



International Journal of Applied Business and Economic Research

ISSN : 0972-7302

available at <http://www.serialsjournal.com>

© Serials Publications Pvt. Ltd.

Volume 15 • Number 15 • 2017

Clustering on The Basis of Regression Equations: The Case of Three Regressions

Jagabandhu Saha¹

¹Department of Economics and Politics, Visva-Bharati University, India, E-mail: saha_jagabandhu@yahoo.com

Abstract: While Chow test is applicable for testing the equality between sets of coefficients in two linear regressions, this paper attempts to construct a test procedure not only to compare the equality between the sets of coefficients in two linear regressions, but also, in case they are not equal, to provide detailed informations about the inequality of the sets. Not only that, in this paper attempt is also made to accomplish all these for not just two linear regressions but for the two linear regressions of all possible pairs of linear regressions out of any number of given linear regressions, and then the results of all these comparisons are used in order to form clusters among the regressions on the basis of some principle stated therein. However, in this paper we consider the case of three regressions only. The procedure is then illustrated through comparison of Engel curves on food for different Socioeconomic groups of Rural India, using NSSO data.

1. INTRODUCTION

Testing the Equality between Sets of Coefficients in two Linear Regressions by Chow Test

(Chow 1960)^[1] is well known. There is no problem if in the Chow Test the null hypothesis of equality between the sets of coefficients is not rejected (as in the examples in his paper). But if rejected, then, naturally, one is probed to the questions:

- at which component/s the sets differ, and
- for each of these components, between the two coefficients of the two regressions concerned, which one is larger/smaller.

Chow test, however, does not provide any answer to these questions at all. These problems can be resolved with some modifications of the model (Saha and Pal 2014)^[2]. Saha and Pal introduced the concept of “component wise complete comparison” (CCC)² in order to overcome this problem. The test procedure

for CCC between every two successive regressions out of any number of given successive regressions was developed. Now, this paper attempts to construct a test procedure for CCC between not only every two successive regressions out of any number of given successive regressions, but between the two linear regressions of *all possible pairs of regressions out of any number of given regressions*, and then uses all these comparisons in order to form clusters among the regressions on the basis of this simple principle: *all those regressions which satisfy the condition that the vectors of coefficients of any two of these regressions do not differ from each other significantly will form a cluster*. But, however, in this paper we consider the case of three regressions only³. The rest of the paper can be outlined as follows. In Section 2, we put the problem in the formal terms, problem of finding test procedure for CCC between the two linear regressions of all possible pairs of regressions out of three given regressions. Section 3 is devoted for the methodology for solving this problem and then for discussion about process of clustering. In Section 4 we consider a numerical example in order to illustrate the methodology and the process of clustering while in Section 5 we present our conclusions.

2. THE MODEL

Consider the problem of finding test procedure for CCC between the two linear regressions of all possible pairs of regressions out of three given regressions as follows:

$$\begin{aligned} y^{(1)} &= a_1^{(1)} + a_2^{(1)} x_2^{(1)} + a_3^{(1)} x_3^{(1)} + \dots + a_k^{(1)} x_k^{(1)} + u^{(1)}, \\ y^{(2)} &= a_1^{(2)} + a_2^{(2)} x_2^{(2)} + a_3^{(2)} x_3^{(2)} + \dots + a_k^{(2)} x_k^{(2)} + u^{(2)} \\ y^{(3)} &= a_1^{(3)} + a_2^{(3)} x_2^{(3)} + a_3^{(3)} x_3^{(3)} + \dots + a_k^{(3)} x_k^{(3)} + u^{(3)}, \end{aligned} \quad \dots (1)$$

where, n_1, n_2, n_3 are the nos. of observations for these regressions. We, assume as in Chow test that the components of each of $u^{(1)}, u^{(2)}, u^{(3)}$ are *iid* $N(0, \sigma^2)$. Also the vectors $u^{(1)}, u^{(2)}, u^{(3)}$ are mutually independent.

Now, in order to accomplish CCC between first and second regressions in (1), also, first and third regressions, one requires to decide whether the differentials:

$$\begin{aligned} c_j^{12} &= a_j^2 - a_j^1 < 0 \text{ or } = 0 \text{ or } > 0, \text{ for all } j = 1, 2, \dots, k, \\ c_j^{13} &= a_j^3 - a_j^1 < 0 \text{ or } = 0 \text{ or } > 0, \text{ for all } j = 1, 2, \dots, k. \end{aligned}$$

Lastly, for second and third regressions, one requires to decide whether:

$$c_j^{23} = a_j^3 - a_j^2 < 0 \text{ or } = 0 \text{ or } > 0, \text{ for all } j = 1, 2, \dots, k.$$

A moment's reflection shows that the desired CCC, *i.e.*, CCC for all pairs of regressions out of the three regressions in (1), will be over when all the decisions enlisted above are completed.

3. THE METHODOLOGY

Methodology consists of two steps as follows:

- (a) combine the three regressions in (1) into a single regression equation model; modify the combined model in such way that the differentials c_j^{12}, c_j^{13} , for all $j = 1, 2, \dots, k$, appear as regression coefficients in the modified model, and then, run regression with the modified model and perform tests on the regression coefficients of this model (particularly, on the differentials concerned) and decide for each

of these differentials whether it is < 0 or $= 0$ or > 0 . This will complete CCC for two pairs of regressions: first and second, first and third.

- (b) do exactly similar as in the step 1). but with the last two regressions in (1). This will cover CCC for one and the last pair of regressions: second and third.

Thus CCC for all possible pairs of regressions (three pairs) is completed.

Let us work out these steps.

- (a) We combine the regressions in (1) as follows:

$$\begin{pmatrix} y_1^{(1)} \\ \vdots \\ y_{n_1}^{(1)} \\ y_1^{(2)} \\ \vdots \\ y_{n_2}^{(2)} \\ y_1^{(3)} \\ \vdots \\ y_{n_3}^{(3)} \end{pmatrix} \begin{pmatrix} 1 & x_{21}^{(1)} & \dots & x_{k_1}^{(1)} & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{2,n_1}^{(1)} & \dots & x_{k,n_1}^{(1)} & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 1 & x_{21}^{(2)} & \dots & x_{k_1}^{(2)} & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 & 1 & x_{2,n_2}^{(2)} & \dots & x_{k,n_2}^{(2)} & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 1 & 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 & 0 & 0 & \dots & 1 & 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} \end{pmatrix} \begin{pmatrix} a_1^{(1)} \\ \vdots \\ a_k^{(1)} \\ a_1^{(2)} \\ \vdots \\ a_k^{(2)} \\ a_1^{(3)} \\ \vdots \\ a_{k_3}^{(3)} \end{pmatrix} + \begin{pmatrix} u_1^{(1)} \\ \vdots \\ u_{n_1}^{(1)} \\ u_1^{(2)} \\ \vdots \\ u_{n_2}^{(2)} \\ u_1^{(3)} \\ \vdots \\ u_{n_3}^{(3)} \end{pmatrix} \quad \dots (2)$$

where, in (2) and subsequently, $n_i = n_j$, for all $i = 1, 2, 3$.

To modify the model (2) in order to get the differentials c_j^{12}, c_j^{13} , for all $j = 1, 2, \dots, k$, as regression coefficients, we rewrite it as:

$$\begin{pmatrix} y_1^{(1)} \\ \vdots \\ y_{n_1}^{(1)} \\ y_1^{(2)} \\ \vdots \\ y_{n_2}^{(2)} \\ y_1^{(3)} \\ \vdots \\ y_{n_3}^{(3)} \end{pmatrix} \begin{pmatrix} 1 & x_{21}^{(1)} & \dots & x_{k_1}^{(1)} & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{2,n_1}^{(1)} & \dots & x_{k,n_1}^{(1)} & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 1 & x_{21}^{(2)} & \dots & x_{k_1}^{(2)} & 1 & x_{21}^{(2)} & \dots & x_{k_1}^{(2)} & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{2,n_2}^{(2)} & \dots & x_{k,n_2}^{(2)} & 1 & x_{2,n_2}^{(2)} & \dots & x_{k,n_2}^{(2)} & 0 & 0 & \dots & 0 \\ 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} & 0 & 0 & \dots & 1 & 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} \\ \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} & 0 & 0 & \dots & 1 & 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} \end{pmatrix} \begin{pmatrix} a_1^{(1)} \\ \vdots \\ a_k^{(1)} \\ c_1^{12} \\ \vdots \\ c_k^{12} \\ c_1^{13} \\ \vdots \\ c_k^{13} \end{pmatrix} + \begin{pmatrix} u_1^{(1)} \\ \vdots \\ u_{n_1}^{(1)} \\ u_1^{(2)} \\ \vdots \\ u_{n_2}^{(2)} \\ u_1^{(3)} \\ \vdots \\ u_{n_3}^{(3)} \end{pmatrix} \quad \dots (3)$$

Now, we run regression with (3) and carry out tests as prescribed above.

(b) We combine the last two regressions in (1) as follows:

$$\begin{pmatrix} y_1^{(2)} \\ \vdots \\ y_{n_2}^{(2)} \\ y_1^{(3)} \\ \vdots \\ y_{n_3}^{(3)} \end{pmatrix} \begin{pmatrix} 1 & x_{21}^{(2)} & \dots & x_{k_1}^{(2)} & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{2,n_2}^{(2)} & \dots & x_{k,n_2}^{(2)} & 0 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 & 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & 0 & 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} \end{pmatrix} \begin{pmatrix} a_1^{(2)} \\ \vdots \\ a_{n_2}^{(2)} \\ a_1^{(3)} \\ \vdots \\ a_k^{(3)} \end{pmatrix} + \begin{pmatrix} u_1^{(2)} \\ \vdots \\ u_{n_2}^{(2)} \\ u_1^{(3)} \\ \vdots \\ u_{n_3}^{(3)} \end{pmatrix} \quad \dots(4)$$

To modify the model (4) in order to get the differentials c_j^{23} , for all $j = 1, 2, \dots, k$, as regression coefficients, we rewrite it as:

$$\begin{pmatrix} y_1^{(2)} \\ \vdots \\ y_{n_2}^{(2)} \\ y_1^{(3)} \\ \vdots \\ y_{n_3}^{(3)} \end{pmatrix} \begin{pmatrix} 1 & x_{21}^{(2)} & \dots & x_{k_1}^{(2)} & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{2,n_2}^{(2)} & \dots & x_{k,n_2}^{(2)} & 0 & 0 & \dots & 0 \\ 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} & 1 & x_{21}^{(3)} & \dots & x_{k_1}^{(3)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} & 1 & x_{2,n_3}^{(3)} & \dots & x_{k,n_3}^{(3)} \end{pmatrix} \begin{pmatrix} a_1^{(2)} \\ \vdots \\ a_k^{(2)} \\ c_1^{23} \\ \vdots \\ c_k^{23} \end{pmatrix} + \begin{pmatrix} u_1^{(2)} \\ \vdots \\ u_{n_2}^{(2)} \\ u_1^{(3)} \\ \vdots \\ u_{n_3}^{(3)} \end{pmatrix} \quad \dots(5)$$

Now, we run regression with (5) and carry out tests as prescribed above.

Thus CCC between the two regressions of all possible pairs of regressions out of the three given regressions in (1) are done.

Having done this one can partition the set of the given regressions into, say, clusters, *every cluster consisting of those regressions which satisfy the condition that the vectors of coefficients of any two of these regressions do not differ from each other significantly*. In order to sort out the clusters easily we first present the outcomes of all the above comparisons neatly by introducing a matrix, let it be called *Indicator Matrix* (IM), as:

$$i_{3 \times 3} = (r_{ij}), \quad \dots(6)$$

where, r_{ij} indicates whether the i th and the j th regressions in (1) differ from each other significantly or not and is defined as follows:

$$r_{ij} = 0,$$

when

$$i = j, 0,$$

when $i \neq j$ and the two regressions concerned do not differ from each another significantly,

$$1,$$

when $i \neq j$ and the two regressions concerned differ from each another significantly,

where,

$$i, j = 1, 2, 3^4.$$

It is this IM which will clearly show the clusters: no. of clusters as well as the regressions in each of these clusters. Needless to say that IM is a symmetric matrix with all the diagonal elements as zeros and it will be a null matrix (of order 3×3) iff all the three regressions coincide.

4. ILLUSTRATION

In the context of Engel curve for food, for each of the Socioeconomic groups ST, SC and OBC, considering regression of “proportion of Monthly Per Capita Expenditure on Food (pMPCEF)” on “Monthly Per Capita Expenditure (MPCE)”, *i.e.*, regression of pMPCEF on MPCE, totally three regression equations are constructed as follows ($m = 3$). State level data on Rural India on MPCE and MPCEF for the different Socioeconomic groups mentioned are used, the source of data being NSSO Report^[3].

Define

$x^{(i)}$ = MPCE of a State for the i^{th} Socioeconomic group, for all $i = 1, 2, 3$.

$y^{(i)}$ = pMPCEF of a State for the i^{th} Socioeconomic group, for all $i = 1, 2, 3$.

Then, the regressions considered are as follows:

$$y^{(i)} = a_1^{(i)} + a_2^{(i)} x_2^{(i)} + u^{(i)}, \quad \dots(7)$$

for all

$$i = 1, 2, 3,$$

and the task is to perform CCC between the two regressions of all pairs of regressions out of these three regressions. (The no. of observations for each regression is thirty (no. of States/UTs in India)⁵ and, so, we have: $m = 3, k = 2, n = 30$.)

Now, the two steps in the Methodology described above, are completed as follows.

1. Firstly, we run the regression (3) *with all the equations in (7)*. The decisions on the basis of the regression results come out as follows:

$$c_1^{12} < 0, c_2^{12} > 0, c_1^{13} = 0, c_2^{13} = 0.$$

This in turn means that as long as the relationship of MPCE with pMPCEF is concerned, the Socioeconomic group ST differs from SC significantly (at both intercept and slope terms) but not from OBC.

2. Secondly and lastly, we run the regression (5) *but with the last two equations in (7)*. The decisions are:

$$c_1^{23} > 0, c_2^{23} < 0.$$

This, in terms of the relationship concerned, means that the Socioeconomic group SC differs from OBC significantly (at both intercept and slope coefficients). This is quite expected because ST differs from SC significantly but not from OBC. So, we can say that SC should differ significantly from both ST and OBC.

Thus CCC between the two regressions of all pairs of regressions out of these three regressions are completed.

Now, in order to get the clusters in the present context, following the discussions above, we first construct the IM defined by (6) which comes out to be as follows:

$$i_{3 \times 3} = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \quad \dots(8)$$

Hence, representing the Socioeconomic groups ST, SC, and OBC as S_1 , S_2 and S_3 respectively, we get here two clusters as follows:

$$C_1 = \{S_1, S_3\} \text{ and} \\ C_2 = \{S_2\}.$$

So, it may be concluded that in respect of the relation of MPCE with proportion of MPCE on food, the behaviour of the SC group do not match with that of the other two groups while the other two groups behave in the same manner.

5. CONCLUSIONS

The above test procedure enables one to perform CCC for the two regressions in each and every pair of regressions possible out of three given regressions. Needless to say, this generalises the Chow test in two directions. This, further, has the following important implications.

Suppose the regressions are now arranged successively in a definite order for investigating existence of structural change. Then, obviously, if this procedure is carried out, CCC between every two successive regressions out of the given regressions, arranged now successively, is automatically done and hence detailed informations regarding all structural changes are obtained, if there is any such at all⁶.

Again, once this procedure is performed, *i.e.*, CCC between the two regressions of all possible pairs of regressions out of the three given regressions are done, one can partition the set of the given regressions into, say, clusters, *every cluster consisting of those regressions which satisfy the condition that, as already stated, the vectors of coefficients of any two of these regressions do not differ from each other significantly*. An example has already been illustrated.

It is to be noted that no. of clusters need not always be two; it may be more than two as well as be one when and only when all the regressions coincide, *i.e.*, the Indicator Matrix becomes a null matrix of an appropriate order.

It may be noted that the clustering introduced here is quite different from that which is carried out by using Mahalanobis D² Statistic^[4] or using k-means Method^[5]—clustering by means of each of the later approaches is based on a single variable/a vector of variables while that by means of our procedure is based entirely on *relationship of variables*.

It is a hunch that the procedure introduced here can be extended for the purpose of clustering on the basis of not only one relationship of variables but more than one relationship together!

REFERENCES

- ^[1]Chow, Gregory C. (1960), Tests of Equality between Sets of Coefficients in Two Linear Regressions, *Econometrica*, Vol. 28, No. 3, pp. 591-605.

- ^[2]Saha, J. and Pal, M. (2014), A Modified Chow Test Approach Towards Testing Differences in The Engel Elasticities, *Asian-African Journal of Economics and Econometrics*, Vol. 14, No. 1, 2014; 57-67.
- ^[3]National Sample Survey Organisation (2007), Household Consumer Expenditure Among Socio-Economic Groups: 2004-2005, NSS 61st Round (July 2004 - June 2005), *NSSO Report No. 514(61/1.0/7)*, Government of India.
- ^[4]Mahalanobis P.C. (1936), On the generalised distance in statistics, *Proceedings of the National Institute of Sciences of India* 2(1): 49–55. Retrieved 2012-05-03.
- ^[5]R in Action, Second Edition, Manning publishing house.

FOOTNOTES

- ²By complete comparison between any two parameters a and b we mean to decide whether $a < b$ or $a = b$ or $a > b$. By component wise complete comparison (CCC) between two vectors of parameters of the same size $(a_1 a_2 \dots a_m)$ and $(b_1 b_2 \dots b_m)$ we mean complete comparison between $(a_1$ and $b_1)$, $(a_2$ and $b_2)$, ... and $(a_m$ and $b_m)$. By CCC between/of/ for two regressions with same no. of parameters we mean CCC between the two vectors of parameters of these regressions. In the paper by Saha and Pal, CCC is done between every two successive regressions out of any number of given successive regressions with same no. of parameters.
- ³The general case (the case of any number of linear regressions) is expected to follow soon as a separate article.
- ⁴Of course, IM provides limited knowledge that for every two regressions, whether they coincide or not, nothing more.
- ⁵The no. of States/UTs, as in the Report, is 35. But there are 5 for each of which some figure/s is/are missing and hence we have excluded those. Also, it may be noted that the numbers of observations for the regressions need not be the same.
- ⁶It is to be noted that if one requires only CCC between every two successive regressions out of given several successive regressions and nothing more, it is not necessary to carry out the procedure described here; for that purpose it is sufficient to work out only the procedure laid for that purpose in the paper by Saha and Paul (2014) referred earlier^[2].