

Iwona Konarzewska*, Włodzisław Milo**

ON SOME CONSEQUENCES OF THE LACK OF INDEPENDENCE
BETWEEN NOISE COMPONENT AND EXPLANATORY VARIABLES

IN LINEAR REGRESSION MODEL

1. Introduction

In this work we consider a following linear model from which we obtain a system of equations:

$$\begin{aligned} Y_{MRL} &= (R^k, S, S_Y, S_{\underline{X}}, S_Z, Y = \beta' \underline{X} + E, P_Z = \mathcal{N}_E(0, \sigma_E^2)) \\ P_{\underline{X}} &= \mathcal{N}_{\underline{X}}(\mu, \Sigma) \end{aligned}$$

where

$$\beta \in R^k, \mu = (\mu_1, \dots, \mu_k)' \in R^k, \Sigma = (\sigma_{ij})_{i,j=1}^k =$$

$$= \text{cov}(X_i, X_j)_{i,j=1}^k \in R^{k \times k}, \sigma_E^2 \in R_+$$

$S = (U, \mathcal{F}, P)$ - elementary events space,

\mathcal{F} - σ -field of subsets of U

P - complete probability measure defined on \mathcal{F} ,

$$Y : (U, \mathcal{F}, P) \rightarrow (R, \mathcal{F}_R, P_Y) = S_Y$$

$$E : (U, \mathcal{F}, P) \rightarrow (R, \mathcal{F}_R, P_E) = S_E$$

$$\underline{X} : (U, \mathcal{F}, P) \rightarrow (R^k, \mathcal{F}_{R^k}, P_{\underline{X}}) = S_{\underline{X}}$$

* Senior Assistant, Institute of Econometrics and Statistics,
University of Łódź.

** Lecturer, Institute of Econometrics and Statistics, University of Łódź.

While analyzing properties of parameter vector β estimators one traditionally accepts an assumption about distributional independence between random variables X_i , $i = 1, k$ and random noise component E . This is, assuming normal distributions of considered variables, equivalent to the assumption that $\text{cov}(X, E) = 0$. We will check the influence of avoiding this assumption on an explanation level of the model measured by the square of correlation coefficient between the variable Y and linear combination of explanatory variables $\beta'X$. The assumption of independence between explanatory variables and a random noise component cannot be accepted if, for example due to multicollinearity, one rejects from the model an important explanatory variable (in such a case this rejected variable acts as a part of random noise component E). In most cases this variable is not independent from the rest explanatory variables - and consequently E is not independent from them either (and often $\text{E}(E) \neq 0$). The existence of such conditions we will call "error in specification of the set of explanatory variables". We will compare explanation levels of models for those well or badly specified (in the above sense).

2. Theoretical explanation coefficient of the model NMRL

We define the theoretical explanation coefficient $\rho^2 = \rho^2(Y, \beta'X)$ of the model NMRL as follows:

$$(1) \quad \rho^2 = \frac{\text{cov}(Y, \beta'X)^2}{\text{var}(Y)\text{var}(\beta'X)}$$

It is equal to the square of multiple correlation coefficient ρ between random variable Y and explanatory variables of the model NMRL. Availing definition of variance and the model NMRL we obtain:

$$\text{var}(\beta'X) = \beta'X\beta,$$

$$\begin{aligned} \text{var}(Y) &= \sigma_Y^2 = \text{var}(\beta'X) + \sigma_E^2 + 2 \text{cov}(\beta'X, E) \\ &= \beta'X\beta + \sigma_E^2 + 2\beta' \text{cov}(X, E). \end{aligned}$$

Because

$$\begin{aligned}\text{cov}(Y, \beta' \underline{X}) &= \mathbb{E}(Y - \mathbb{E}(Y)) (\beta' \underline{X} - \beta' \mu)' = \\ &= \mathbb{E}(\beta' \underline{X} + \Sigma - \beta' \mu - \mu_{\Sigma}) (\underline{X} - \mu)' \beta = \\ &= \mathbb{E}[\beta' (\underline{X} - \mu) + (\Sigma - \mu_{\Sigma})] (\underline{X} - \mu)' \beta = \\ &= \beta' \Sigma \beta + [\text{cov}(\underline{X}, \Sigma)]' \beta\end{aligned}$$

then (1) is of the form

$$(2) \quad \rho^2 = \frac{(\beta' \Sigma \beta + [\text{cov}(\underline{X}, \Sigma)]' \beta)^2}{(\beta' \Sigma \beta + \sigma_{\Sigma}^2 + 2\beta' \text{cov}(\underline{X}, \Sigma)) \beta' \Sigma \beta}$$

The formula (2) can be introduced also in the form

$$(3) \quad \rho^2 = 1 - \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} + \frac{\beta' \text{cov}(\underline{X}, \Sigma) \text{cov}(\underline{X}', \Sigma) \beta}{\sigma_Y^2 \beta' \Sigma \beta}$$

In the case when $\text{cov}(\underline{X}, \Sigma) : 0 = 0$, ρ^2 takes the form

$$(4) \quad \rho_{(0)}^2 = \frac{\beta' \Sigma \beta}{\beta' \Sigma \beta + \sigma_{\Sigma}^2}, \quad \rho_{(0)}^2 = \rho \text{cov}(\underline{X}, \Sigma) = 0$$

or

$$(4') \quad \rho_{(0)}^2 = 1 - \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} \quad \sigma_{Y(0)}^2 = \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} \text{cov}(\underline{X}, \Sigma) = 0$$

The difference between ρ^2 and $\rho_{(0)}^2$ is expressed by a value

$$\begin{aligned}\rho^2 - \rho_{(0)}^2 &= \frac{\beta' \text{cov}(\underline{X}, \Sigma) \beta}{\sigma_Y^2 \beta' \Sigma \beta} - \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} + \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} \\ &= \frac{\beta' \text{cov}(\underline{X}, \Sigma) \beta}{(\beta' \Sigma \beta + \sigma_{\Sigma}^2 + 2\beta' \text{cov}(\underline{X}, \Sigma)) \beta' \Sigma \beta} + \frac{\sigma_{\Sigma}^2}{\sigma_Y^2} \frac{2\beta' \text{cov}(\underline{X}, \Sigma) \beta}{(\sigma_Y^2 - 2\beta' \text{cov}(\underline{X}, \Sigma)) \beta' \Sigma \beta}.\end{aligned}$$

Denoting $\beta' \text{cov}(\underline{X}, \Sigma) \beta = a$, $\beta' \Sigma \beta = \gamma$ we can rewrite the above difference in a following way

$$\rho^2 - \rho_{(0)}^2 = \frac{\alpha^2}{\gamma(\gamma + \sigma_{\Sigma}^2 + 2\alpha)} + \frac{2\sigma_{\Sigma}^2 \alpha}{(\gamma + \sigma_{\Sigma}^2 + 2\alpha)(\gamma + \sigma_{\Sigma}^2)},$$

$$\rho^2 - \rho_{(0)}^2 = \frac{\alpha^2(\gamma + \sigma_{\Sigma}^2) + 2\sigma_{\Sigma}^2 \alpha \gamma}{\gamma(\gamma + \sigma_{\Sigma}^2)(\gamma + \sigma_{\Sigma}^2 + 2\alpha)}.$$

Because of the fact that factors γ , $\gamma + \sigma_{\Sigma}^2$, $\gamma + \sigma_{\Sigma}^2 + 2\alpha$ denote variances of variables $\beta' \underline{X}$, $\text{cov}(\underline{X}, \Sigma) = 0$, γ - the denominator of the above fraction is nonnegative (to secure this condition it

is necessary that $\alpha > \frac{-\gamma - \sigma_{\Sigma}^2}{2}$). We assume that this denominator is positive. Now we check for which values of α the considered difference takes the nonnegative values.

$$\begin{aligned} \rho^2 - \rho_{(0)}^2 \geq 0 &\iff \alpha^2(\gamma + \sigma_{\Sigma}^2) + \\ &+ 2\sigma_{\Sigma}^2 \alpha \gamma \geq 0 \iff \alpha [\alpha(\gamma + \sigma_{\Sigma}^2) + 2\sigma_{\Sigma}^2 \gamma] \geq 0 \end{aligned}$$

The above inequality occurs for

$$\alpha \in \left\{ (-\infty, \frac{-2\sigma_{\Sigma}^2 \gamma}{\gamma + \sigma_{\Sigma}^2}) \cup (0, +\infty) \right\}.$$

Including condition of σ_{Σ}^2 to be positive, that is $\alpha > \frac{-\gamma - \sigma_{\Sigma}^2}{2}$

and noticing that $\forall \gamma, \sigma_{\Sigma}^2 : \frac{-\gamma - \sigma_{\Sigma}^2}{2} \leq \frac{-2\sigma_{\Sigma}^2 \gamma}{\gamma + \sigma_{\Sigma}^2}$, we obtain

$$(6) \quad \rho^2 - \rho_{(0)}^2 \geq 0 \iff \alpha \in \left\{ \left(\frac{-(\gamma + \sigma_{\Sigma}^2)}{2}, \frac{-2\sigma_{\Sigma}^2 \gamma}{\gamma + \sigma_{\Sigma}^2} \right) \cup (0, +\infty) \right\}$$

The equality $\rho^2 = \rho_{(0)}^2$ occurs in the case $\alpha = \frac{-2\sigma_{\Sigma}^2 \gamma}{\gamma + \sigma_{\Sigma}^2}$ or $\alpha = 0$

(which, for $\beta \neq 0$, means that $c = 0$). Among models where identical linear dependencies exist (that is $\mathcal{D}(\underline{X})$ for these models

is constant and fixed and σ_{Σ}^2 is constant and fixed) the model in which $c \neq 0$ has greater value of theoretical explanation level (in the sense $\rho^2 - \rho_{(0)}^2 \geq 0$), except the case when $\beta'c \in (\alpha^*, 0)$

where $\alpha^* = \frac{-2\sigma_{\Sigma}^2\gamma}{\gamma + 6_{\Sigma}^2}$. The following example illustrates this conclusion.

Example. Let $\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$, $\beta' = [1 \ -2]$, $\sigma_{\Sigma}^2 = 1$. From this we obtain: $\gamma = \beta'\Sigma\beta = 1.4$,

$$\sigma_Y^2 = \gamma + \sigma_{\Sigma}^2 = 2.4,$$

$$\rho_{(0)}^2 = 1 - \frac{1}{2.4} \approx 0.5833.$$

1. Let $c^{(1)} = [0.5 \ 0.2]$. Then $\alpha^{(1)} = \beta'c^{(1)} = 0.1$,

$$\sigma_Y^{2(1)} = \gamma + \sigma_{\Sigma}^2 + 2\alpha^{(1)} = 2.6,$$

$$\rho^{2(1)} = 1 - \frac{1}{2.6} + \frac{0.01}{2.6 \cdot 1.4} \approx 0.6177.$$

Consequently $\rho^{2(1)} > \rho_{(0)}^2$.

2. Let $c^{(2)} = [-0.2 \ 0.5]$. Then $\alpha^{(2)} = -1.2$,

$$\sigma_Y^{2(2)} = 1.4 + 1 - 2.4 = 0 \implies \rho^{2(2)} \text{ does not exist.}$$

3. Let $c^{(3)} = [0.2 \ 0.5]$. Then $\alpha^{(3)} = -0.8$,

$$\sigma_Y^{2(3)} = 1.4 + 1 - 1.6 = 0.8,$$

$$\rho^{2(3)} = 1 - \frac{1}{0.8} + \frac{0.64}{0.8 \cdot 1.4} \approx 1 - 1.25 + 0.57 = 0.32.$$

Consequently $\rho^{2(3)} < \rho_{(0)}^2$.

Assuming that $\frac{\sigma_E^2}{\sigma_Y^2} < 1$, the measure $\rho_{(0)}^2$ is normalized in the range $(0, 1)$. Now, we find the adequate assumptions for ρ^2 when $\text{cov}(X, E) \neq 0$. We check¹ whether there exist such vectors c for which $\rho^2 < 0$ or $\rho^2 > 1$.

At first we assume $\alpha = \beta'c > \frac{-(\gamma + \frac{\sigma_E^2}{\sigma_X^2})}{2}$ (necessary condition for positivity of σ_Y^2).

$$\begin{aligned}\rho^2 < 0 &\iff 1 - \frac{\frac{\sigma_E^2}{\sigma_X^2}}{\frac{\sigma_E^2}{\sigma_X^2} + \gamma + 2\alpha} + \frac{\alpha^2}{(\frac{\sigma_E^2}{\sigma_X^2} + \gamma + 2\alpha)\gamma} < 0 \iff \\ &\iff \frac{\gamma(\frac{\sigma_E^2}{\sigma_X^2} + \gamma + 2\alpha) - \frac{\sigma_E^2}{\sigma_X^2}\gamma + \alpha^2}{\gamma(\frac{\sigma_E^2}{\sigma_X^2} + \gamma + 2\alpha)} < 0 \iff \\ &\iff \gamma^2 + 2\gamma\alpha + \alpha^2 < 0 \iff (\gamma + \alpha)^2 < 0.\end{aligned}$$

Therefore $\forall \alpha, \gamma, \beta : \rho^2 > 0$.

$$\begin{aligned}\rho^2 > 1 &\iff \frac{-\frac{\sigma_E^2}{\sigma_X^2}\gamma + \alpha^2}{(\frac{\sigma_E^2}{\sigma_X^2} + \gamma + 2\alpha)\gamma} > 0 \iff \\ &\iff \alpha^2 \geq \frac{\sigma_E^2}{\sigma_X^2}\gamma.\end{aligned}$$

¹ From definition of correlation coefficient ρ it should be that $\rho^2 \leq 1$. This inequality is obvious for $\rho_{(0)}^2$ (see the relation (4)). However, due to relation (2) and its additive complex form it is not obvious why it should be that $\rho^2 \leq 1$. From further findings (formulas, relations and examples) it is seen that there are points of the parameter space of the MURL model in which the relation $\rho^2 \leq 1$ does not hold. Selecting these points is very useful in planning such experiments in which the inequality $\rho^2 \leq 1$ holds (sensible experiments which are coherent with statistical and econometric theory) and for pre-testing purposes.

Hence for $\alpha \in \{(-\infty, -\sqrt{\frac{6^2}{E}\gamma}) \cup (\sqrt{\frac{6^2}{E}\gamma}, +\infty) \} \wedge (\alpha > \frac{-(\gamma + \frac{6^2}{E})}{2}) \}$
 the value of ρ^2 exceeds 1.

Summing up, $\rho^2 \in (0, 1)$ for $\beta'c = \alpha \in \{(-\sqrt{\frac{6^2}{E}\gamma}, \sqrt{\frac{6^2}{E}\gamma}) \wedge \alpha > \frac{-(\gamma + \frac{6^2}{E})}{2}\}$.
 In the example from this paragraph the normalizing condition was
 fulfilled, except $\alpha^{(2)}$.

3. Consequences of wrong specification of explanatory variables

As we noticed in the introduction, the wrong specification of the set of explanatory variables is the one of the causes of the dependence between X_i , $i = 1, k$ and E . This dependence is an essential difficulty in an estimation process of the model MML parameters. The least squares estimator is then biased and inconsistent.

To investigate the influence of the wrong specification of the set of explanatory variables we consider two models m_1 , m_2 . We assume about m_1 and m_2 that

$$(7) \quad \begin{aligned} m_1: Y &= \beta'_1 X_1 + E_1, \\ m_2: Y &= \beta'_2 X_1 + \beta_{2,k+1} X_{k+1} + E_2, \\ P_{X_1} &= \mathcal{N}_{X_1}(\mu, \Sigma), \quad P_{X_{k+1}} = \mathcal{N}_{X_{k+1}}(\mu_{k+1}, \Sigma_{k+1}), \\ \text{cov}(X_{k+1}, X_1) &= d, \quad P_{E_2} = \mathcal{N}_{E_2}(0, \Sigma_{E_2}), \\ \beta_1, \beta_2 &\in \mathbb{R}^k, \quad \beta_{2,k+1} \in \mathbb{R}^1, \quad \text{cov}(E_2, X_1) = 0, \\ \text{cov}(E_2, X_{k+1}) &= 0. \end{aligned}$$

Let us suppose that

1) m_1 is badly specified because the essential explanatory variable X_{k+1} is not included. We analyze the consequences induced by this fact.

2) μ_2 is well specified (because variable X_{k+1} is included).
 3) X_1, \dots, X_k, X_{k+1} and ε_2 are independently distributed.
 From the definitions of μ_1 and μ_2 we obtain the following properties of the noise component ε_1 of the model μ_1 .

$$(8) \quad \varepsilon_1 = (\beta_2 - \beta_1)' \underline{X}_1 + \beta_{2,k+1} X_{k+1} + \varepsilon_2$$

$$(9) \quad \mathbb{E}(\varepsilon_1) = \mu_{\varepsilon_1} = (\beta_2 - \beta_1)' \mu + \beta_{2,k+1} \mu_{k+1}$$

$$(10) \quad \text{var}(\varepsilon_1) = \sigma_{\varepsilon_1}^2 = (\beta_2 - \beta_1)' \Sigma (\beta_2 - \beta_1) + \beta_{2,k+1}^2 \sigma_{k+1}^2 + \\ + \sigma_{\varepsilon_2}^2 + 2 \beta_{2,k+1} (\beta_2 - \beta_1)' d$$

$$(11) \quad \text{cov}(\underline{X}_1, \varepsilon_1) = c = (\beta_2 - \beta_1)' \Sigma + \beta_{2,k+1} d'$$

It follows that, in general, wrong specification implies dependence (11) between \underline{X}_1 and ε_1 , that is, $c \neq 0$ except the case when $d = 0$ and $\beta_1 = \beta_2$ (X_{k+1} is independent from \underline{X}_1) - variance-covariance matrix Σ is assumed to be nonsingular and positive definite.

4. Wrong specification of explanatory variables and theoretical explanation coefficient of a model

Now, we compare the explanation coefficients of μ_1 and μ_2 .
 From (3) we have:

$$(12) \quad \rho_1^2 = \rho^2(\mu_1) = 1 - \frac{\sigma_{\varepsilon_1}^2}{\sigma_Y^2} + \frac{\beta_1' c c' \beta_1}{\sigma_Y^2 \beta_1' \Sigma \beta_1}.$$

Substituting for $\sigma_{\varepsilon_1}^2$ an expression (10) we obtain

$$\rho_1^2 = 1 - \frac{\sigma_{\varepsilon_2}^2}{\sigma_Y^2} -$$

$$- (\beta_2 - \beta_1)' \ntriangleleft (\beta_2 - \beta_1) + \sigma_{k+1}^2 \beta_{2,k+1}^2 + 2\beta_{2,k+1}(\beta_2 - \beta_1)' d + \sigma_Y^2$$

$$(13) + \frac{\beta_1' c e' \beta_1}{\sigma_Y^2 \beta_1' \ntriangleleft \beta_1}.$$

The explanation coefficient of d_{k2} can be expressed by the following formula

$$(14) \quad \rho_2^2 = 1 - \frac{\sigma_{\Sigma}^2}{\sigma_Y^2}.$$

From (13) and (14) it follows that the condition for $\rho_2^2 > \rho_1^2$ (that is, the explanation level of the well specified model d_{k2} to be greater than adequate value for badly specified model d_k) is a relation

$$(15) \quad (\beta_2 - \beta_1)' \ntriangleleft (\beta_2 - \beta_1) + \sigma_{k+1}^2 \beta_{2,k+1}^2 + 2\beta_{2,k+1}(\beta_2 - \beta_1)' d > \frac{\beta_1' c e' \beta_1}{\beta_1' \ntriangleleft \beta_1}.$$

Because of positive definiteness of matrix ∇ both sides of the inequality (15) can be, without changing the sign, multiplied by $\beta_1' \ntriangleleft \beta_1$. Substituting for c the expression (11) we obtain

$$\begin{aligned} & (\beta_2 - \beta_1)' \nabla (\beta_2 - \beta_1) \beta_1' \ntriangleleft \beta_1 + \beta_{k+1}^2 \beta_{2,k+1}^2 \beta_1' \ntriangleleft \beta_1 + \\ & + 2\beta_{2,k+1}(\beta_2 - \beta_1)' d \beta_1' \ntriangleleft \beta_1 > \\ & > \beta_1' \ntriangleleft (\beta_2 - \beta_1)(\beta_2 - \beta_1)' \ntriangleleft \beta_1 + \beta_{2,k+1}^2 \beta_1' d d' \beta_1 + \\ & + 2\beta_{2,k+1} \beta_1' \ntriangleleft (\beta_2 - \beta_1) d' \beta_1. \end{aligned}$$

After simplification and utilization of the following substitutions

$$\beta'_2 \Delta \beta_2 = \gamma_{22}, \quad \beta'_2 \Delta \beta_1 = \gamma_{21}, \quad \beta'_1 \Delta \beta_1 = \gamma_{11}$$

we obtain a relation

$$(16) \quad \begin{aligned} & \gamma_{22} \cdot \gamma_{11} - \gamma_{21}^2 + \beta_{2,k+1}^2 \beta_1' (\delta_{k+1}^2 \Delta - dd') \beta_1 + \\ & + 2\beta_{2,k+1} (\beta_2' d \gamma_{11} - \gamma_{21} d' \beta_1) > 0. \end{aligned}$$

Fulfilment of the relation (16) depends on concrete values of vector d elements - that is, on the values of covariances between X_{k+1} and other explanatory variables, and on the degree of dependence between X_{k+1} and Y . Therefore it is not obvious whether some well specified models have always greater coefficients (level) of explanation than some badly specified ones.

5. Concluding remarks

In the paper we formulated a measure of theoretical explanation level in linear regression model as the multiple correlation coefficient of explained variable with the set of explanatory variables. We showed that a value of this measure is closely connected with a lack of existence of dependencies between explanatory variables \underline{X} and noise component E which expresses the model error. We found the necessary conditions for this measure to be normalized in the range $(0, 1)$ when there exist non-zero covariances between \underline{X} and E . In the badly specified model M_1 , a set of explanatory variables is such that considered dependence of the form $\text{cov}(\underline{X}_1, E_1) \neq 0$ exists, except the situation when the variables not included are independent from the rest of explanatory variables. This fact is of great importance in the case when we use mechanical techniques for variable selection (for instance rejection of the variable because of strong multicollinearity although from economical theory we know about its important

influence on explained variable). It can imply that least squares estimates are biased and inconsistent.

It is not, however, obvious in what way good specification implies the level of explanation coefficient of the considered model. It was shown that inclusion of all important variables does not necessarily leads to an increase in the value of explanation coefficient of our linear model.

Our results can be especially useful in designing of simulation experiments - they make possible a serious reduction of experimental parameter space. The authors plan to make similar analysis for empirical models substituting theoretical (population) measures by their estimators.

Iwona Konarzewska, Włodzimierz Mikołaj

O PEWNYCH SKUTKACH BRAKU NIEZALEŻNOŚCI
MIĘDZY SKŁADNIKIEM ZABURZAJĄCYM
A ZMIENNYMI OBJASNIAJĄCYMI W LINIOWYM MODELU REGRESJI

Przy analizie własności estymatorów wektora parametrów β liniowego modelu regresji przyjmuje się, tradycyjnie, założenie o niezależności rozkładów zmiennych objaśniających X_i , $i = 1, k$ i zaburzenia losowego ε (oznacza to, że w przypadku rozkładów normalnych, $\text{cov}(X, \varepsilon) = 0$). W pracy badamy wpływ uchylenia tego założenia na stopień objaśnienia modelu, zdefiniowany jako kwadrat współczynnika korelacji między zmienią objaśnianą Y a liniową kombinacją zmiennych objaśniających $\beta'X$. Rozważamy także przypadek, nazywany błędem specyfikacji zbioru zmiennych objaśniających, polegający na nieuwzględnieniu w tym zbiorze ważnej zmiennej, ze względu na jej wysokie skorelowanie z pozostałymi zmiennymi objaśniającymi. Pokazano także, że nie jest oczywiste, jak "dobra" specyfikacja wpływa na wartość stopnia objaśnienia modelu.

* Lecturer, Institute of Econometrics and Statistics, University of Łódź.