*Janusz Wywiał\**

# ON ESTIMATION OF DOMINANT
# OF MULTIDIMENSIONAL RANDOM VARIABLE

## Abstract

The problem of estimation of the mode of a continuous distribution function of multi-dimensional random variable is considered. The estimator of the mode is the vector of means from appropriately truncated sample. The truncation sample is obtained through rejecting the observation in such a way that the measure of skewnees of multidimensional variable takes value as close zero as possible. We can expect that through successful truncation of the sample the vector of sample means approach to vector of modes of multidimensional variable. The estimator constructed in such a way is usually biased estimators of the mode. Moreover, the biased estimators of values of modal regressions are proposed. The well known "jackknife" procedure is proposed to evaluate the mean square errors of the estimators.

**Key words:** moments of truncated distribution, estimation of distribution mode.

### I. ESTIMATION OF MODE IN ONE DIMENSIONAL CASE

Let $X_{(1)} \leqslant X_{(2)} \leqslant ... \leqslant X_{(n)}$ be the sequence of the order statistics from a simple sample drawn from one-dimensional distribution. The third central moment from the right hand truncated simple sample is defined by the expression:

$$M_3(X_{(k)}) = \frac{1}{k} \sum_{i=1}^{k} (X_{(i)} - \overline{X}_k)^3, \quad 1 < k \leqslant n$$

where: $\overline{X}_k = \frac{1}{k} \sum_{i=1}^{k} X_{(i)}$. Particularly, $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i = \overline{X}_n$.

---

\* Professor, Department of Statistics, University of Economics in Katowice.

Let us assume that $M_{(3)}(X_{(n)}) = M_3 > 0$. The sample quasi-mode $D_q$ is defined as follows:

$G_e = \overline{X}_e$, if $M_3(X_{(e)}) \leqslant 0$ and $M_3(X_{(k)}) > 0$ for $k = e + 1, ..., n$.

The statistic $Dq$ will be usually a biased estimator of the mode $\gamma$.

**Theorem 1** (Johnson and Rogers, 1951). Let the density function of a random variable $X$ be one-modal end it is concave from the left side of the dominant $\gamma$ and it is convex from the right side of the dominant $\gamma$. Moreover, let $D^2(X) > 0$. There exists one dimensional probability distribution if and only if $(E(X) - \gamma)^2 \leqslant 3 \ D^2(X)$.

Hence, the well known Pearson coefficient of skewnees is bounded.

Let us assume that probability distribution of $X$ is right side truncated in the point a and $E(X|a) = \gamma(a)$, $D^2(X|a)$ and $\eta_3(X|a)$ *be the expected value, the variance and the third central moment of the truncated distribution, respectively*. The theorem 1 lead to the following. If $\eta_3(X|a) \to 0$ and $D^2(X|a)$ do not increases then the $(E(X|a) - \gamma)^2$ do not increase. Hence, the appropriate truncation of the probability distribution can lead to decreasing the distance between mode and the expected value of the truncated probability distribution. This lead to the conclusion that the estimator $G_e$ is closer to the mode than the common mean $\overline{X}$ from non-truncated simple sample. That is why the mean from the truncated sample can be used to estimation the mode.

## II. ESTIMATION OF MODE IN MULTI-DIMENSIONAL CASE

Let $f(x_1, ..., x_k) = f(\mathbf{x})$ where $\mathbf{x} = [x_1...x_k]$, be a density function of an $k$-dimensional random variable defined in $R^k$. The vectors of expected values and variances are denoted by $\boldsymbol{\mu} = [\mu_1...\mu_k]$ and $\boldsymbol{\sigma}^2 = [\sigma_1^2...\sigma_k^2]$. The mode of the $k$-dimensional random variable is denoted by $\gamma = [\gamma_1...\gamma_k]$ and $f(\gamma) = $ maximum. The central moments of the order 3 of the $r$-dimensional variable are denoted by:

$$\eta_{uwz}(X_i, X_j, X_g) = \int_{-\infty}^{\infty} ... \int_{-\infty}^{\infty} (x_i - \mu_i)^u (x_j - \mu_j)^v (x_g - \mu_g)^g f(x_1, ..., x_k)dx_1...dx_k,$$

$$\eta_u(X_i) = \eta_{u00}(X_i, X_j, X_g)$$

Let us consider the truncated multidimensional variable of the following density function:

$$f(x_1, \ldots, x_k | A) = \frac{f(x_1, \ldots, x_k)}{P([X_1, \ldots, X_k] \in A)}$$

where A is a convex region and $A \subseteq R^k$. Hence, $f(x_1, \ldots, x_k | A)$ is density function of truncated distribution. The vectors of expected values and variances of the truncated distributions are as follows: $\mu(A) = [\mu_1(A) \ldots \mu_k(A)]$, $\sigma^2(A) = [\sigma_1^2(A) \ldots \sigma_k^2(A)]$, where:

$$\mu_i(A) = \int_{-\infty}^{\infty} \ldots \int_{-\infty}^{\infty} x_i f(x_1, \ldots, x_k | A) dx_1 \ldots dx_k,$$

$$\sigma_i^2(A) = \int_{-\infty}^{\infty} \ldots \int_{-\infty}^{\infty} (x_i - \mu_i(A))^2 f(x_i, \ldots, x_k | A) dx_1 \ldots dx_k,$$

$$\eta_{uwz}(X_i, X_j, X_g | A) =$$

$$= \int_{-\infty}^{\infty} \ldots \int_{-\infty}^{\infty} (x_i - \mu_i(A))^u (x_j - \mu_j(A))^v (x_g - \mu_g(A))^g f(x_1, \ldots, x_k | A) dx_1 \ldots dx_k.$$

Let us introduce the following notation:

$$\theta = [\theta_1 \; \theta_2 \; \theta_3] \text{ is of dimensions } 1 \times (k^2 + k(k-1)(k-2)/6),$$

where

$$\theta_1 = [\eta_3(X_1) \; \eta_3(X_1) \; \ldots \; \eta_3(X_k)] \text{ of dimensions } 1 \times k,$$

$$\theta_2 = [\eta_{12}(X_1, X_2) \; \eta_{21}(X_1, X_2) \; \ldots \; \eta_{21}(X_{k-1}, X_k) \; \eta_{21}(X_k, X_{k-1})]$$
$$\text{of dimensions } 1 \times k(k-1),$$

$$\theta_3 = [\eta_{111}(X_1, X_2, X_3) \; \eta_{111}(X_1, X_2, X_4) \; \ldots \; \eta_{111}(X_{k-3}, X_{k-1}, X_k)$$
$$\eta_{111}(X_{k-2}, X_{k-1}, X_k)] \text{ of dimensions } 1 \times \frac{1}{6} k(k-1)(k-2),$$

The vector can be estimated by the vector $L = [L_1 \, L_2 \, L_3]$, where:

$$L_1 = [C_{30}(X_1, X_2) \; C_{03}(X_2, X_1) \; \ldots \; C_{03}(X_{k-1}, X_k)],$$

$$L_2 = [C_{12}(X_1, X_2) \; C_{21}(X_1, X_2) \; \ldots \; C_{21}(X_{k-1}, X_k) \; C_{21}(X_k, X_{k-1})],$$

$$L_3 = [C_{111}(X_1, X_2, X_3) \; C_{111}(X_1, X_2, X_4) \; \ldots \; C_{111}(X_{k-3}, X_{k-1}, X_k)$$
$$C_{111}(X_{k-2}, X_{k-1}, X_k)].$$

Similarly, we define the vectors of moments of truncated distribution:

$$\theta(A) = [\theta_1(A) \; \theta_2(A) \; \theta_3(A)],$$

where

$$\theta_1(A) = [\eta_i(X_1|A)], \quad \theta_2(A) = [\eta_{12}(X_i,X_j|A)]$$

and

$$\theta_3(A) = [\eta_{111}(X_i,X_j,X_t|A)].$$

It is well known that if a distribution function is symmetric, then all central moments of the order 3 of the marginal one or two dimensional distributions are equal to zero and $\theta = 0$. In the case when $\theta \neq 0$ we can find such set $A = \{A: A \subseteq R^k$ and $\theta(A) = 0\}$. The vector of mean values $\mu(A) = [\mu_1(A)...\mu_k(A)]$ from truncated distribution we define as the A-mode of this distribution and it will be denoted by $\gamma(A)$. It is obvious that the A-mode can not necessary be equal to the mode of the entire distribution. The conditions of such equality can be stated quite simple only in the case of one-dimensional distribution, see Wywiał (2000a–2000b). That is why the parameter $\gamma(A)$ can be called a quasi-mode as in one-dimensional case.

The simple sample of size $n$ is denoted by $X = [X_{*1}, ..., X_{*n}]$, where:

$$X_{*t} = \begin{bmatrix} X_{1t} \\ ... \\ X_{kt} \end{bmatrix}, \quad t = 1, ..., n.$$ The sample moments are as follows:

$$C_{uvz}(X_i,X_j,X_g) = \frac{1}{n} \sum_{t=1}^{n} (X_{it} - \overline{X}_i)^u (X_{jt} - \overline{X}_j)^v (X_{gt} - \overline{X}_g)^z, \quad \overline{X}_i = \frac{1}{n} \sum_{t=1}^{n} X_{it}$$

Let $X_{*(h)}$ $(h = n, ..., 1)$ be such vector that

$$d(X_{*(h)}, \overline{X}_{*(h)}) = \underset{i \in F_{(h)}}{\text{maximum}} \{Q_{i,h}\}$$

where:

$$F_{(h)} = \{t : t = 1, ..., n \quad \text{and} \quad X_{*t} \neq X_{*(p)}, \quad \text{for} \quad p > h]$$

$$Q_{i(h)} = (X_{*1} - \overline{X}_{*(h)})^T (X_{*i} - \overline{X}_{*(h)}),$$

$$\overline{X}_{*(h)} = \frac{1}{n_{(h)}} \sum_{t \in F_{(h)}} X_{*t}, \quad n_{(h)} = \text{Card}\{F_{(h)}\}$$

$$\overline{X}_{*(h)} = [\overline{X}_{1(h)}...\overline{X}_{k(h)}], \quad \overline{X}_{*(n)} = \overline{X} = \frac{1}{n} \sum_{t=1}^{n} X_{*t}.$$

A truncated sample is identified by the set $F_{(h)}$.

Let

$$C_{uv}(X_i, X_j, X_g | h) = \frac{1}{n_{(h)}} \sum_{t \in F_{(h)}} (X_{it} - \overline{X}_{i(h)})^u (X_{jt} - \overline{X}_{j(h)})^v (X_{gt} - \overline{X}_{g(h)})^z,$$

$$L_{(h)} = [L_{1(h)} \ L_{2(h)} \ L_{3(h)}] = [L_{1(h)} \ L_{2(h)} \ ... \ L_{w(h)}], \text{ where}$$

$$w = 1 \times (k^2 + k(k-1)(k-2)/6) \quad \text{and}$$

$$L_{1(h)} = [C_{30}(X_1, X_2 | h) \ C_{03}(X_2, X_1 | h) \ ... \ C_{03}(X_{k-1}, X_k | h)],$$

$$L_{2(h)} = [C_{12}(X_1, X_2 | h) \ C_{21}(X_1, X_2 | h) \ ... \ C_{21}(X_{k-1}, X_k | h) \ C_{21}(X_k, X_{k-1} | h)],$$

$$L_{3(h)} = [C_{111}(X_1, X_2, X_3 | h) \ C_{111}(X_1, X_2, X_4 | h) \ ...$$

$$C_{111}(X_{k-3}, X_{k-1}, X_k | h) \ C_{111}(X_{k-2}, X_{k-1}, X_k | h)].$$

The vector parameter $\gamma$ can be estimated bz means of the statistic $G_i(u)$ determined bz the following procedure:.

Let

$$Z_{(h)} = [L^2_{1(h)} \ L^2_{2(h)} ... L^2_{w(h)}]$$

$$L*_{(h)} = \underset{i=1,...,w}{\text{maximum}} \{L^2_{i(h)}\}$$

$$G_{1(u)} = X*_{(u)}: \quad u = \max\{h: h = n,..., 1 \quad \text{and} \quad L*_{(h)} \leqslant e\}$$

where $e \geqslant 0$. It seems that $e$ should be assigned in such a way that the size $n_{(u)}$ of the truncated sample is sufficiently large.

The next estimator is based on measure of multivariate skewness defined by Mardia (1970). Let $\sigma^{ij}(A)$ be an $(i, j)$ element of the matrix $\Sigma^{-1}(A)$ where $\Sigma(A) = [\sigma_{ij}(A)]$, $\sigma_{ij}(A) = \eta_{11}(X_i, X_j | A)$, is the variance-covariance matrix of the truncated distribution. In the case of truncated distribution the skewness coefficient can be written in the following way.

$$\beta(A) = \sum_{\{i_1, i_2, i_3\}} \sum_{\{j_1, j_2, j_3\}} \sigma^{i_1 j_1}(A)\sigma^{i_2 j_2}(A)\sigma^{i_3 j_3}(A)\eta_{111}(X_{i_1}, X_{i_2}, X_{i_3} | A)\eta_{111}$$

$$(X_{j_1}, X_{j_2}, X_{j_3} | A)$$

where $\{i_1, i_2, i_3\}$ and $\{j_1, j_2, j_3\}$ are variations with replications. Each variation is determined on the basis of sequence: 1, 2, ..., $k$.

The coefficient of skewness from the truncated sample $F_{(h)}$ is as follows:

$$B_{(h)} = \frac{1}{n_{(h)}^2} \sum_{i,j \in F_{(h)}} ((X_{*i} - \overline{X}_{*(h)})^T S_{(h)}^{-1} (X_{*i} - \overline{X}_{*(h)}))^3$$

where:

$$S_{(h)} = [C_{11(h)}(X_i, X_j)], \quad C_{11(h)}(X_i, X_j) = \frac{1}{n_{(h)}} \sum_{t \in F_{(h)}} (X_{it} - \overline{X}_{i(h)})(X_{jt} - \overline{X}_{j(h)}).$$

The second estimator of quasi-mode is as follows:

$$G_{2(v)} = X_{*(v)}: \quad v = \max\{h: h = n, ..., 1 \quad \text{and} \quad B_{(h)} \leqslant e\}$$

The statistics $G_{1(u)}$ and $G_{1(u)}$ can be biased estimators of the mode $\gamma$. Their variances can be estimated using the well known method of jackknife.

Let us assume that $\theta_1 = [\eta_3(X_1) \; \eta_3(X_1) \; ... \; \eta_3(X_k)] > 0$ for the values $x \in B \subseteq R^k$ of a multidimensional random variable $X$. Let $A_\# \subset B$ be such the sup-set that $\theta_1(A_\#) = 0$. If $D^2(X_i) \geqslant D^2(X_i|A_\#)$ for each $i = 1, ..., k$ and $D^2(X_j) > D^2(X_j|A_\#)$ for at least one index $j$ where $j = 1, ..., k$ then

$$\sum_{i=1}^{k} (E(X_i|A_\#) - \gamma_i)^2 \leqslant \sum_{i=1}^{k} (E(X_i) - \gamma_i)^2.$$

It means that the appropriate truncation of the multivariate probability distribution can lead to decreasing the distance between mode and the expected value of the truncated probability distribution. This leads to the conclusion that in special cases the expected values of the above introduced estimators can be closer to the mode than the vector of means from non-truncated simple sample.

Some simulation methods should be used to analysis of the accuracy of the proposed estimators. The similar estimators can be constructed on the basis of other coefficients of multivariate skewness e.g. like those proposed by Wywiał (1983, 1985).

### III. MODAL REGRESSION

Let us consider the following regression model:

$$Y = XB + U$$

where $Y^T = [Y_1 \; ... \; Y_n]$ is the vector of independent random variables. The matrix of non-random values of explanatory variables is denoted by $X$ and it is of dimension $n \times r$. The vector of parameters $B$ has dimensions $r \times 1$.

Moreover, $E(U) = 0$ and variance-covariariance matrix of $U$ is diagonal and $D^2(U_i) = \sigma^2$ and $\eta_3(U_i) > 0$ for each $i = 1, ..., n$. The modal values of the variables $Y^T = [Y_1 ... Y_n]$ will be denoted by $\gamma^T = [\gamma_1 ... \gamma_n]$ and the modal values of the variables $U^T = [U_1 ... U_n]$ are the same and equal to $\kappa$. The modal regression is defined in the following way:

$$\gamma = XB + \kappa J_n$$

where $J_n$ is the vector of dimensions $n \times 1$.

Under these assumption $E(Y) = XB$ can be estimated by means of the statistic:

$$\hat{Y} = XB$$

where $\hat{B}$ is unbiased estimator of $B$ obtained by the well known method of least squares and

$$\hat{B} = (X^T X)^{-1} X^T Y$$

It is obvious that $E(\hat{Y}) = E(Y)$

The vector of residual is as follows:

$$\hat{U} = Y - \hat{Y} = MY$$

where

$$M = I_n - X(X^T X)^{-1} X^T, \ I_n \ \text{is unit matrix of degree } n.$$

The problem is determining the vector $\gamma$ of modal values. Firslty, as suggested Pawłowski (1973), the mode $\kappa$ should be estimated by means of a statistic $\hat{G}$ which is a function of the residuals $\hat{U}$. This leads to the following estimator of the regression

$$\tilde{Y} = \hat{Y} + \hat{G}J_n = X\hat{B} + \hat{G}J_n$$

where $J_n$ is the vector of dimensions $n \times 1$. In order to simplify the consideration let us assume that $\kappa > 0$. The statistic $\hat{G}$ can take form of the estimator of mode proposed in the first. Let $\hat{U}_{(1)} \leqslant \hat{U}_{(2)} \leqslant ... \leqslant \hat{U}_{(n)}$ be the sequence of the order statistics section. The third central moment from the right hand truncated sample is defined by the expression:

$$M_3(\hat{U}_{(k)}) = \frac{1}{k} \sum_{i=1}^{k} (\hat{U}_{(i)} - \hat{U}_{\#k})^3, \quad 1 < k \leq n$$

where:

$$\hat{U}_{\#k} = \frac{1}{k} \sum_{i=1}^{k} \hat{U}_{(i)}$$

When we assume that $M_3(\hat{U}_{(n)}) = M_3(\hat{U}) > 0$, the sample quasi-mode $\hat{G}$ is defined as follows:

$$\hat{G} = \hat{U}_{\#c}, \text{ if } M_3(\hat{U}_{\#c}) \leq 0 \text{ and } M_3(\hat{U}_{\#c}) > 0 \text{ for } k = c + 1, ..., n$$

The statistic $\tilde{Y}$ is biased estimator of $\gamma$ because $\hat{G}$ is usually biased estimator of $\kappa$ as it was discussed in the previous paragraph.

The idea of construction of the next estimator is based on the truncated least-square method proposed by Ruppert and Carrolla (1980). Let $Y_{(e)}$ is consisted of those variables of the vector $Y$ which indexes are equal to appropriate indexes of residuals in the set $\{\hat{U}_i : i = 1, ..., n; \hat{U}_i \leq \hat{U}_{(e)}\}$ where evaluation of $\hat{U}_{(e)}$ was showed above. In the same way the submatrix $X_{(e)}$ of the observations matrix $X$ is determined. Next, we again evaluate the parameters of the linear regression:

$$B_{(e)} = (X_{(e)}^T X_{(e)})^{-1} X_{(e)}^T Y_{(e)}$$

If $e > r$, the next estimator of the modal regression is as follows:

$$\tilde{Y}_{(e)} = X_{(e)} B_{(e)}$$

In general case the statistic $\tilde{Y}_{(e)}$ is not unbiased estimator of the modal regression.

Let us note that estimators $\tilde{Y}$ and $\tilde{Y}_{(e)}$ can be easy adopted to evaluation conditional prognosis of the variable under study for some fixed values of auxiliary variables. Let $x$ be row vector (of dimensions $r \times 1$) of values of auxiliary variables. The two predictors of the value $\gamma(x)$ of modal regression for given vector x are as follows:

$$\tilde{Y}(x) = x\hat{B} + \hat{G} \quad \text{or} \quad \tilde{Y}_{(e)}(x) = xB_{(e)}$$

The mean square error of the estimators or predictors of the modal regression values and the parameters can be evaluated on the basis of the jackknife method. The accuracy of the both estimators or predictors of the modal regression can be analysed on the basis of appropriate simulation studies.

REFERENCES

Johnson N.L., Rogers C.A. (1951), The moment problem for unimodal distributions, *Annals of Mathematical Statistics*, **22**, 433–439.
Mardia K.V. (1970), Measures of multivariate skewneess and kurtosis with application, *Biometrika*, **57**, 3, 519–530.
Pawłowski Z. (1973), *Prognozy ekonometryczne*, PWN, Warszawa.
Ruppert D., Carroll R.J. (1980), Trimmed least squares estimation in the linear model, *Journal of the American Statistical Association*, **75**, 828–838.
Wywiał J. (1983), Normalized coefficients of deviation from multinormal distribution (in Polish). *Przegląd Statystyczny*, **30**, 77–83.
Wywiał J. (1985), *Test normalności dla wielowymiarowej zmiennej losowej dla dużych prób*, (A test for normality of a multidimensional random variable in the large sample case; in Polish), *Przegląd Statystyczny*, **32**, 355–364.
Wywiał J. (2000a), Estimation of distribution function mode on the basis of sample moment or sample median, *Badania Operacyjne i Decyzje*, no 2, 89–98.
Wywiał J. (2000b), Estimation of mode on the basis of a truncated sample, *Acta Universitatis Lodziensis. Folia Oeconomica*, **152**, 73–81.

*Janusz Wywiał*

O ESTYMACJI DOMINANTY WIELOWYMIAROWEJ ZMIENNEJ LOSOWEJ

Streszczenie

Praca dotyczy problemu estymacji dominanty rozkładu prawdopodobieństwa wielowymiarowej zmiennej losowej. Zajmowano się problemem estymacji dominanty zmiennej losowej ciągłej. Analizowano jednomodalne rozkłady prawdopodobieństwa.

W niniejszej pracy analizowano głównie klasę rozkładów prawdopodobieństwa zmiennych losowych charakteryzującą się tym, że zachodzą dla nich pewne nierówności między wartościami oczekiwanymi i dominantą rozkładu funkcją jego trzecich momentów centralnych. Dla takiej klasy rozkładów są wprowadzone dwa estymatory dominanty, których wyznaczanie w praktyce ma charakter iteracyjny. Z grubsza rzecz biorąc, wyliczenie wartości estymatora dominanty wiąże się z sukcesywnym obcinaniem obserwacji próby do chwili, gdy zostaną spełnione pewne warunki, a w szczególności, że pewna funkcja trzech mieszanych momentów centralnych rozkładu uciętego osiągnie wartość zero. Wówczas oceną punktu będącego dominantą rozkładu wielo-wymiarowego są właśnie średnie z uciętych rozkładów brzegowych z próby. Proponowane estymatory mogą być obciążone. W niektórych przypadkach dało się oszacować maksymalny poziom takiego obciążenia. Proponowane są również dwa estymatory regresji modalnej.