*Janusz Wywiał*[*]

# PERFORMING QUANTILES IN MULTIPLE REGRESSION SAMPLING STRATEGY

**ABSTRACT.** Estimation of the population average in a finite population by means of sampling strategy dependent on the sample quantile of an auxiliary variables is considered. The sampling design is proportionate to the determinant of the matrix dependent on some quantiles of an auxiliary variables. The sampling scheme implementing the sampling design is proposed. The derived inclusion probabilities are applied to estimation the population mean using the well known Horvitz-Thompson estimator. Moreover, the regression estimator is defined as the function of the coefficient dependent on the quantiles of the auxiliary variables. The properties of this estimator under the above defined sampling design are studied. The considerations are supported by empirical examples.

**Key words**: sampling design, order statistic, sample quantile, auxiliary variable, Horvitz-Thompson statistic, inclussion probabilities, sampling scheme, regression estimator.

## I. INTRODUCTION

We are going to consider a finite and fixed population of the size N. A variable under study will be denoted by *y*, an auxiliary variable by x. Let $((y_1,x_1)$, $(y_2,x_2),\ldots,(y_i,x_i),\ldots,(y_n,x_n))$ be the observations of the variable (y, x). The sample means of the variables *y* and *x* will be denoted by $\bar{y}$ and $\bar{x}$, respectively.

Let $s_R = \{i : x_i \le \bar{x}\}$, $s_L = \{i : x_i > \bar{x}\}$, $n_R = Card\{s_R\}$, $n_L = n - s_R$ and

$$\bar{x}_R = \frac{1}{n_R}\sum_{i \in s_R} x_i, \quad \bar{x}_L = \frac{1}{n_L}\sum_{i \in s_L} x_i, \quad \bar{y}_R = \frac{1}{n_R}\sum_{i \in s_R} y_i, \quad \bar{y}_L = \frac{1}{n_L}\sum_{i \in s_L} y_i.$$

[*]Professor, Department of Statistics, Katowice University of Economics, Katowice.

So, the $\bar{x}_R$ and $\bar{x}_L$, statistics are the sample means of the variable $x$ from the right truncated and left truncated sample, respectively, in the point $\bar{x}$.

The estimator of the slope coefficient of the linear regression of the variable $y$ on $x$ considered by Wald (1940), Kendall and Stuart (1961), pp. 399-400 and Hellwig (1956, 1963), pp.138-155 is as follows

$$b_s = \frac{\bar{y}_R - \bar{y}_L}{\bar{x}_R - \bar{x}_L}.$$

In this paper we are going to generalize those results into the multidimensional case when the linear regression function depends on at least two explanatory variables.

## II. BASIC DEFINITIONS AND NOTATIONS

We are going to consider the finite population U of size $N>k$. The observation of a variable under study and an auxiliary variables are denoted by $y_i$ and $x_{i,j}$, respectively $i=1,...,N$, $j=1,...,k$. Let $[y \ x]$ where $x=[x_1 \ x_2 \ ... \ x_n]$ be the vector which values are coordinates of a point in a $(k+1)$ -dimensional Euclidean space. Elements of the vector $y^T=[y_1 \ y_2 \ ... \ y_N]$ be values of a variable under study observed in the population. The vector $y_k$ of dimensions $k \times 1$ consists of k-elements of the vector $y$, $k \le N$. The observed in the population values of k-auxiliary variables are the elements of the matrix $x=[x_{ij}]$ of dimensions $N \times k$ and $i=1,2,...,N$, $j=1,2,...,k$. Moreover, $x=[x_{\bullet 1} \ x_{\bullet 2} \ ... \ x_{\bullet k}]$ where $(x_{\bullet j})^T=[x_{1j} \ x_{2j} \ ... \ x_{kj}]$ and

$$x = \begin{bmatrix} x_{1\bullet} \\ ... \\ x_{N\bullet} \end{bmatrix}$$ where $x_{i\bullet}=[x_{i1} \ x_{i2} \ ... \ x_{ik}]$. Let $x_k$ $(k<N)$ be the submatrix of the matrix

$x$. The matrix $x_k$ is obtained after removing (N-k) rows from the matrix $x$. The column vector consisted of $k$-th values each equal to one will be denoted by $J_k$. The column vector consisted of $k$-th values each equal to zero will be denoted by $0_k$.

The well known equation of a $k$ dimensional hiper-plain spanned on $(k+1)$ points $z_i=[y_i \ x_{i\bullet}]$, $(i=0,1,2,...,k)$ in the $(k+1)$-dimensional Euclidean space is as follows (see. e.g. Borsuk (1969)):

$$\begin{vmatrix} 1 & y & x \\ 1 & y_0 & x_0 \\ J_k & y_k & x_k \end{vmatrix} = 0,$$

If we subtract the second row of the above matrix from the remainder rows we get:

$$\begin{vmatrix} 0 & y - y_0 & x - x_0 \\ 1 & y_0 & x_0 \\ 0_k & d_y(y_0) & d(x_k, x_0) \end{vmatrix} = 0 \tag{1}$$

or

$$\begin{vmatrix} y - y_0 & x - x_0 \\ d_y(y_0) & d(x_k, x_0) \end{vmatrix} = 0, \tag{2}$$

where

$$c(y_0) = y_k - y_0 J_k, \tag{3}$$

$$d(x, x_0) = x_k - J_k x_0. \tag{4}$$

Let $d^{(.,j)} = d^{(.,j)}(x_k, x_0)$ be the matrix obtained through removing a $j$-*th* column of the matrix $d(x_k, x_0)$, $j = 1, \ldots, k$. Similarly, let $d^{(i,.)} = d^{(i,.)}(x_k, x_0)$ be the matrix obtained through removing an $i$-*th* row of the matrix $d(x_k, x_0)$, $i = 1, \ldots, k$. Finally, let $d^{(i,j)} = d^{(i,j)}(x_k, x_0)$ be the matrix obtained through removing an i-th row and a $j$-*th* column of the matrix $d(x_k, x_0)$, $i, j = 1, \ldots, k$. Moreover, let $D(x_k, x_0) = [d^{(i,j)}(x_k, x_0)] = [d^{(i,j)}]$. This notation let us rewrite the determinant (1) or (2) in the following ways.

$$(y - y_0) |d(x_k, x_0)| + \sum_{j=1}^{k} (-1)^j (x_j - x_{0j}) |c(y_0) \, d^{(.,j)}(x_k, x_0)| = 0$$

or

$$(y - y_0) |d(x_k, x_0)| + \sum_{i=1}^{k} (-1)^i (y_i - y_0) |d^{(i,.)}(x_k, x_0)| = 0$$

or

$$(y-y_0)\left|\mathbf{d}(\mathbf{x}_k,\mathbf{x}_0)\right|+\sum_{i=1}^{k}\sum_{j=1}^{k}(-1)^{i+j}\left(\mathbf{x}_j-\mathbf{x}_{0j}\right)\left(y_i-y_0\right)\left|\mathbf{d}^{(i,j)}(\mathbf{x}_k,\mathbf{x}_0)\right|=0$$

or

$$(y-y_0)\left|\mathbf{d}^{-1}(\mathbf{x}_k,\mathbf{x}_0)\right|-(x-\mathbf{x}_0)\mathbf{D}(\mathbf{x}_k,\mathbf{x}_0)\mathbf{c}(y_0)=0$$

or

$$y-y_0-(x-\mathbf{x}_0)\mathbf{d}^{-1}(\mathbf{x}_k,\mathbf{x}_0)\mathbf{c}(y_0)=0.$$

These equations are equivalent and can be rewritten in the following ways:

$$y=y_0-\frac{1}{\left|\mathbf{d}(\mathbf{x}_k,\mathbf{x}_0)\right|}\sum_{j=1}^{k}(-1)^j\left(\mathbf{x}_j-\mathbf{x}_{0j}\right)\left|\mathbf{c}(y_0)\ \mathbf{d}^{(\cdot,j)}(\mathbf{x}_k,\mathbf{x}_0)\right|$$

or

$$y=y_0-\frac{1}{\left|\mathbf{d}(\mathbf{x}_k,\mathbf{x}_0)\right|}\sum_{i=1}^{k}(-1)^i\left(y_i-y_0\right)\left|\mathbf{d}^{(i,\cdot)}(\mathbf{x}_k,\mathbf{x}_0)\right|$$

or

$$y=y_0-\frac{1}{\left|\mathbf{d}(\mathbf{x}_k,\mathbf{x}_0)\right|}\sum_{i=1}^{k}\sum_{j=1}^{k}(-1)^{i+j}\left(\mathbf{x}_j-\mathbf{x}_{0j}\right)\left(y_i-y_0\right)\left|\mathbf{d}^{(i,j)}(\mathbf{x}_k,\mathbf{x}_0)\right|$$

or

$$y=y_0+\frac{1}{\left|\mathbf{d}(\mathbf{x}_k,\mathbf{x}_0)\right|}(x-\mathbf{x}_0)\mathbf{D}(\mathbf{x}_k,\mathbf{x}_0)\mathbf{c}(y_0)$$

or

$$y=y_0+(x-\mathbf{x}_0)\mathbf{d}^{-1}(\mathbf{x}_k,\mathbf{x}_0)\mathbf{c}(y_0). \tag{5}$$

Let U be a fixed population of size $N$. Moreover, let $x_{i,1} < x_{i+1,1}$, $i=1,...,N-1$. Our problem is estimation of the population average $\bar{y}=\sum_{i=1}^{N}y_i/N$. Let $s$ be a simple sample of the fixed size $n>k+1$. The sample mean of the variable under study is defined by $\bar{y}_s=\sum_{i\in s}y_i/n$. The population mean vector of auxiliary variables we denote by $\bar{\mathbf{x}}=\begin{bmatrix}\bar{x}_1 & \bar{x}_2 & ... & \bar{x}_k\end{bmatrix}$ where $\bar{x}_j=\sum_{i=1}^{N}x_{i,j}/N$ and the sample mean vector of auxiliary variables will be denoted by

$\overline{\mathbf{x}}_s = \begin{bmatrix} \overline{x}_{1s} & \overline{x}_{2s} & \dots & \overline{x}_{ks} \end{bmatrix}$ where, $j=1,\dots,k$. Let $\mathbf{x}_n$ ($n<N$) be such a submatrix of the matrix $\mathbf{x}$ that the rows of the matrix $\mathbf{x}_n$ are observations of the auxiliary variables in the sample s of size n.

Let us consider the estimator which we obtain through changing $\overline{y}_s$ for $y_0$ and $(\overline{\mathbf{x}} - \mathbf{x}_0)$ for $(x - \mathbf{x}_0)$ only at the first row of the matrix given by the expression (2) and $\overline{\mathbf{x}}$ for $\mathbf{x}_k$. This leads to the following estimator of the population mean:

$$\tilde{y} = \overline{y}_s + (\overline{\mathbf{x}} - \mathbf{x}_0) \mathbf{d}^{-1}(\mathbf{x}_k, \mathbf{x}_0) \mathbf{c}(y_0) \tag{6}$$

or

$$\tilde{y} = \overline{y}_s + \frac{1}{\left|\mathbf{d}(\mathbf{x}_k, \mathbf{x}_0)\right|}(\overline{\mathbf{x}} - \mathbf{x}_0) \mathbf{D}(\mathbf{x}_k, \mathbf{x}_0) \mathbf{c}(y_0). \tag{7}$$

Particularly, if $k=1$ then

$$\tilde{y}_s = \overline{y}_s + \frac{y_1 - y_0}{x_1 - x_0}(\overline{x} - \overline{x}_s).$$

The next estimator can be constructed in the following way. At the right side of the equation (5) let us change: the mean $\overline{y}_s$ for $y_0$, $\overline{\mathbf{x}}$ for $x$ and $\overline{\mathbf{x}}_s$ for $x_0$. This leads to the following estimator of the population mean:

$$\hat{y} = \overline{y}_s + (\overline{\mathbf{x}} - \overline{\mathbf{x}}_s) \mathbf{d}^{-1}(\mathbf{x}_k, \overline{\mathbf{x}}_s) \mathbf{c}(\overline{y}_s). \tag{8}$$

Let us note that particularly, if $k=1$, then

$$\hat{y}_s = \overline{y}_s + \frac{y_1 - \overline{y}_s}{x_1 - \overline{x}_s}(\overline{x} - \overline{x}_s).$$

Our problem is how to determine the observation of the auxiliary variables consisting the rows of the matrix $\mathbf{x}$. Firstly, let us note that they should be evaluate in such a way that $\mathbf{d}(\mathbf{x}_k, \overline{\mathbf{x}}_s) \neq 0$ or $\mathbf{d}(\mathbf{x}_k, \mathbf{x}_0) \neq 0$. The several ways of the matrix $\mathbf{x}_k$ determining can be proposed. The elements of $\mathbf{x}_k$ can be truncated

means or quantiles. Some proposition of determining the elements of the matrix will be presented in the next sections.

## III. QUANTILES OF AUXILIARY VARIABLES

The sample space of the samples $s$ we denote by $S$. The sample is of the fixed effective size $1 < n < N$. The sampling design is denoted by $P(s)$. We assume that $P(s) > 0$ for all $s \in S$ and $\sum_{s \in S} P(s) = 1$.

Let us assume that observations of the first auxiliary variable increase when their indexes increase. Let $(X_{(j)1}) = (X_{(1)1}, X_{(2)1}, \ldots, X_{(n)1})$ be the sequence of the order statistics of observations of the first auxiliary variable in the sample $s$. The sample quantile of order $\alpha$ is defined, see e.g. Fisz (1963), as follows:

$$Q_{s,\alpha} = X_{(r)1}, \tag{9}$$

where $r = [n\alpha] + 1$, the function $[n\alpha]$ means the integer part of the value $n\alpha$, $r = 1, 2, \ldots, n$. Let us note that $X_{(r)1} = Q_{s,\alpha}$ for $\dfrac{r-1}{n} \leq \alpha \leq \dfrac{r}{n}$. In this paper it will be more conveniently to consider the order statistic than the quantile.

Let $G(r_i, t_i, i = 1, \ldots, k+1) = \{s : X_{(r_i)1} = x_{j_i, 1}, \ i = 1, \ldots k+1\}$ be the set of all samples which $r_i$ -th order statistics $(i = 1, \ldots, k+1)$ of the first auxiliary variable are equal to $x_{t_i, 1}$, respectively where $r_1 \leq t_1 < t_2 < \ldots < t_{k+1} \leq N - n + r_k$. The size of the set $G(r_i, t_i, \ i = 1, \ldots, k+1)$ is denoted by $g(r_i, t_i, \ i = 1, \ldots, k+1)$ and

$$g(r_i, t_i, \ i = 1, \ldots, k+1) = \prod_{i=1}^{k+2} \binom{t_i - t_{i-1} - 1}{r_i - r_{i-1} - 1}, \tag{10}$$

where $r_0 = 0$, $t_0 = 0$, $r_{k+2} = n+1$, $t_{k+2} = N+1$.

The sets $G(r_i, t_i, \ i = 1, \ldots, k+1)$ and $G(r_e, t_e, \ e = 1, \ldots, k+1)$ are disjoint for $i \neq e$. This and the expression (10) lead to the following ones:

$$\bigcup_{t_1=1}^{N-n+r_1} \bigcup_{t_2=t_1+r_2-r_1}^{N-n+r_2} \cdots \bigcup_{t_k=t_{k-1}+r_k-r_{k-1}}^{N-n+r_k} G(r_i, t_i, \ i = 1, \ldots, k+1) = S,$$

$$\sum_{t_1=1}^{N-n+r_1} \sum_{t_2=t_1+r_2-r_1}^{N-n+r_2} \cdots \sum_{t_k=t_{k-1}+r_k-r_{k-1}}^{N-n+r_k} g(r_i,t_i, \ i=1,...,k+1) = \binom{N}{n}. \tag{11}$$

Hence, we have the following probability distribution of the order statistics from the simple sample.

$$P\left(X_{(r_1)1}=x_{t_1,1}, \ X_{(r_2)1}=x_{t_2,1},...,X_{(r_{k+1})1}=x_{t_{k+1},1}\right)=\frac{g(r_i,t_i,i=1,...,k+1)}{\binom{N}{n}}. \tag{12}$$

In order to simplifying the notation let $X_{(r)1}^T = \begin{bmatrix} X_{(r_1)1} & X_{(r_2)1} & \cdots & X_{(r_{k+1})1} \end{bmatrix}$ and

$x_{(t)1}^T = \begin{bmatrix} x_{t_1,1} & x_{t_2,1} & \cdots & x_{t_{k+1},1} \end{bmatrix}$. Now we have: $G\left(X_{(r)1}=x_{(t)1}\right)=G(r_i,t_i,i=1,...,k+1)$

and $g\left(X_{(r)1}=x_{(t)1}\right)=g(r_i,t_i,i=1,...,k+1)$. Moreover, the equation (12) take

the following shorter form: $P\left(X_{(r)1}=x_{(t)1}\right)=\dfrac{g\left(X_{(r)1}=x_{(t)1}\right)}{\binom{N}{n}}$.

## III. SAMPLING DESIGN AND ESTIMATION STRATEGY

Let $\begin{bmatrix} X_{(t)} \\ X_{t_{k+1}*} \end{bmatrix}$ where $x_{t_{k+1}*} = \begin{bmatrix} x_{t_{k+1},1} & x_{t_{k+1},2} & \cdots & x_{t_{k+1},k} \end{bmatrix}$,

$$X_{(t)} = \begin{bmatrix} x_{t_1,1} & x_{t_1,2} & \cdots & x_{t_1,k} \\ x_{t_2,1} & x_{t_2,2} & \cdots & x_{t_2,k} \\ \cdots & \cdots & \cdots & \cdots \\ x_{t_k,1} & x_{t_k,2} & \cdots & x_{t_k,k} \end{bmatrix}.$$

The first column $X_{(t)1}$ of the matrix $\begin{bmatrix} X_{(t)} \\ X_{t_{k+1}*} \end{bmatrix}$ is the observation of the vector

of the order statistic $X_{(r)1}$. Now let us determine the matrix $d\left(x_{(t)},x_{t_{k+1}*}\right)$ in the following way

$$\mathbf{d}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right) = \mathbf{x}_{(t)} - \mathbf{J}_k \mathbf{x}_{t_{k+1}*} .$$ (13)

Moreover,

$$\mathbf{c}(y_0) = y_k - y_0 \mathbf{J}_k,$$ (14)

Let us define the following expression.

$$z\left(\mathbf{X}_{(r)1}, \mathbf{x}_{t_{k+1}*}\right) = \sum_{t_1=1}^{N-n+r_1} \sum_{t_1=t_2+r_2-r_1}^{N-n+r_2} \dots \sum_{t_k=t_{k-1}+r_k-r_{k-1}}^{N-n+r_k} \left|\mathbf{d}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right)\right| g\left(\mathbf{X}_{(r)1} = \mathbf{x}_{(t)1}\right).$$ (15)

**Definition 1.** The sampling design proportional to the determinant $\left|\mathbf{d}\left(\mathbf{x}_{(t)1}, \mathbf{x}_{j_{k+1}*}\right)\right|$ is as follows.

$$P\left(s \middle| \mathbf{X}_{(r)1}, \mathbf{x}_{t_{k+1}*}\right) = \frac{\left|\mathbf{d}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right)\right|}{z\left(\mathbf{X}_{(r)1}, \mathbf{x}_{t_{k+1}*}\right)}$$ (16)

for $s \in G\left(\mathbf{X}_{(r)1} = \mathbf{x}_{(t)1}\right) = G(r_i, t_i, i = 1, \dots, k+1)$.

On the basis of the expression (7) we construct the following regression estimator:

$$\tilde{y} = \overline{y}_s + \frac{1}{\left|\mathbf{d}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right)\right|} \left(\overline{\mathbf{x}} - \mathbf{x}_{t_{k+1}*}\right) \mathbf{D}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right) \mathbf{c}(y_{k+1})$$ (17)

or

$$\tilde{y} = \overline{y}_s + \left(\overline{\mathbf{x}} - \mathbf{x}_{t_{k+1}*}\right) \mathbf{d}^{-1}\left(\mathbf{x}_{(t)}, \mathbf{x}_{t_{k+1}*}\right) \mathbf{c}(y_{k+1})$$ (18)

Particularly, Wywial (2004) considered the case if $k=1$:

$$\tilde{y}_s = \overline{y}_s + \frac{y_2 - y_1}{x_{(r_2)} - x_{(r_1)}} \left( \overline{x} - \overline{x}_s \right). \tag{19}$$

He show that the strategy $\left( \tilde{y}_s, P\left( s \middle| \mathbf{X}_{(r)1}, \mathbf{X}_{t_{k+1}*} \right) \right)$ is not unbiased for population mean $\overline{y}$. So, in the considered case of multidimensional auxiliary variable the strategy is biased, too. From the other point of view the strategy can be useful when some outliers exists or the some observations of variable under study are censored.

## IV. SUPERPOPULATION APPROACH

Let us consider the following regression model:

$$Y = xb + \xi, \tag{19}$$

where $\mathbf{x}$ is the matrix of dimension N×k of non-random auxiliary variables observations, the observations of the random variables $\mathbf{Y}^T = [Y_1 \dots Y_N]$ are values of the variable under study, $\mathbf{b}$ is the column vector of non-random parameters and $\xi^T = [\xi_1 \dots \xi_N]$ is the vector of residuals and $E(\xi) = 0$, variance covariance matrix $V(e) = \mathbf{I}_N \sigma^2$, where $\mathbf{I}_N$ is the unit matrix of the degree N. Hence, $Y_i = x_i \cdot b + \xi_i$ and $E(Y_i) = x_i \cdot b$.

On the basis of the expression (8) we construct the following predictor of the mean value $\overline{Y} = \frac{1}{N} \sum_{i=1}^{N} Y_i$ :

$$\hat{Y} = \overline{Y}_s + \left( \overline{x} - \overline{x}_s \right) \mathbf{d}^{-1} \left( x_k, \overline{x}_s \right) \mathbf{c}(\overline{Y}_s). \tag{20}$$

The similar operations lead to the following result:

$$E_\xi \left( \hat{Y}_s \right) = \overline{x} \mathbf{b} \tag{21}$$

So, the predictor $\hat{Y}_s$ is $\xi$-unbiased. Its variance is as follows:

$$D_\xi^2\left(\hat{Y}\right) = D_\xi^2\left(\overline{Y}_s\right) + D_\xi^2\left(\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{Y}_s)\right) +$$

$$+2Cov_\xi^2\left(\overline{Y}_s, \left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{Y}_s)\right) =$$

$$= D_\xi^2\left(\overline{\xi}_s\right) + D_\xi^2\left(\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right) +$$

$$+\_Cov_\xi^2\left(Y\overline{\xi}_s, \left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right) =$$

$$= \frac{\sigma^2}{n} + E_\xi\left(\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right)^2 + E_\xi\left(\overline{\xi}_s\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right),$$

(22)

$$E_\xi\left(\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right)^2 =$$

$$= \left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)E_\xi\left(\mathbf{c}(\overline{\xi}_s)\mathbf{c}^T(\overline{\xi}_s)\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\left(\overline{x} - \overline{x}_s\right)^T$$

$$E_\xi\left(\mathbf{c}(\overline{\xi}_s)\mathbf{c}^T(\overline{\xi}_s)\right) = E_\xi\left(\xi_k - \overline{\xi}_s\mathbf{J}_k\right)\left(\xi_k^T - \overline{\xi}_s\mathbf{J}_k^T\right) = \sigma^2\left(\mathbf{I}_k - \frac{2}{n}\mathbf{J}_k\mathbf{J}_k^T\right)$$

$$E_\xi\left(\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right)^2 =$$

$$= \sigma^2\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{M}_k\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\left(\overline{x} - \overline{x}_s\right)^T,$$

(23)

where

$$\mathbf{M}_k = \mathbf{I}_k - \frac{1}{n}\mathbf{J}_k\mathbf{J}_k,$$

(24)

$$E_\xi\left(\overline{\xi}_s\left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{c}(\overline{\xi}_s)\right) = \left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)E_\xi\left(\mathbf{c}(\overline{\xi}_s)\overline{\xi}_s\right) = 0$$

(25)

because

$$E_\xi\left(\mathbf{c}(\overline{\xi}_s)\overline{\xi}_s\right) = E_\xi\left(\xi_k - \overline{\xi}_s\mathbf{J}_k\right) = \frac{\sigma^2}{n}\left(\mathbf{J}_k - \mathbf{J}_k\right) = \mathbf{0}_k.$$

The expressions (12)-(25) lead to the following one

$$D_\xi^2\left(\hat{Y}\right) = \frac{\sigma^2}{n}\left(1 + \left(\overline{x} - \overline{x}_s\right)\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\mathbf{M}_k\mathbf{d}^{-1}\left(x_k, \overline{x}_s\right)\left(\overline{x} - \overline{x}_s\right)^T\right)$$

(26)

Particularly, if k=1 then:

$$\hat{Y}_s = \overline{Y}_s + \frac{Y_1 - \overline{Y}_s}{x_1 - \overline{x}_s}\left(\overline{x} - \overline{x}_s\right).$$

We can show that $E_m\left(\hat{Y}_s\right) = E_m\left(\overline{Y}\right)$ and

$$D_\xi^2\left(\hat{Y}\right) = \sigma^2\left(\frac{1}{n} + \left(1 - \frac{1}{n}\right)\frac{\left(\overline{x} - \overline{x}_s\right)^2}{\left(x_1 - \overline{x}_s\right)^2}\right) \geq \frac{\sigma^2}{n}.$$

So, when a purposive sample $s_1$ is such a one that $\overline{x} = \overline{x}_{s_1}$ the above variance takes minimal value. The sample $s_1$ is called the balanced one.

The next particular case of the predictor given by the expression (20) is as follows

$$\hat{Y}_s = \overline{Y}_s + \frac{Y_2 - Y_1}{x_{(r_2)} - x_{(r_1)}}\left(\overline{x} - \overline{x}_s\right).$$

We can show that $E_\xi\left(\hat{Y}_s\right) = E_\xi\left(\overline{Y}\right)$ and

$$D_\xi^2\left(\hat{Y}\right) = \sigma^2\left(\frac{1}{n} + \frac{2\left(\overline{x} - \overline{x}_s\right)^2}{\left(x_{(r_2)} - x_{(r_1)}\right)^2}\right) \geq \frac{\sigma^2}{n}.$$

Similarly like in the previous case when a purposive sample $s_1$ is such a one that $\overline{x} - \overline{x}_{s_1} = 0$ the above variance takes minimal value.

The more general predictor is as follows:

$$\breve{Y} = Y_0 + \left(\overline{x} - x_s\right)d^{-1}\left(x_k, x_s\right)c(Y_0) \tag{27}$$

where $Y_0$ is an observed value of variable under study and it is not the element of the vector $\mathbf{Y}_k$, so $Cov(Y_0, Y_i)=0$ for $Y_i \in \mathbf{Y}_k$. We can show that

$$E_{\xi}\left(\breve{Y}_s\right) = \overline{x}b. \tag{28}$$

So, the predictor $\breve{Y}_s$ is $\xi$-unbiased. Its variance can be derived similarly as the parameter $D_{\xi}^2\left(\hat{Y}\right)$ and is as follows:

$$D_{\xi}^2\left(\breve{Y}\right) = \sigma^2\left(1 + \left(\overline{x} - \overline{x}_s\right)d^{-1}\left(x_k, \overline{x}_s\right)L_k d^{-1}\left(x_k, \overline{x}_s\right)\left(\overline{x} - \overline{x}_s\right)^T\right), \tag{29}$$

where

$$L_k = I_k + J_k J_k. \tag{30}$$

Hence, the statistic $\breve{Y}_s$ is not consistent predictor of $\overline{Y}$.

Finally let us consider the following predictor:

$$\tilde{Y}_s = \overline{Y}_s + \left(\overline{x} - x_0\right)d^{-1}\left(x_{(t)}, x_0\right)c(Y_0), \tag{31}$$

where $\overline{Y}_s = \dfrac{1}{N}\sum_{i \in s} Y_i$, $c(Y_0)=Y_k-Y_0J_k$ and $Y_0$ is an observed value of variable under study and it is not the element of the vector $Y_k$. Under these assumptions we evaluate the $\xi$-expected value:

$$E_{\xi}\left(c(Y_0)\right) = \left[E_{\xi}\left(Y_i\right) - E_{\xi}\left(Y_0\right)\right] = \left[x_i b - x_0 b\right] = \left(x_{(t)} - x_0 J_k\right)b = d\left(x_{(t)}, x_0\right)b$$

$$= \overline{x}_s b + \left(\overline{x} - x_0\right)d^{-1}\left(x_{(j)}, x_0\right)d\left(x_{(j)}, x_0\right)b$$

$$E_{\xi}\left(\tilde{Y}_s\right) = \overline{x}_s b + \left(\overline{x} - x_0\right)b. \tag{32}$$

So, the predictor $\tilde{Y}_s$ is $\xi$-biased. Its variance is derived in the following way:

$$D_{\xi}^2\left(\tilde{Y}\right) = D_{\xi}^2\left(\overline{Y}_s\right) + D_{\xi}^2\left(\left(\overline{x} - x_0\right)d^{-1}\left(x_k, x_0\right)c(Y_0)\right) +$$

$$+2Cov_{\xi}\left(\overline{Y}_s, \left(\overline{x} - x_0\right)d^{-1}\left(x_k, x_0\right)c(Y_0)\right) =$$

$$= D_\xi^2\left(\overline{Y}_s\right) + D_\xi^2\left(\left(\overline{x}-x_0\right)d^{-1}\left(x_k,x_0\right)c(Y_0)\right) +$$

$$+2Cov_\xi\left(\overline{Y}_s,\left(\overline{x}-\overline{x}_s\right)d^{-1}\left(x_k,x_0\right)c(Y_0)\right) =$$

$$= D_\xi^2\left(\overline{Y}_s\right) + D_\xi^2\left(\left(\overline{x}-x_0\right)d^{-1}\left(x_k,x_0\right)c(Y_0)\right) + 0$$

$$D_\xi^2\left(\tilde{Y}_s\right) = \sigma^2\left(\frac{1}{n} + \left(\overline{x}-\overline{x}_s\right)d^{-1}\left(x_k,\overline{x}_s\right)L_k d^{-1}\left(x_k,\overline{x}_s\right)\left(\overline{x}-\overline{x}_s\right)^T\right) \qquad (33)$$

Hence, the statistic $\tilde{Y}_s$ is not consistent predictor of $\overline{Y}$.

The particular case of the above predictor can be as follows:

$$\tilde{Y}_s = \overline{Y}_s + \frac{Y_2-Y_1}{x_{(r_2)}-x_{(r_1)}}\left(\overline{x}-x_0\right).$$

We can show that $E_m\left(\tilde{Y}_s\right) = E_m\left(\overline{Y}\right)$ and

$$D_\xi^2\left(\tilde{Y}_s\right) = \sigma^2\left(\frac{1}{n} + \frac{2\left(\overline{x}-x_0\right)^2}{\left(x_{(r_2)}-x_{(r_1)}\right)^2}\right) \geq \frac{\sigma^2}{n}.$$

Let us remind the bias of the predictor is $E_\xi\left(\tilde{Y}_s\right) = \overline{x}_s b + \left(\overline{x}-x_0\right)b$. So, when a purposive sample $s_2$ is such a one that $\overline{x} = x_0$ the above variance takes minimal value. For instance $s_2$ can be such a balanced sample that sample median $x_0$ is equal to the population mean $\overline{x}$.

### Acknowledgement

## REFERENCES

Borsuk K. (1969): *Multidimensional Analytic Geometry*. PWN, Warsaw.

Fisz, M. (1963). *Probability Theory and Mathematical Statistics*. Wiley and Sons Inc., New York.

Hellwig Z. (1963). *Linear Regression and Its Application in Economy (in Polish)*. PWN, Warszawa.

Kendall, M. G., Stuart, A. (1961). *The Advanced Theory of Statistics. Vol. II: Inference and Relationship*. Charles Griffin & Company Limited, London.

Wald, A. (1940): The fitting the stright lines if both variables are subject to errors. *Annals of Mathematical Statistics*, 11, pp. 284

Wywiał. J. (2004). Quqantile regression sampling strategy. In: *Metoda Reprezentacyjna w Badaniach Ekonomiczno-Społecznych (Survey Sampling in Economical and Social Research)*. (Edited by J. Wywiał ) Katowice University of Economics, Katowice.

*Janusz Wywiał*

## OCENA WARTOŚCI PRZECIĘTNEJ ZA POMOCĄ REGRESYJNEJ STRATEGII LOSOWANIA WYKORZYSTUJĄCEJ KWANTYLE ZMIENNEJ POMOCNICZEJ

Problem oceny wartości średniej z wykorzystaniem danych o wszystkich wartościach cech pomocniczych jest rozważany. W tym celu znany estymator regresyjny zależny od wielu zmiennych pomocniczych jest wykorzystywany. W odróżnieniu od zwykłego podejścia znanego w metodzie reprezentacyjnej do oceny parametrów regresji są wykorzystywane kwantyle jednej ze zmiennych dodatkowych. Otrzymane na tym polu wyniki są adoptowane do konstrukcji predytorów wartości średniej w nadpopulacji. Wyprowadzono również wariancje różnych odmian proponowanych predykatorów.