



Munich Personal RePEc Archive

Autocorrelation - Prevalence of identification of collinearity cause

Merce, Emilian and Merce, Cristian Calin and Pocol,
Cristina Bianca

USAMV Cluj-Napoca, USAMV Cluj-Napoca, USAMV Cluj-Napoca

16 November 2017

Online at <https://mpra.ub.uni-muenchen.de/85090/>
MPRA Paper No. 85090, posted 11 Mar 2018 22:24 UTC

AUTOCORRELATION - PREVALENCE OF IDENTIFICATION OF COLLINEARITY CAUSE-

Merce E., C. C. Merce, Cristina B. Pocol*
*Corresponding author

Abstract: *The paper demonstrates that autocorrelation is an accidental statistical phenomenon, whose origin is the incomplete data base. It also shows that the attempts to redistribute factors interactions have focused on the development of methods of solving the effect rather than identifying the cause that generates collinearity. Three possible methods for collinearity removal are analysed comparatively. The premise for two of these methods is autocorrelation redistribution, and the third reveals the cause of collinearity and, implicitly, its cancellation. The three methods are named as follows:*

1. Classic method [1,7];
2. Method of Merce E., Merce C.C.[6];
3. Method of Merce E., Merce C.C.[5];

It is demonstrated that the first two methods are conventional approximations on the distribution of factors' interaction, with possible subjective consequences.

The ideal solution is the use of a complete data base. If this is not possible, as is often the case with databases of economic or sociological research, solving can be the completion of information with theoretical values, obtained by adjusting the causal relationship, in the hypothesis of a certain regression model, a procedure that represents, in fact and implicitly, a way of redistributing the interaction on the influence factors included in the causal model.

INTRODUCTION

Collinearity is an objective reality in the research of complex causal relationships, which externalises, as it will be demonstrated, whenever information about the causal complex is incomplete. The presence of collinearity alters the accuracy of numerical determinations between factors, on the one hand, and the studied effect, on the other. The phenomenon of collinearity cannot, however, always be avoided. This is primarily about research in economics, sociology, psychology. Therefore, it seems natural to evaluate the collinearity and then correct the determination relationship between factors and effect. For this purpose, methods of individualizing the influence of each factor have been outlined, respectively by calculating the partial correlation coefficients [1,6,7]. It will be emphasized that such attempts, although rigorous from a methodological point of view, are working conventions and that neither of these methods leads to the actual numerical determination ratios between factors and effects, ratios which can only be obtained in the case in which there are complete information on the causal complex. If the specificity of the researches necessary requires the use of an incomplete data base, then, in our opinion, the way of redistribution of the factors' interaction must be solved through the integration of the data base.

MATERIAL AND METHOD

We appreciate that, in the construction of methods of separating factors' influence, in cases of incomplete information, a principle mistake was made by which enthusiasm pushes us to combat or adjust the effects and not to explain the causes that produce them. This explains the presence in specialty literature of many methods which, with a higher or lower dose of conventionalism, offer the possibility of deciphering collinearity and collinearity redistribution by factors. All these methods, however, fall under the scope of conventional or, even more severely, of approximation of research results in violation of scientific rigor. That is, in the search for causality, according to these methods, the cause of collinearity was not identified. What is, therefore, the cause that generates collinearity (interdependence) among factors?

Studies, observations and concrete processing are the grounds that lead us to the conclusion that the source of collinearity is the incomplete information on the way of the exteriorization of the effect under the influence of the investigated factors. In such a case, the effect

of collinearity, respectively autocorrelation does not occur if all the states of the resultant variable (y_{ij}) are known for all possible combinations of states comprised of the factorial variables (x_{1j}, x_{2j}). Any deviation from this imperative generates collinearity. Compliance with this requirement means complete experimental plans, including all possible combinations of predetermined factors variants.

In the case of socio-economic phenomena, where the experiment is often impossible, the only alternative is to fill in the information with data adjusted in the hypothesis of a certain regression model, based on incomplete data in the experiment.

And in the case of agricultural experiments, it happens often to encounter situations that only contain some of the possible combinations of influence factors variants. In this regard, it was assumed the following experimental plan for corn crops, which is aimed at the evolution of average production according to NP doses (Table 1).

Table 1 The evolution of average corn production according to NP doses (conventional data)

Dose	Kg/ha	Dose	Kg/ha	Dose	Kg/ha	Dose	Kg/ha
N ₀ P ₀	4600	N ₅₀ P ₈₀	5865	N ₁₀₀ P ₁₂₀	7725	N ₁₅₀ P ₁₆₀	7920
N ₀ P ₄₀	4945	N ₁₀₀ P ₄₀	6095	N ₁₅₀ P ₈₀	7820	N ₂₀₀ P ₁₂₀	8050
N ₅₀ P ₄₀	5980	N ₁₀₀ P ₈₀	7590	N ₁₅₀ P ₁₂₀	7935	N ₂₀₀ P ₁₆₀	7915

RESULTS AND DISCUSSIONS

The picture of the possible combinations, respectively the corresponding average production, is shown in Table 2.

Table 2 The range of possible combinations of the five variants of each factor

X ₁ \ X ₂	0	50	100	150	200
0	4600	?	?	?	?
40	4945	5980	6095	?	?
80	?	5865	7590	7820	?
120	?	?	7725	7935	8050
160	?	?	?	7920	7915

This is a typical example of incomplete information, which generates collinearity and all shortcomings related to redistribution. Correspondences between the levels of factors allocated and the average production obtained for data processing are presented in Table 3.

Table 3

X ₁	X ₂	Y	X ₁	X ₂	Y	X ₁	X ₂	Y
0	0	4600	100	40	6095	150	120	7935
0	40	4945	100	80	7590	150	160	7920
50	40	5980	100	120	7725	200	120	8050
50	80	5865	150	80	7820	200	160	7915

In the case of the first two methods, from those mentioned, for autocorrelation redistribution, it is necessary to determine the correlation coefficients in the hypothesis of a certain theoretical regression model. To express the causal relation between the two factors and the average production, a linear bifactorial model was used. The bifactorial model is, at the same time, the starting point for calculating the adjusted values for completing the baseline data for the third method. Based on the hypothesis that the link could be expressed by a bifactorial, respectively a mono-factorial linear model, by processing the database, the following concrete forms of the models mentioned were obtained:

$$\bar{y}(x_1x_2) = 4985,1 + 11,83x_1 + 8,56x_2; R_{y_{x_1x_2}} = 0,934; D_{y_{x_1x_2}} = 87,2 \%$$

$$\bar{y}(x_1) = 5097,1 + 17,02x_1; r_{y_{x_1}} = 0,914; ;$$

The calculation relations, respectively the calculations made according to the judgments presented in Figure 1, are as follows:

a. The general case:

The coefficient of partial correlation represents the square root of the average of determinations average explained step by step (iterative) in the context of a certain causal complex, in all possible successions, calculating according to the relationship:

$$r_{01 \bullet 23 \dots n} = \sqrt{\frac{A}{n}}$$

$$A = R_{01}^2 + \frac{(R_{012}^2 - R_{02}^2) + \dots + (R_{01n}^2 - R_{0n}^2)}{C_{n-1}^1} +$$

$$+ \frac{(R_{0123}^2 - R_{023}^2) + \dots + [R_{01(n-1)n}^2 - R_{0(n-1)n}^2]}{C_{n-1}^2} +$$

$$+ \frac{(R_{01234}^2 - R_{0234}^2) + \dots + [R_{01(n-2)(n-1)n}^2 - R_{0(n-2)(n-1)n}^2]}{C_{n-1}^3} + \dots$$

$$\dots + \frac{[R_{012 \dots (n-3)}^2 - R_{02 \dots (n-3)}^2] + \dots + [R_{015 \dots n}^2 - R_{05 \dots n}^2]}{C_{n-1}^{n-4}} +$$

$$+ \frac{[R_{012 \dots (n-2)}^2 - R_{02 \dots (n-2)}^2] + \dots + [R_{014 \dots n}^2 - R_{04 \dots n}^2]}{C_{n-1}^{n-3}} +$$

$$+ \frac{[R_{012 \dots (n-1)}^2 - R_{02 \dots (n-1)}^2] + \dots + [R_{013 \dots n}^2 - R_{03 \dots n}^2]}{C_{n-1}^{n-2}} +$$

$$+ [R_{012 \dots n}^2 - R_{02 \dots n}^2]$$

b. The three-factors case:

$$r_{01 \bullet 23} = \sqrt{\frac{R_{01}^2 + \frac{(R_{012}^2 - R_{01}^2) + (R_{013}^2 - R_{03}^2)}{2} + (R_{0123}^2 - R_{023}^2)}{3}}$$

$$r_{02 \bullet 13} = \sqrt{\frac{R_{02}^2 + \frac{(R_{012}^2 - R_{01}^2) + (R_{023}^2 - R_{03}^2)}{2} + (R_{0123}^2 - R_{013}^2)}{3}}$$

$$r_{03 \bullet 12} = \sqrt{\frac{R_{03}^2 + \frac{(R_{013}^2 - R_{01}^2) + (R_{023}^2 - R_{02}^2)}{2} + (R_{0123}^2 - R_{012}^2)}{3}}$$

c. The two-factors case and related processing:

$$r_{y_{x_1} \bullet x_2} = \sqrt{\frac{r_{01}^2 + (R_{012}^2 - r_{02}^2)}{2}} = \sqrt{\frac{(0,914)^2 + [(0,934)^2 - (0,862^2)]}{2}} = 0,694$$

$$d_{y_{x_1} \bullet x_2} = (0,694)^2 \cdot 100 = 48,2 \%$$

$$r_{y_{x_2} \cdot x_1} = \sqrt{\frac{r_{02}^2 + (R_{012}^2 - r_{01}^2)}{2}} = \sqrt{\frac{(0,862)^2 + [(0,934)^2 - (0,914)^2]}{2}} = 0,624$$

$$d_{y_{x_2} \cdot x_1} = (0,624) \cdot 100 = 39,00 \%$$

Method 3

As in many other areas, scientists remain stuck in efforts to combat the effects, neglecting the decipherment of causes that produce unwanted effects. This is the case with collinearity. As a result of many applications and statistical processing by authors, there was a suspicion that autocorrelation could be caused by the incomplete data base. Remaining in the field of scientific speculation, it has been shown that interaction distribution could be done by filling in the missing information with the adjusted values of the regression model used. By generating the adjusted values, using the elaborated bifactorial model, the complete database is as shown in Table 4.

Table 4

X ₂ \ X ₁	0	50	100	150	200
0	4600	5487	6078	6670	7261
40	4945	5919	6095	7012	7604
80	5580	6261	7590	7820	7946
120	5922	6514	7725	7935	8050
160	6265	6856	7448	7920	7915

Through data processing, the following concrete forms of the bifactorial model and of the mono-factorial models were obtained:

$$\bar{y}(x_1, x_2) = 4919,3 + 11,70x_1 + 8,59x_2; R_{y_{x_1}, x_2} = 0,955; D_{y_{x_1}, x_2} = 91,20 \%$$

$$\bar{y}(x_1) = 5606,8 + 11,70x_1; r_{y_{x_1}} = 0,823; d_{y_{x_1}} = 67,73 \%$$

$$\bar{y}(x_2) = 6089,1 + 8,59x_2; r_{y_{x_2}} = 0,4845; d_{y_{x_2}} = 23,47 \%$$

$$\bar{x}_1(x_2) = 100 + 0x_2; r_{y_{x_2}} = 0; d_{y_{x_2}} = 0 \%$$

It can be noticed that, for the third method, the interaction does not operate, and the coefficients of the simple correlation are at the same time coefficients of the partial correlation, respectively reflecting the pure influence of each factor.

Synthetically, the aggregate influence and the separate influences of the two factors for the three methods are presented in Table 5.

Table 5 The comparative situation of total determination and by factors (%)

Factor's influence	Incomplete data base		Complete data base
	Method 1	Method 2	Method 3
X ₁	46,2	48,2	67,73
X ₂	41,0	39,0	23,47
X ₁ , X ₂	87,2	87,2	91,20

For all three methods, the assignment of the total determination by factors is complete, but not unique. Moreover, the total determination is the same for the first two methods, but different for the third.

Method three confirms the truth that autocorrelation is generated by incomplete data bases, but, even in this case, it is assumed that total determination and true partial determinations can only be obtained in the case of the complete data base, obtained through the experimental plan.

BIBLIOGRAPHY

1. Merce E. – Statistică – aplicații practice, Babeș-Bolyai University, Faculty of Economics, Cluj-Napoca, 1986, p.92
2. Merce E. – Cu privire la calculul coeficientului corelației parțiale (I). Definiere, conținut. Studia 1/1989, Babeș-Bolyai University, Cluj-Napoca, pag. 51-60.
3. Merce E., V. Pârv – O nouă metodă de determinare a coeficientului corelației parțiale. Revista de Statistică nr. 3/1991, pag. 32-39.
4. Merce E., B. Pârv, Flavia Laun - Calculul coeficienților determinației parțiale. Metodă și program. Romanian Journal of Statistics no. 6, 1992, pag. 29-37.
5. Merce E., C. C. Merce (2003) – Eliminarea autocorelației dintre variabilele cauzale prin întregirea bazei de date, in the volume Specializare, dezvoltare și integrare; ISBN; 973-86547-4-2; pag,275-281, Babeș-Bolyai University, Faculty of Economics, Cluj-Napoca
6. Merce E., C. C. Merce – Statistică – paradigme consacrate și paradigme întregitoare, AcademicPres Publishing House, Cluj-Napoca, 2009, p.311;
7. Moineagu C. - Modelarea corelațiilor în economie, Scientific Publishing House, Bucharest, 1974.