

ECO341: ASSIGNMENT REPORT

ACKNOWLEDGEMENT:

We would like to thank **Dr. Deep Mukherjee** for giving us this opportunity to work on this project which helped us gain a deep insight into the concepts by working on a real-life problem. We would also like to thank **Mr. Debarun Sengupta** for teaching us the R language, without which we wouldn't have been able to solve the problem. It was indeed a great learning experience.

PROBLEM DESCRIPTION:

The problem statement suggests a “cure” for multicollinearity between the columns of \mathbf{X}_1 and \mathbf{X}_2 , where \mathbf{X}_1 and \mathbf{X}_2 are two sets of explanatory variables partitioned from the same dataset. First, we regress each variable of \mathbf{X}_2 on \mathbf{X}_1 and compute the residuals matrix \mathbf{Z}_2 . Then we regress the dependent variable \mathbf{y} on $(\mathbf{X}_1 \& \mathbf{Z}_2)$ and $(\mathbf{X}_1 \& \mathbf{X}_2)$. $\mathbf{c} = (\mathbf{c}_1, \mathbf{c}_2)$ are the coefficients of the regression of \mathbf{y} on $(\mathbf{X}_1 \& \mathbf{Z}_2)$ and $\mathbf{b} = (\mathbf{b}_1, \mathbf{b}_2)$ are the coefficients of the regression of \mathbf{y} on $(\mathbf{X}_1 \& \mathbf{X}_2)$. First, we algebraically prove that $\mathbf{c}_2 = \mathbf{b}_2$ and \mathbf{c}_1 is biased and \mathbf{c}_2 unbiased. After that, we verify our claims using real-life data (on gasoline) by using R for computation.

THEORETICAL RESULT:

We have,

$$\mathbf{Z}_2 = \mathbf{M}_1 \mathbf{X}_2 \text{ (from the regression of } \mathbf{X}_2 \text{ on } \mathbf{X}_1 \text{)}$$

where

(i) $\mathbf{M}_1 = \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$ from the regression of \mathbf{X}_2 on \mathbf{X}_1

(ii) \mathbf{Z}_2 denotes residuals from the regression of each variable of \mathbf{X}_2 on \mathbf{X}_1

Now,

$$\mathbf{X}_1^T \mathbf{M}_1 = \mathbf{X}_1^T - (\mathbf{X}_1^T \mathbf{X}_1)(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T = \mathbf{0}$$

Thus,

$$\mathbf{X}_1^T \mathbf{Z}_2 = \mathbf{X}_1^T \mathbf{M}_1 \mathbf{X}_2 = \mathbf{0} \dots (3)$$

So, when we regress \mathbf{y} on $(\mathbf{X}_1, \mathbf{Z}_2)$, we can instead regress separately on \mathbf{X}_1 and \mathbf{Z}_2 . This will hold because the moment matrix is a block diagonal, with blocks $\mathbf{X}_1^T \mathbf{X}_1$ and $\mathbf{Z}_2^T \mathbf{Z}_2$, its inverse will likewise be a block diagonal and hence regression of \mathbf{y} on \mathbf{X}_1 and \mathbf{Z}_2 is the same as separate regressions on the two data sets (as \mathbf{X}_1 and \mathbf{Z}_2 are orthogonal using (3)).

Now,

$$\mathbf{c}_1 = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y} \text{ ---- (1) (regressing } \mathbf{y} \text{ on } \mathbf{X}_1 \text{ separately)}$$

$$\mathbf{c}_2 = (\mathbf{Z}_2^T \mathbf{Z}_2)^{-1} \mathbf{Z}_2^T \mathbf{y} \text{ ---- (2) (regressing } \mathbf{y} \text{ on } \mathbf{Z}_2 \text{ separately)}$$

Since \mathbf{M}_1 is symmetric and idempotent and $\mathbf{Z}_2 = \mathbf{M}_1 \mathbf{X}_2$, using both of these in (2)

$$\mathbf{c}_2 = (\mathbf{X}_2^T \mathbf{M}_1 \mathbf{X}_2)^{-1} (\mathbf{X}_2^T \mathbf{M}_1 \mathbf{y})$$

In the original regression of \mathbf{y} on \mathbf{X}_1 and \mathbf{X}_2 , \mathbf{b}_2 can be calculated using Frisch–Waugh–Lovell Theorem. Using that, it comes out to be the same as \mathbf{c}_2 , implying $\mathbf{c}_2 = \mathbf{b}_2$.

But,

$$\begin{aligned} \mathbf{c}_1 &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y} \\ &= (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T (\mathbf{X}_1 \mathbf{b}_1 + \mathbf{X}_2 \mathbf{b}_2 + \mathbf{e}) \\ &= \mathbf{b}_1 + (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \mathbf{b}_2 \quad (\mathbf{X}_1^T \mathbf{e} = \mathbf{0} \text{ as OLS residuals}) \end{aligned}$$

From the above expression, it is clear that \mathbf{c}_1 is not equal to \mathbf{b}_1 unless $\mathbf{X}_1^T \mathbf{X}_2 = \mathbf{0}$

Hence we conclude **\mathbf{c}_1 is biased whereas \mathbf{c}_2 is unbiased**

EMPIRICAL RESULTS:

We have taken

X_1 matrix as [1(constant), GASP, and PCINCOME]

X_2 matrix as [PD, PN, and PS]

y = GASEXP

Z_2 = [Ed, En, Es], the residuals matrix from the regression of X_2 on X_1 .

where PCINCOME is taken to be INCOME/POP (Per capita Income)

Models		
=====		
	Dependent variable:	

	GASEXP)	
	x1-x2	x1-z2
	(1)	(2)

INTERCEPT	-22.343*** (4.466)	-64.568*** (3.461)
GASP	1.231*** (0.063)	1.537*** (0.026)
PCINCOME	452.034*** (97.353)	768.261*** (63.805)
PD	-0.697*** (0.068)	
PN	-0.125 (0.214)	
PS	0.684*** (0.094)	
Ed		-0.697*** (0.068)
En		-0.125 (0.214)
Es		0.684*** (0.094)

Observations	52	52
R2	0.999	0.999
Adjusted R2	0.999	0.999
Residual Std. Error (df = 46)	2.319	2.319
F Statistic (df = 6; 46)	13,139.820***	13,139.820***

CONCLUSION:

If we compare the coefficients of both the regressions, we find that the coefficient of X_2 and Z_2 in both the cases is the same showing that $b_2 = c_2$ or c_2 is unbiased whereas c_1 is not.

Hence theoretical results match with the empirical results.

TEAM DETAILS:

Danish Ahmad(160219)

Sakib Malik(180652)

Shivam Goel(180714)