*Anna Budka\*, Wiesław Wagner\*\**

## ANALYSIS OF LINEAR REGRESSION MODEL AT DIVIDED SYSTEM MATRIX

**Abstract.** In this study problems connected with the detection of influential observations are investigated in the linear regression model using the least squares estimation of structural parameters. This issue has been presented in three cuts: the complete model, 1-cut model and *m*-cut model.

**Key words**: linear regression model, influential observations, hat matrix, complete, 1-cut and *m*-cut model.

### 1. INTRODUCTION

Methods of causality analysis are commonly used in research. Since the moment of its formulation until recently the least squares method has been developed in terms of its theoretical foundations by the so-called diagnostic support. These are advanced statistical methods making it possible to determine which units have a significantly large effect on the quality of estimated parameters of the linear regression model. In consequence this leads to the determination of the so-called influential observations. Frequently they take the form of outliers or leverage observations.

Influential observations may occur separately or in groups. While in the former case numerous methods have been developed to detect them, statistical methods are still being developed for the detection of many such observations occurring simultaneously. Separation of such observations one by one, even at the applied step procedure, does not necessarily provide

---
\* Ph.D., Agricultural University of Poznań.
\*\* Professor, University of Information Technology and Management in Rzeszów.

correct solutions. This is so due to the fact that in the set of observable explanatory variables the so-called masking effect may occur, i.e. single observations may individually be influential observations, although their cluster may not confirm it.

In this study problems connected with the detection of influential observations are investigated in the linear regression model using the least squares estimation of structural parameters. This issue has been presented in three cuts: the complete model, 1-cut model and $m$-cut model. Detailed methods to investigate influential observations are presented in each case. For this purpose the primary statistics are diagonal elements of the so-called orthogonal projection matrices. Their high values, while all belonging to the (0, 1) interval, and exceeding set threshold values, make it possible to indicate the occurrence of influential observations. Obviously, various possible statistics being to some extent functions of elements of the above mentioned matrix will provide diagnostic information of varying importance concerning influential observations.

## 2. DETERMINATION OF THE LINEAR REGRESSION MODEL

It is assumed that the investigations covered "causality" $x_1, x_2, ..., x_{p-1} \rightarrow y$ of set $p-1$ of established causes defined by the system $x_1, x_2, ..., x_{p-1}$ "of independent variables" (explanatory, clarifying, regression variables) and effect $y$ being a "dependent variable" (explained or clarified variable, predictor). Moreover, it is assumed that the set effect, apart from the mentioned controlled causes, is affected also by random causes expressed by variable $e$. In relation to it stochastic assumptions are made:

(E1) $\mathrm{E}(e) = 0$, expected value equals 0,

(E2) $\mathrm{D}(e) = \sigma$, standard deviation equals a certain positive constant $\sigma > 0$,

(E3) $e \sim \mathrm{N}(0, \sigma)$, random variable $e$ has a "normal distribution" (Gaussian distribution) with indicated parameters.

Assumptions (E1) and (E2) concern moments of random variable $e$, whereas assumption (E3) defines its distribution type. This means that random variable $e$ is of the continuous type and its distribution belongs to the class of normal distributions $\mathcal{N}\{(\mu, \sigma): \mu = 0, \sigma > 0\}$. The connection between variable $y$ and variables $x_1, x_2, ..., x_{p-1}$, $e$ is presented by an "additive linear model"

$$y = \beta_0 + \beta_1 x_1 + ... + \beta_{p-1} x_{p-1} + e \tag{1}$$

also called the "multiple rectilinear regression model". Constants $\beta_0, \beta_1, ..., \beta_{p-1}$ are called "structural parameters or regression coefficients", where paameter $\beta_0$ is a free term, whereas $\beta_1, ..., \beta_{p-1}$ are slopes expressing unitary changes of variables $x_1, x_2, ..., x_{p-1}$ on variable $y$.

Variables $y$ and $e$ are treated as certain "random variables" with set "probability distributions", whereas $x_1, x_2, ..., x_{p-1}$ are set "real variables". Random variable $e$ in model (1) is called "random error" (component).

Random variable $y$ due to the set equation (1) and adopted assumptions (E1), (E2), (E3) for random variable $e$ assumes the following parameters and type of distribution:

$$(Y1)\, E(y) = \sum_{j=0}^{p-1} \beta_j x_j, \quad (Y2)\, D^2(y) = \sigma^2, \quad (Y3)\, y \sim N\left(\sum_{j=0}^{p-1} \beta_j x_j, \sigma\right).$$

Thus assumption (Y3) means that random variable $y$ belongs to the class of normal distributions $\mathcal{N}\left\{(\mu, \sigma):\ \mu = \sum_{j=0}^{p-1} \beta_j x_j \in R,\ \sigma > 0\right\}$. For the purpose of estimation of unknown structural parameters $\beta_0, \beta_1, ..., \beta_{p-1}$ and standard deviation $\sigma$ of the random component, statistical testing is conducted on a finite set $n$ of "units" (cases) $J_1, J_2, ..., J_n$. These units constitute a "random sample" selected from a certain "general population" according to a set "sampling pattern". It is assumed that for each unit $J_i$ a $(p-1)$-dimensional "vector of observation" $x_i' = (x_{i1}, x_{i2}, ..., x_{i, p-1})$ is known and $y_i$, $i = 1, 2, ..., n$ on independent variables and the dependent variable. This system of $n$ vectors of observation constitutes a "multidimensional sample" with size $n$. This sample makes it possible to present model (1) in the form of a vector-matrix linear model

$$y = X\beta + e \tag{2}$$

where $X = (x_0', x_1', x_2', ..., x_n')' : n \times p$ – the system matrix at $x_0 : n \times 1$ unit vector, $y : n \times 1$, $\beta = (\beta_0, \beta_1, ..., \beta_{p-1}) : p \times 1$ – vector of structural parameters, $e = (e_1, e_2, ..., e_n) : n \times 1$ – vector of random errors.

From the adopted stochastic assumptions for random errors in model (1), we obtain the following assumptions for the vector of random errors $e$ in model (2):

(WE1) $E(e) = 0$, vector of expected values is equal to the zero vector,

(WE2) $D^2(e) = \sigma^2 I$, variance-covariance matrix is equal to the scalar diagonal matrix, where $\sigma^2 > 0$ is the variance of random errors. This assumption states also that components of random vector $e$ are not correlated, i.e. "covariance" $cov(e_i, e_{i'}) = 0$ for $i \neq i'$; $i, i' = 1, 2, ..., n$,

(WE3) $e \sim N_n(0, \sigma^2 I)$, random vector $e$ has an $n$-dimensional normal distribution with a zero vector of expected value $\mu = 0$ and covariance matrix $\Sigma = \sigma^2 I$, i.e. belongs to the class of normal distributions $\mathcal{N}_n\{(\mu, \Sigma) : \mu = 0, \Sigma = \sigma^2 I\}$.

Assumption (WE3) implies that vector of observation $y$ has also an $n$-dimensional normal distribution belonging to class $\mathcal{N}_n\{(\mu, \Sigma) : \mu = X\beta, \Sigma = \sigma^2 I\}$, which means that $E(y) = X\beta$ and $D^2(y) = \sigma^2 I$.

### 3. SELECTED PROPERTIES OF LINEAR REGRESSION MODEL

For the purposes of investigations of the influence of influential observations on the quality of evaluation of the structural parameter vector, numerous analytical results, connected with a complete linear regression model, are derived. These results may be found in many publications concerning the theory of linear models. Some of them have been presented in new versions using orthogonal projection matrices. Proofs of many of them are available in literature. While presenting results we are using their specification in the set problem groups.

(S1) Estimation using the least squares method:

a) $X'X\beta = X'y$ – the system of standard equations, where $X'X$ – matrix of moments on explanatory variables, $X'y$ – vector of moments of explanatory variables and the explained variable,

b) $\hat{\beta} = (X'X)^{-1}X'y = GX'y = Ay$ – estimator of structural parameters vector, where $G = (X'X)^{-1}$ and $A = GX'$ – associated matrix,

c) $E(\hat{\beta}) = AE(y) = AX\beta = \beta$ – property of unbiasedness,

d) $D^2(\hat{\beta}) = AD^2(y)A' = \sigma^2 AA' = \sigma^2 G$ – a variance-covariance matrix,

e) $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 G)$ – distribution of estimator $\hat{\beta}$.

(S2) Vector of estimated observations:

a) $\hat{y} = X\hat{\beta} = XAy = XGX'y = Hy$ – evaluation of vector $\hat{y}$, where $H = XGX' = X(X'X)^{-1}X'$ – orthogonal projection matrix,

b) $E(\hat{y}) = HE(y) = HX\beta = X\beta$ – expected value of vector $\hat{y}$,

c) $D^2(\hat{y}) = D^2(Hy) = HD^2(y)H' = H(\sigma^2 I)H = \sigma^2 HH = \sigma^2 H$ – variance-covariance matrix.

(S3) Residuls

a) $r = y - \hat{y} = (I - H)y = My$ – vector of residuals, where $M = I - H$,

b) $E(r) = ME(y) = MX\beta = 0$ – expected value of vector of residuals,

c) $D^2(r) = D(My) = MD(y)M' = \sigma^2MM = \sigma^2M$ – variance-covariance matrix of vector of residuals.

(S4) Sum of squares for SSE error of the least squares method:

a) $SSE = r'r = y'My$ – sum of squares for error,

b) $SSE = y'y - y'Hy = \dfrac{\begin{vmatrix} y'y & y'X \\ X'y & X'X \end{vmatrix}}{|X'X|}$ – SSE expression as a quotient of two determinants of an augmented matrix of the $(X, y)$ system, and the matrix of the $X$ system, which results from formula $\begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = |A_{22}||A_{11} - A_{12}A_{22}^{-1}A_{21}|$, where $A_{11}$ and $A_{22}$ are square matrices and $A_{22}$ is non-singular,

c) $E(SSE) = E(y'My) = E(y')ME(y) + tr[MD(y)] = \sigma^2tr(M) = (n-p)\sigma^2$,
d) $SSE \sim \chi^2_{n-p}$ – the chi-square distribution of sum of squares for error,
e) $E(SSE) = (n-p)\sigma^2$ and $D^2(SSE) = 2(n-p)\sigma^4$ – moments of random variable SSE resulting directly from the distribution given in point c),
f) $SSY = SSR + SSE$ – the factorization of sum of squares of deviations into two summands: sum of squares for regression and sum of squares for error, where $SSY = y'(I - 11'/n)y$ and $SSR = y'(H - 11'/n)y$.

(S5) Estimation of parameter $\sigma^2$:

(a) $s^2 = \dfrac{SSE}{n-p} = \dfrac{y'My}{r(M)}$ – estimator of parameter $\sigma^2$,
(b) $E(s^2) = \sigma^2$ – property of unbiasedness of estimator $s^2$,
(c) $D^2(s^2) = D^2\left(\dfrac{SSE}{n-p}\right) = \dfrac{2(n-p)\sigma^2}{(n-p)^2} = \dfrac{2\sigma^2}{n-p}$ – variance of estimator $s^2$.

(S6) Augmented matrix of orthogonal projection

a) $H^* = Z(Z'Z)^{-1}Z'$, $H^*: n \times n$, $Z = (X, y): n \times (p+1)$,
b) $H^*$ possesses all properties of matrix $H$,
c) $H^* = H + \dfrac{rr'}{r'r} = H + \dfrac{rr'}{SSE}$,
d) $h^*_{ij} = h_{ij} + \dfrac{r_ir_j}{SSE}$, $i, j = 1, 2, ..., n$,
e) $h^*_{ii} = h_{ii} + \dfrac{r_i^2}{SSE}$,

while properties c) and d) are given after J. B. G r a y and R. F. L i n g (1984).

## 4. THE CONCEPT OF INFLUENTIAL OBSERVATIONS

For the purpose of formal determination of influential observations let us introduce the notation of a sample of $n(p+1)$-dimensional observations as a sequence of row vectors of matrix $X$ and of vector $y$ in the form of $(X, y) = ((x_1', y_1),\ (x_2', y_2), ..., (x_n', y_n))' = ((x_i', y_i);\ i = 1, 2, ..., n)' = P_n^{p+1}$, while $(x_i', y_i) \in R^{p+1}$. It is necessary to use the transposition sign, as vectors are always treated as column vectors.

The concept of influential observations is given in the following definition.

**Definition.** *The system of $m\ (m \geqslant 1)$ vector observations $\{(x_{i1}', y_{i1}), ..., (x_{im}', y_{im})\}$ in sample $P_n^{p+1}$ indexed by the discriminanted set of $m$ indexes $\{i_1, i_2, ..., i_m\} \in \{1, 2, ..., n\}$ is called influential observations (points) if they significantly contribute to changes in the values of analyzed numerical characteristics referring to the investigated model of linear regression.*

The above definition indicates that among numerical data there may be one $(m = 1)$ or more observable vectors $(m > 1)$, which will constitute outliers in the direction of the $x$-axis or in the direction of the $y$-axis, or both at the same time. In the first case such influential points may be detected by inspecting diagonal elements of matrix $H$, whereas in the second case – by investigating Studentized residuals. There are numerous solutions to this problem.

While determining influential observations diagonal elements of matrix $H$ are used. These elements are expressed by vector values referring to explanatory variables. This makes it possible to investigate the behaviour of atypicality of these vectors manifested in their distance from the regression cluster. Excessive concentration is understood as a homogenous set corresponding to the regression dependency characteristic. In case of a regression model with one explanatory variable this is equivalent to the configuration of points on a plane arranged along a certain straight line, whereas in case of a multiple regression model such a configuration constitutes a generalized ellipsoid with an intersecting hyperplane. This means that the vector cases will vary in their number in the estimated linear regression model. In connection with the above remarks, we are going to introduce the following definition.

**Definition.** *Observations of cases, i.e. row vector of matrix $X$ corresponding to diagonal values of matrix $H$ are called levearage points.*

Values of diagonal elements of matrix $H$ fall within the interval of $(1/n, 1)$, i.e. they are normalized numbers independent of the number of cases $n$ and the number of properties $p - 1$. Thus, assuming the *a priori* set threshold value, let's say $h_0$ for these diagonal elements, "leverage points" may be distinguished among them, exceeding this value. They are commonly referred to as "high-leverage points". Such points are interesting in terms of their effect on the estimated linear regression model.

## 5. 1-CUT LINEAR REGRESSION MODEL

In case of investigations of regression models, it is interesting to study the dependency between the complete system matrix and its submatrix divided with the use of vectors. Such a division is most frequently connected with the fact that it is necessary to investigate the submatrix distinguished from the matrix of system $X$ in the context of estimating model parameters excluding any of the observed vectors of matrix $X$.

The row division of matrix $X$ is understood as follows. Let us say set $\{1, 2, ..., n\}$ denotes successive numbers of row vectors of matrix $X$. Let us assume that in this set an $i$-th vector was distinguished, which is transposed with the $n$-th vector. Vectors with numbers $i + 1, i + 2, ..., n - 1, n$ will be transposed by one position to loci with numbers $i, i + 1, ..., n - 1$. Such a procedure is called the operation of re-numbering and translocation of row vectors of matrix $X$.

Let us denote with the symbol $X_{(i)}$ a submatrix formed from matrix $X$ without the $i$-th row vector $x_i'$. According to the operation presented above, this means the division of the system matrix into $X = \begin{bmatrix} X_{(i)} \\ x_i' \end{bmatrix}$, where $X : n \times p$, $X_{(i)} : (n - 1) \times p$ is a cut matrix and $x_i' : 1 \times p$ is a distinguished vector, and moreover $X'X = [X_{(i)}'\, x_i] \begin{bmatrix} X_{(i)} \\ x_i' \end{bmatrix} = X_{(i)}'X_{(i)} + x_i x_i'$. To the above division we make a permanent assumption of rank $r(X_{(i)}) = p$ and we introduce the following denotations:

a) $G_{(i)} = (X_{(i)}'X_{(i)})^{-1}$, $G_{(i)} : p \times p$,
b) $H_{(i)} = X_{(i)}G_{(i)}X_{(i)}'$, $H_{(i)} : (n - 1) \times (n - 1)$,
c) $v_{(i)} = X_{(i)}G_{(i)}x_i$, $v_{(i)} : (n - 1) \times 1$,
d) $c_i = x_i'G_{(i)}x_i$,
e) $d_i = (1 + c_i)^{-1}$,
f) $h_{kl(i)} = x_k'G_{(i)}x_l$, $k, l = 1, 2, ..., n$.

In correspondence to the given division of matrix $\mathbf{X}$ with one row vector, the division is conducted of the vector of observable random variables $\mathbf{y} = \begin{bmatrix} \mathbf{y}_{(i)} \\ y_i \end{bmatrix}$, where $\mathbf{y}_{(i)} : (n-1) \times 1$, while $y_i$ is the $i$-th component of vector $\mathbf{y}$.

We will further use the denotation "1-cut" to emphasize that it refers to a linear model investigated for vector $\mathbf{y}$ and the matrix of system $\mathbf{X}$ after elimination of the $i$-th component of vector $\mathbf{y}$ and the $i$-th row vector in matrix $\mathbf{X}$. Let us denote a 1-cut model in the ternary form

$$\{\mathbf{y}_{(i)}, \mathbf{X}_{(i)} \boldsymbol{\beta}_{(i)}, \sigma^2 \mathbf{I}\} \tag{3}$$

where here the identity matrix $\mathbf{I}$ is of the $(n-1)$-th degree. We will add to the previously given characteristics $\hat{\boldsymbol{\beta}}$, $\hat{\mathbf{y}}$, $\mathbf{r}$, $SSE$ their equivalents $\boldsymbol{\beta}_{(i)}$, $\hat{\mathbf{y}}_{(i)}$, $\mathbf{r}_{(i)}$, $SSE_{(i)}$, after the application of the row division of the system matrix $\mathbf{X}$ and vector $\mathbf{y}$.

Stochastic properties for the 1-cut model will be noted according to the same principle as for the complete model, but denotation "-1" is additionally placed for emphasize the fact that the cut model is used.

(S1-1) Least squares estimation:

a) $\hat{\boldsymbol{\beta}}_{(i)} = \mathbf{G}_{(i)} \mathbf{X}'_{(i)} \mathbf{y}_{(i)} = \mathbf{A}_{(i)} \mathbf{y}_{(i)}$ – estimated vector of structural parameters,

b) $\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_{(i)} = \mathbf{G}_{(i)} \mathbf{x}_i r_i = \dfrac{r_i}{m_{ii}} \mathbf{G} \mathbf{x}_i$ – the difference between estimated vectors

of structural parameters of the complete and 1-cut models, where $r_i$ denote residuals,

(S2-1) Forecasting vector $\mathbf{y}$ in the 1-cut model:

a) $\hat{\mathbf{y}}_{(i)} = \mathbf{X}_{(i)} \hat{\boldsymbol{\beta}}_{(i)} = \theta_{(i)}$ – estimation (prediction) of vector $\mathbf{y}$,

b) $\hat{y}_{j(i)} = \mathbf{x}'_j \hat{\boldsymbol{\beta}}_{(i)} = \theta_{j(i)}$ – $j$-th component of vector $\hat{\mathbf{y}}_{(i)}$, i.e. forecast of the $j$-th observation after the elimination of the $i$-th observation,

c) $\hat{y}_j - \hat{y}_{j(i)} = \dfrac{r_i}{m_{ii}} \mathbf{x}'_j \mathbf{G} \mathbf{x}_i$ – the difference of the $j$-th component estimated in the complete and 1-cut models,

d) $\hat{y}_i = h_{ii} y_i + m_{ii} \theta_{(i)}$ – results from (S1-1) b) – linear combination of observations $y_i$ and $\theta_{(i)} = \mathbf{x}'_i \hat{\boldsymbol{\beta}}_{(i)}$ at weights expressed by the $i$-th diagonal element of matrix $\mathbf{X}$,

(S3-1) Residuals and sum squares for error:

a) $r_{(i)} = y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}_{(i)} = y_i - \theta_{(i)} = \dfrac{r_i}{m_i}$ – residuals in the 1-cut model,

b) $SSE_{(i)} = y'_{(i)}(I - H_{(i)})y_{(i)}$ – sum of squares for error in the 1-cut model, where $H_{(i)}$ is the submatrix of matrix $H$,

c) $SSE = SSE_{(i)} + d_i r^2_{(i)}$ – connection between sum of squares for error in the complete and 1-cut models,

d) $s^2_{(i)} = \dfrac{1}{\nu - 1} \sum\limits_{j=1}^{n} (y_j - x'_j \hat{\beta}_{(i)})^2 = \dfrac{SSE_{(i)}}{\nu - 1} = \dfrac{\nu}{\nu - 1} s^2 - \dfrac{r^2_i}{(\nu - 1)m_{j \neq i}}$ – estimator of parameter $\sigma^2$ in the 1-cut model,

e) $D^2(r_{(ii)}) = s_{(i)}[1 + x'_i G_{(i)} x_i]^{1/2} = \dfrac{S_{(i)}}{\sqrt{m_{ii}}}$ – standard error for residuals in the 1-cut model (H a o g l i n, W e l s c h 1978),

f) $\nu s^2 = (\nu - 1)s^2_{(i)} + w_{ii}r^2_{(i)}$ – equation b) expressed by estimators $s^2$ and $s^2_{(i)}$ of parameter $\sigma^2$ of the complete and 1-cut models, where $w_{ii} = \dfrac{h_{ii}}{m_{ii}}$,

g) $s^2_{(i)} \sim \dfrac{\sigma^2}{\nu - 1} \chi_{\nu - 1}$ – chi-square distribution of estimator $s^2_{(i)}$,

h) squares of residuals $r^2_i$ and estimates $s^2_{(i)}$ are statistically independent (L a   M o t t e 1994),

(S4-1) Studentized residuals:

a) $t_{(i)} = \dfrac{r_{(i)}}{D(r_{(i)})} = \dfrac{r_i}{s_{(i)}\sqrt{m_{ii}}}$ – expression of parameter $\sigma^2$ by estimator $s^2_{(i)}$ in the 1-cut model (the so-called "external Studentized residuals of least squares"),

b) $t_{(i)} = t_i \sqrt{\dfrac{\nu - 1}{\nu - t^2_i}}$ – expression by standardized residuals and $t$-Student distribution with $\nu - 1$ degrees of freedom.

## 6. M-CUT LINEAR REGRESSION MODEL

Apart from the investigations of influential observations in the 1-cut model, the problem of the $m$-cut model is also studied. Let us denote submatrix $X_{(I)}$ formed without the subset of $I$ row vectors of matrix $X$ and $X_I$ the system of such discriminated vectors. Matrix $X$ will be denoted by the form: $X = \begin{bmatrix} X_{(I)} \\ X_I \end{bmatrix}$, where $X : n \times p$, $X_{(I)} : (n - m) \times p$ – $m$-cut matrix and $X'_I : m \times p$ – discriminated matrix (discriminated system of $m$ row vectors). The following ratio occurs for it $X'X = [X'_{(I)} X_I]\begin{bmatrix} X_{(I)} \\ X_I \end{bmatrix} = X'_{(I)} X_{(I)} + X_I X'_I$, and

moreover for the given division we adopt a permanent assumption with rank $r(X_{(I)}) = p$ and we introduce the following denotations:

a) $G_{(I)} = (X'_{(I)} X_{(I)})^{-1}$, $G_{(I)} : p \times p$,

b) $H_{(I)} = X_{(I)} G_{(I)} X'_{(I)}$, $H_{(I)} : (n - m) \times (n - m)$,

c) $F_I = X'_I G_{(I)} X_I$, $F_I : m \times m$,

d) $V_{(I)} = X_{(I)} G_{(I)} X_I$, $V_{(I)} : (n - m) \times m$,

e) $E_I = [I + F_I]^{-1}$, $E_I : m \times m$.

The given division of matrix of system $X$ is applied to the vector of observable random variables $y = \begin{bmatrix} y_{(I)} \\ y_I \end{bmatrix}$, where $y_{(I)} : (n - m) \times 1$, while $y_I : m \times 1$ is a column vector containing, without the loss of generality, the last $m$ components of vector $y$. Further we apply the "$m$-cut" denotation to emphasize that we are considering a linear model of vector $y$ and system matrix $X$ after the elimination from the above the indicated subvector $y_I$ and matrix $X_I$. Let us present the $m$-cut model in the ternary form

$$\{y_{(I)}, X_{(I)} \beta_{(I)}, \sigma^2 I\} \tag{4}$$

where here the identity matrix $I$ is a matrix of $(n - m)$-th degree. Let us also stress that symbol $\beta_{(I)}$ does not refer to the cutting of the vector of structural parameters in model (4), but that we apply such a denotation to emphasize that a $p$-dimensional vector of parameters $\beta$ is estimated from the $m$-cut model.

To stochastic characteristics $\hat{\beta}, \hat{y}, r, SSE$ for the complete model, we will give their equivalents $\beta_{(I)}, \hat{y}_{(I)}, r_{(I)}, SSE_{(I)}$ after the application of row division of the system matrix $X$ and vector $y$. Stochastic properties of the $m$-cut model will be denoted according to the same principle as for the 1-cut model, but additionally the denotation "-$m$" will be placed to emphasize the application of the $m$-cut model. Earlier we will give several denotations:

- $X = \begin{bmatrix} X_{(I)} \\ X'_I \end{bmatrix}$, $X : n \times p$, $X_{(I)} : (n - m) \times p$, $X'_I : m \times p$,

- $r(X'_{(I)} X_{(I)}) = p$, $G_{(I)} = (X'_{(I)} X_{(I)})^{-1}$.

(S1-$m$) Estimation of least squares:

a) $X'_{(I)} X_{(I)} \hat{\beta}_{(I)} = X'_{(I)} y_{(I)}$ – system of normal equations,

b) $\hat{\beta}_{(I)} = G_{(I)} X'_{(I)} y_{(I)}$ – estimator of vector of structural parameters $\beta$,

c) $E(\hat{\beta}_{(I)}) = \beta_{(I)}$ and $D^2(\hat{\beta}_{(I)}) = \sigma^2 G_{(I)}$ at the assumptions of model (4),

d) $\hat{\beta}_{(I)} = \hat{\beta} + GX_I(I - H_I)^{-1} r_I$, where matrix $H_I$ is the submatrix from

$$X = \begin{bmatrix} X_{(I)} \\ X_I' \end{bmatrix} G[X_{(I)}' \ X_I] = \begin{bmatrix} X_{(I)} GX_{(I)}' & X_{(I)} GX_I \\ X_I' GX_{(I)}' & X_I' GX_I \end{bmatrix} = \begin{bmatrix} H_{(I)} & H_{I(I)} \\ H_{I(I)} & H_I \end{bmatrix},$$

and $r_I = y_I - X_I'\hat{\beta}$,

e) $\hat{y}_{(I)} = H_{(I)} y_{(I)}$ – estimated vector $y_{(I)}$ from model (4),

f) $X_I' \hat{\beta}_{(I)} = X_I' \hat{\beta} + H_{(I)} r_I$, which results from expression d) and formula

$$G_{(I)} X_I = GX_I + GX_I(I - H_I)^{-1} H_I = GX_I(I - H_I)^{-1}.$$

(S2-$m$) Vectors of residuals $r_{(I)}$ and $r_I$:

a) $r_{(I)} - y_{(I)} - X_{(I)}\hat{\beta} = (I - H_{(I)})y_{(I)} - H_{(I)I} y_I$,

b) $r_I - y_I - X_I'\hat{\beta} = (I - H_I)y_I - H_{I(I)} y_{(I)}$,

c) $E(r_{(I)}) = 0$, $D^2(r_{(I)}) = \sigma^2(I - H_{(I)})$,

$$E(r_{(I)}) = (I - H_{(I)})E(y_{(I)}) - H_{(I)I} E(y_I) = [(I - H_{(I)})X_{(I)} - H_{(I)I} X_I']\beta =$$

$$= [H_{(I)I} X_I' - H_{(I)I} X_I']\beta = 0,$$

$$D^2(r_{(I)}) = (I - H_{(I)})D^2(y_{(I)})(I - H_{(I)}) + H_{(I)I}D^2(y_I)H_{(I)I} =$$

$$= \sigma^2[(I - H_{(I)})(I - H_{(I)}) + H_{(I)I} H_{I(I)}] =$$

$$= \sigma^2[I - H_{(I)} - H_{(I)I}H_{I(I)} + H_{(I)I}H_{I(I)}] = \sigma^2(I - H_{(I)}),$$

d) $E(r_I) = 0$, $D^2(r_I) - = \sigma^2(I - H_I)$, as is shown analogically as in c),

e) $r_{(I)} = [(I - H_{(I)}) - H_{I(I)}(I - H_I)^{-1}H_{I(I)}]y_{(I)}$ – vector of residuals expressed by the vector of observations in the $m$-cut model,

f) $Q_m = r_I'(I - H_I)^{-1}r_I$ – "outlier sum of squares" connected with the $m$-dimensional vector of parameters $\gamma$ in the extended model $\begin{bmatrix} E(y_{(I)}) \\ E(y_I) \end{bmatrix} = \begin{bmatrix} X_{(I)} & 0 \\ X_I & I \end{bmatrix}\begin{bmatrix} \beta \\ \gamma \end{bmatrix}$ in comparison to model $\begin{bmatrix} E(y_{(I)}) \\ E(y_I) \end{bmatrix} = \begin{bmatrix} X_{(I)} \\ X_I \end{bmatrix}\beta$ (Gen-tleman, Wilk 1975; Draper, John 1981),

g) $Q_m = r_I'r_I + r_{(I)}'H_{(I)}r_{(I)}$ – outlier sum of squares split into sum of squares of direct residuals indexed with the vector of observations $y_I$ of residuals from the $m$-cut model (4).

(S3-$m$) Sum of squares for error and estimation of variance:

a) $SSE_{(I)} = y_{(I)}'(I - H_{(I)})y_{(I)} = SSE - r_I'(I - H_I)^{-1}r_I = SSE - Q_m$, where $H_{(I)}$ is the matrix of orthogonal projection for the system matrix in model (4),

matrix $H_I$ was given and vector of residuals $r_I$ is given by formula S2-$m$ b), whereas $Q_m$ was given in S2-$m$ f),

b) $s_{(I)}^2 = \dfrac{I}{v-m}(vs^2 - Q_m)$ – estimation of variance for error expressed by the variance of the complete model $s^2$ corrected by the outlier sum of squares,

c) $F_{(I)} = \dfrac{v-m}{m} \cdot \dfrac{Q_m}{SSE_{(I)}}$ – the $F$ distribution with $v-m$ and $m$ degrees of freedom at stochastic assumptions of model (4).

In connection with these stochastic considerations for the $m$-model of linear regression, let us supply two more important results (C o o k, W e i s - b e r g 1980):

a) $r_I'(I - H_I)^{-1}r_I = \sum\limits_{i \in I} h_{ii}^2$ – generalized Studentized residual for the set of cases $I$, i.e. it is directly expressed by the sum of squares of diagonal elements of matrix $H$, indexed by set $I$,

b) $\mathrm{tr}[H_I(I - H_I)^{-1}]$ – generalized leverage point for set $I$ of discriminated cases.

## 7. CONCLUSIONS

The paper presents theoretical results referring to the complete, 1-cut and $m$-cut models in linear regression. Individual observations or the $m$--system are investigated in the above mentioned cut models. As it was shown in the last two chapters, this investigation may be conducted on units from the complete model. This makes it possible to considerably simplify numerical operations. Various testing statistics are derived from the data determined for the 1-cut and $m$-cut models in order to analyze the occurrence of influential observations. They are used for the purpose of practical identification of influential observations found in the linear regression model. A list of such statistics will be supplied y the authors in another study.

## REFERENCES

C o o k  R. D., W e i s b e r g  S. (1980), *Characterizations of an Empirical Influence Function for Detecting Influential Cases in Regression*, "Technometrics", 22, 495–507.
C o o k  R. D., E i s b e r g  S. (1982), *Residuals and Influence in Regression*, Chapman and Hall, New York.

Draper N. R., John J. A. (1981), *Influential Observations and Outliers in Regression*, "Technometrics", **23**, 21–26.

Gentleman J. F., Wilk M. B. (1975), *Detecting Outliers II: Supplementing the Direct Analysis of Residuals*, "Biometrics", **31**, 387–410.

Gray J. B., Ling R. F. (1984), *K – Clustering as a Detection Tool for Influential Subsets in Regression*, "Technometrics", **26**, 305–318.

Haoglin D. C., Welsch R. E. (1978), *The Hat Matrix in Regression and ANOVA*, "American Statistician", **32**, 17–22.

La Motte L. (1994), *A Note on the Role of Independence in T Statistics Constructed from Linear Statistics in Regression Model*, Amer. Statist., **48**, 238–240.

*Anna Budka, Wiesław Wagner*

## ANALIZA MODELU REGRESJI LINIOWEJ
## PRZY PODZIELONEJ MACIERZY UKŁADU

W pracy przedstawiono zagadnienia związane z wykrywaniem obserwacji wpływowych w modelu regresji liniowej przy zastosowaniu estymacji parametrów strukturalnych za pomocą MNK. Temat ten jest ujęty w trzech przekrojach: model pełny, 1-ucięty oraz model $m$-ucięty. W każdym przypadku prezentowane są szczegółowe metody badania obserwacji wpływowych. Podstawowymi dla tych celów statystykami są elementy diagonalne tzw. macierzy ortogonalnego rzutu. Ich duże wartości, przy czym wszystkie należą do przedziału (0, 1), przekraczające zadane wartości progowe pozwalają na wskazanie istnienia obserwacji wpływowych. Oczywiście różne możliwe statystyki będące w jakimś stopniu funkcjami elementów wspomnianej macierzy będą dostarczały informacji diagnostycznych o różnym znaczeniu, dotyczącym obserwacji wpływowych.