correctly specified.

specification bias comes from
 (1) add redundant variable(s)

6. No exact collinearity between the X variables or

or (2) omit important variable(s)
 (Exclusion Errors)

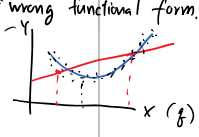
Example of exact linear relationship

suppose $X_{3i} = 2X_{2i}$

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

$$= \beta_1 + \beta_2 X_{2i} + \beta_3 [2X_{2i}] + u_i$$

$$= \beta_1 + [\beta_2 + 2\beta_3] X_{2i} + u_i$$



Partial Effect of X_{2i} on Y and of X_{3i} only cannot be estimated when X_{2i} and X_{3i} have an exact linear relationship.

EXACT LINEAR RELATIONSHIP

X_{2i}	X_{3i}
2	10
4	20
6	30
8	40
$X_{3i} = 5X_{2i}$	

IV EXACT LINEAR RELATIONSHIP

X_{2i}	X_{3i}	$X_{3i} = 5X_{2i} + u_i$
2	10	$10 = 5 \cdot 2 + 0$
4	22	$22 = 5 \cdot 4 + 2$
6	28	$28 = 5 \cdot 6 - 2$
8	40	$40 = 5 \cdot 8 + 0$

be estimated when x_{2i} and x_{3i} have an exact linear relationship.

so called "Perfect Multicollinearity"

4	20
6	30
8	40

$x_{3i} = 5x_{2i}$

2	10	10 = 5*2 + 0
4	22	22 = 5*4 + 2
6	28	28 = 5*6 - 2
8	40	40 = 5*8 + 0

7.1 OLS Estimation of the Partial Regression Coefficients 117

By the above assumptions, we can find out the conditional expectation of Y_i : $E[Y_i | x_{2i}, x_{3i}]$

Take $E[\cdot | x_{2i}, x_{3i}]$ to the PRF:

$$E[Y_i | x_{2i}, x_{3i}] = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + E[u_i | x_{2i}, x_{3i}] = 0$$

The meaning of partial coefficients:

β_2

$$\frac{d E[Y_i | x_{2i}, x_{3i}]}{d x_{2i}} = \beta_2 \Rightarrow \text{on average, when } x_{2i} \text{ changes by 1 unit, } E[Y_i | x_{2i}, x_{3i}] \text{ will change by } \beta_2 \text{ units, holding } x_{3i} \text{ constant.}$$

So $\beta_2 =$ partial effect of x_2 on $E[Y_i | x_{2i}, x_{3i}]$, holding x_{3i} constant.

β_1

$$D - I - Y$$

In general $\lambda_2 x_{2i} + \lambda_3 x_{3i} = 0$

If λ_2, λ_3 are found to make the above equation true, x_{2i} and x_{3i} has "EXACT" linear relationship.

w/ inexact linear relationship (Imperfect multicollinearity) β_2 and β_3 can still be estimated.

7.1 OLS Estimation of the Partial Regression Coefficients

In order to find the OLS estimators, we need to write down the sample regression function (SRF) corresponding to the PRF:

ACTUAL Y ESTIMATED Y RESIDUALS

$$y_i = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + u_i$$

GOAL: FIND $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$ (ALSO $\hat{\beta}_i$)

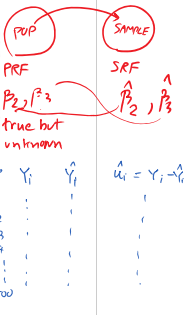
OLD TRICK: MDV $\sum \hat{u}_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})^2$

F.O.C $\frac{\partial \sum \hat{u}_i^2}{\partial \hat{\beta}_1} = 2 \sum (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-1) = 0 \Rightarrow \sum \hat{u}_i = 0 \quad \text{--- (1)}$

$\frac{\partial \sum \hat{u}_i^2}{\partial \hat{\beta}_2} = 2 \sum (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-x_{2i}) = 0 \Rightarrow \sum \hat{u}_i x_{2i} = 0 \quad \text{--- (2)}$

$\frac{\partial \sum \hat{u}_i^2}{\partial \hat{\beta}_3} = 2 \sum (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})(-x_{3i}) = 0 \Rightarrow \sum \hat{u}_i x_{3i} = 0 \quad \text{--- (3)}$

3 Equations, 3 unknown, we can solve for $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \dots$



Chapter 7. Multiple Regression Analysis: The Problem of Analysis

From the FOC, we then get the normal equations:

$$\begin{aligned} \sum y_i x_{2i} &= \beta_1 \sum x_{2i} + \beta_2 \sum x_{2i}^2 + \beta_3 \sum x_{2i} x_{3i} \\ \sum y_i x_{3i} &= \beta_1 \sum x_{3i} + \beta_2 \sum x_{2i} x_{3i} + \beta_3 \sum x_{3i}^2 \end{aligned}$$

We therefore get:

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum y_i \sum x_{2i} x_{3i} - \sum x_{2i} \sum y_i \sum x_{3i}}{(\sum x_{2i}^2)(\sum x_{3i}^2) - (\sum x_{2i} x_{3i})^2} \\ \hat{\beta}_2 &= \frac{(\sum y_i x_{2i})(\sum x_{3i}^2) - (\sum y_i x_{3i})(\sum x_{2i} x_{3i})}{(\sum x_{2i}^2)(\sum x_{3i}^2) - (\sum x_{2i} x_{3i})^2} \\ \hat{\beta}_3 &= \frac{(\sum y_i x_{3i})(\sum x_{2i}^2) - (\sum y_i x_{2i})(\sum x_{2i} x_{3i})}{(\sum x_{2i}^2)(\sum x_{3i}^2) - (\sum x_{2i} x_{3i})^2} \end{aligned}$$

solutions

Variance and Standard Errors of OLS Estimators

$$\text{var}(\hat{\beta}_1) = \left[\frac{1}{n} + \frac{\sum x_{2i}^2 \sum x_{3i}^2 - 2 \sum x_{2i} x_{3i} \sum x_{2i} x_{3i}}{\sum x_{2i}^2 \sum x_{3i}^2 - (\sum x_{2i} x_{3i})^2} \right] \sigma^2$$

$$se(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)}$$

$$\text{var}(\hat{\beta}_2) = \frac{\sum x_{3i}^2}{(\sum x_{2i}^2)(\sum x_{3i}^2) - (\sum x_{2i} x_{3i})^2} \sigma^2$$

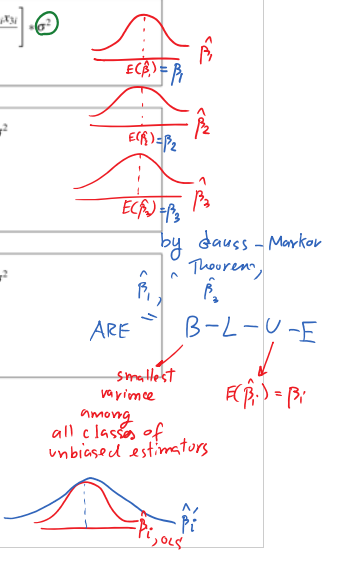
$$\text{var}(\hat{\beta}_3) = \frac{\sigma^2}{\sum x_{2i}^2 (1 - r_{23}^2)}$$

$$se(\hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_2)}$$

$$\text{var}(\hat{\beta}_1) = \frac{\sum x_{2i}^2}{(\sum x_{2i}^2)(\sum x_{3i}^2) - (\sum x_{2i} x_{3i})^2} \sigma^2$$

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

$$se(\hat{\beta}_3) = \sqrt{\text{var}(\hat{\beta}_3)}$$



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

$$\hat{\sigma}^2 = \frac{\sum \hat{u}_i^2}{n-3} \rightarrow \text{unbiased estimator of } \sigma^2 !!!$$

Here, $k=3 \rightarrow$ NUMBER OF EXPLANATORY VARIABLES + INTERCEPT.

7.2 Properties of OLS Estimators: $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$

① The three-variable regression plane passes through the Means $\bar{Y}, \bar{X}_2, \bar{X}_3$

So, $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{X}_2 + \hat{\beta}_3 \bar{X}_3$

② The mean value of the estimated Y_i (\hat{Y}_i) is equal to the mean value of actual Y_i .

Proof: $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i}$

we have $\hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X}_2 - \hat{\beta}_3 \bar{X}_3$

So $\hat{Y}_i = (\bar{Y} - \hat{\beta}_2 \bar{X}_2 - \hat{\beta}_3 \bar{X}_3) + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i}$

$$\hat{Y}_i = \bar{Y} + \hat{\beta}_2 x_{2i} - \hat{\beta}_2 \bar{X}_2 + \hat{\beta}_3 x_{3i} - \hat{\beta}_3 \bar{X}_3$$

$$= \bar{Y} + \hat{\beta}_2 (x_{2i} - \bar{X}_2) + \hat{\beta}_3 (x_{3i} - \bar{X}_3)$$

$$\hat{Y}_i = \bar{Y} + \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} \quad \text{where } x_{2i} = x_{2i} - \bar{X}_2, x_{3i} = x_{3i} - \bar{X}_3$$

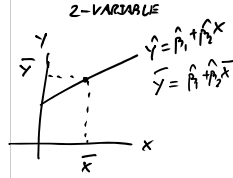
Take summation operators throughout and then divide w/ N:

$$\sum \hat{Y}_i = \sum \bar{Y} \quad \#$$

③ $\hat{Y}_i = \bar{Y} + \hat{\beta}_3 x_{3i}$

actual $y_i = \hat{y}_i + \hat{u}_i$ residuals

$y_i = \hat{\beta}_2 x_{2i} + \hat{\beta}_3 x_{3i} + \hat{u}_i \rightarrow$ SRF in Deviation Form.



$$\hat{u}_i = y_i - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i}$$

$$\sum \hat{u}_i^2 = \sum \hat{u}_i (y_i - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i})$$

Properties of OLS Estimators (Cont):

$$\sum \hat{u}_i^2 = \sum \hat{u}_i y_i - \hat{\beta}_2 \sum \hat{u}_i x_{2i} - \hat{\beta}_3 \sum \hat{u}_i x_{3i}$$

$$\sum \hat{u}_i^2 = \sum (y_i - \hat{\beta}_2 x_{2i} - \hat{\beta}_3 x_{3i}) \cdot y_i$$

RSS Residual Sum of Squares

$$= \sum y_i^2 - \hat{\beta}_2 \sum x_{2i} y_i - \hat{\beta}_3 \sum x_{3i} y_i$$

\rightarrow we will make use of this soon...

④ $\sum \hat{u}_i = 0$ and $\sum \hat{u}_i = \bar{\hat{u}}_i = 0$

⑤ \hat{u}_i are uncorrelated with X_{2i} : $\sum \hat{u}_i X_{2i} = 0$

\hat{u}_i are uncorrelated with X_{3i} : $\sum \hat{u}_i X_{3i} = 0$

⑥ \hat{u}_i are uncorrelated with \hat{Y}_i : $\sum \hat{u}_i \hat{Y}_i = 0$

Proof: $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i}$
 $\sum \hat{u}_i \hat{Y}_i = \sum \hat{\beta}_1 \hat{u}_i + \sum \hat{\beta}_2 X_{2i} \hat{u}_i + \sum \hat{\beta}_3 X_{3i} \hat{u}_i$
 $= 0 + 0 + 0$

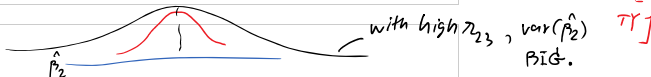
So $\sum \hat{u}_i \hat{Y}_i = 0$

*** ⑦ Recall that $\text{var}(\hat{\beta}_2) = \frac{1}{\sum x_{2i}^2} \cdot \frac{\sigma^2}{(1-r_{23}^2)}$
 $\text{var}(\hat{\beta}_3) = \frac{1}{\sum x_{3i}^2} \cdot \frac{\sigma^2}{(1-r_{23}^2)}$

Given $\sum x_{2i}^2$ and σ^2 , if $r_{23} \uparrow \rightarrow \text{var}(\hat{\beta}_2) \uparrow \rightarrow$ Reduce Accuracy of Estimation!

If $r_{23} = 1$, then $\text{var}(\hat{\beta}_2)$ becomes INFINITE!

In this case, $\hat{\beta}_2$ cannot be estimated [PERFECT MULTICOLLINEARITY]



Properties of OLS Estimators (Cont):

Properties of OLS Estimators (Cont.)

- Given λ_{23} and σ^2 , If $\sum x_{2i}^2 \uparrow$, $\text{var}(\beta_2)$ will become lower

Advice: get more observations as much as possible!
This will help to increase power of estimation!

- Given λ_{23} and $\sum x_{2i}^2$, if $\sigma^2 \uparrow$, $\text{var}(\beta_2) \uparrow$ ☹️

⑧ $\hat{\beta}_2$ and $\hat{\beta}_3$ are Best Linear Unbiased Estimator.

with λ_{23} , $\text{var}(\beta_2)$ BIG.

improve accuracy of estimation.

	low x_{3i}	high x_{3i}
low x_{2i}	?	?
high x_{2i}	?	?

(wealth) x_{3i}
(income) x_{2i}

The Multiple Coefficient of Determination R^2 and the Multiple Coefficient of Correlation R

In this section, we will study how to measure the proportion of the variation in Y explained by the variables X_2 and X_3 jointly. This is the same concept of r^2 that we have learned before.

The quantity that gives this information is known as the **multiple coefficient of determination** and is denoted by R^2 .

To derive R^2 , we firstly write down the following equation:

$$Y_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \hat{u}_i \rightarrow \text{STANDARD FORM OF SRF}$$

$$= \hat{Y}_i + \hat{u}_i \quad (7.1)$$

where \hat{Y}_i is the estimated value of Y_i from the fitted regression line and is an estimator of true $E(Y_i|X_{2i}, X_{3i})$.

7.1 may be written as

$$y_i = \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \hat{u}_i \rightarrow \text{DEVIATION FORM OF SRF}$$

$$= \hat{y}_i + \hat{u}_i \quad (7.2)$$

Squaring 7.2 on both sides and summing over the sample values, we obtain

$$\sum \hat{y}_i^2 = \sum \hat{y}_i^2 + \sum \hat{u}_i^2 + 2 \sum \hat{y}_i \hat{u}_i$$

$$\sum \hat{y}_i^2 = \sum \hat{y}_i^2 + \sum \hat{u}_i^2 \quad (7.3)$$

TSS = $\sum \hat{y}_i^2$ (variation in y_i that x_2 and x_3 cannot help to explain)
ESS = $\sum \hat{y}_i^2$ (variation in y_i that could be explained by the help of adding x_2 and x_3)
RSS = $\sum \hat{u}_i^2$

Baseline errors when we simply use sample mean to estimate y_i

Recall that $TSS = ESS + RSS$

$$1 = \frac{ESS}{TSS} + \frac{RSS}{TSS}$$

R^2 → Proportion Variation in Y that can be explained by the regression (by using x_2 and x_3)

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2} = 1 - \frac{\sum (Y_i - \hat{Y}_i)^2}{\sum (Y_i - \bar{Y})^2}$$

RESIDUAL SUM OF SQUARES (RSS)
TOTAL SUM OF SQUARES (TSS)

$$ESS \leq TSS \leq ESS + RSS$$

FROM $\sum \hat{u}_i^2 = \sum y_i^2 - \hat{\beta}_2 \sum x_{2i} y_i - \hat{\beta}_3 \sum x_{3i} y_i$,
 $\sum y_i^2 = \sum \hat{y}_i^2 + \sum \hat{u}_i^2$
 $\sum y_i^2 - \sum \hat{y}_i^2 = \sum \hat{u}_i^2 = \hat{\beta}_2 \sum x_{2i} y_i + \hat{\beta}_3 \sum x_{3i} y_i$
 $\therefore \sum \hat{y}_i^2 = \hat{\beta}_2 \sum x_{2i} y_i + \hat{\beta}_3 \sum x_{3i} y_i$
 ESS

$R^2 = \frac{ESS}{TSS} = \frac{\hat{\beta}_2 \sum y_i x_{2i} + \hat{\beta}_3 \sum y_i x_{3i}}{\sum y_i^2} \rightarrow R^2$ in another form. (7.4)

$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki}$
 # Regressors = $k-1$ variables
 R_j^2 is the R^2 in the regression of X_j on the remaining regressors ($k-2$).

$x_2 = f(x_3)$

The three or more variable analogue of r is the coefficient of multiple correlation, denoted by R , and it is a measure of the degree of association between Y and all the explanatory variables jointly. Although r can be positive or negative, R is always taken to be positive.

$Var(\hat{\beta}_j) = \frac{\sigma^2}{\sum y_i^2} \left(\frac{1}{1-R^2} \right)$

7.2.1 R^2 and the Adjusted R^2

It should be noted that R^2 is a nondecreasing function of the number of explanatory variables. Thus, when the number of regressors increases, R^2 almost invariably increases and never decreases. **In other words, an additional X variable will not decrease R^2 !**

To explain this fact, let us write down the definition of R^2 again:

$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2}$ (7.5)

Therefore, in comparing two regression models with the same dependent variable but differing number of X variables, one should be very wary of choosing the model with the highest R^2 .

In light of comparing two R^2 terms, we have to take into account the number of X variables present in the model. To achieve this goal, we can consider the alternative coefficient of determination, which is as follows:

To compare between the two models, it is advised to compare adjusted r^2 instead of R^2 .
 $\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2 / (n-3)}{\sum y_i^2 / (n-1)}$ for $k=3$
 General form:
 $\hat{Y}_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + u_i$
 $\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2 / (n-k)}{\sum y_i^2 / (n-1)}$

k = the number of parameters in the model including the intercept term. \rightarrow
 n = the number of observations in the sample data.

The above equation is known as the adjusted R^2 , denoted by \bar{R}^2 . The term adjusted means adjusted for the df associated with the sums of squares entering into 7.5.

We can rewrite the the adjusted R^2 as:

$$\bar{R}^2 = 1 - \frac{\hat{\sigma}^2}{S_Y^2}$$

Where $\hat{\sigma}^2 = \frac{\sum \hat{u}_i^2}{n-k} \rightarrow$ unbiased estimator of σ^2 residual variance

$S_Y^2 = \frac{\sum y_i^2}{n-1} \rightarrow$ sample variance of Y

We can also get the equation which shows the relationship between \bar{R}^2 and R^2 :

$$\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2} \cdot \frac{(n-1)}{(n-k)}$$

As $R > 1 \rightarrow \frac{(n-1)}{(n-k)} > 1$

$$\bar{R}^2 = 1 - [1 - R^2] \cdot \frac{(n-1)}{(n-k)}$$

$\bar{R}^2 < R^2$ ***

Note: \bar{R}^2 CAN BE NEGATIVE.
Ex: When $R^2 = 0$. (verify this)

Besides R^2 and \bar{R}^2 as goodness of fit measures, other criteria are often used to judge the adequacy of a regression model. Two of these are Akaike's Information criterion and Amemiya's Prediction criteria, which are used to select between competing models. We will discuss these criteria in greater detail later.

* When comparing two R^2 values, make sure that

- ① Dependent variables of the two models are the same,
 - ② sample size must be the same.
- ① $\ln Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i}$
 ② $Y_i = \alpha_1 + \alpha_2 X_{2i} + \alpha_3 X_{3i}$ } ① AND ② CANNOT BE COMPARED DIRECTLY BY USING R^2 AS THE INDEPENDENT VARIABLES ARE NOT THE SAME.